

# Feature-based Transfer Learning via Kernel Embedding of Distributions

[Sinno Jialin Pan \(Ph.D.\)](#)

Nanyang Assistant Professor

School of Computer Science and Engineering  
Nanyang Technological University, Singapore

# Transfer of Learning

A psychological point of view

- The study of dependency of human conduct, learning or performance on prior experience.
  - [Thorndike and Woodworth, 1901] explored how individuals would transfer in one context to another context that share similar characteristics



# Transfer Learning

In the machine learning community

- The ability of a system to recognize and apply knowledge and skills learned in previous domains/tasks to novel tasks/domains, which share some commonality
- Given a target domain/task, how to identify the commonality between the domain/task and previous domains/tasks, and transfer knowledge from the previous domains/tasks to the target one?

Transfer learning for classification, and regression problems.

- [**Pan** and Yang, A Survey on Transfer Learning, IEEE TKDE 2010]
- [**Pan**, Transfer learning, Book Chapter 2014]

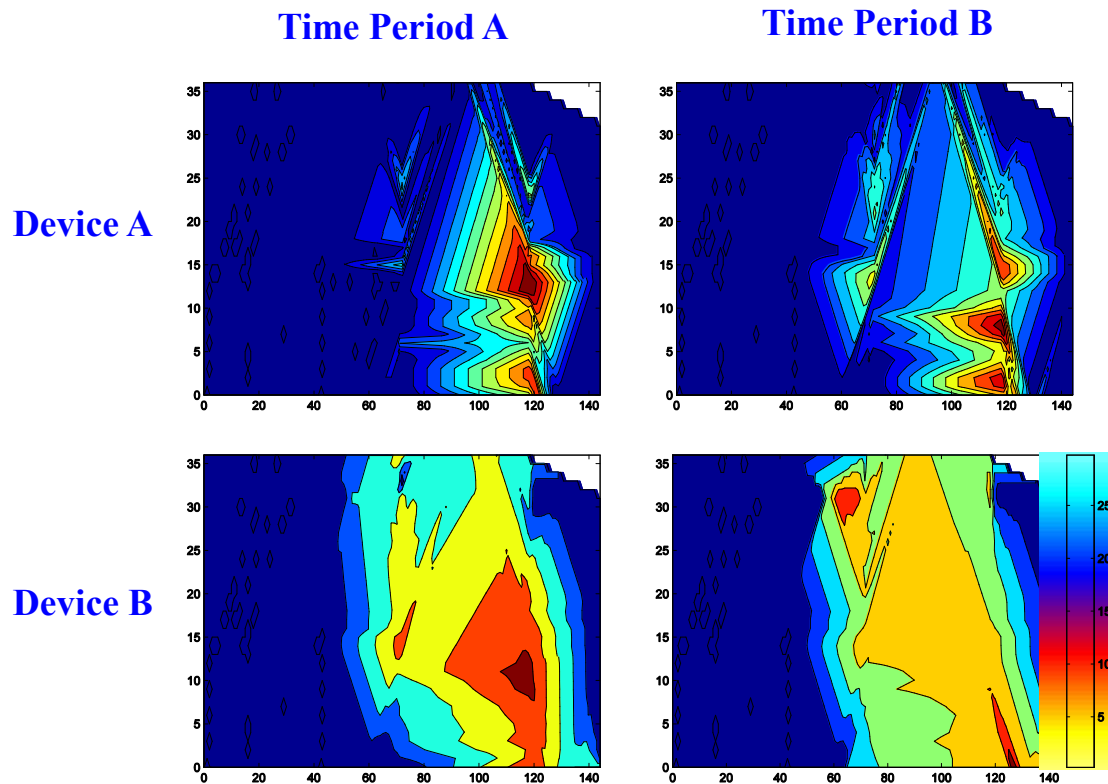
Transfer learning for reinforcement learning problems.

- [Taylor and Stone, Transfer Learning for Reinforcement Learning Domains: A Survey, JMLR 2009]

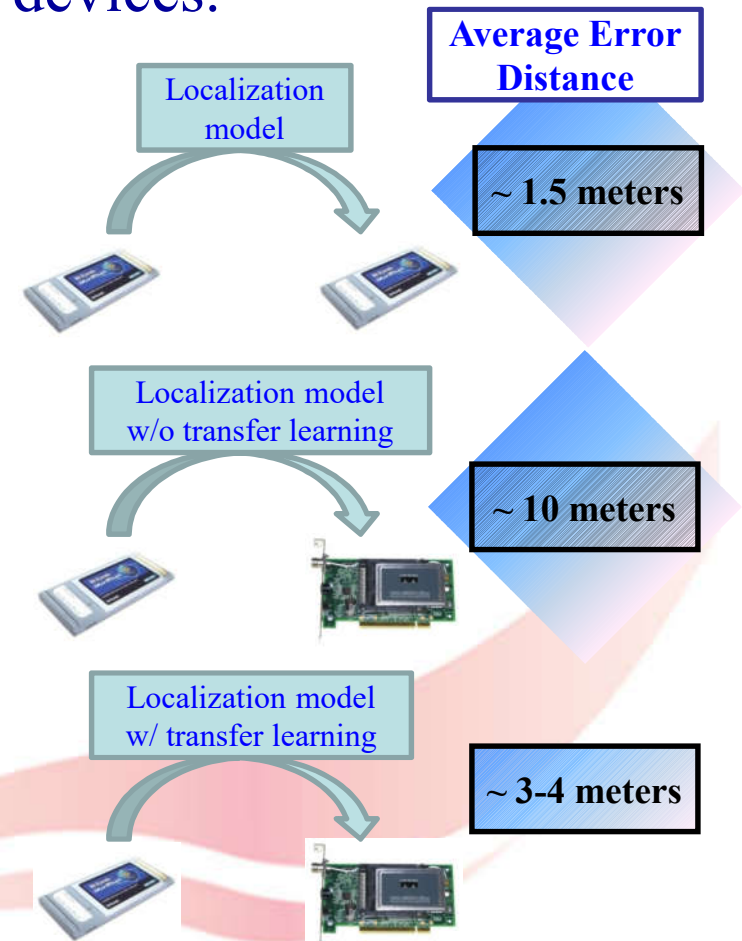
*Focus!*

# Applications

- WiFi localization: signal strength changes a lot over different time periods, or across different mobile devices.



Contour of signal strength values

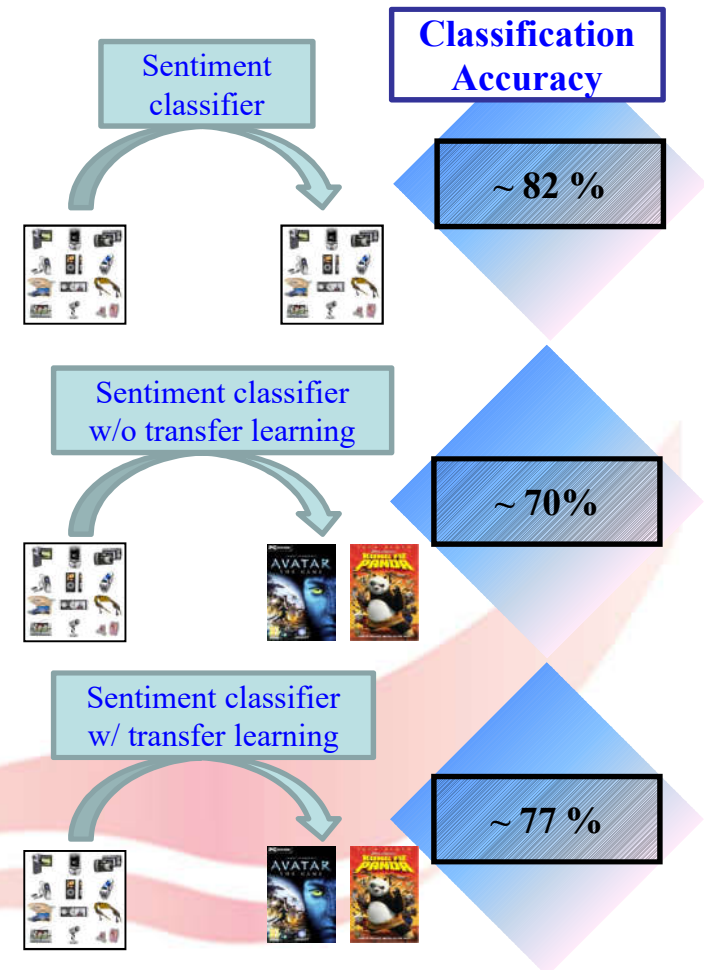


# Applications (cont.)

- Sentiment analysis: users may use different sentiment words across different domains.

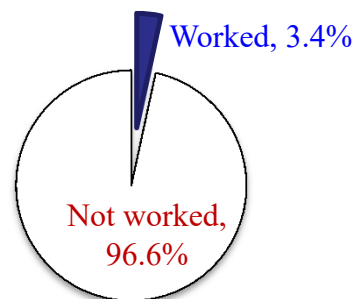
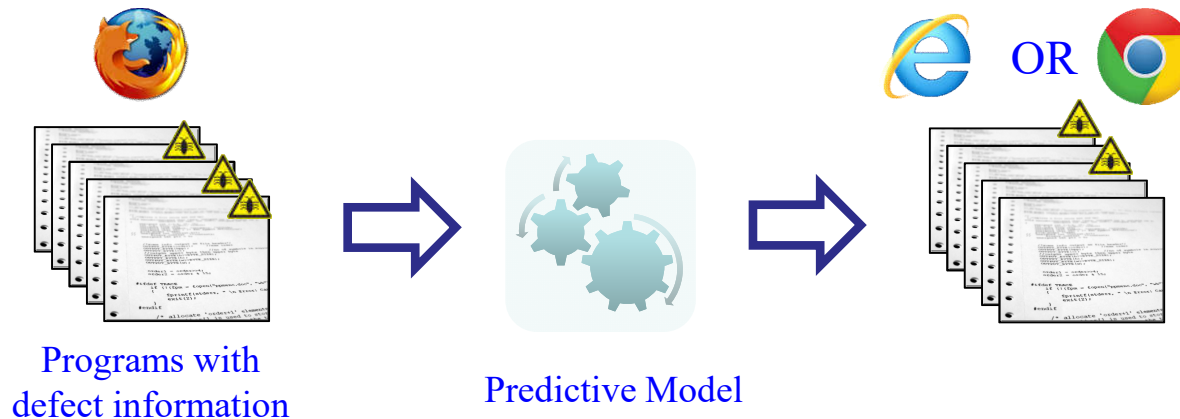
Electronics	Video Games
(1) <b>Compact</b> ; easy to operate; very good picture quality; looks <b>sharp</b> !	(2) A very good game! It is action packed and full of excitement. I am very much <b>hooked</b> on this game.
(3) I purchased this unit from Circuit City and I was very excited about the quality of the picture. It is really nice and <b>sharp</b> .	(4) Very <b>realistic</b> shooting action and good plots. We played this and were <b>hooked</b> .
(5) It is also quite <b>blurry</b> in very dark settings. I will never buy HP again.	(6) The game is so <b>boring</b> . I am extremely unhappy and will probably never buy UbiSoft again.

Product reviews on different domains

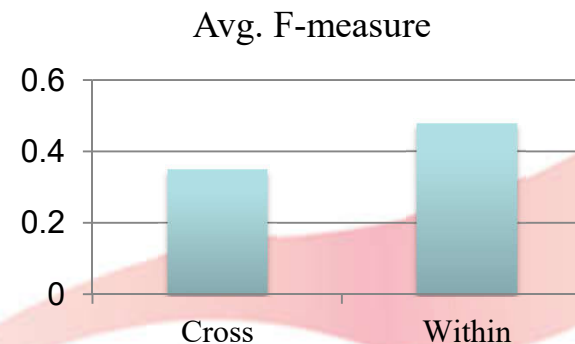


# Applications (cont.)

- **Defect prediction**: development processes can be very different across different projects



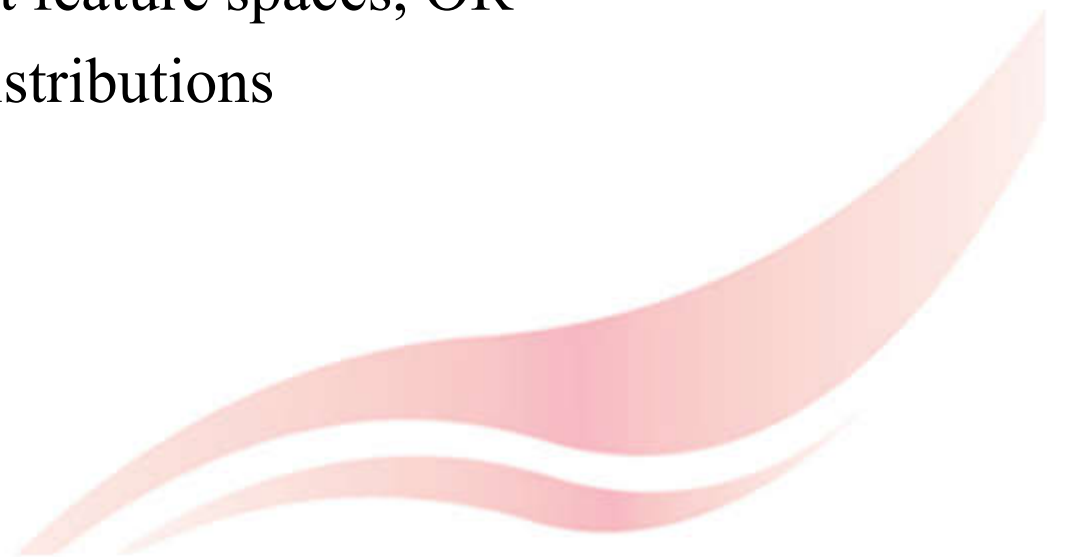
[Zimmerman et al. FSE-09] “We ran 622 cross-project predictions and found only 3.4% actually worked.”



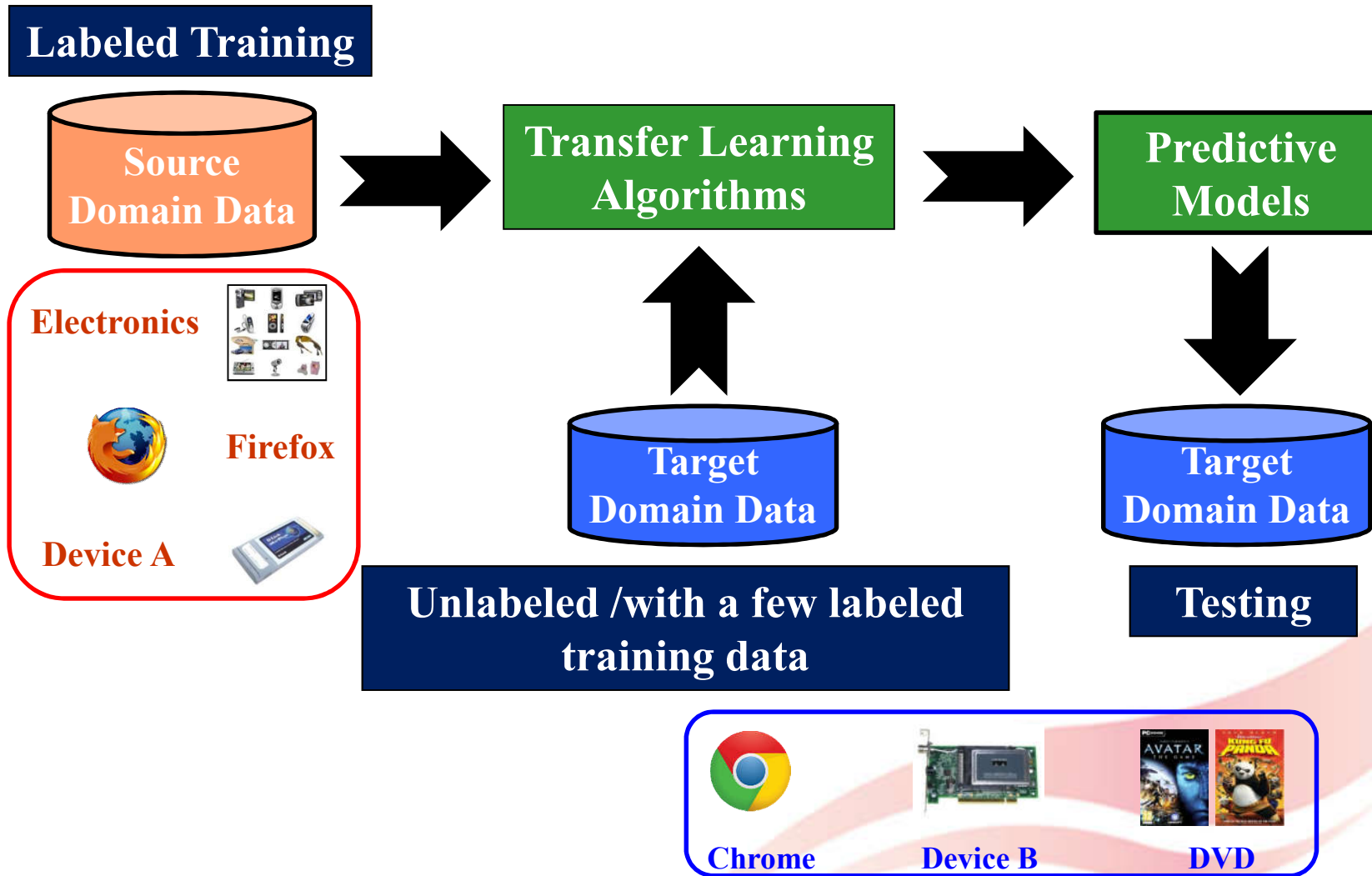
[Rahman, Posnett, and Devanbu. FSE-12]

# Why Models Perform Poor across Domains

- Fundamental assumption in machine learning: training and test data are assumed to be
  - Represented in the same feature space, AND
  - Follow the same data distribution
- Training and test data from different domains may be
  - Represented in different feature spaces, OR
  - Follow different data distributions

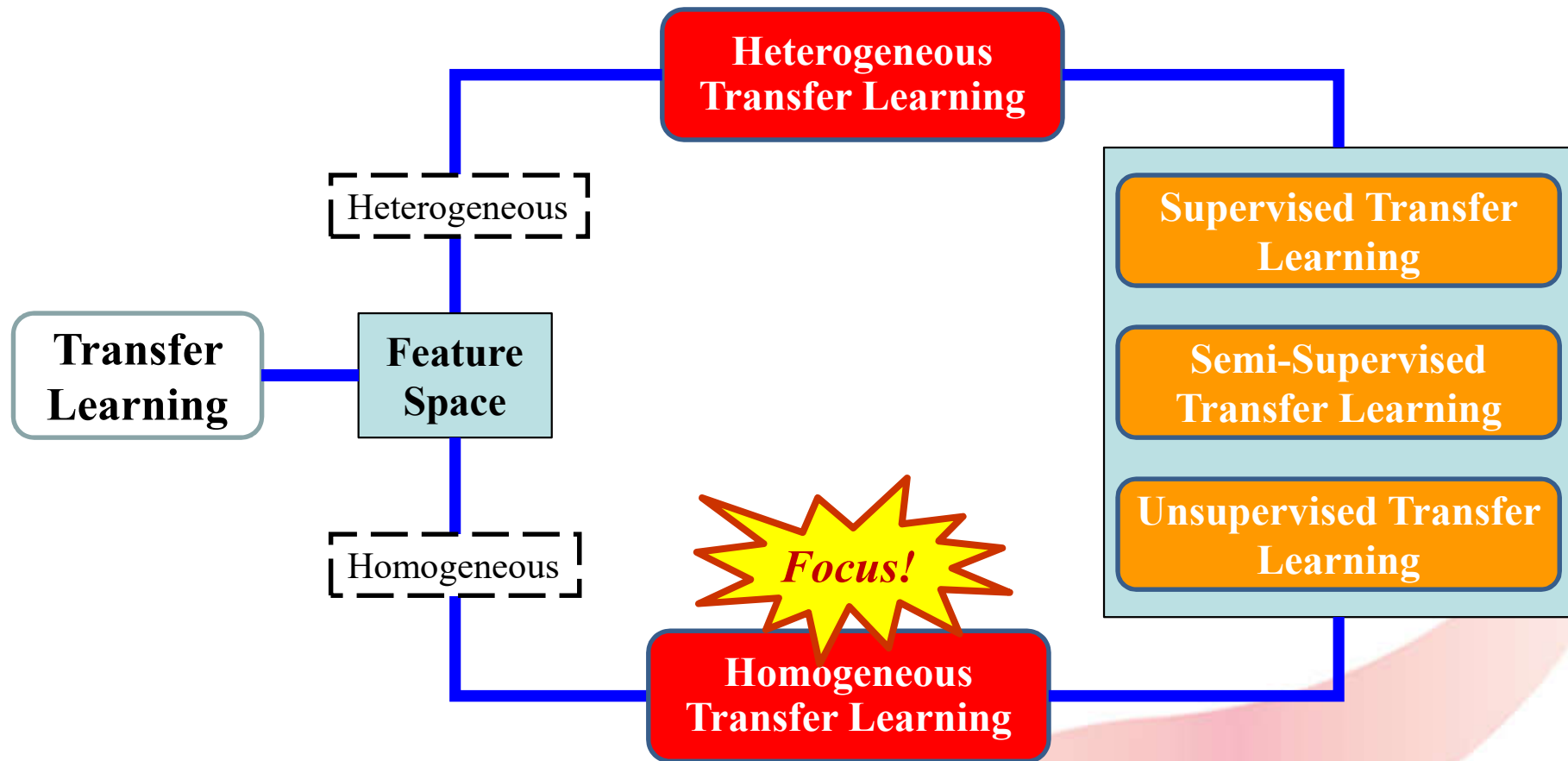


# The Goal of Transfer Learning

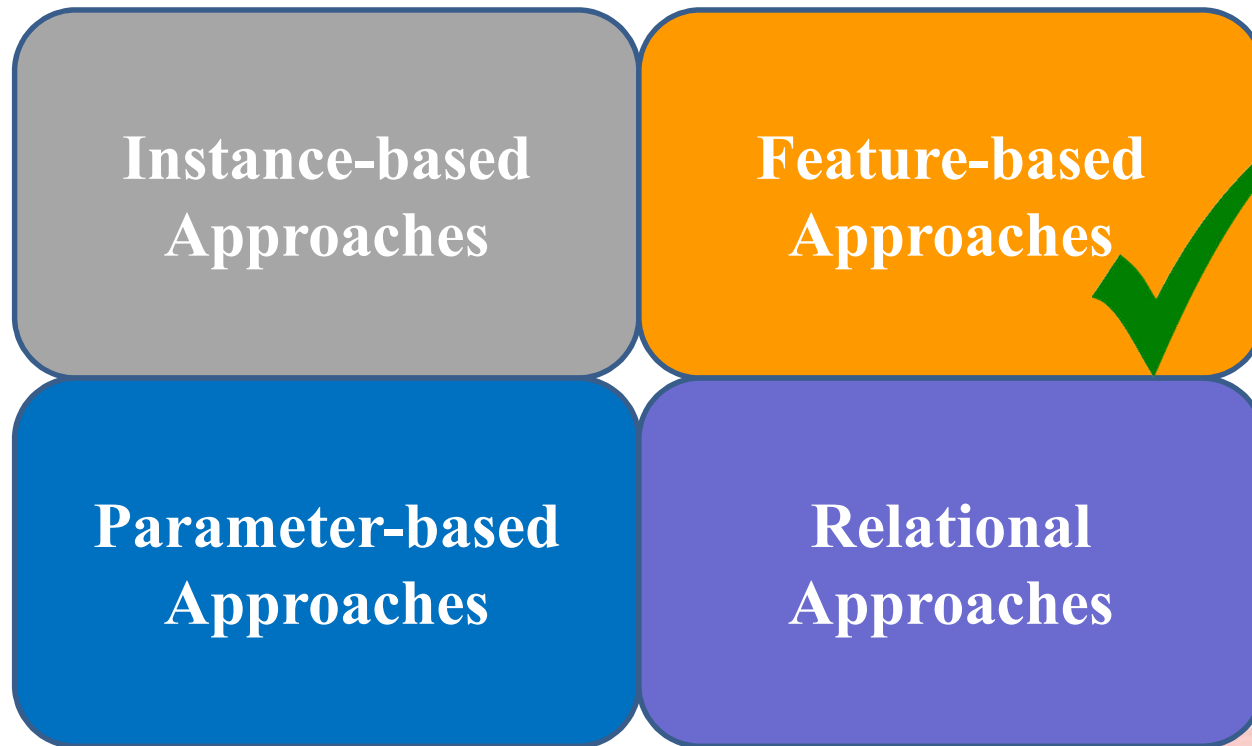




# Transfer Learning Settings



# Transfer Learning Approaches



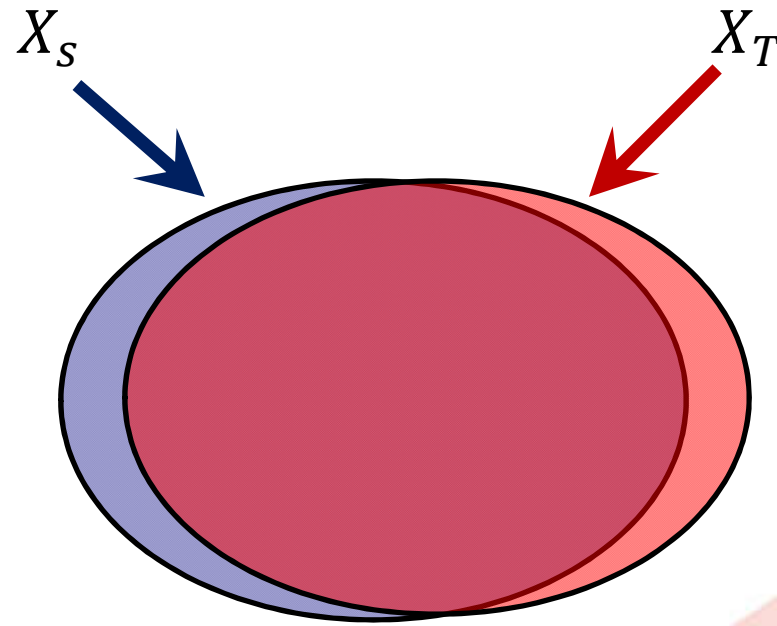
# Instance-based Approaches

## General Assumption

Source and target domains have a lot of overlapping features (domains share the same/similar support)

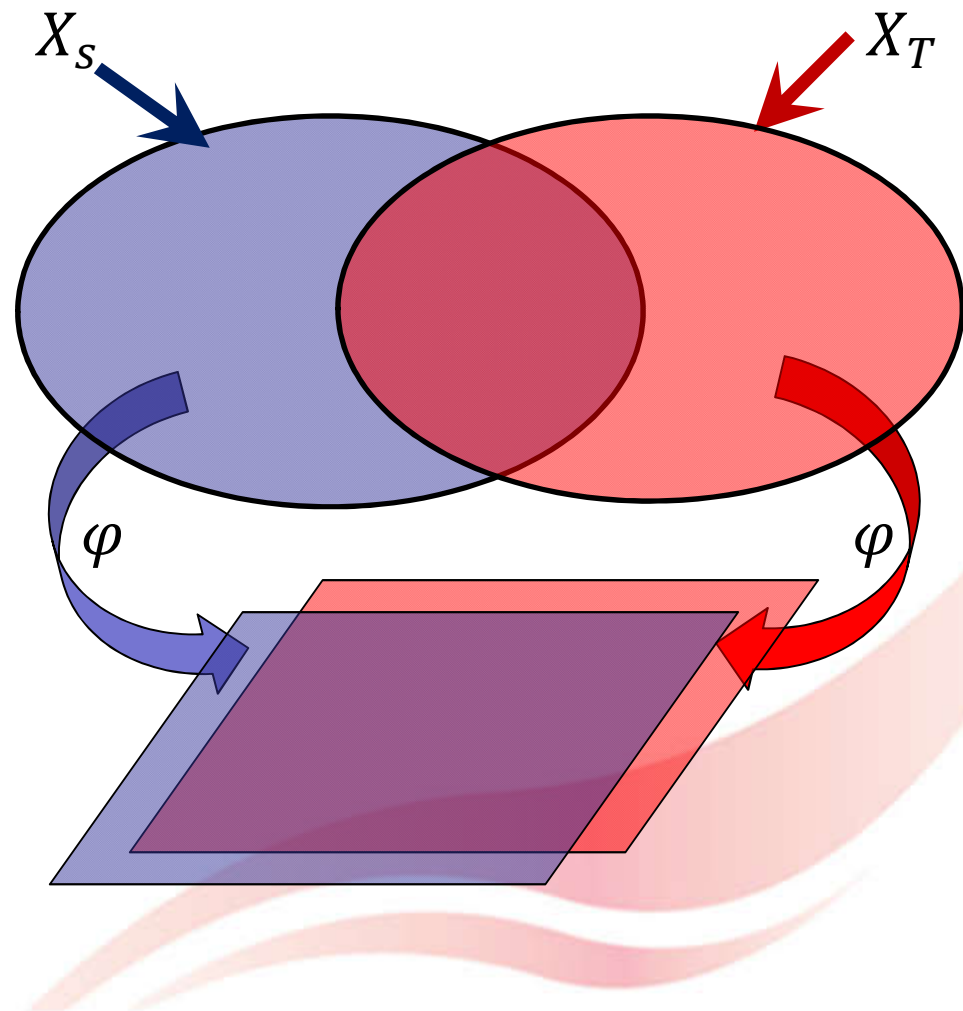
## Motivation

Reweight source-domain labeled data to be reused for the target domain



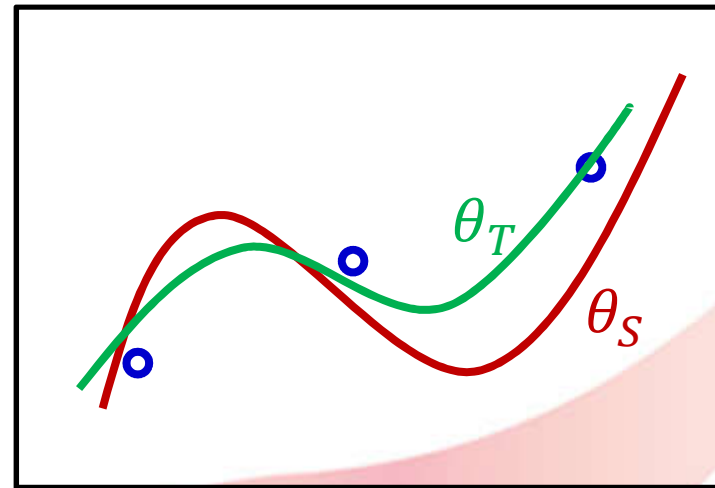
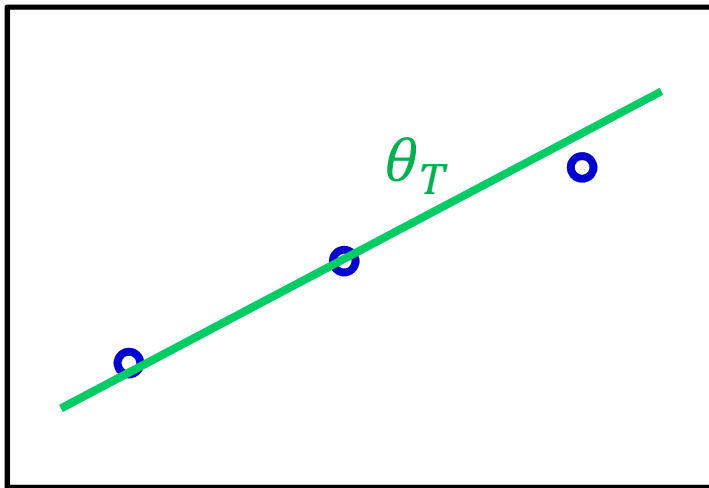
# Feature-based Approaches

When source and target domains only have some overlapping features. (lots of features only have support in either the source or the target domain)



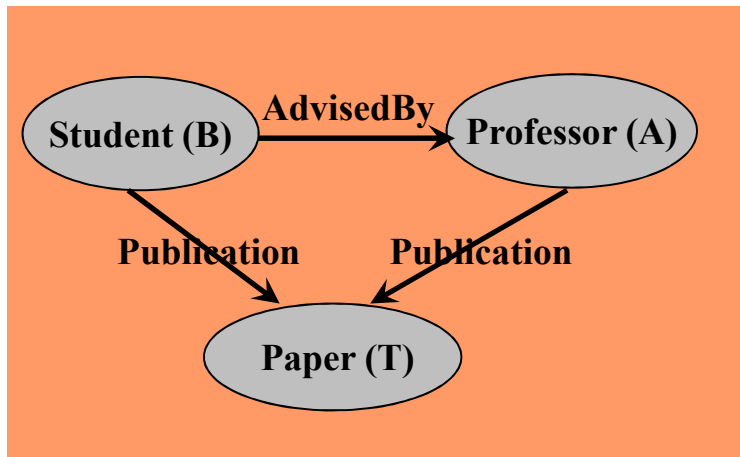
# Parameter-based Approaches

- **Motivation:** A well-trained source model  $\theta_S$  has captured a lot of structure from data. If two tasks are related, this structure can be transferred to learn a more precise target model  $\theta_T$  with a few labeled data in the target domain

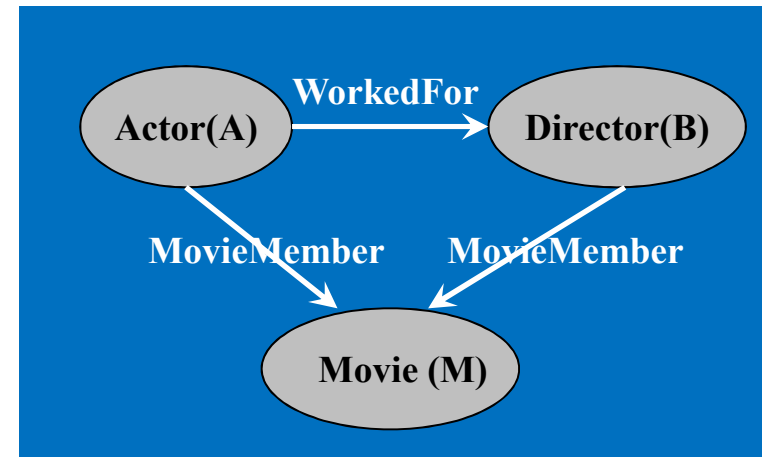


# Relational Approaches

Academic domain (source)



Movie domain (target)



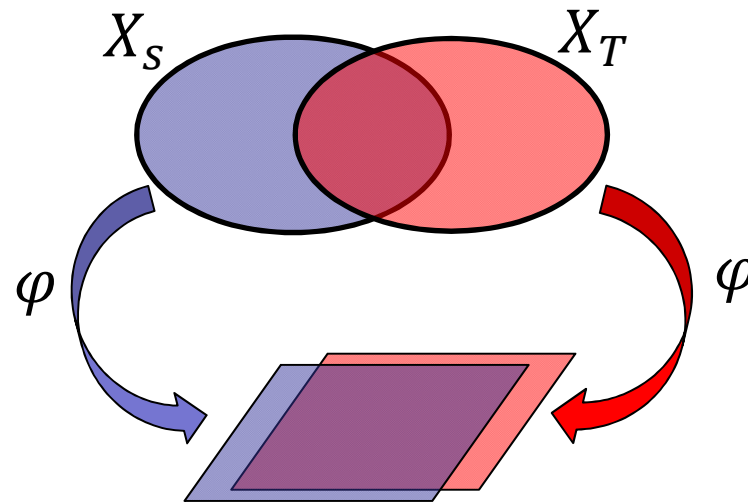
$\text{AdvisedBy}(B, A) \wedge \text{Publication}(B, T) \Rightarrow \text{Publication}(A, T)$

$\text{WorkedFor}(A, B) \wedge \text{MovieMember}(A, M) \Rightarrow \text{MovieMember}(B, M)$

$P1(x, y) \wedge P2(x, z) \Rightarrow P2(y, z)$

# Feature-based Approaches (cont.)

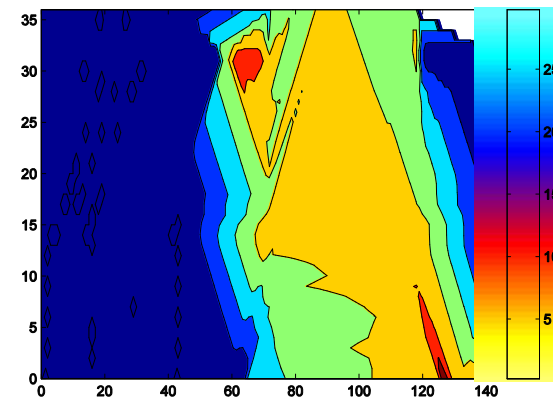
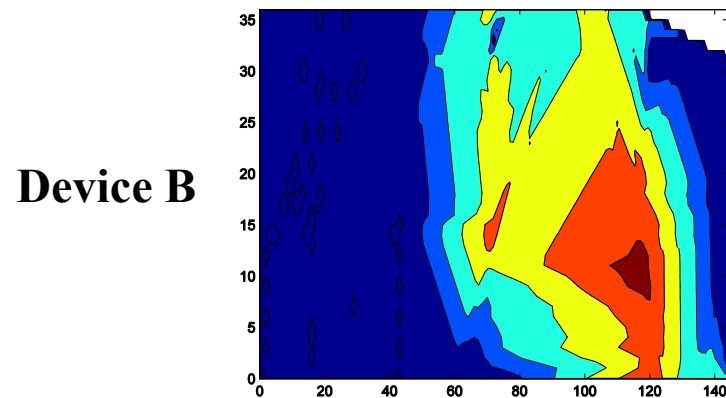
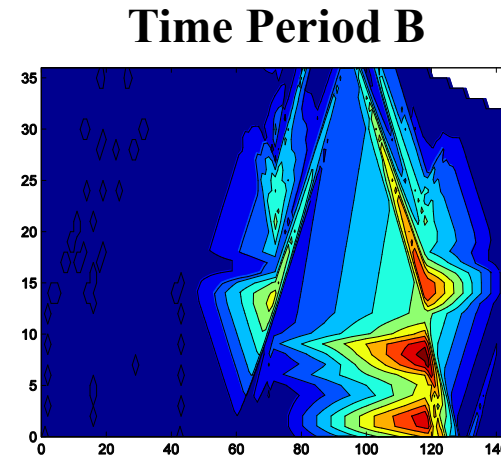
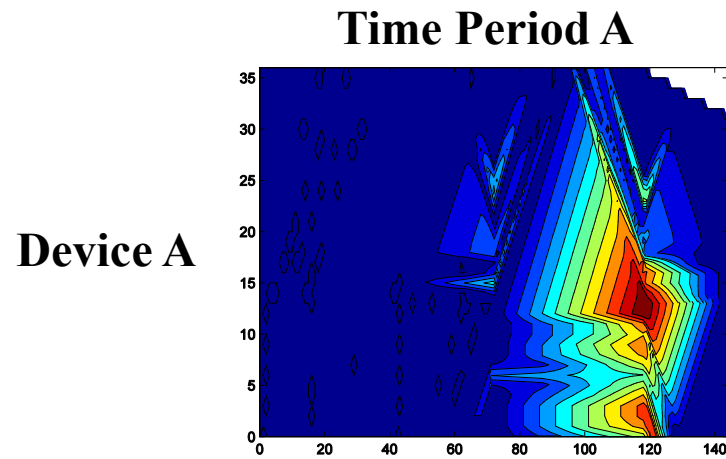
How to learn  $\varphi$ ?



- Solution 1: Encode application-specific knowledge to learn the transformation, e.g., sentiment analysis
- Solution 2: General approaches to learning the transformation

# Developing General Approaches

## An illustrating Example

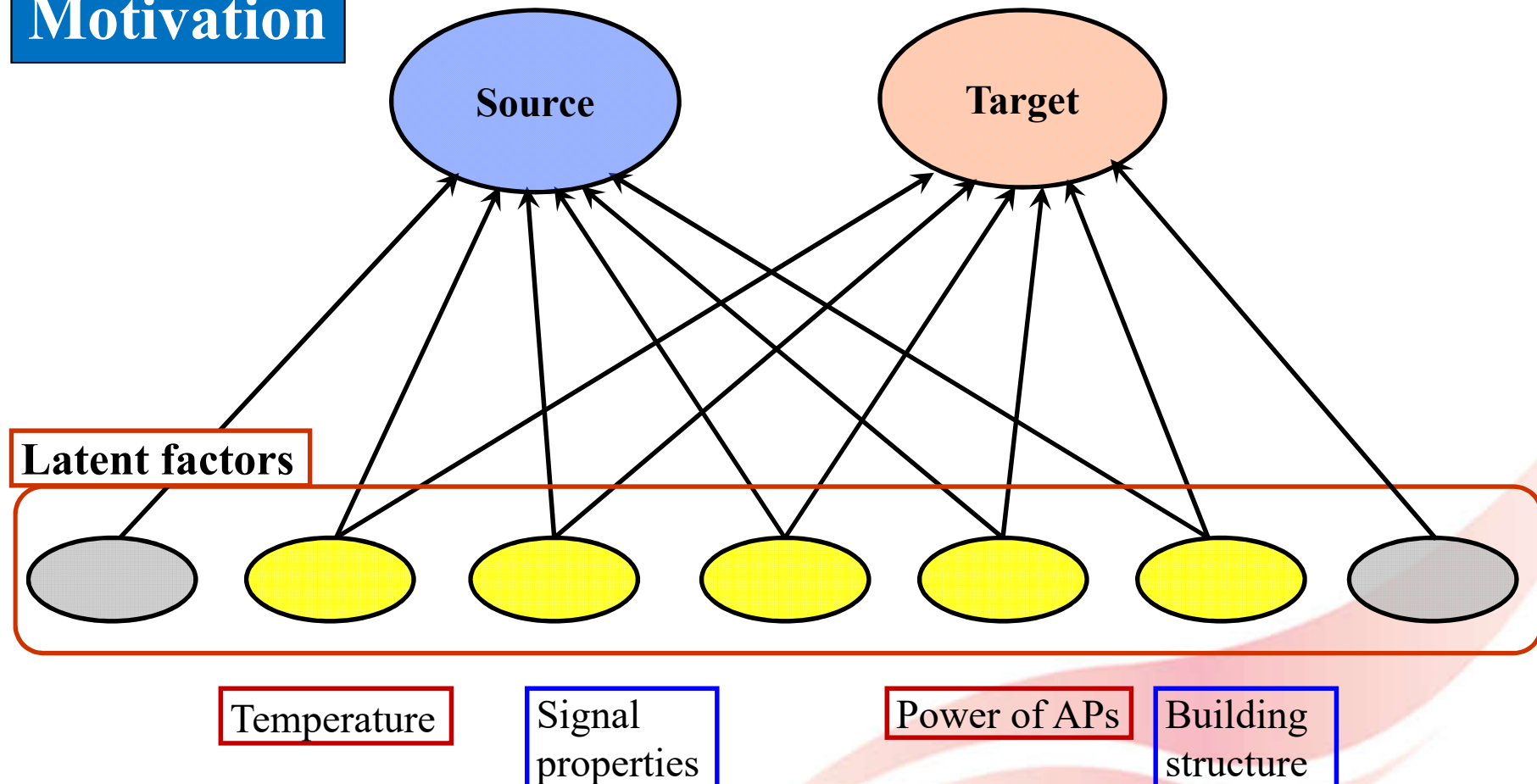




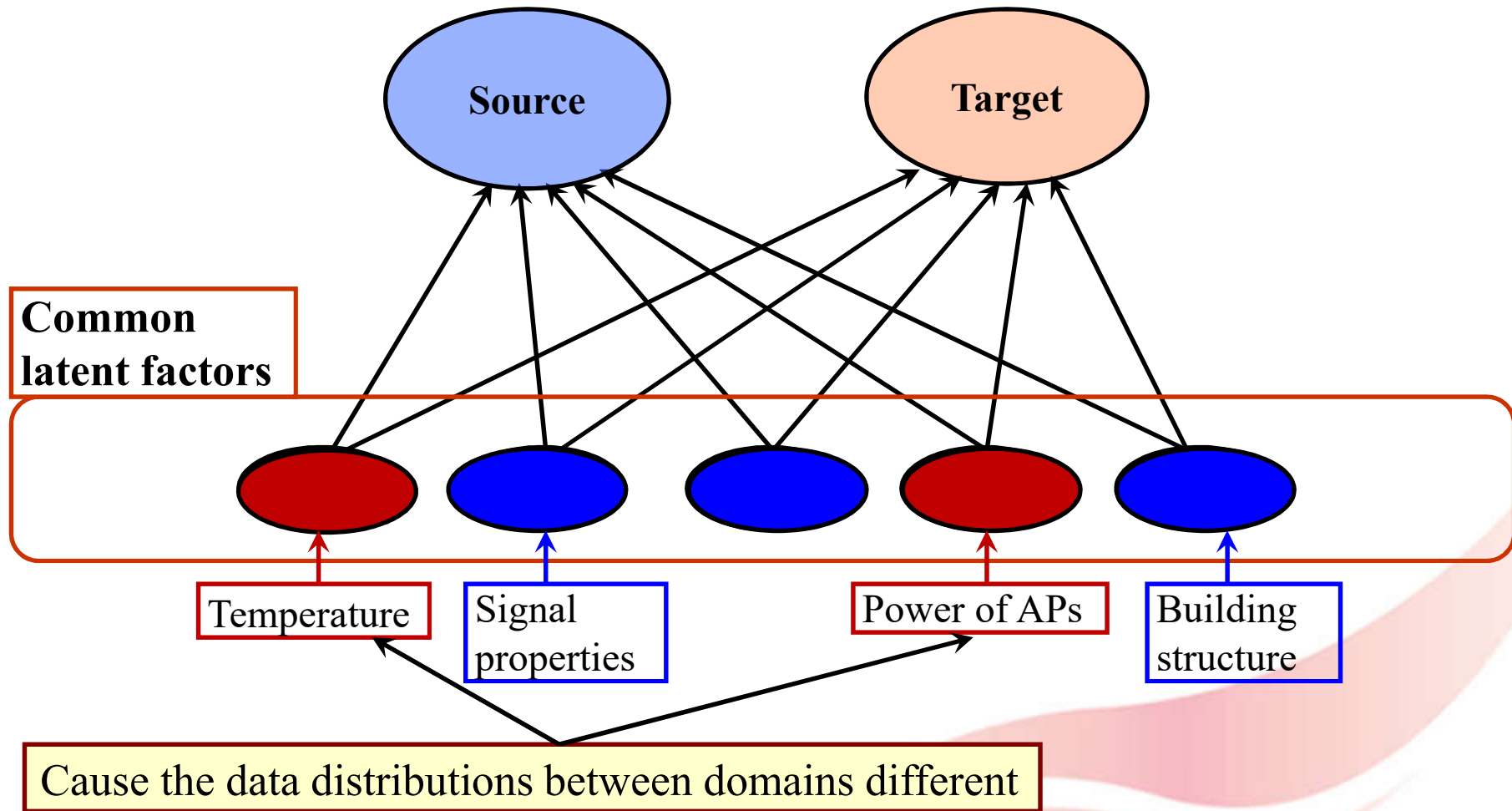
# Learning Features via Kernel Embedding of Distributions

Transfer Component Analysis (TCA) [Pan et al., IJCAI-09, TNN-11]

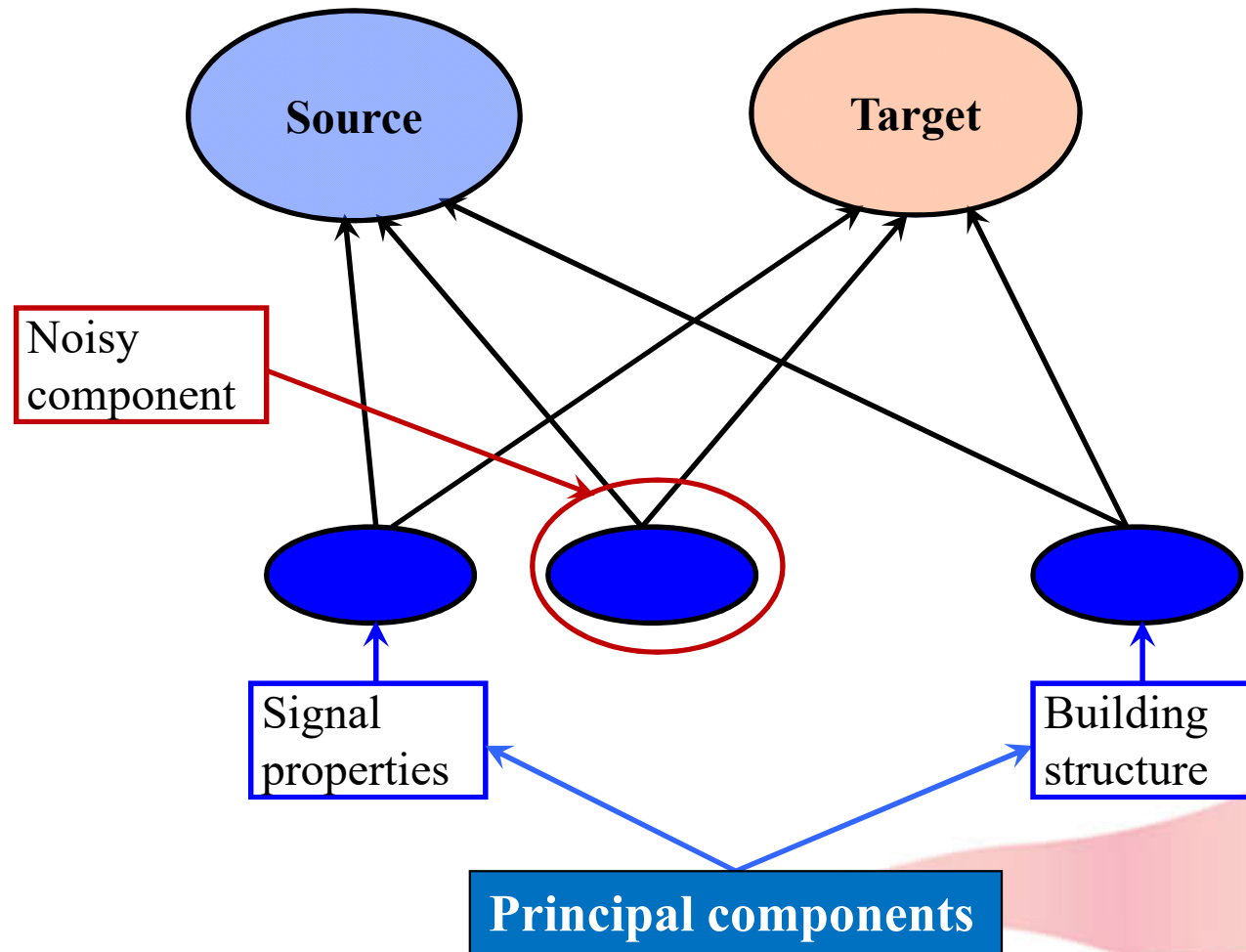
## Motivation



# Transfer Component Analysis (cont.)

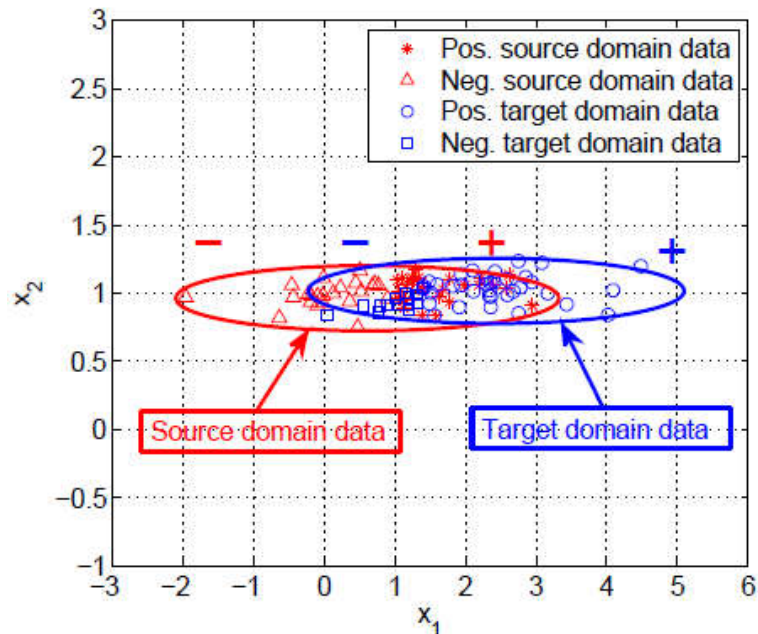


# Transfer Component Analysis (cont.)



# Transfer Component Analysis (cont.)

- Learning  $\varphi$  by only minimizing distance between distributions may map the data onto noisy factors



# Transfer Component Analysis (cont.)

- Main idea: the learned  $\varphi$  should map the source domain and target domain data to a latent space spanned by the factors that reduce domain distance as well as preserve data structure
- High level optimization problem

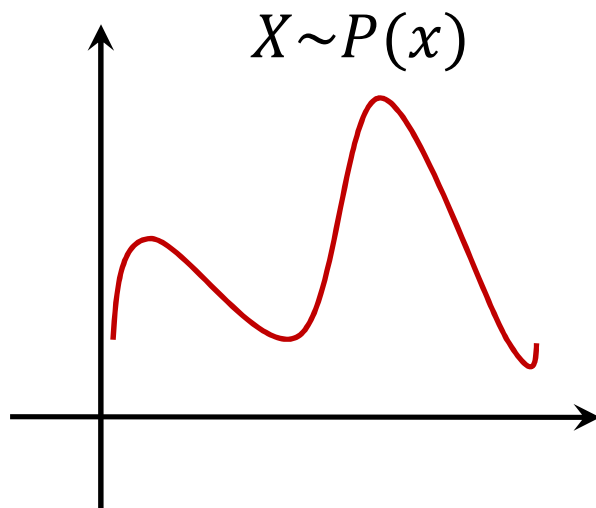
$$\begin{array}{l} \min_{\varphi} \quad \boxed{\text{Dist}(\varphi(X_S), \varphi(X_T))} + \lambda\Omega(\varphi) \\ \text{s.t.} \quad \text{constraints on } \varphi(X_S) \text{ and } \varphi(X_T) \end{array}$$

Maximum Mean Discrepancy (MMD)



# Representing Distributions in RKHS

Can we use a vector  $\mu_X$  to represent the distribution?



What about  $\mu_X \in \mathcal{H}$  (RKHS)?

Kernel trick can be applied!

$$\mu_X = (\mathbb{E}[X])$$



$$\mu_X = \begin{pmatrix} \mathbb{E}[X] \\ \mathbb{E}[X^2] \end{pmatrix}$$



$$\mu_X = \begin{pmatrix} \mathbb{E}[X] \\ \mathbb{E}[X^2] \\ \mathbb{E}[X^3] \end{pmatrix}$$



$$\mu_X = \begin{pmatrix} \mathbb{E}[X] \\ \mathbb{E}[X^2] \\ \mathbb{E}[X^3] \\ \dots \\ \dots \end{pmatrix}$$

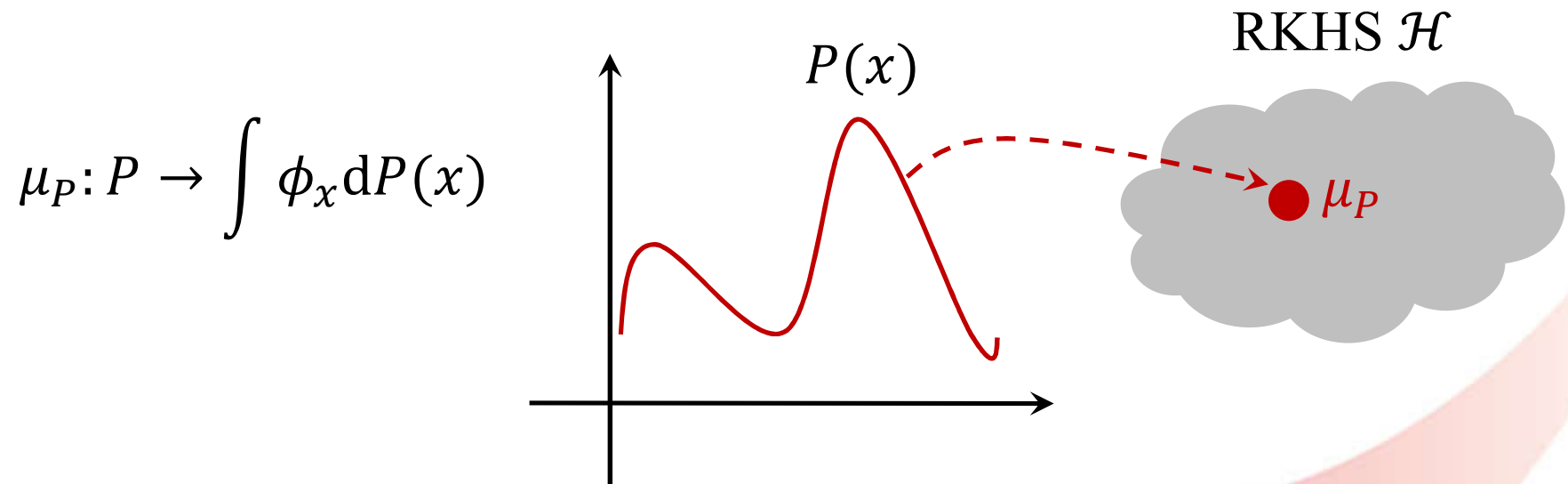


But, infinite,  
cannot be  
explicitly  
computed!

# Mean Map in RHKS

Suppose  $X \sim P(x)$ , and denote  $k(x, \cdot) = \phi_x$

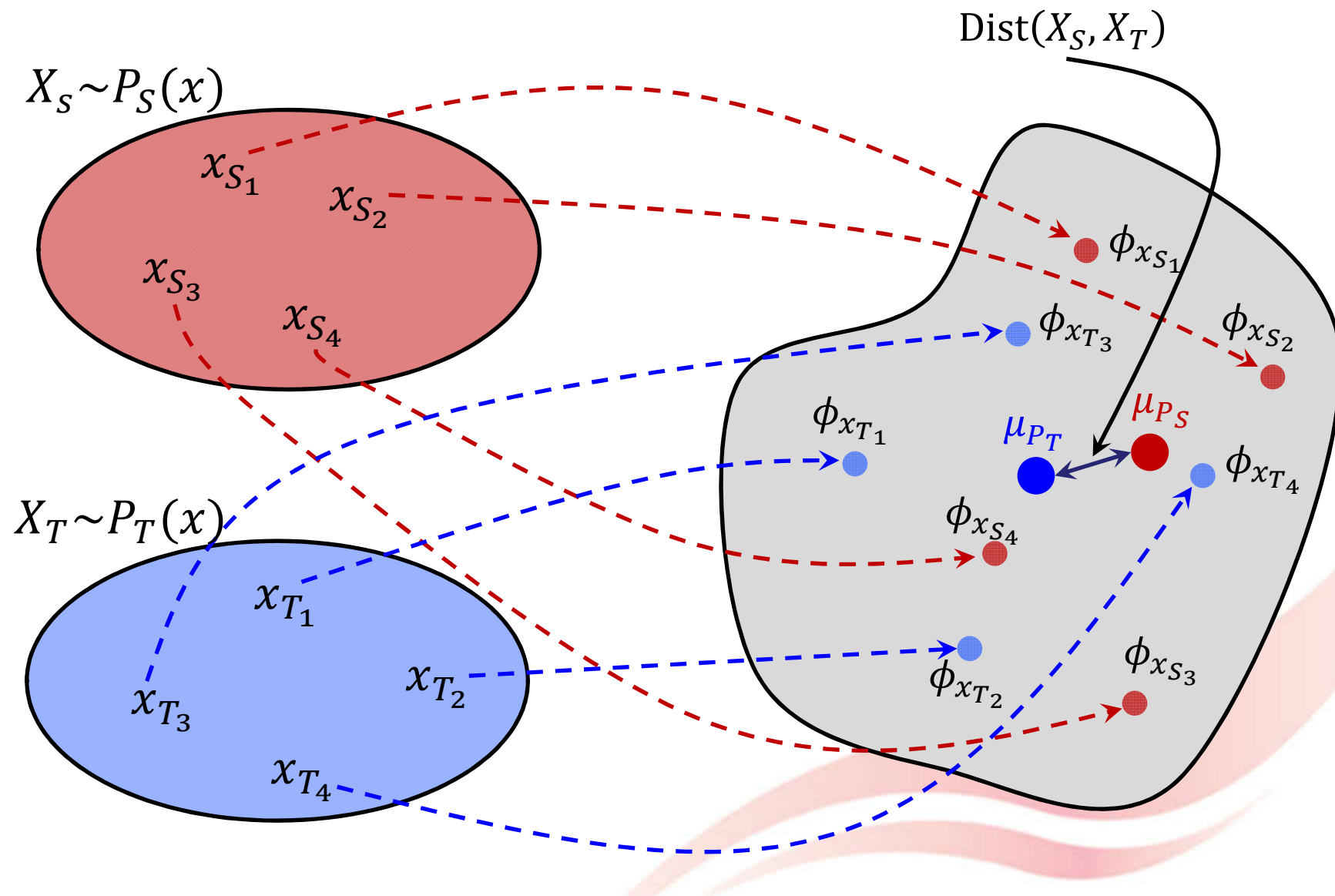
Mean map:  $\mu_P = \mathbb{E}_{x \sim P(x)}[\phi_x]$



Empirical mean map for  $\{x_i\}_{i=1}^n$ :  $\hat{\mu}_P = \frac{1}{n} \sum_{i=1}^n \phi_{x_i} = \frac{1}{n} \sum_{i=1}^n k(x_i, \cdot)$

[Berlinet and Thomas-Agnan 2004; Smola et al. ALT-07]

# Mean Map in RHKS (cont.)





# Distance Measure via MMD

$$\begin{aligned}\text{Dist}(\varphi(X_S), \varphi(X_T)) &= \left\| \mathbb{E}_{x \sim P_T(x)} [\phi(\varphi(x))] - \mathbb{E}_{x \sim P_S(x)} [\phi(\varphi(x))] \right\|_{\mathcal{H}} \\ &\approx \left\| \frac{1}{n_T} \sum_{i=1}^{n_T} \phi(\varphi(x_{T_i})) - \frac{1}{n_S} \sum_{i=1}^{n_S} \phi(\varphi(x_{S_i})) \right\|_{\mathcal{H}}\end{aligned}$$

Assume  $\psi = \phi \circ \varphi$  be a RKHS with kernel  $k(x_i, x_j) = \psi(x_i)^T \psi(x_j)$

$$\text{Dist}(\varphi(X_S), \varphi(X_T))^2 = \left\| \frac{1}{n_T} \sum_{i=1}^{n_T} \psi(x_{T_i}) - \frac{1}{n_S} \sum_{i=1}^{n_S} \psi(x_{S_i}) \right\|_{\mathcal{H}}^2 = \text{tr}(KL)$$

$$K = \begin{bmatrix} K_{S,S} & K_{S,T} \\ K_{T,S} & K_{T,T} \end{bmatrix} \quad L_{ij} = \begin{cases} \frac{1}{n_S^2} & x_i, x_j \in X_S \\ \frac{1}{n_T^2} & x_i, x_j \in X_T \\ -\frac{1}{n_S n_T} & \text{otherwise} \end{cases}$$

# Transfer Component Analysis (cont.)

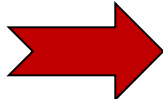
$$\begin{aligned} \min_{\varphi} \quad & \text{Dist}(\varphi(X_S), \varphi(X_T)) + \lambda\Omega(\varphi) \\ \text{s.t.} \quad & \text{constraints on } \varphi(X_S) \text{ and } \varphi(X_T) \end{aligned}$$



$$\begin{aligned} \min_{\varphi} \quad & \text{tr}(KL) + \lambda\Omega(\varphi) \\ \text{s.t.} \quad & \text{constraints on } \varphi(X_S) \text{ and } \varphi(X_T) \end{aligned}$$

- In general, the kernel function  $k(\varphi(x_i), \varphi(x_j))$  can be a highly nonlinear function of  $\varphi$  that is unknown
- A direct optimization of minimizing the quantity w.r.t.  $\varphi$  may get stuck in poor local minima

# Solution I [Pan *et al.*, AAAI-08]

Learning  $\varphi$  

Minimize the distance  
between domains

- 1) Learning  $K$
- 2) Low-dimensional reconstructions of  $X_S$  and  $X_T$  based on  $K$

1)  $\min_{K \succeq 0} \text{tr}(KL) + \lambda \text{tr}(K)$  ← Maximize data variance

s.t.  $K_{ii} + K_{jj} - 2K_{ij} = d_{ij}^2, \forall (i, j) \in \mathcal{N}$ , ← Preserve the local geometric structure  
 $K\mathbf{1} = 0.$

2) Perform PCA on  $K$

- It is a SDP problem, expensive!
- It is transductive, cannot generalize on unseen instances!
- PCA is post-processed on the learned kernel matrix, which may potentially discard useful information

# Solution II [Pan *et al.*, IJCAI-09, IEEE TNN-11]

Assume  $K$  be low-rank, then  $K = \bar{K}W W^T \bar{K}$  Known, given by user

$$W \in \mathbb{R}^{(n_S+n_T) \times m} \text{ and } m \ll n_S + n_T$$

Learning  $K$   Learning a low-rank matrix  $W$

Minimize distance  
between domains

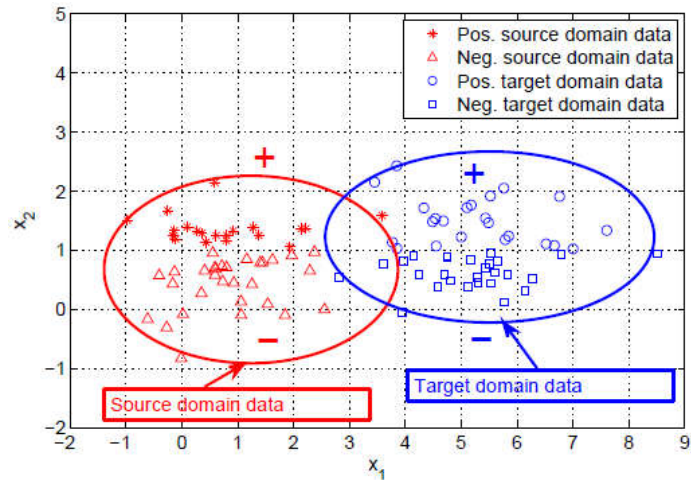
$$\begin{aligned} & \text{tr}(W^T \bar{K} L \bar{K} W) \\ \min_W & \text{tr}(\bar{K} W W^T \bar{K} L) + \lambda \text{tr}(W^T W) \quad \leftarrow \text{Regularization term on } W \\ \text{s.t.} & \quad W^T \bar{K} H \bar{K} W = I \quad \leftarrow \text{Maximize data variance} \end{aligned}$$

Closed form solution for  $W^*$ :

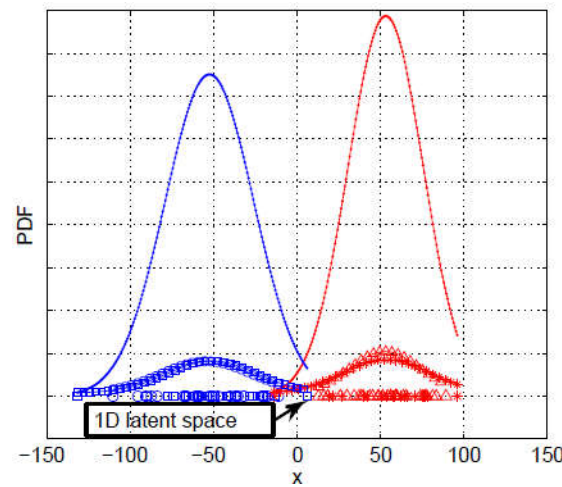
$m$  leading eigenvectors of  $(\bar{K} L \bar{K} + \lambda I)^{-1} \bar{K} H \bar{K}$

# Transfer Component Analysis (cont.)

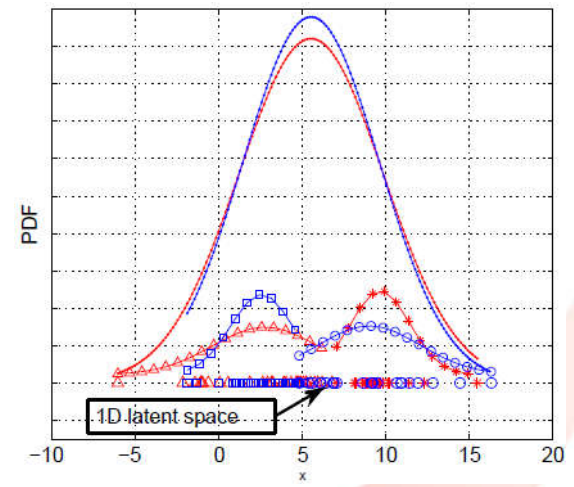
**An illustrative example**  
*Latent features learned by PCA and TCA*



Original feature space



PCA



TCA

# Future Direction

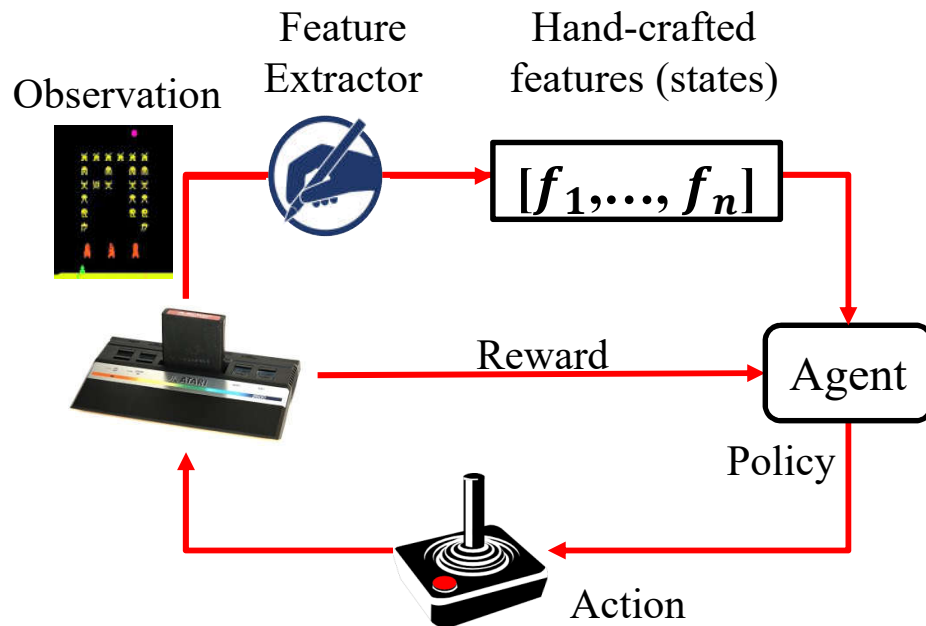
- Theoretical study beyond generalization error bound
  - Given a source domain and a target domain, determine whether transfer learning should be performed
  - For a specific transfer learning method, given a source and a target domain, determine whether the method should be used for knowledge transfer



# Future Direction (cont.)

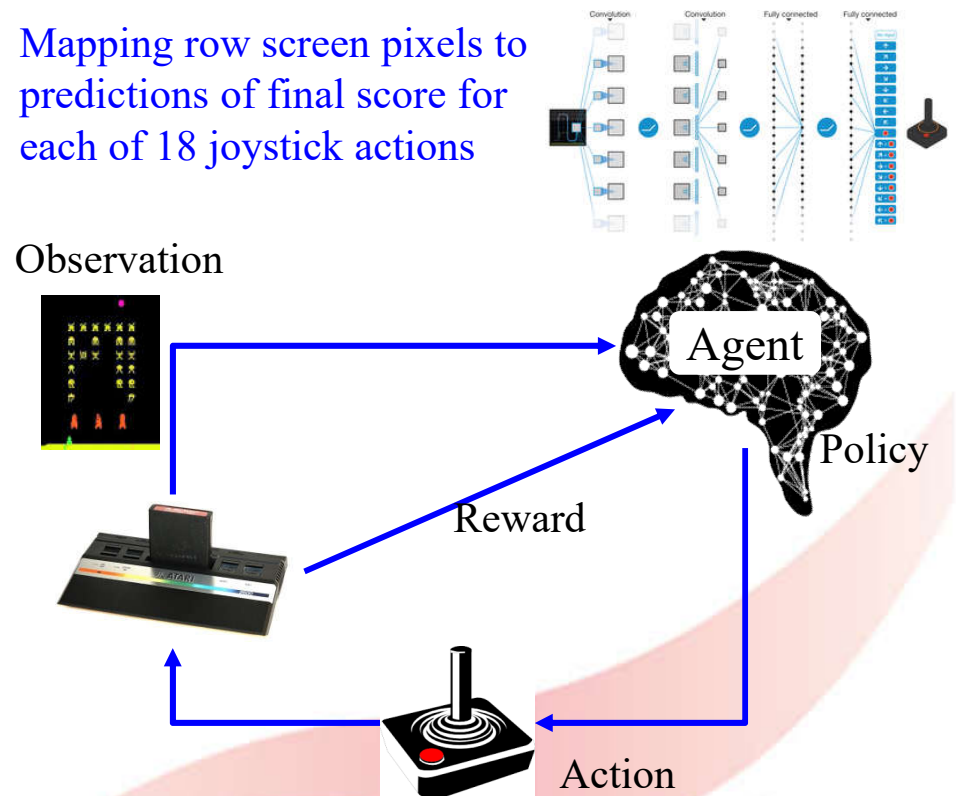
- Transfer learning for deep reinforcement learning

Mapping hand-crafted features (states) to final score for each of 18 joystick actions



**Traditional RL**

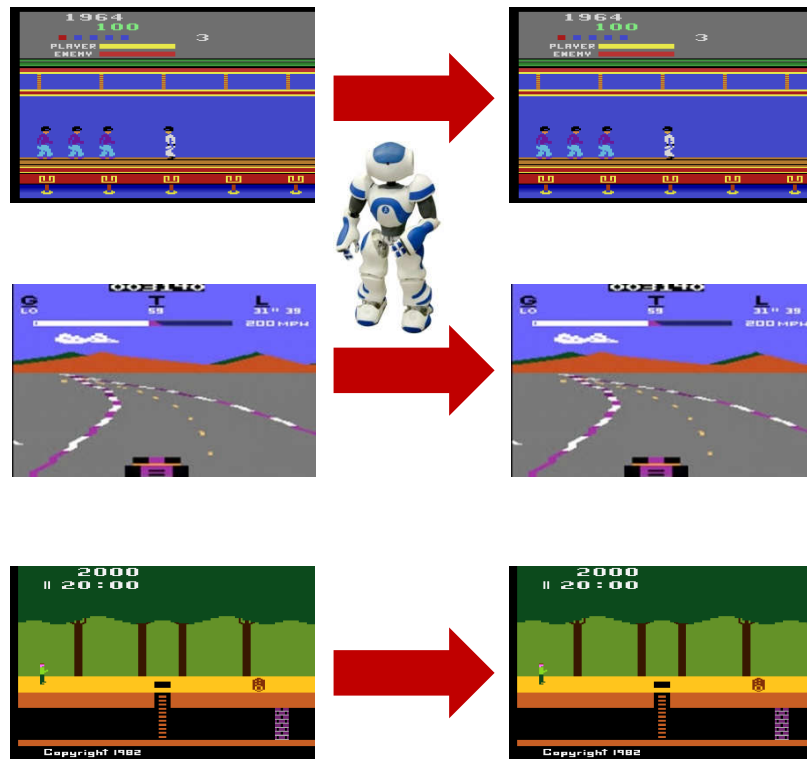
Mapping raw screen pixels to predictions of final score for each of 18 joystick actions



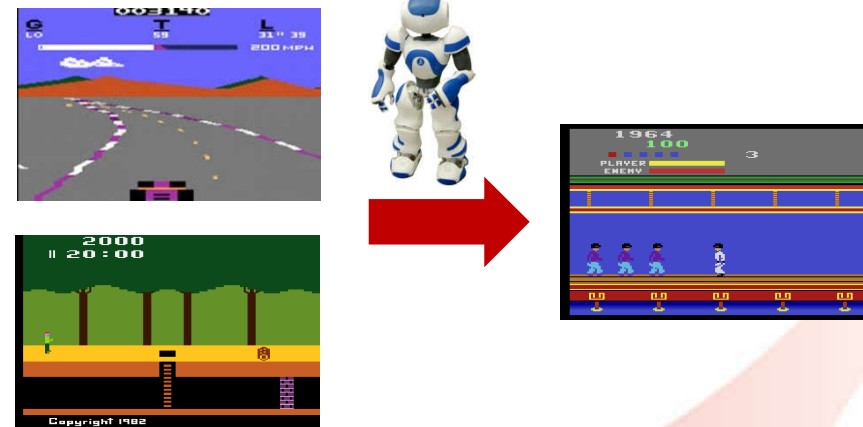
**Deep RL** [Google DeepMind 2015]

# Future Direction (cont.)

- Transfer learning for deep reinforcement learning




**Deep RL**



**Transfer Learning for Deep RL**



# Reference

- Pan and Yang, A Survey on Transfer Learning, IEEE TKDE 2010
  - Pan, Transfer learning, Data Classification: Algorithms and Applications (Chapter 21), 2014
  - Taylor and Stone, Transfer Learning for Reinforcement Learning Domains: A Survey, JMLR 2009
  - Pan, Kwok, and Yang, Transfer Learning via Dimensionality Reduction, AAAI 2008
  - Pan, Tsang, Kwok, and Yang, Domain Adaptation via Transfer Component Analysis, IJCAI 2009, IEEE TNN 2011
  - Alex Smola, Arthur Gretton, Le Song, and Bernhard Scholkopf, A Hilbert Space Embedding for Distributions, ALT 2007
  - A. Berlinet and C. Thomas-Agnan, Reproducing Kernel Hilbert Spaces in Probability and Statistics. Kluwer Academic Publisher 2004
- 

**Thank You!**

