

# Reshaping Datasets for Unsupervised Domain Adaptation

**Boqing Gong**



*Joint work with Kristen Grauman and Fei Sha*

# Data-centric era



# Experiments, observations, and simulations in science



# Internet of things

## Sensors everywhere



# 140 billion images, *12M hourly*



# 300 hour new video every minute




# 200B tweets yearly, 500M daily



# Great sources of discovery and knowledge

 predicted flu outbreak two weeks before CDC, and now they collaborate.

 correctly predicted 2012 presidential election.

 GPS provides real-time traffic information.

 matches cancer patients to cord-blood donors in real-time.

# Challenges

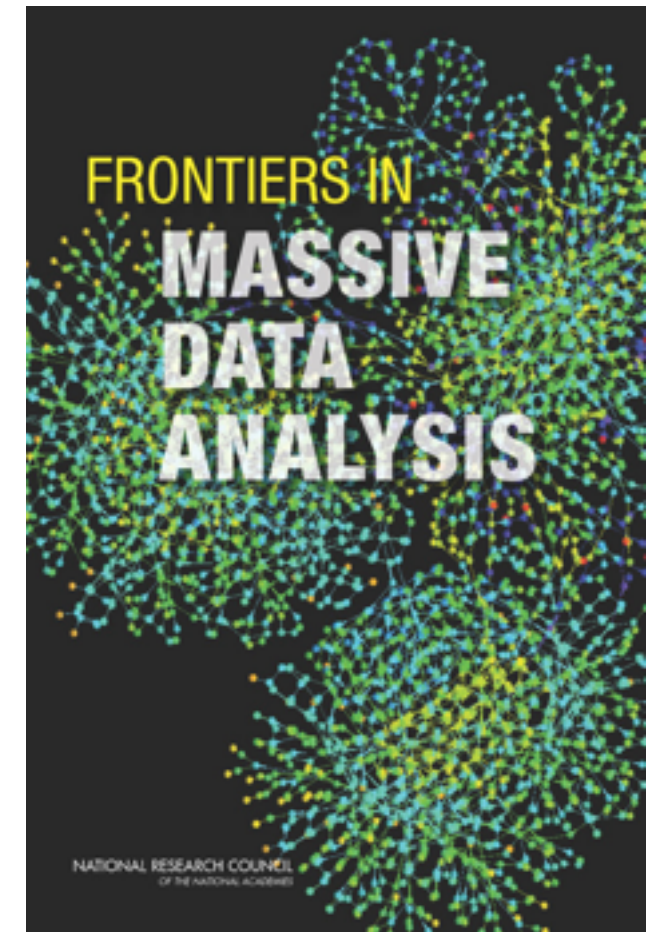
Dealing with highly distributed data

Coping with sampling biases and heterogeneity

Exploiting parallel and distributed architectures

Data visualization, integration, validation, security, sharing, etc.

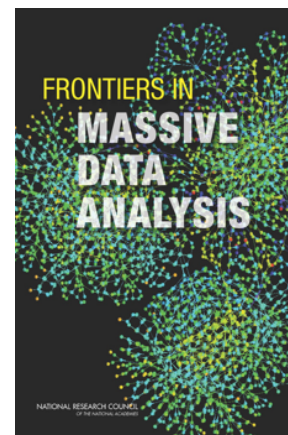
... ..



National  
Academies Report

# Sampling bias & heterogeneity

“(training) Data may have been collected according to a certain criterion ..., but (testing) the inferences and decisions may refer to a different sampling criterion.”



National  
Academies Report



# Self-driving car: a case study



# Self-driving car: a case study

## Pedestrian detection and avoidance system



Sampling bias → Performance significantly degrades [Dollár et al.'09]

# The perils of mismatched domains

**Cause:** standard assumption in machine learning

Same underlying distribution for training and testing



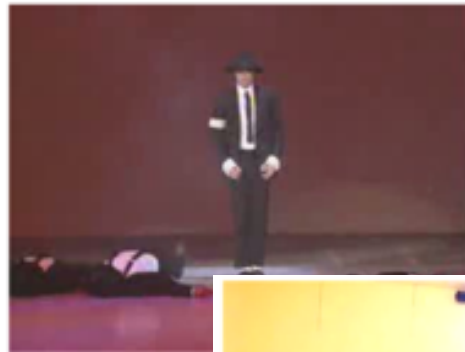
# This is a realistic obstacle for autonomous systems

Systems often deployed to new environment, not  
lab reproducible

Expensive to collect training data from each type  
of target environment

Systems naturally degrade; environment  
dynamically evolves

# Mismatches are common to many areas



## The New York Times

NEW YORK, SATURDAY, FEBRUARY 25, 1995

### Robert Maffery, 79, Junk Artist Behind the Growth of Sculpture

By ROBERTA SHERRY  
Robert Maffery, a junk artist who contributed to the flowering of American postmodernism through his art and work, died last week at the age of 79. He was also a pioneer in the use of plastic and other materials in his work.

The artist was born in New York City, and his work was characterized by a sense of humor and a love of the absurd. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

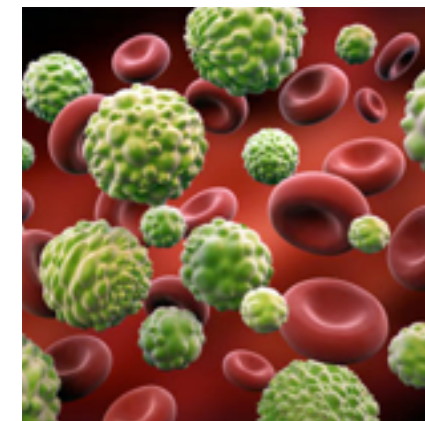
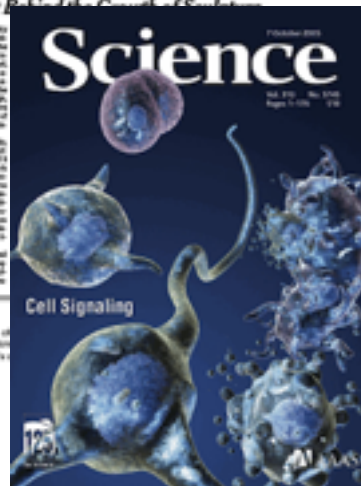
He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.

He was a pioneer in the use of plastic and other materials in his work. He was a pioneer in the use of plastic and other materials in his work.



Biology:  
different  
subjects

Click to LOOK INSIDE!

### Statistical Learning Theory

Most Helpful Customer Reviews

30 of 31 people found the following review helpful

★★★★★ **statistical learning based on the VC class** January 23, 2008

By Michael R. Chernick

Format: Hardcover

Vapnik and Chernovenkis extended the Glivenko-Cantelli Theorem in their book of nonparametric statistical inference based on approximating functions and in an earlier book published by Springer-Verlag he develops the basics of machine learning. In an earlier book published by Springer-Verlag he develops the basics of machine learning. In an earlier book published by Springer-Verlag he develops the basics of machine learning.

Click to open book

### A Step Closer

Most Helpful Customer Reviews

7,580 of 7,707 people found the following review helpful

★★★★★ **A Step Closer** March 15, 2011

By Craig Whisenhunt

Color Name: Black | Item Shape: WiFi + Verizon 3G | Size Name: 16GB

For anyone out there who is considering whether or not to make the iPad 1 and the iPad 2 check out my review of the first generation iPad of people commenting (both positively and negatively) over the past

# Abstract form: *unsupervised* domain adaptation (DA)

## Setup

**Source** domain (with labeled data)

$$D_S = \{(x_m, y_m)\}_{m=1}^M \sim P_S(X, Y)$$

**Target** domain (no labels for training)

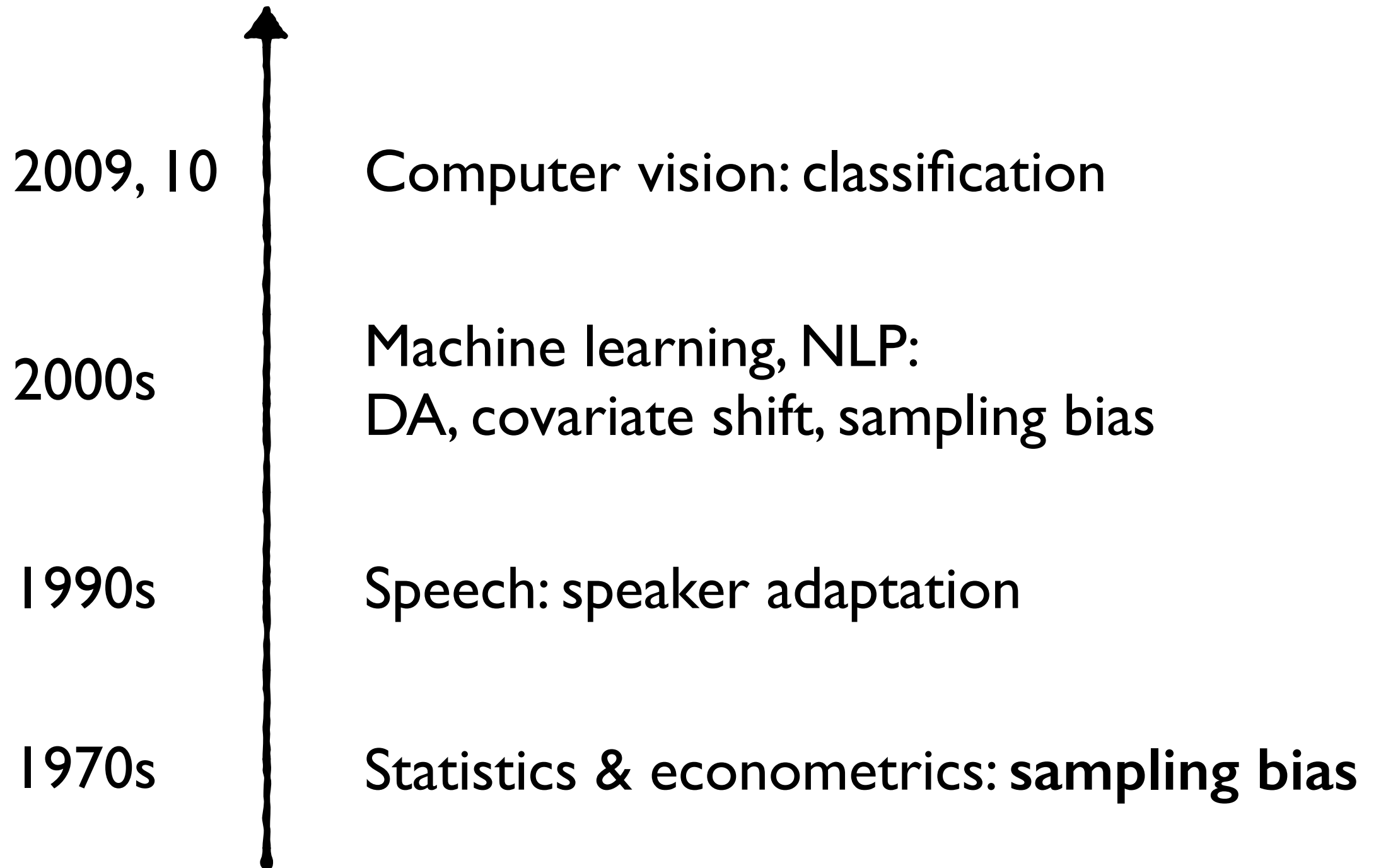
$$D_T = \{(x_n, ?)\}_{n=1}^N \sim P_T(X, Y)$$

**Different distributions**

## Objective

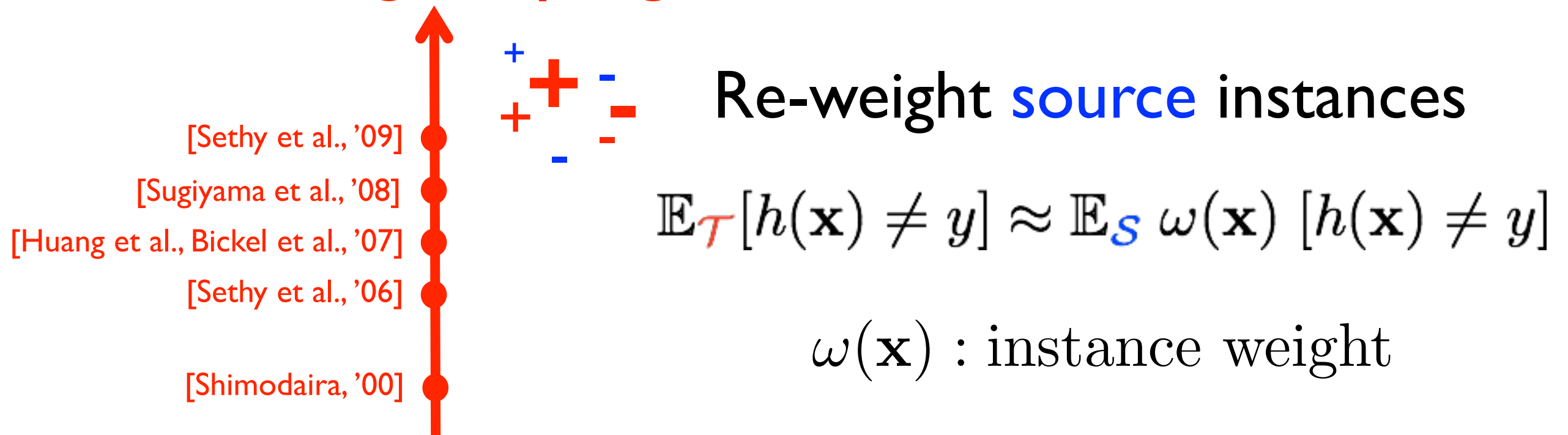
Learn models to work well on **target**

# Background on DA



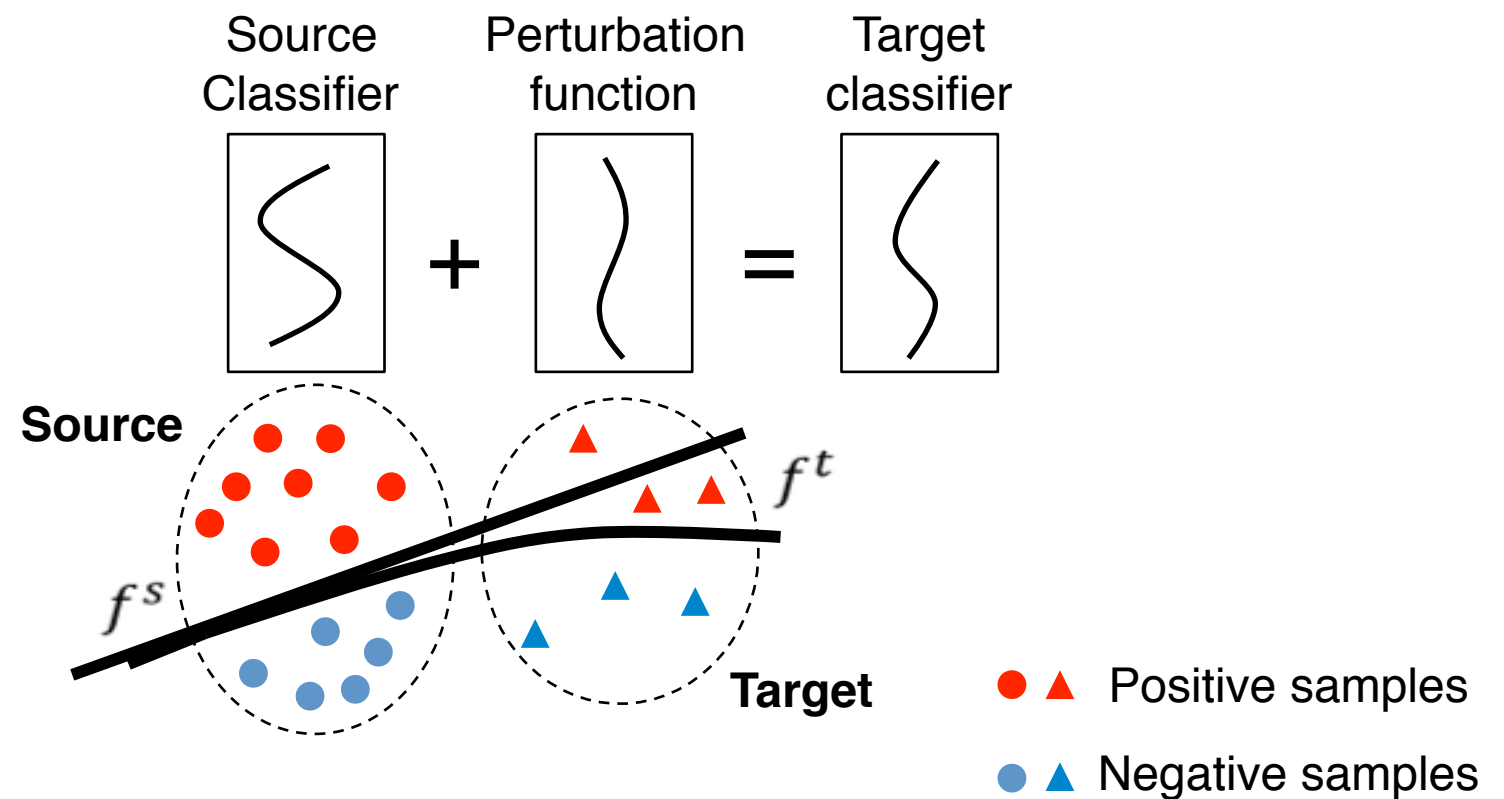
# Background - brief review

## Correcting *sampling* bias





# Background - brief review



[Evgeniou and Pontil, '05]

[Duan et al., '09]

[Duan et al., Daumé III et al., Saenko et al., '10]

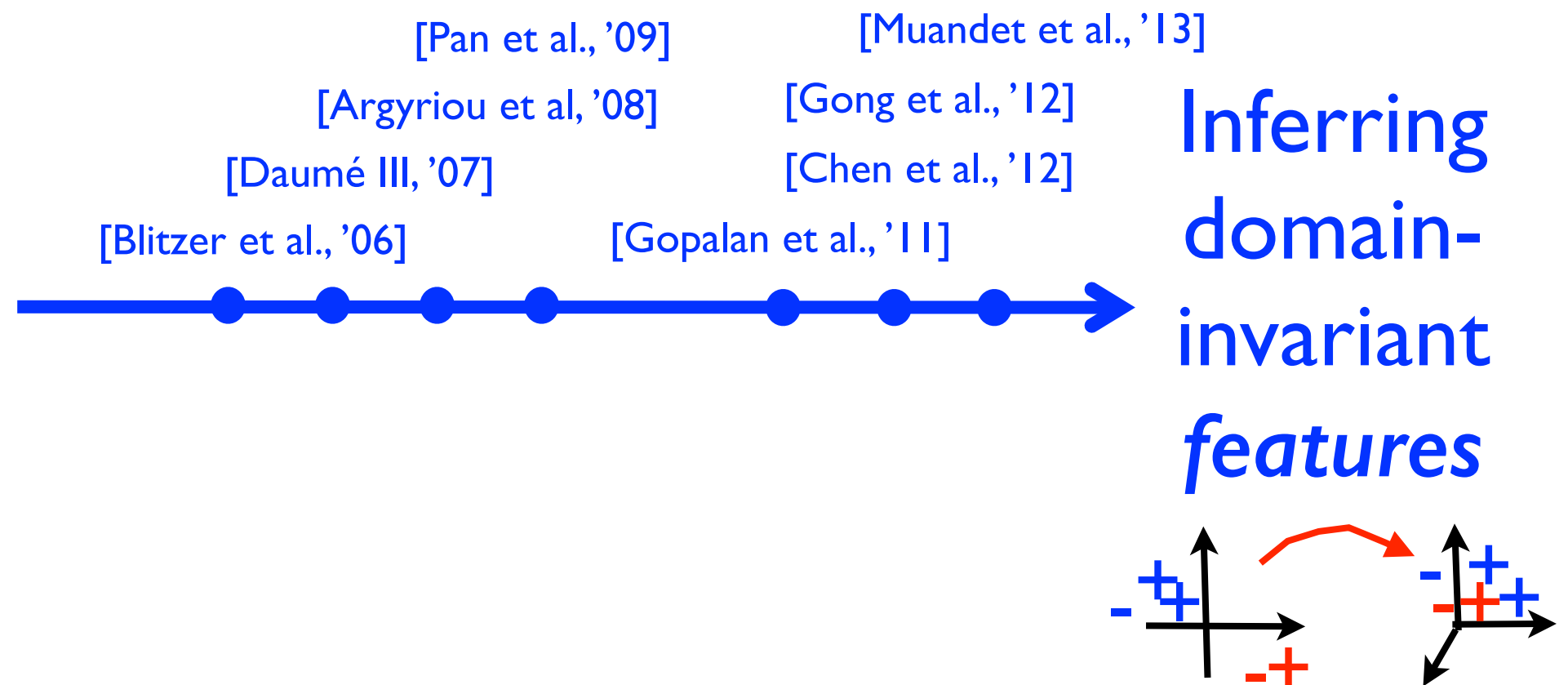
[Kulis et al., Chen et al., '11]

Adjusting mismatched *models*

# Background - brief review

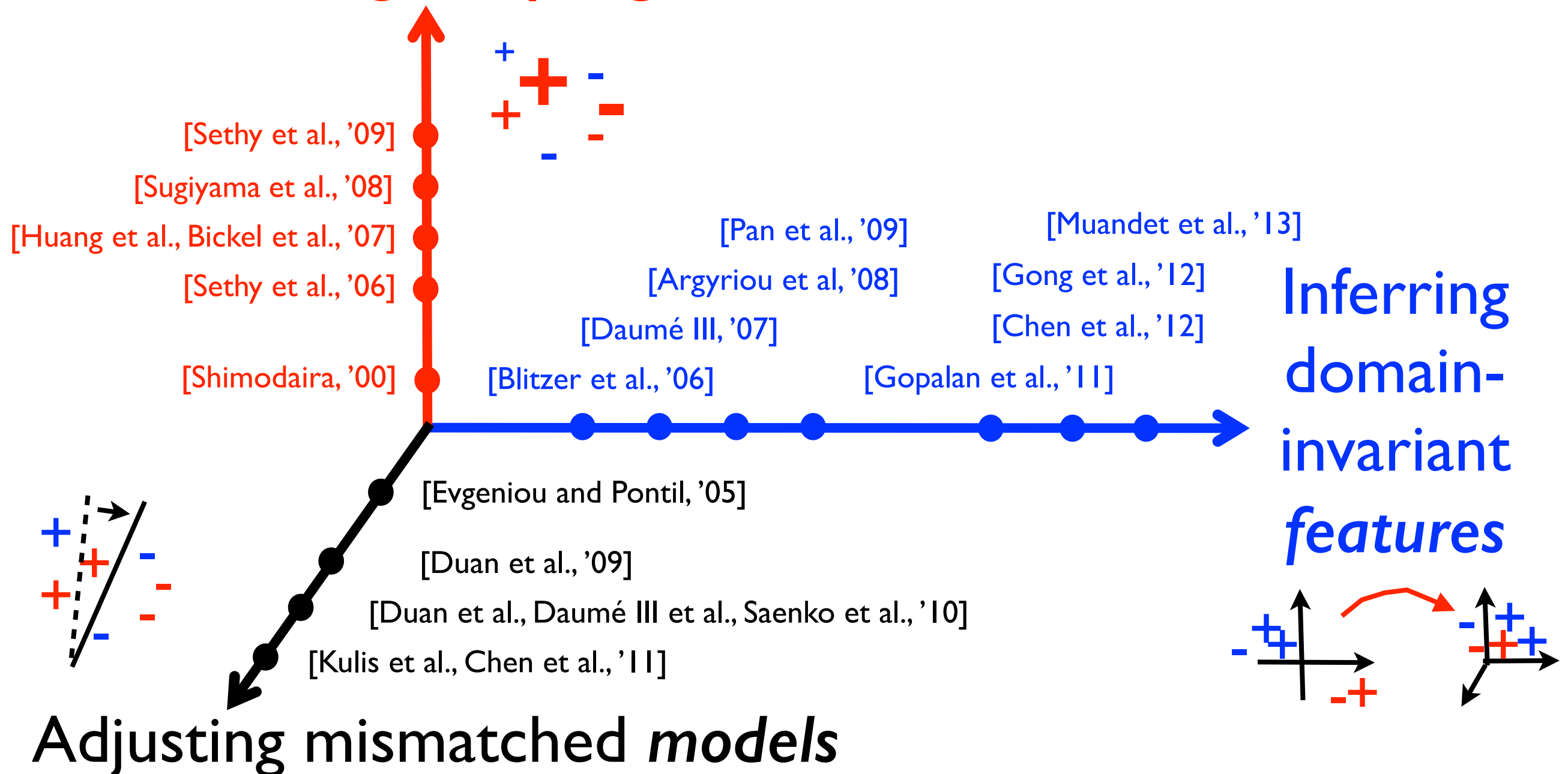
$$\mathbf{x} \mapsto \mathbf{z}, \quad \text{s.t.}$$

$$P_{\mathcal{S}}(z, y) \approx P_{\mathcal{T}}(z, y)$$

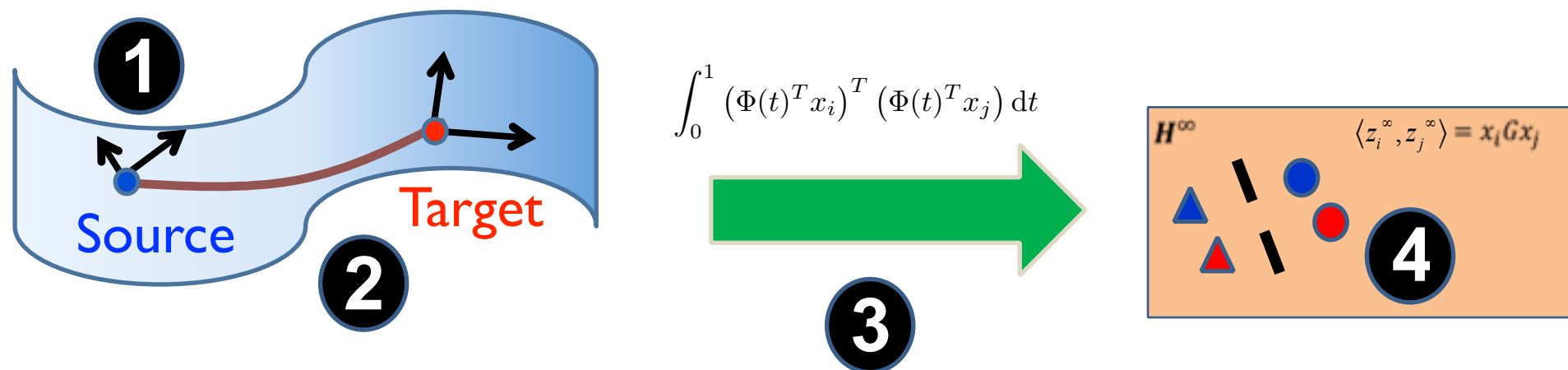


# Background - quick review

## Correcting *sampling* bias



# GFK: inferring a domain-invariant feature space



1. Exploit subspace structure in data
2. Model domain shift with geodesic flow
3. **Derive a domain-invariant kernel**
4. Classify target data in the kernel space

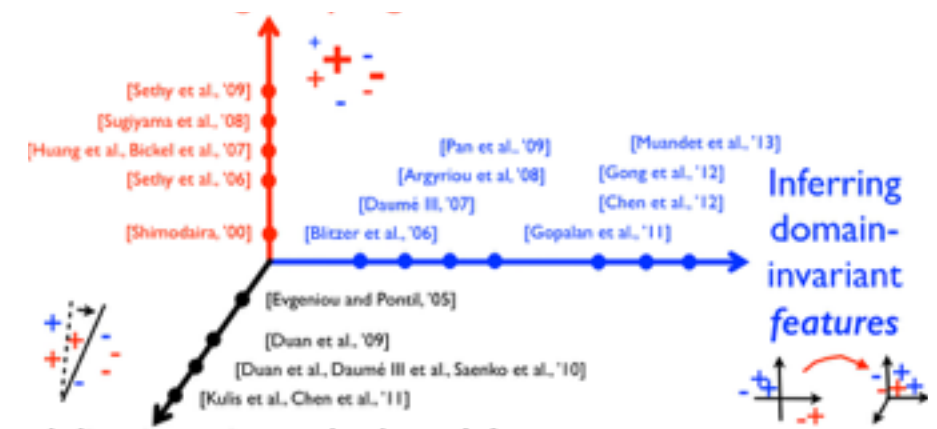
[Gong *et al.*, CVPR'12]

# Key to domain adaptation

*“to reduce **source-target** discrepancy”*



# Snags in previous methods



## Forced adaptation

Attempting to adapt all **source** instances, including “hard” ones

## Implicit discrimination

Learning discrimination biased to **source**, rather than optimized w.r.t. **target**

# Key to domain adaptation

“to reduce **source-target** domain *discrepancy*”

What is a **source** domain?

Is it always fixed?

Can we *reshape* it?

# What constitutes a domain?

In speech and NLP:

Speakers

Languages

Article topics

...other factors

In computer vision:

**Factors?**



Many factors  
overlap & interact

**Pose**

Lighting

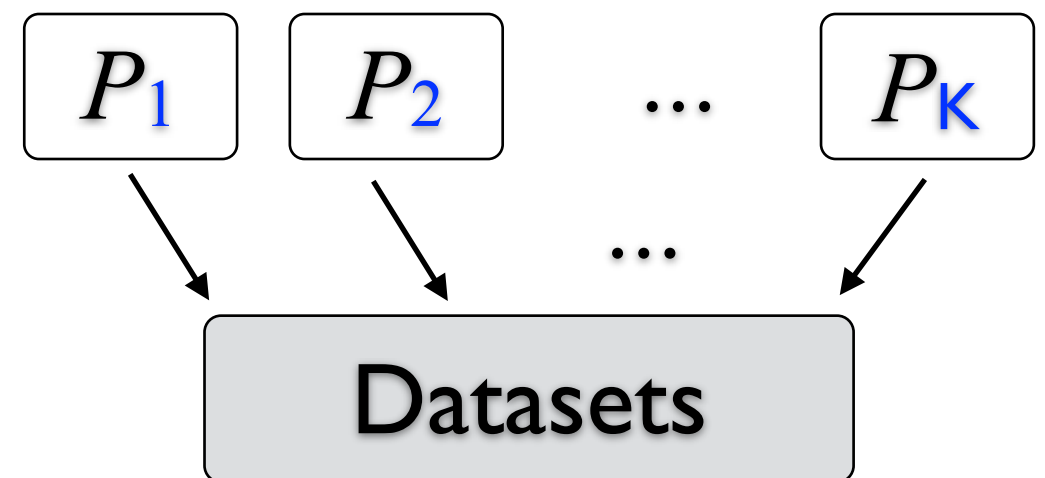
Fore/  
Background

Occlusion

# Some questions revolving around “domain”



Adapt-abilities  
of different domains  
[Gong et al., IJCV'14, CVPR'12]



What is a domain?  
Reshaping data according to  
domains from which they come?  
[Gong et al., NIPS'13]

# Our key insights

Forced adaptation from a prefixed source domain

→ Select the best instances for adaptation

Implicit discrimination

→ Approximate discriminative loss on **target**



# Selecting most adaptable **source** instances

*Landmarks* are labeled **source** instances distributed similarly to the **target** domain.



Source



Target

[Gong et al., ICML'13]

# Selecting most adaptable source instances

*Landmarks* are labeled source instances distributed similarly to the target domain.

Identifying landmarks:

$$\min_{\text{landmarks}} P_{\mathcal{L}}(\text{landmarks}) \approx P_{\mathcal{T}}(\text{target})$$
$$d(P_{\mathcal{L}}, P_{\mathcal{T}})$$



[Gong et al., ICML'13]



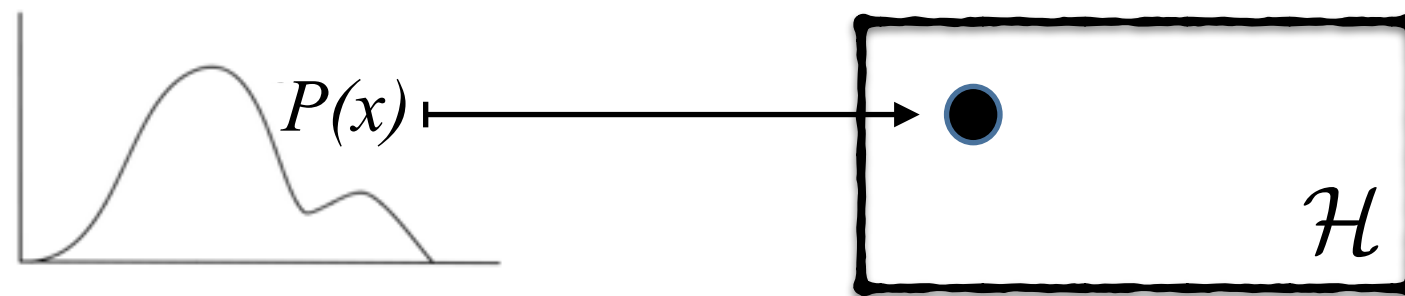
Source



Target

# Kernel embedding of distributions

$$\mu[P] \triangleq \mathbb{E}_x[\phi(x)]$$



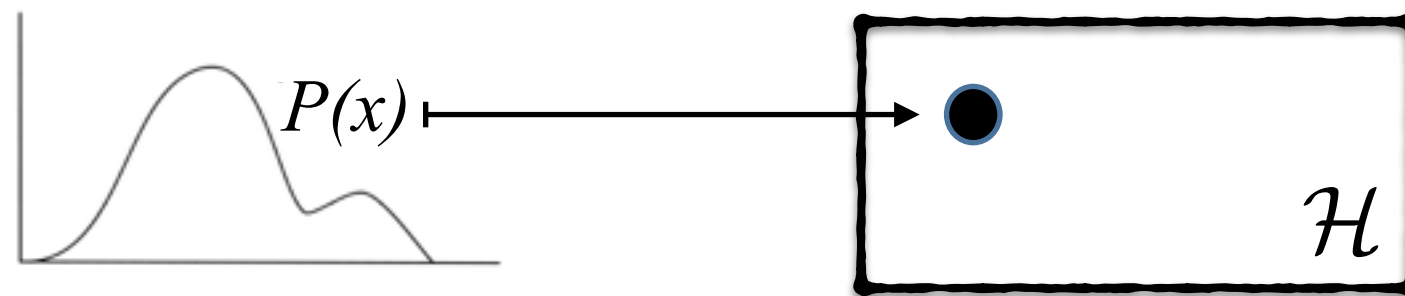
$\mu$  maps distribution  $P$  to Reproducing Kernel Hilbert Space

$\mu$  is injective if  $\phi(\cdot)$  is characteristic

[Müller'97, Gretton et al.'07, Sriperumbudur et al.'10]

# Kernel embedding of distributions

$$\mu[P] \triangleq \mathbb{E}_x[\phi(x)]$$



Empirical kernel embedding:

$$\hat{\mu}[P] = \frac{1}{n} \sum_{i=1}^n \phi(x_i), \quad x_i \sim P$$

# Identifying landmarks by matching kernel embeddings

## Integer programming

$$\min_{\{\alpha_m\}} \left\| \frac{1}{\sum_i \alpha_i} \sum_{m=1}^M \alpha_m \phi(x_m) - \frac{1}{N} \sum_{n=1}^N \phi(x_n) \right\|_{\mathcal{H}}^2$$

where

$$\alpha_m = \begin{cases} 1 & \text{if } x_m \text{ is a landmark wrt target} \\ 0 & \text{else} \end{cases}$$

$$m = 1, 2, \dots, M$$



# Solving by relaxation

## Convex relaxation

$$\min_{\{\alpha_m\}} \left\| \frac{1}{\sum_i \alpha_i} \sum_{m=1}^M \alpha_m \phi(x_m) - \frac{1}{N} \sum_{n=1}^N \phi(x_n) \right\|_{\mathcal{H}}^2$$

$$\beta_m = \frac{\alpha_m}{\sum_i \alpha_i} \rightarrow \text{Quadratic programming}$$

$$\min_{\beta} \quad \beta^T K^s \beta - \frac{2}{N} \beta^T K^{st} \mathbf{1}$$

# How to choose the kernel functions?

$$\min_{\beta} \quad \beta^T K^s \beta - \frac{2}{N} \beta^T K^{st} \mathbf{1}$$

## Gaussian kernels

Plus: universal (characteristic)

Minus: how to choose the bandwidth?

## Our solution: bandwidth---granularity

Examining distributions at multiple granularities

Multiple bandwidths, multiple sets of landmarks

# Other details

Class balance constraint

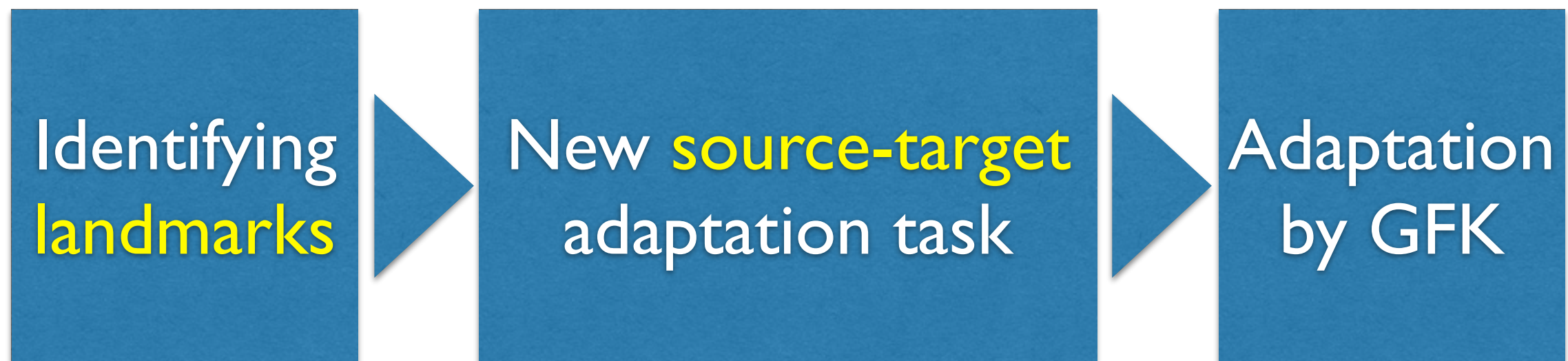
Recovering  $\alpha_m^*$  from  $\beta_m^* (= \frac{\alpha_m}{\sum_i \alpha_i})$

(See [Gong et al., ICML'13, IJCV'14] for details)

# What do landmarks look like?



# Landmark based domain adaptation



# Experimental study

Four vision datasets/domains on visual **object recognition**

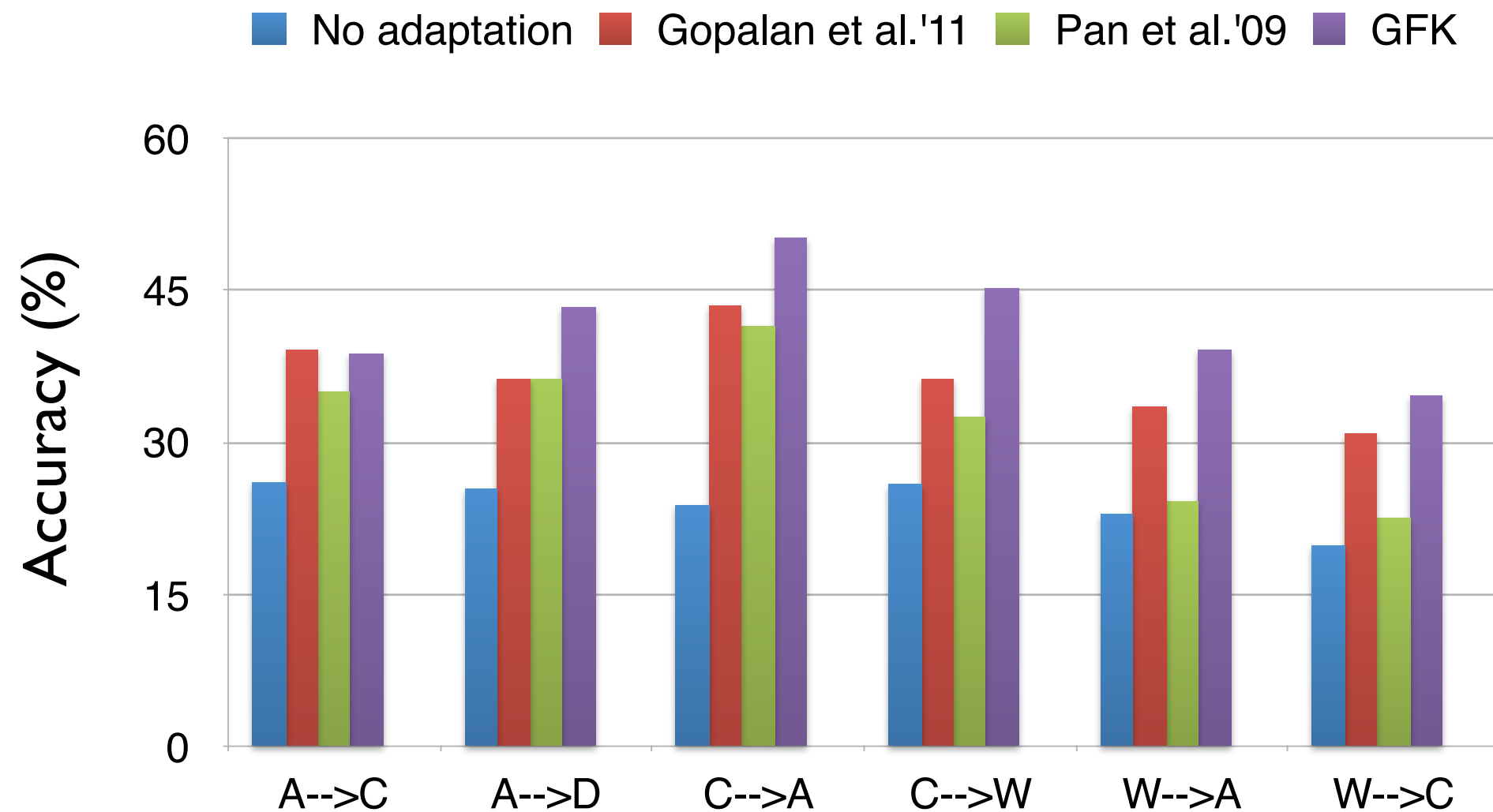
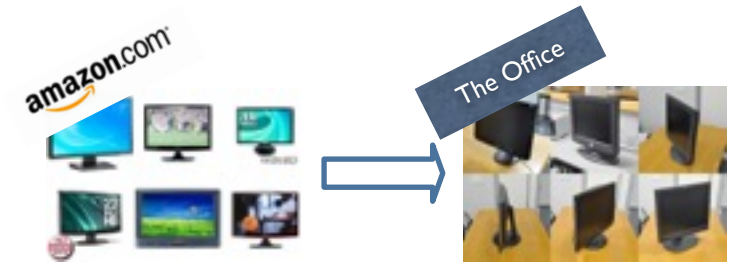
[Griffin et al. '07, Saenko et al. 10']

Four types of product reviews on **sentiment analysis**

Books, DVD, electronics, kitchen appliances [Biltzer et al. '07]

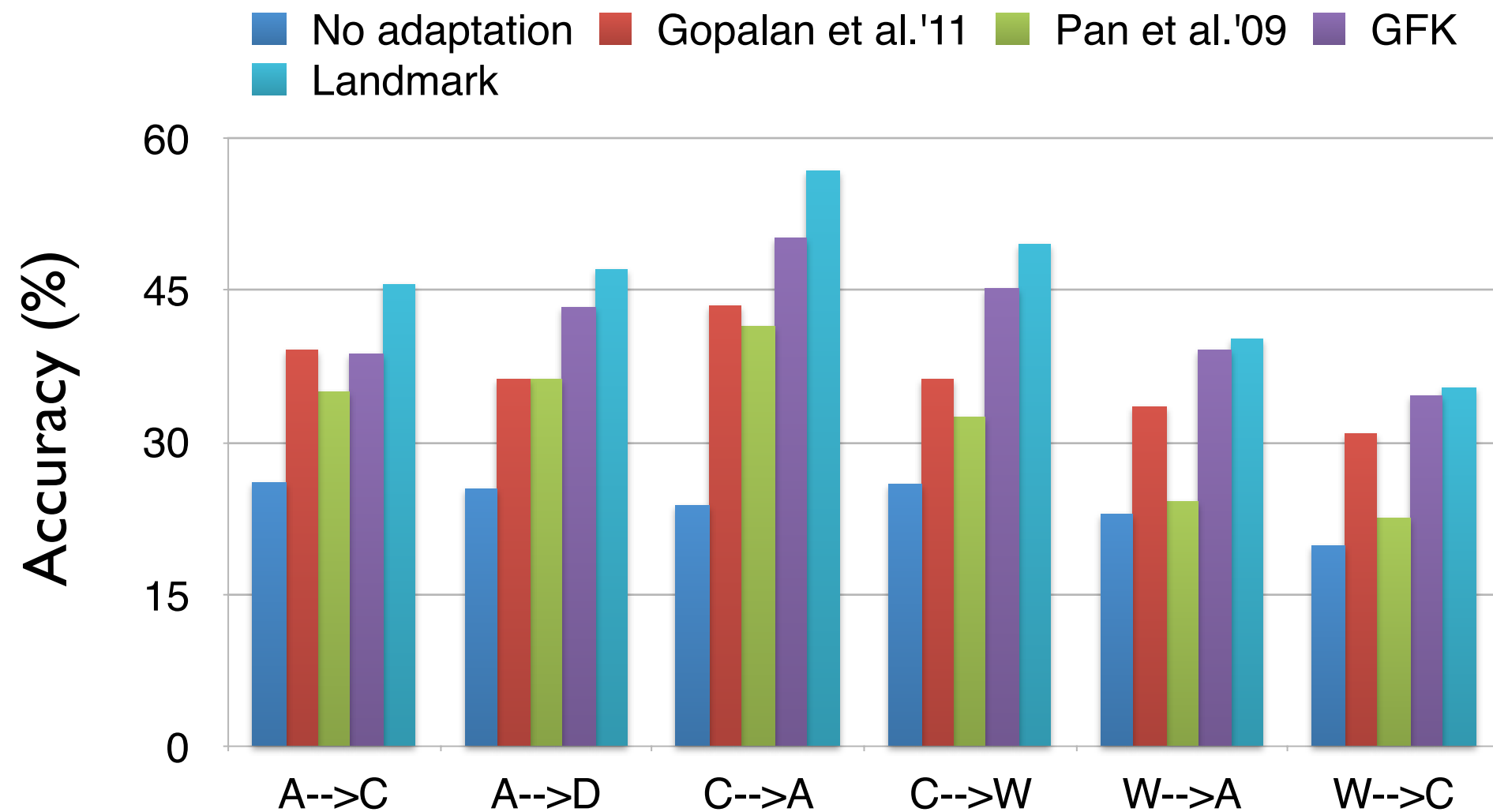
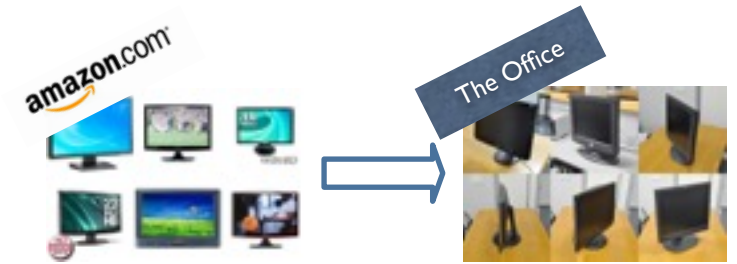


# Comparison results: object recognition



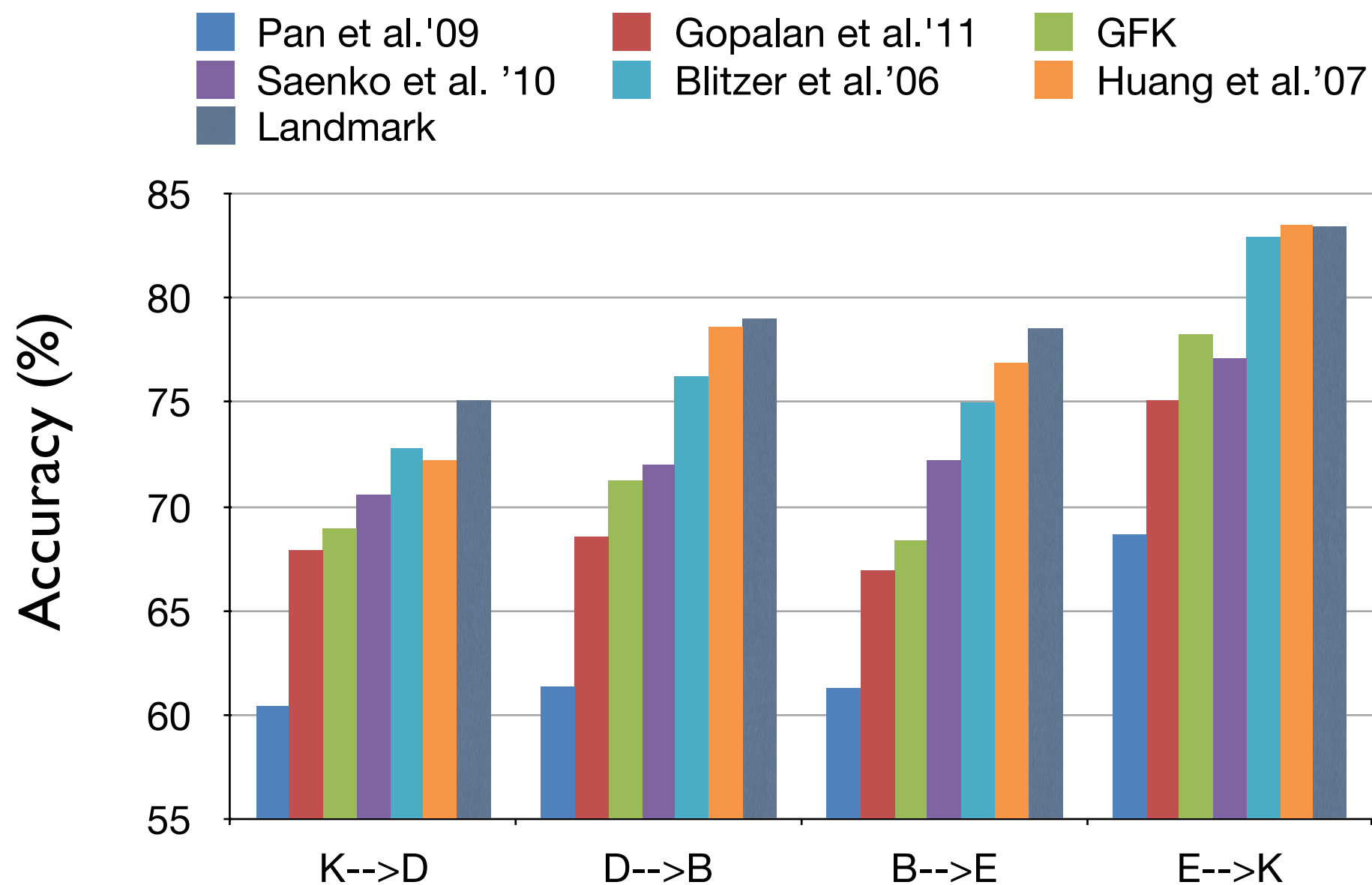


# Comparison results: object recognition





# Comparison results: sentiment analysis



# Summary - Landmarks



Landmarks

[Gong et al., ICML'13]

- Labeled **source** instances, distributed similarly to **target**
- Better approximation of discriminative loss of target
- Automatically identifying landmarks
- Benefiting other adaptation methods

# Key to domain adaptation

“to reduce **source-target** domain *discrepancy*”

What is a **source** domain?

**Landmarks:** reshaped **target**-oriented **source**

What if no a priori knowledge about **target**?



# What constitutes a domain?



# What constitutes a domain?



**Domain I**



**Domain II**

# Two axiomatic properties for latent domains

## I. Maximum **distinctiveness**:

Identifying distinct domains maximally different in **distribution** from each other

## II. Maximum **learnability**

Being able to derive strong discriminative models from the identified domains

[Gong *et al.*, NIPS'13]

# I. Maximum distinctiveness

Domains maximally different in distribution from each other

$$\max_{\{z_{mk}\}} \sum_{k \neq k'} \hat{d}(P_k, P_{k'}; \{z_{mk}\})$$

$$z_{mk} = \begin{cases} 1 & \text{if } x_m \in \text{the } k\text{-th domain} \\ 0 & \text{else} \end{cases}$$

$$m = 1, 2, \dots, M, \quad k = 1, 2, \dots, K$$



## II. Maximum learnability

Able to learn strong classifiers from domains

Within-domain cross-validation

$$\text{Accuracy}(\mathbf{K}) = \sum_{k=1}^{\mathbf{K}} \frac{M_k}{M} \text{Accuracy}_k$$

- Determining the number of domains  $\mathbf{K}$



# Hard to manually define discrete domains



# Our “reshaped” domains

*Adapting from discovered domains > from datasets*



**Domain I**

**Domain II**

# Summary - latent domains



Many factors  
overlap & interact

What constitutes  
domains?

[Gong et al., NIPS'13]

- *Dataset  $\neq$  domain*
- *Suboptimal to use DA methods for cross-dataset problem*
- *Discovering latent domains:*
  - *maximum distinctiveness*
  - *maximum learnability*

# Key to domain adaptation

*“to reduce **source-target** discrepancy”*

What is a **source** domain?

**Landmarks:** reshaped **target**-oriented **source**

Discovering latent domains without **target** *a priori*

*“to define domains / to reshape data well”*

# *Thanks!*

*“to reduce **source-target** discrepancy”*

What is a **source** domain?

**Landmarks:** reshaped **target**-oriented **source**

Discovering latent domains without **target** *a priori*

*“to define domains / to reshape data well”*