

Adaptation for Objects and Attributes

Kristen Grauman

Department of Computer Science

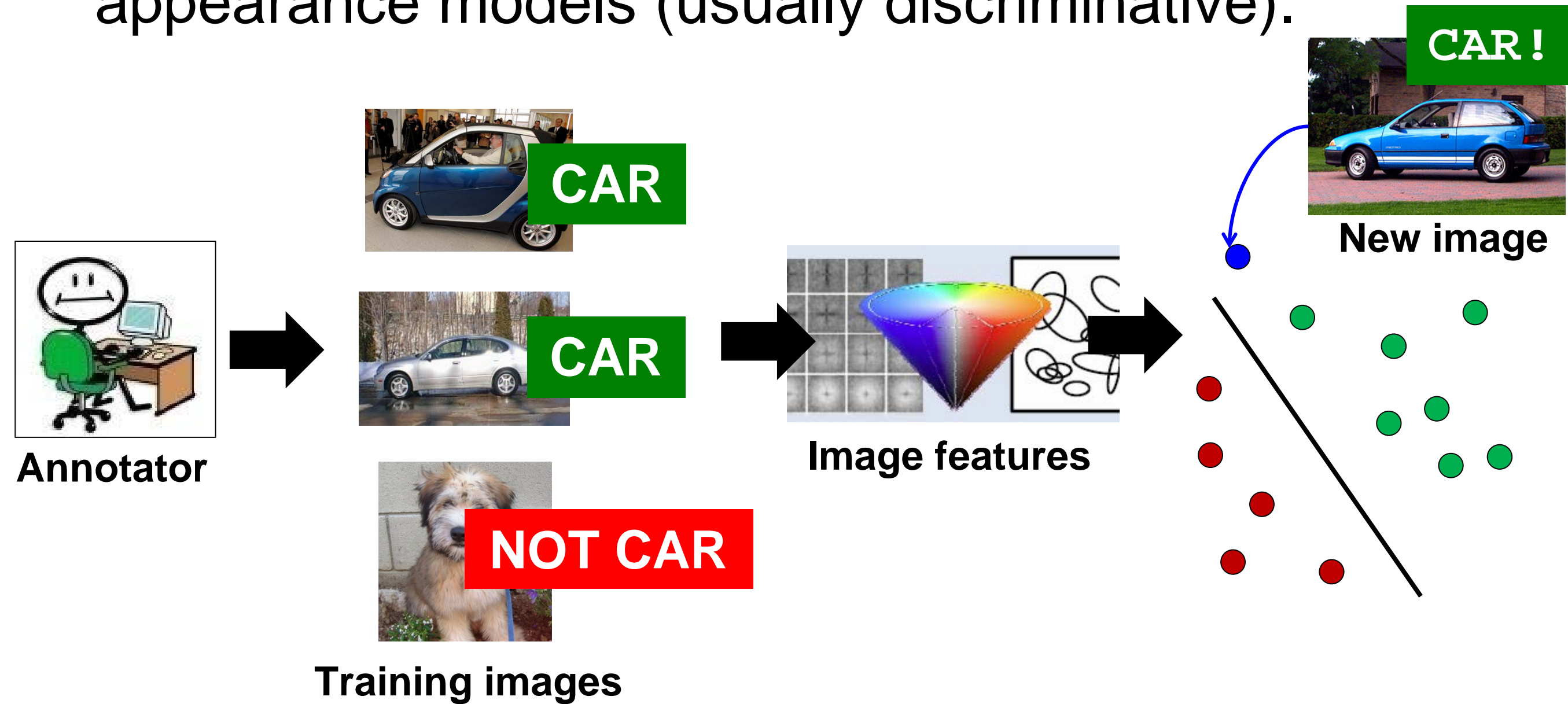
University of Texas at Austin

With Adriana Kovashka (UT Austin),
Boqing Gong (USC), and Fei Sha (USC)



Learning-based visual recognition

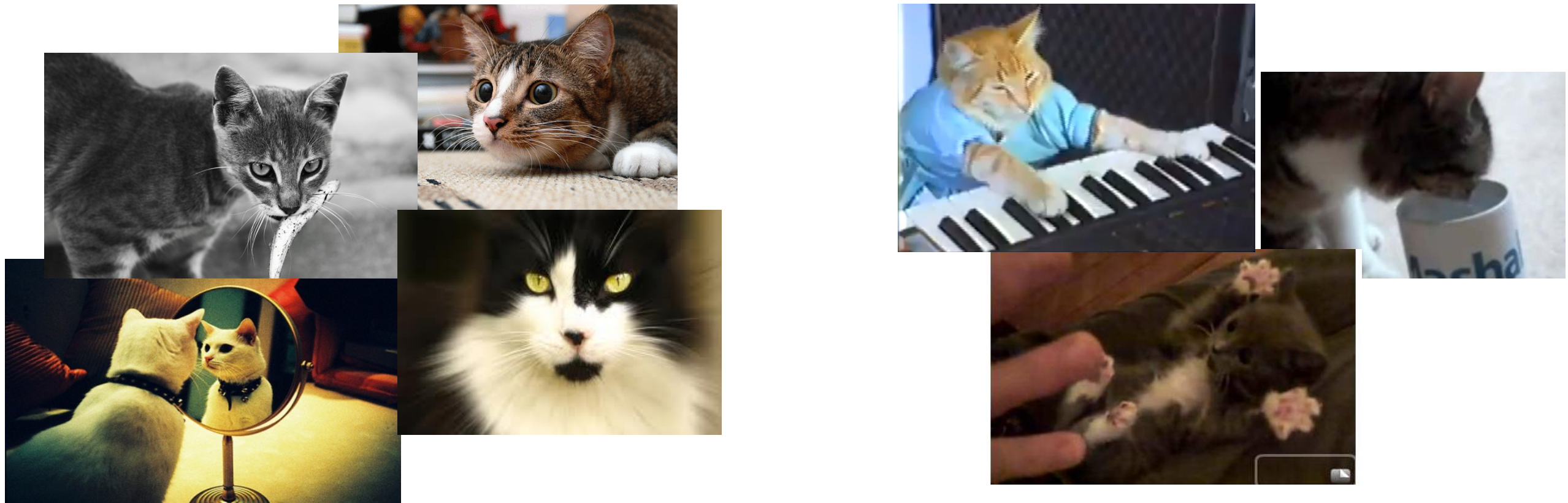
Last 10+ years: impressive strides by *learning* appearance models (usually discriminative).



Typical assumptions

1. Test set will look like the training set.
2. Human labelers “see” the same thing.

Mismatched domains

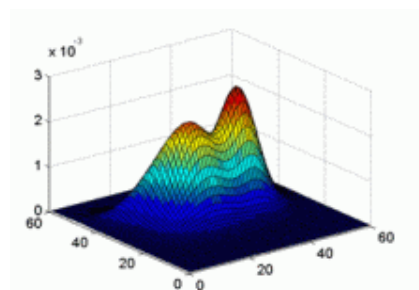


TRAIN

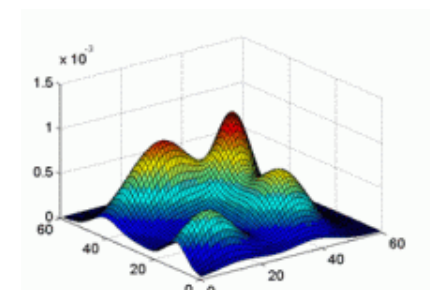


TEST

Flickr



YouTube



Mismatched domains

amazon.com



The Office

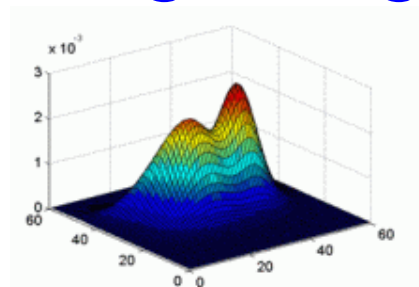


TRAIN

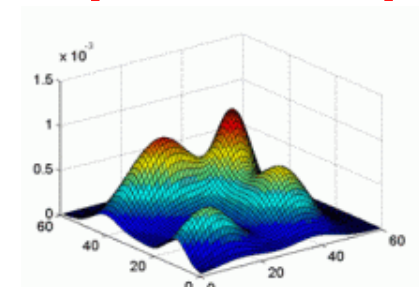


TEST

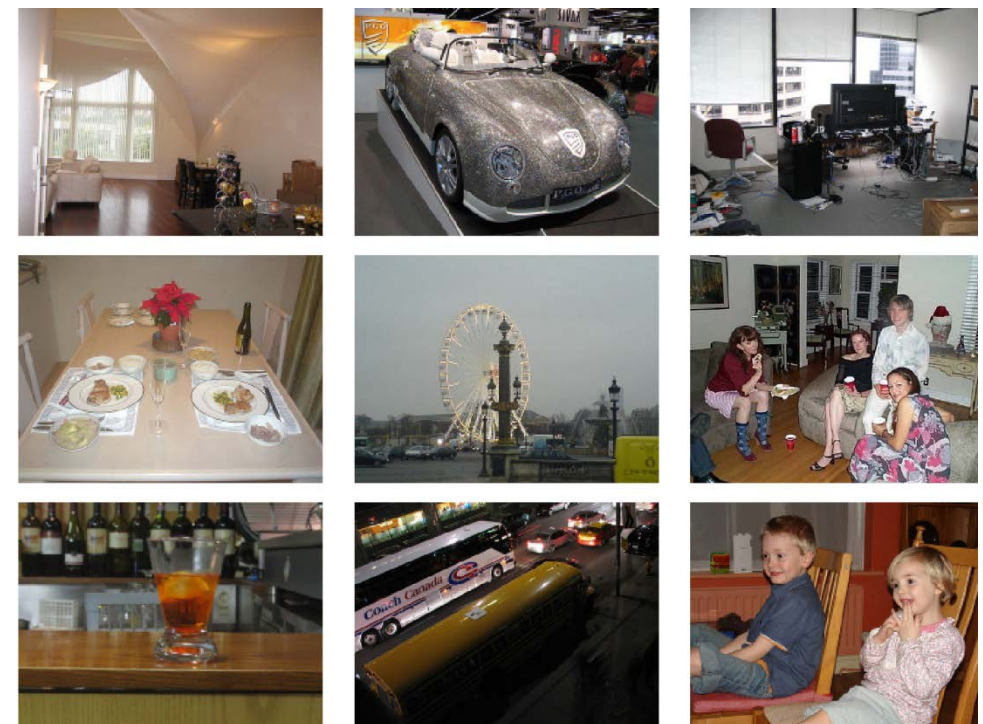
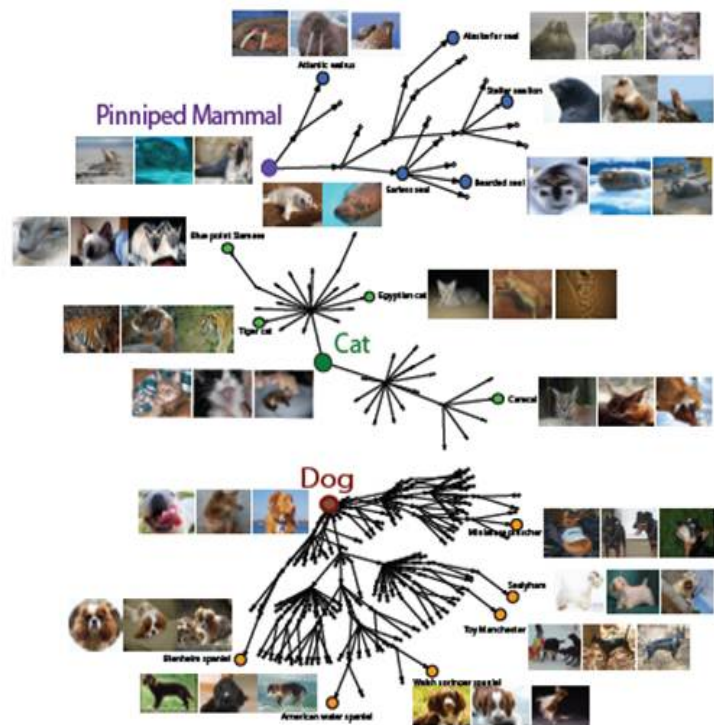
Catalog images



Mobile phone photos



Mismatched domains

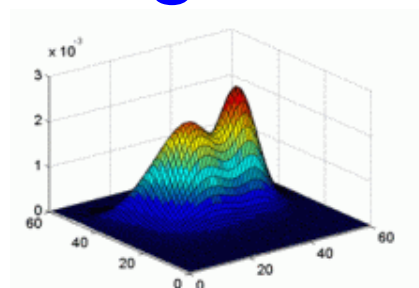


TRAIN

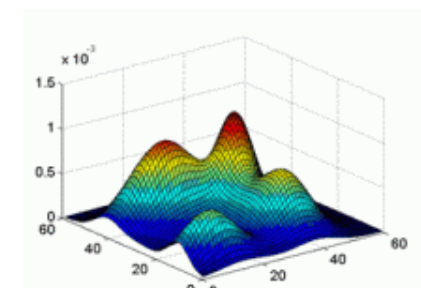


TEST

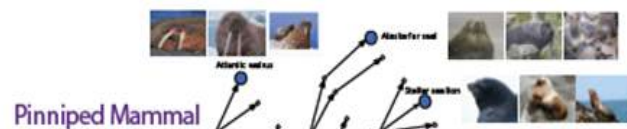
ImageNet



PASCAL VOC



Mismatched domains



“It is worthwhile to note that, **even with 140K training ImageNet images**, we do not perform as well as with 5K PASCAL VOC training images.”

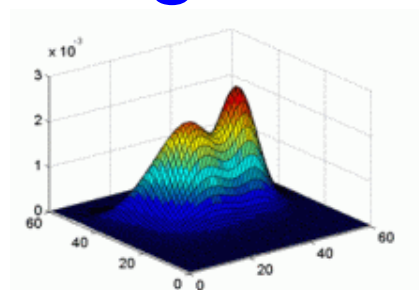
– Perronnin et al. CVPR 2010

TRAIN

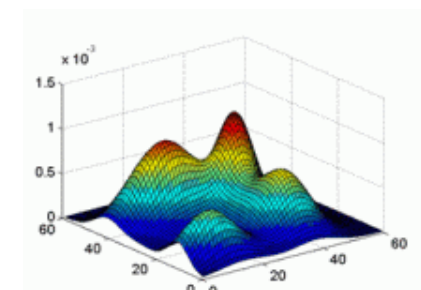


TEST

ImageNet



PASCAL VOC



Mismatched domains

Problem: Poor cross-domain generalization

- Different underlying distributions
- Overfit to datasets' idiosyncrasies

Possible solution:

Unsupervised domain adaptation

Unsupervised domain adaptation

Setup

Source domain (with labeled data)

$$D_{\mathcal{S}} = \{(x_m, y_m)\}_{m=1}^M \sim P_{\mathcal{S}}(X, Y)$$

Target domain (no labels for training)

$$D_{\mathcal{T}} = \{(x_n, y_n)\}_{n=1}^N \sim P_{\mathcal{T}}(X, Y)$$

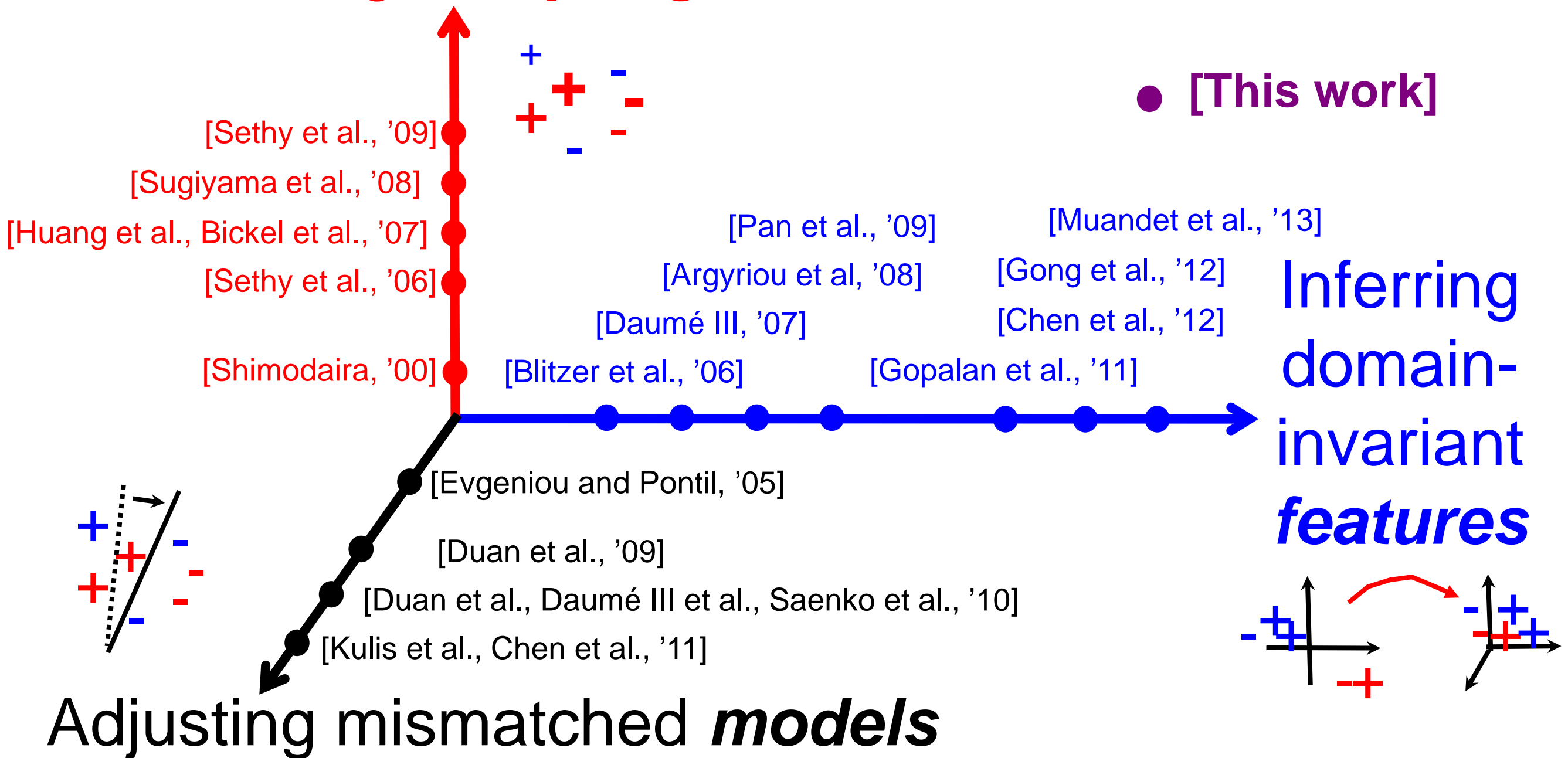
Different distributions

Objective

Learn classifier to work well on the **target**

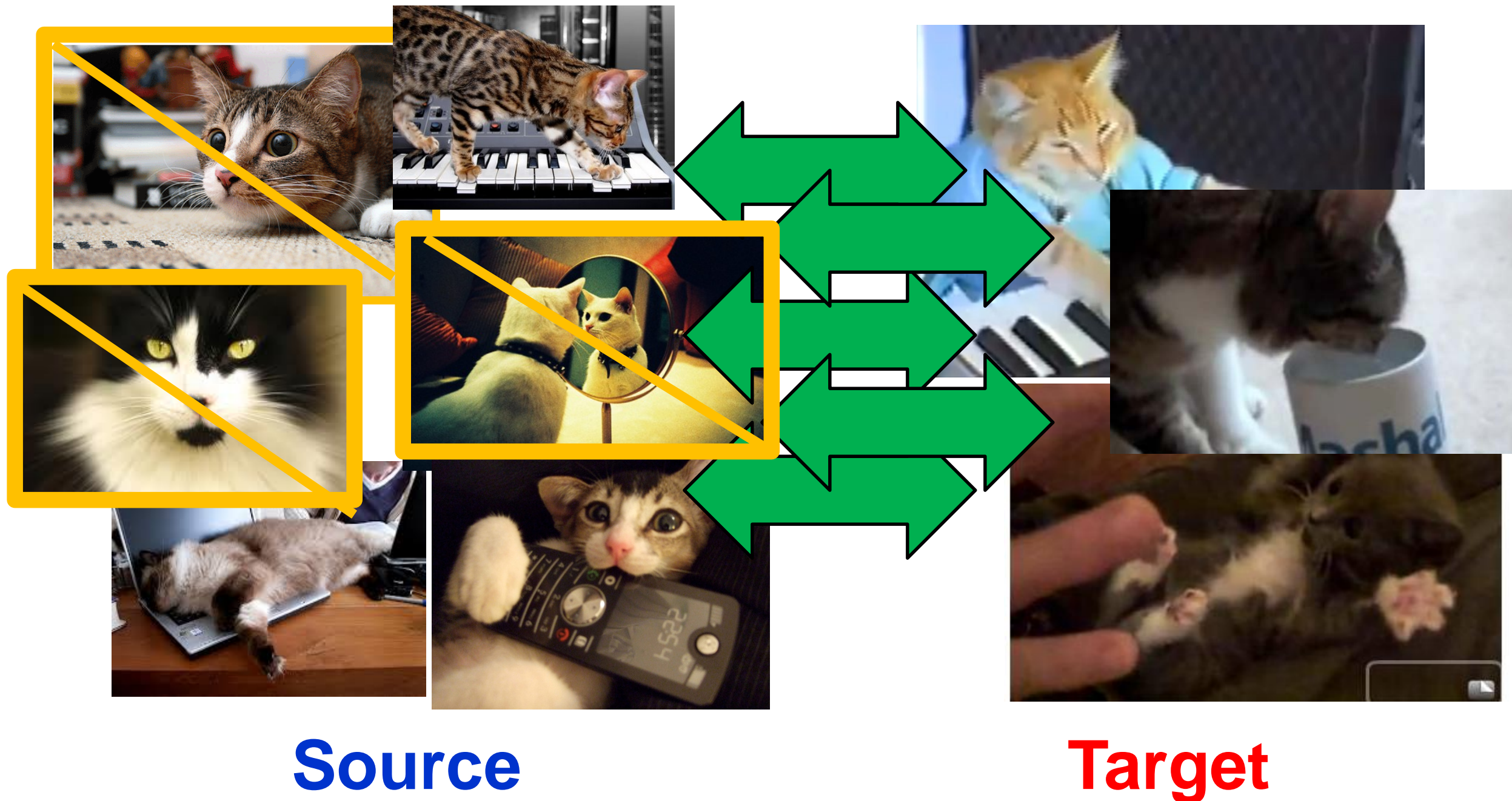
Much recent research

Correcting *sampling* bias



Problem

Existing methods attempt to adapt *all* source data points, including “hard” ones.



Problem

Existing methods attempt to adapt *all* source data points, including “hard” ones.

Our idea

Automatically identify the “most adaptable” instances

Use them to create series of easier auxiliary domain adaptation tasks

Landmarks

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



Source



Target

Landmarks

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

Roles:

Ease adaptation difficulty

Provide discrimination
(biased to **target**)

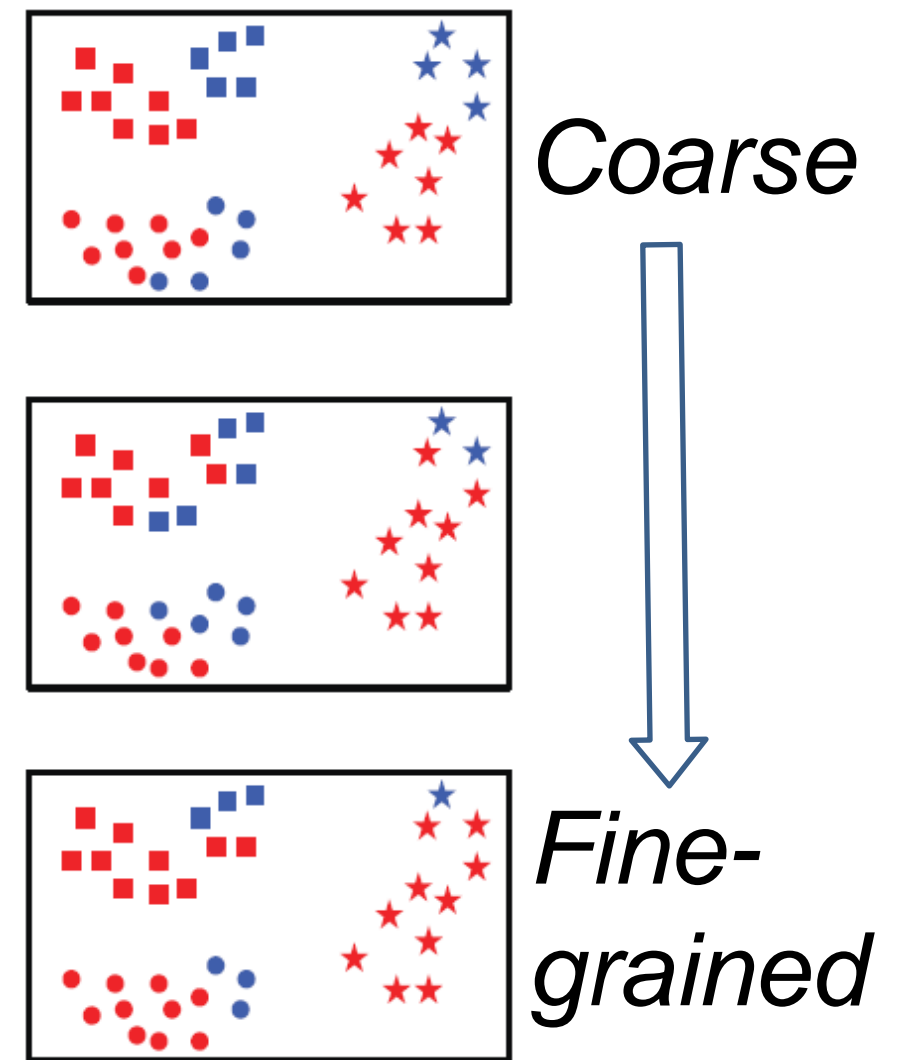
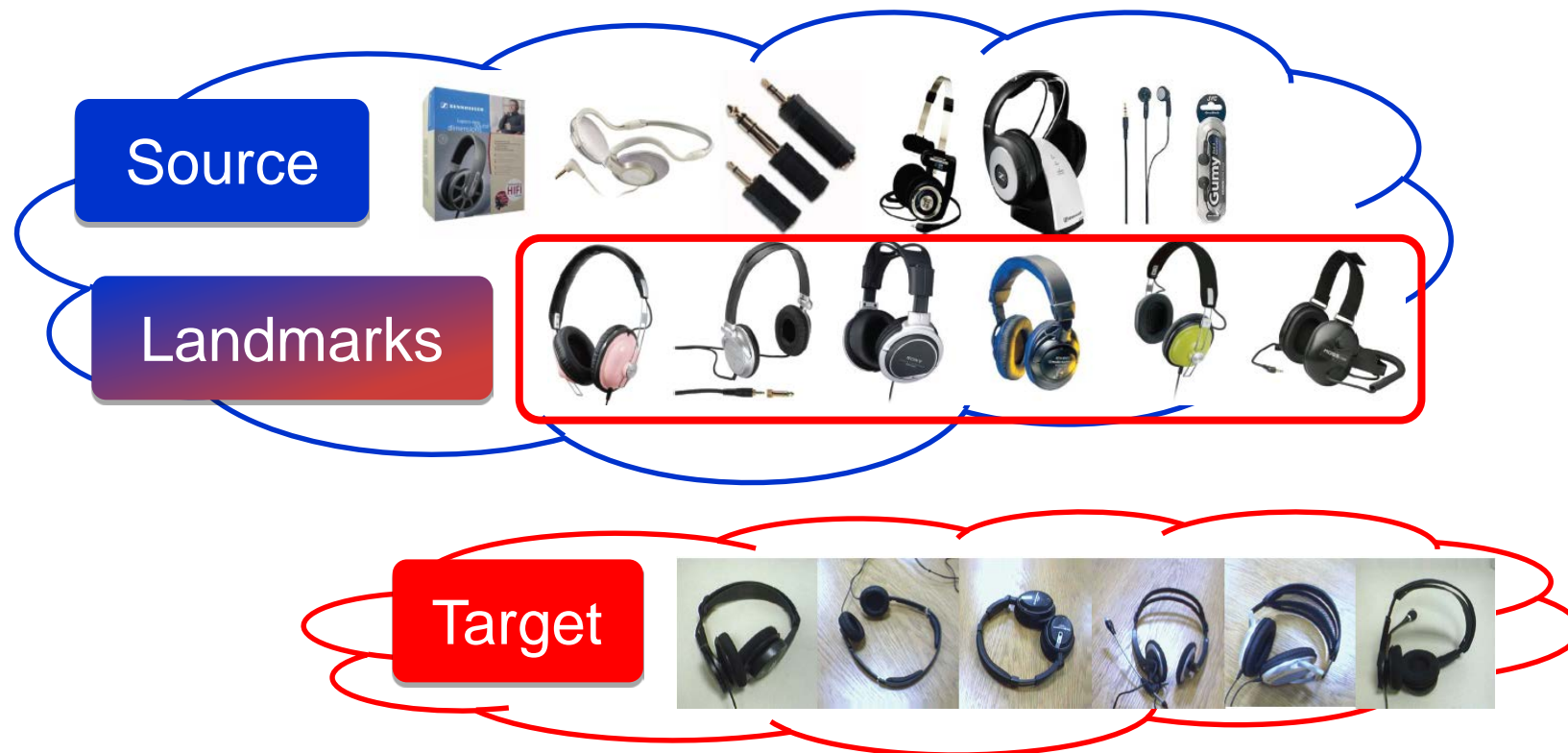


Source



Target

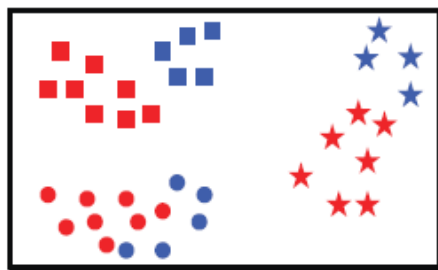
Key steps



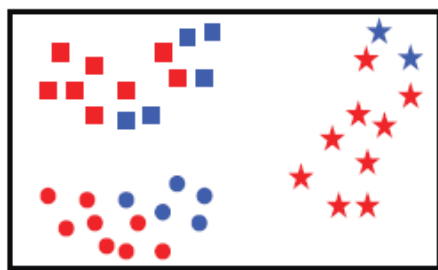
1 Identify landmarks

at multiple scales.

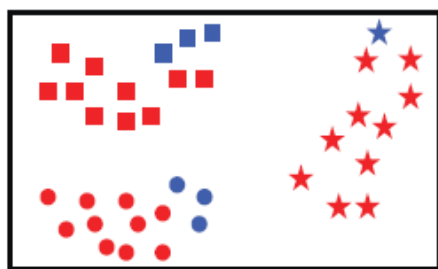
Key steps



⇒ $\Phi_1(\mathbf{x})$



⇒ $\Phi_2(\mathbf{x})$

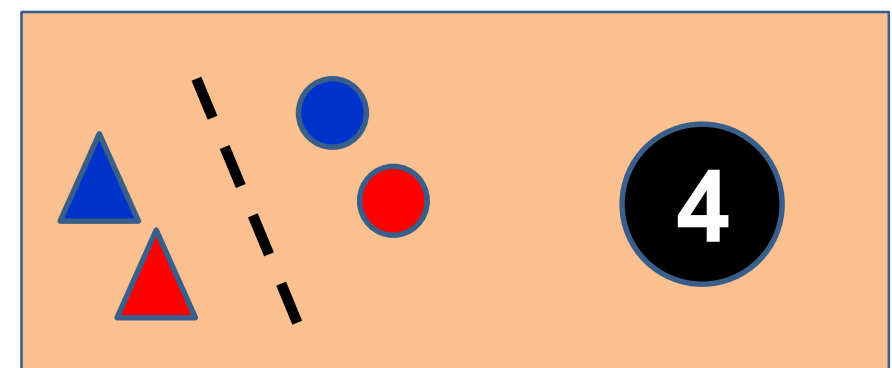


⇒ $\Phi_3(\mathbf{x})$

2 Construct auxiliary domain adaptation tasks

$$\Phi(\mathbf{x}) = \begin{bmatrix} \Phi_1(\mathbf{x}) \cdot w_1 \\ \Phi_2(\mathbf{x}) \cdot w_2 \\ \Phi_3(\mathbf{x}) \cdot w_3 \end{bmatrix} \quad \text{3}$$

Obtain domain-invariant features



Predict target labels

Identifying landmarks

Objective

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

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



Source



Target

Maximum mean discrepancy (MMD)

Empirical estimate [Gretton et al. '06]

$$d(P_{\mathcal{L}}, P_{\mathcal{T}}) = \left\| \frac{1}{L} \sum_{l=1}^L \phi(x_l) - \frac{1}{N} \sum_{n=1}^N \phi(x_n) \right\|_{\mathcal{H}}$$

\mathcal{H} a universal RKHS

$\phi(\cdot)$ kernel function induced by \mathcal{H}

x_l the l -th **landmark** (from the **source** domain)

Method for identifying landmarks

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

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

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

Method for identifying landmarks

Convex relaxation

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

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

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

Scale for landmark similarity?

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

Gaussian kernels

How to choose the bandwidth?

Our solution:

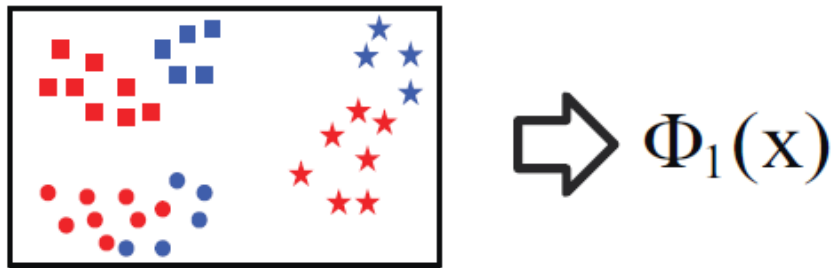
Examine distributions at multiple granularities

Multiple bandwidths \rightarrow multiple sets of landmarks

Landmarks at multiple scales



Key steps



2 Construct auxiliary domain adaptation tasks

Constructing easier auxiliary tasks



At each scale σ

New source = **Source** \setminus Landmarks

New target = **Target** \cup Landmarks

Intuition: distributions are closer (cf. Theorem 1)

Constructing easier auxiliary tasks



At each scale σ

New source = **Source** \setminus Landmarks

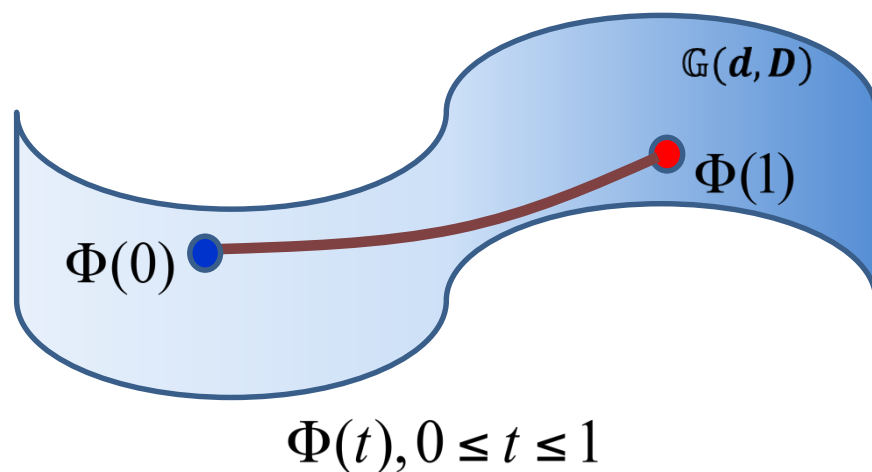
New target = **Target** \cup Landmarks

Intuition: distributions are closer (cf. Theorem 1)

Constructing easier auxiliary tasks

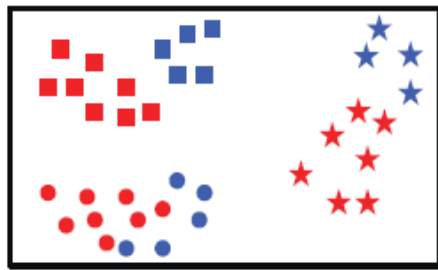
Each task provides new basis of features via geodesic flow kernel (GFK):

$$K_{\sigma}(x_i, x_j) = \int_0^1 (\Phi_{\sigma}(t)'x_i)'(\Phi_{\sigma}(t)'x_j)dt = x_i G_{\sigma} x_j$$

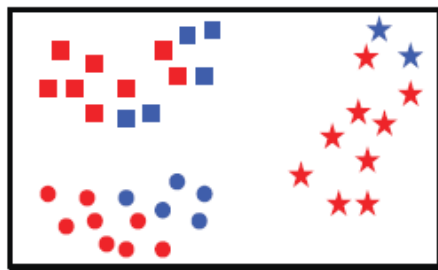


- Integrate out domain changes
- Obtain domain-invariant representation [Gong, et al. '12]

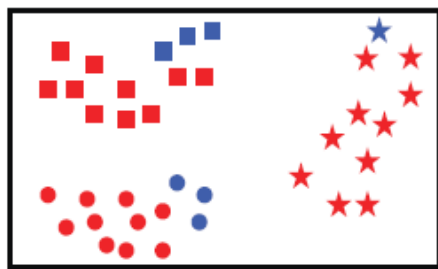
Key steps



⇒ $\Phi_1(\mathbf{x})$



⇒ $\Phi_2(\mathbf{x})$



⇒ $\Phi_3(\mathbf{x})$

$$\Phi(\mathbf{x}) = \begin{bmatrix} \Phi_1(\mathbf{x}) \cdot w_1 \\ \Phi_2(\mathbf{x}) \cdot w_2 \\ \Phi_3(\mathbf{x}) \cdot w_3 \end{bmatrix} \quad \text{3 MKL}$$

Obtain domain-invariant features

2 Construct auxiliary domain adaptation tasks

Combining features discriminatively

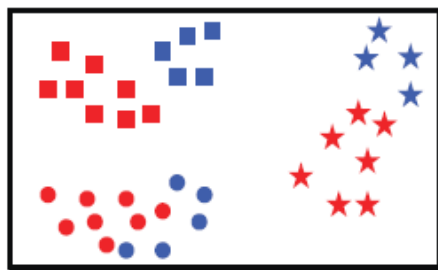
Multiple kernel learning on the **labeled landmarks**

$$F = \sum_{\sigma} w_{\sigma} G_{\sigma}, \quad \text{s.t.} \quad w_{\sigma} \geq 0, \quad \sum_{\sigma} w_{\sigma} = 1$$

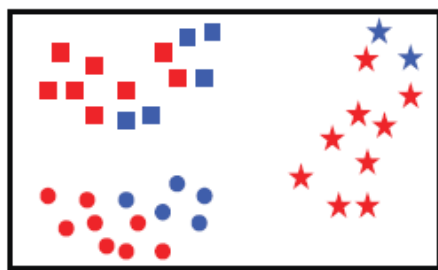
Arriving at domain-invariant feature space

Discriminative loss biased to the **target**

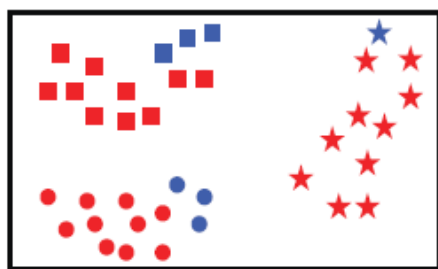
Key steps



$\Rightarrow \Phi_1(x)$



$\Rightarrow \Phi_2(x)$

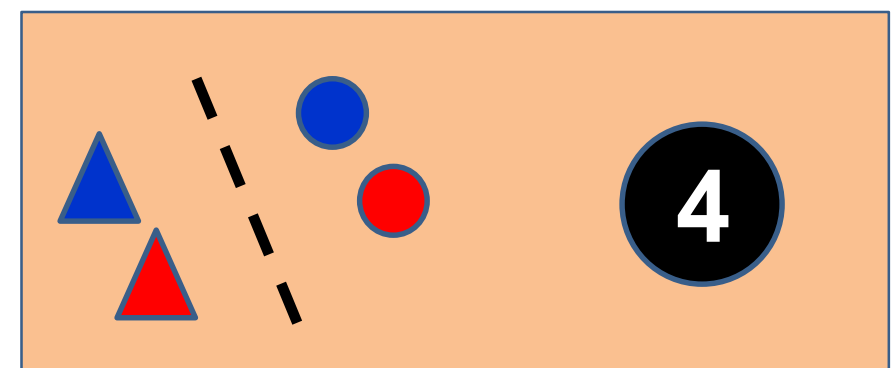


$\Rightarrow \Phi_3(x)$

$$\Phi(x) = \begin{bmatrix} \Phi_1(x) \cdot w_1 \\ \Phi_2(x) \cdot w_2 \\ \Phi_3(x) \cdot w_3 \end{bmatrix} \quad \text{3}$$

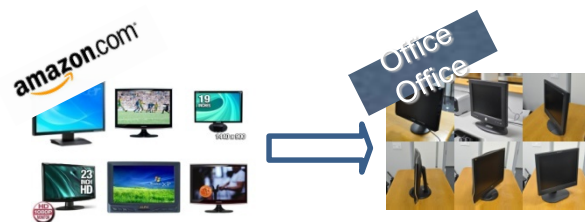
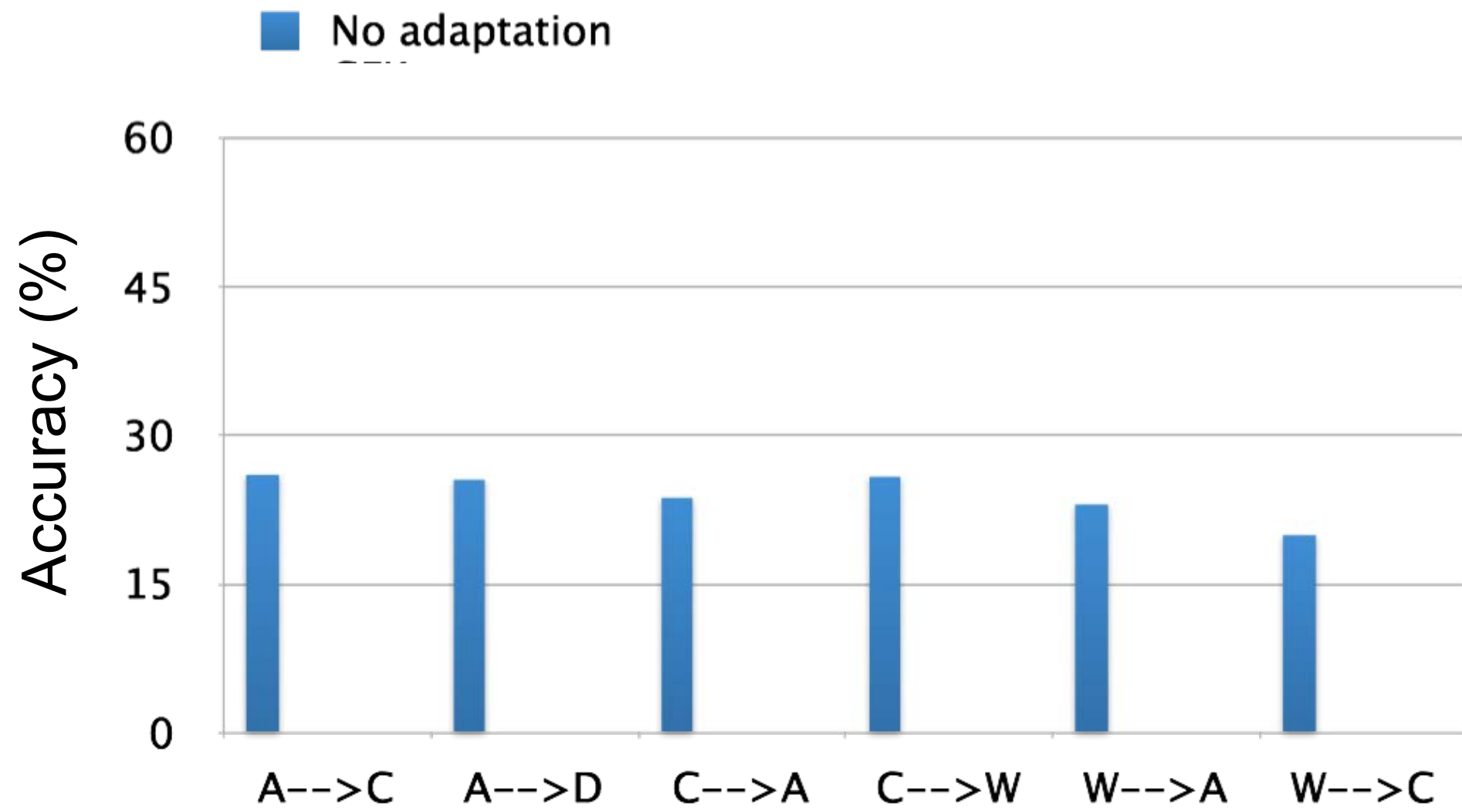
Obtain domain-invariant features

2 Construct auxiliary domain adaptation tasks

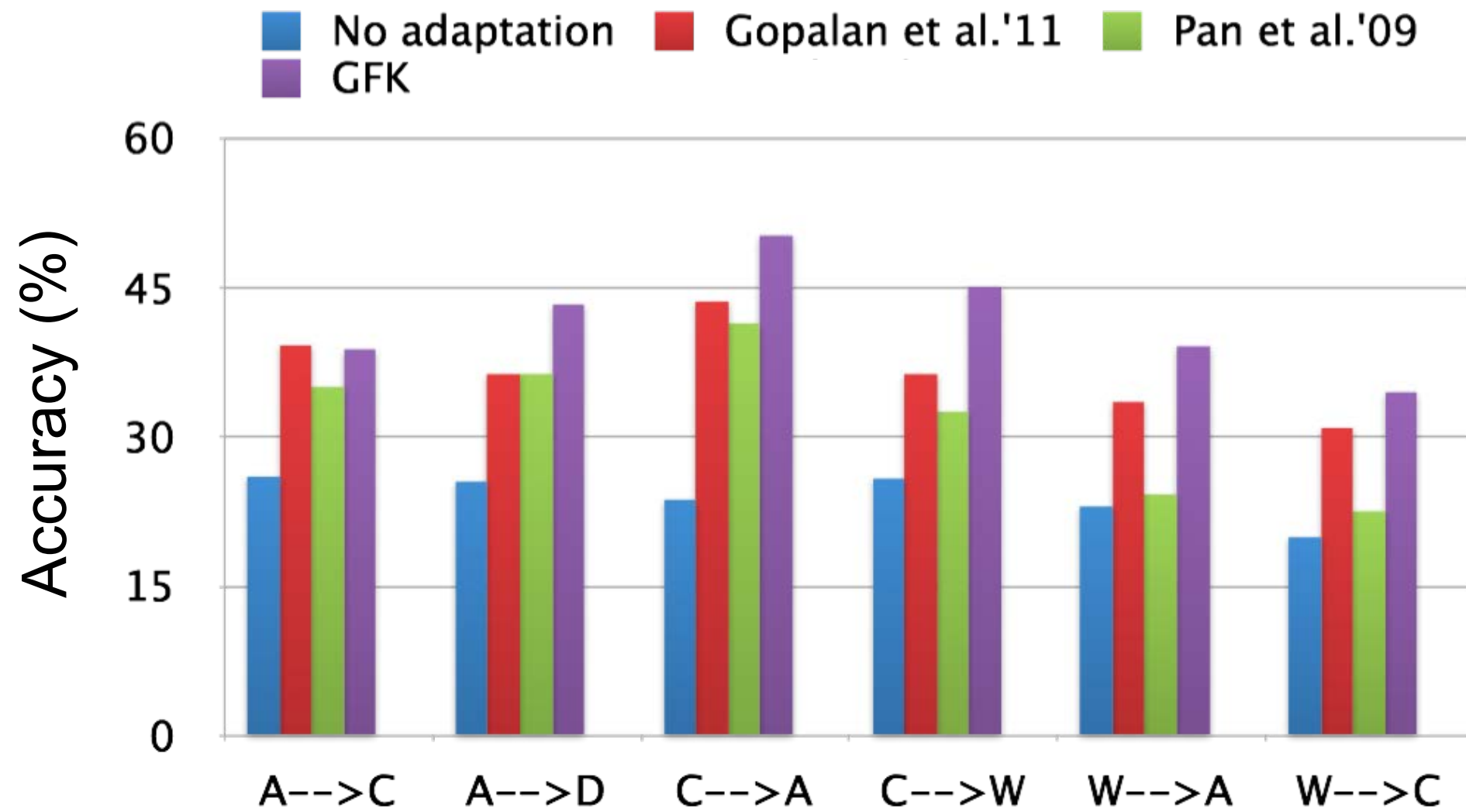


Predict target labels

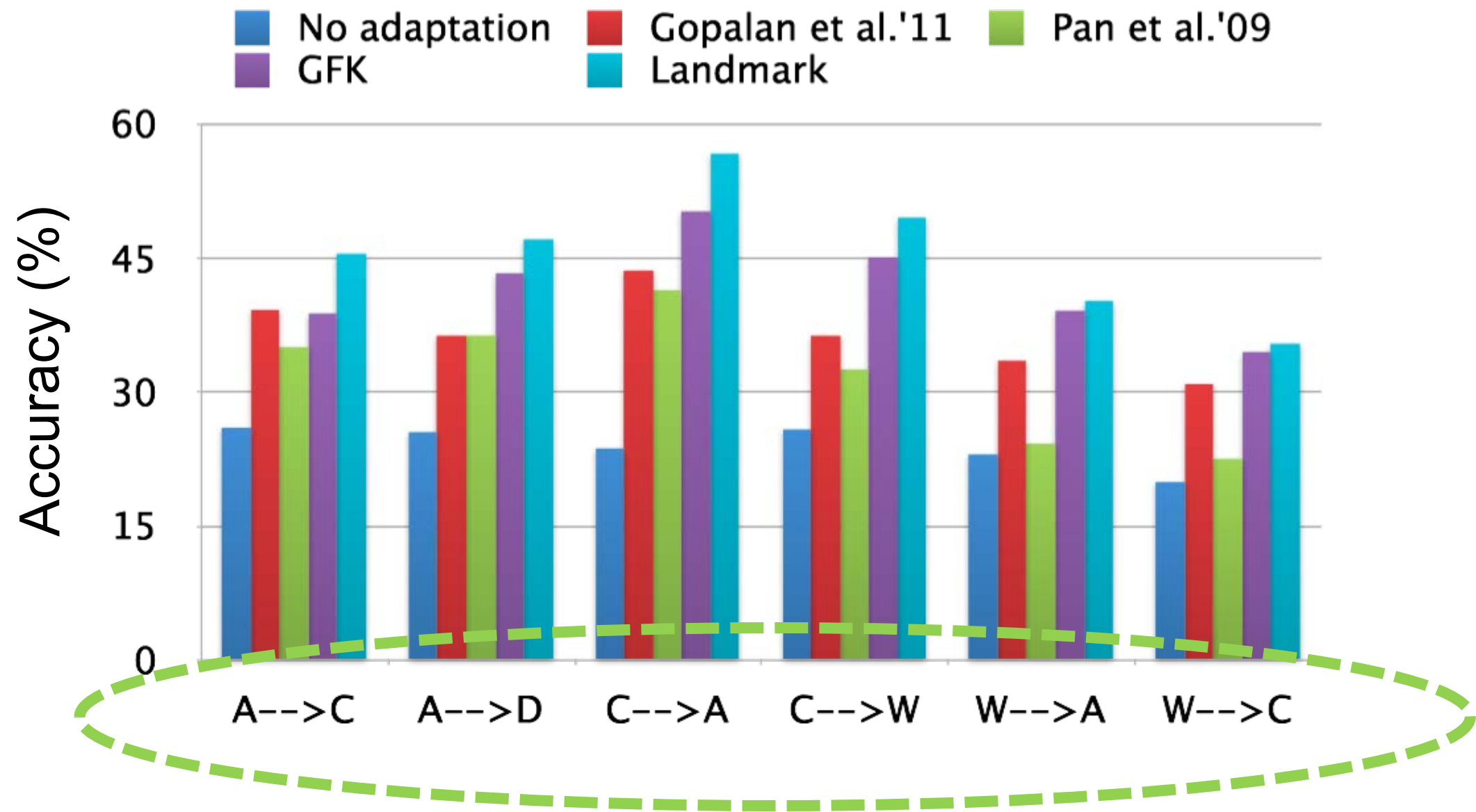
Cross-dataset object recognition



Cross-dataset object recognition



Cross-dataset object recognition



Datasets as domains?

Domain 1

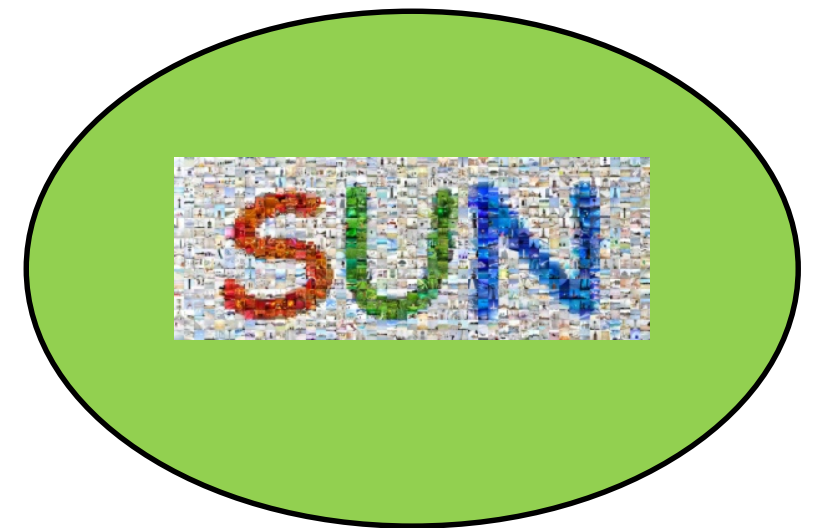


Domain 2



ASSUMED

Domain 3

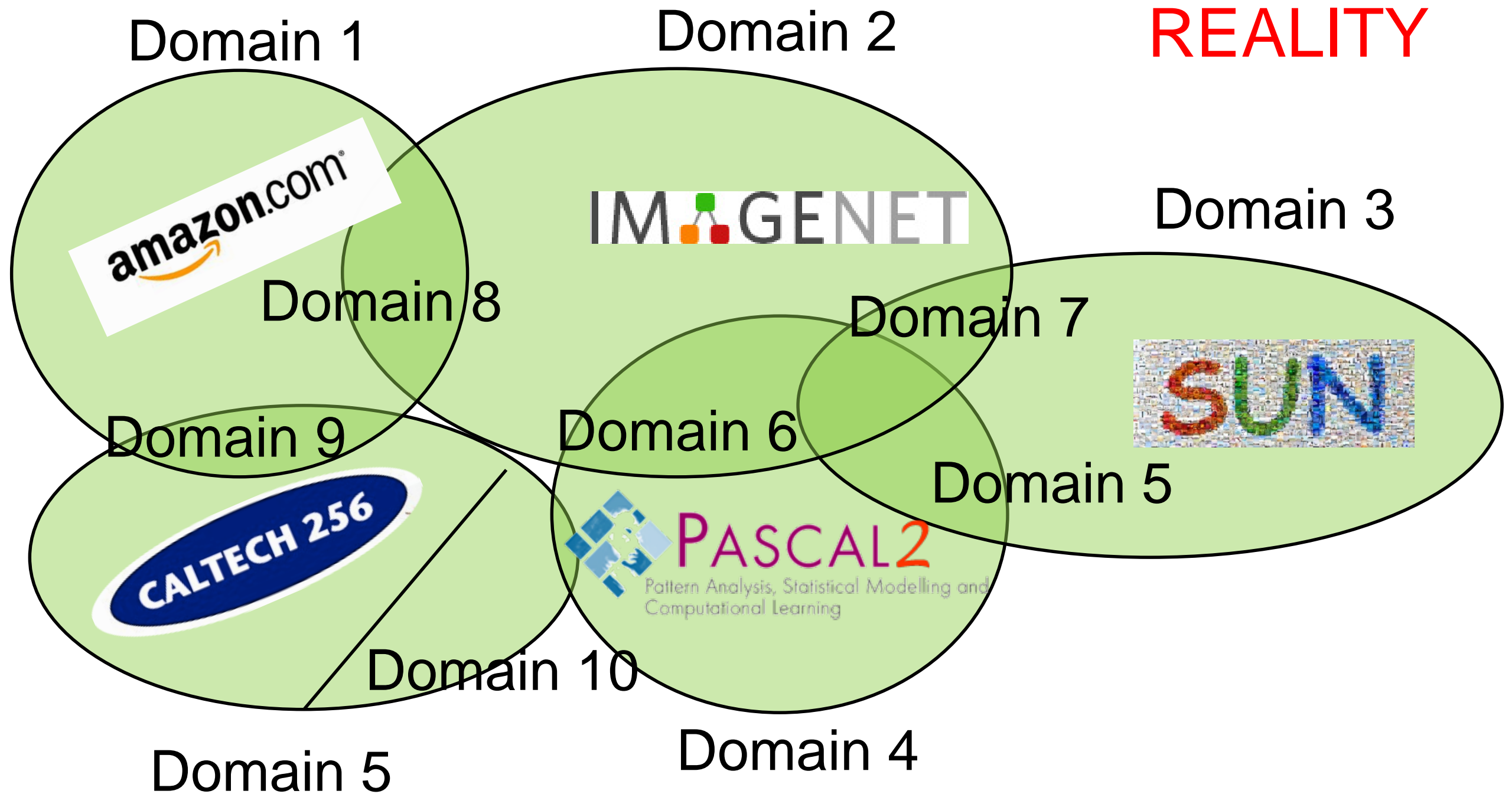


Domain 5



Domain 4

Datasets as domains?



Datasets as domains?

Domain 1

Domain 2

REALITY

Dataset \neq Domain

Cross-dataset adaptation is suboptimal

Domain 5

Domain 10

Domain 4

CALTECH

Pattern Analysis, Statistical Modelling and
Computational Learning

How to define a domain?

NLP: *Language-specific domains* ✓

Speech: *Speaker-specific domains* ✓

Vision: ??

pose-specific? illumination-specific?

occlusion? image resolution? background?

Challenges:

Many continuous factors vs. few discrete

Factors **overlap** and interact

Discovering latent visual domains

We propose to **discover** domains – “reshaping” them to cross dataset boundaries

Maximum distinctiveness

$$\max_{\{z_{mk}\}} \sum_{k \neq k'} d(P(k), P(k')) \longrightarrow \text{MMD}$$

$$\text{where } z_{mk} = \begin{cases} 1 & \text{if } \mathbf{x}_m \text{ belongs to domain } k \\ 0 & \text{else} \end{cases}$$

Maximum learnability

Determine K with domain-wise cross-validation

Results: discovering domains

amazon.com[®]



CALTECH 256

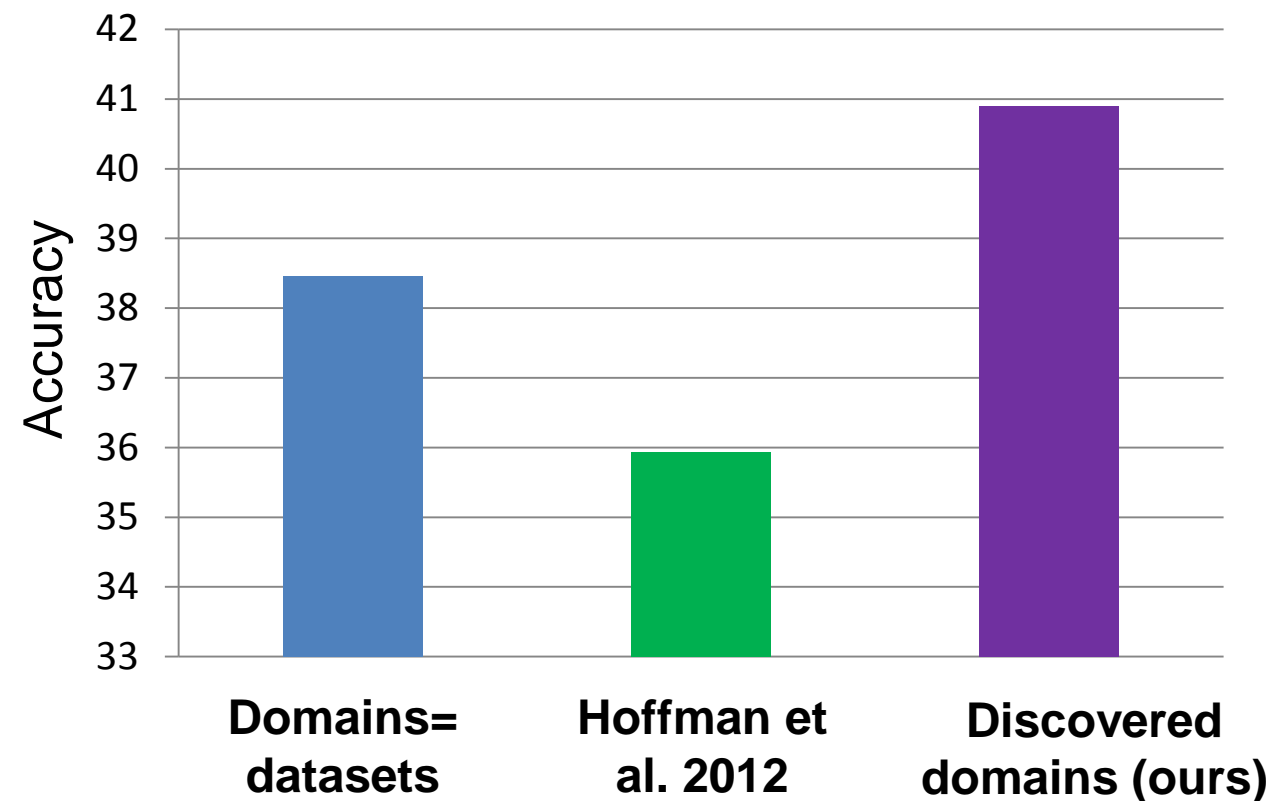


**Discovered
domain I**

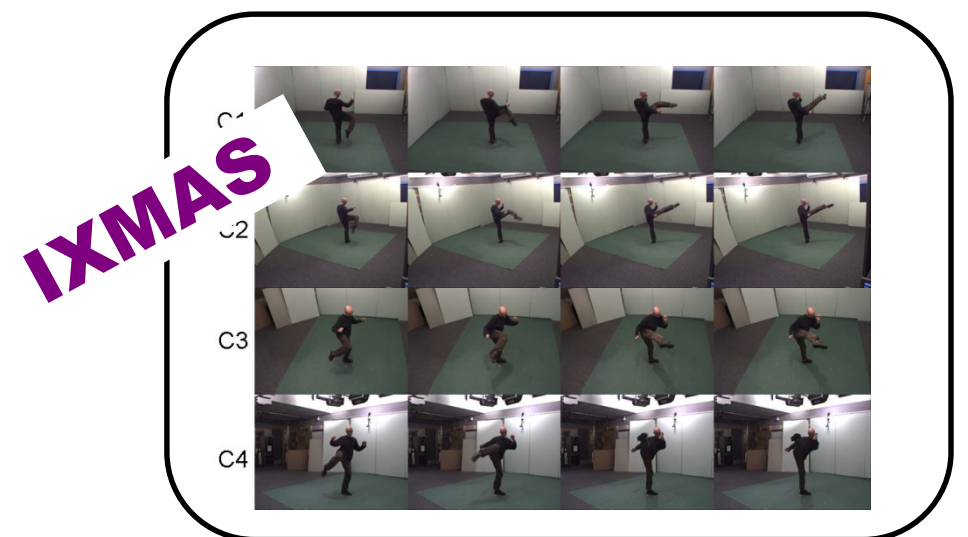
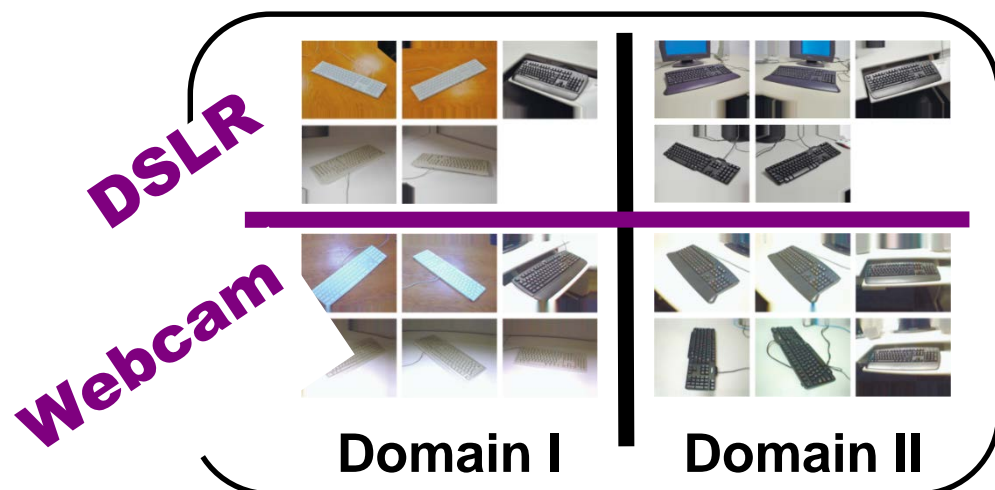
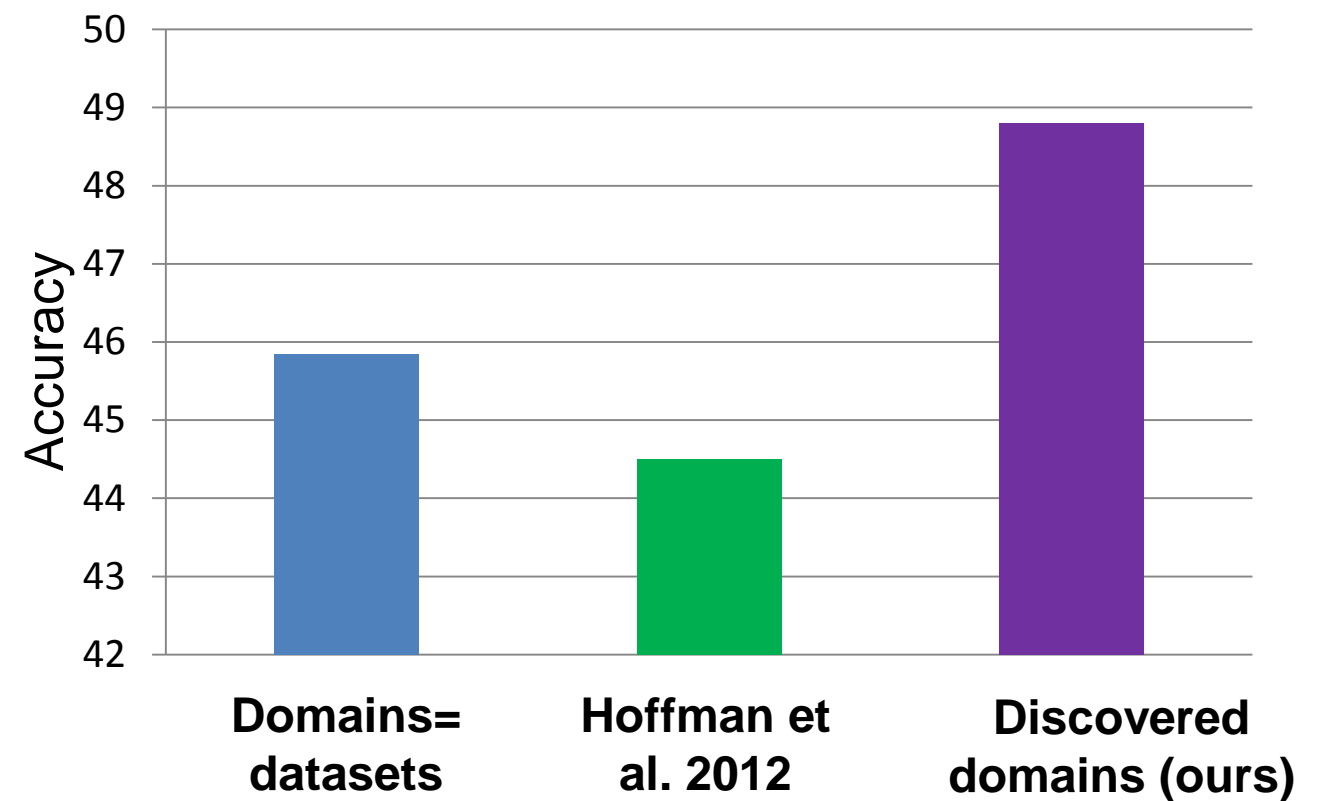
**Discovered
domain II**

Results: discovering domains

Cross-dataset object recognition



Cross-viewpoint action recognition



Summary so far

landmarks

labeled **source** instances

distributed similarly to the **target**

auxiliary tasks provably easier to solve

discriminative loss despite unlabeled **target**

reshaping datasets to latent domains

discover cross-dataset domains

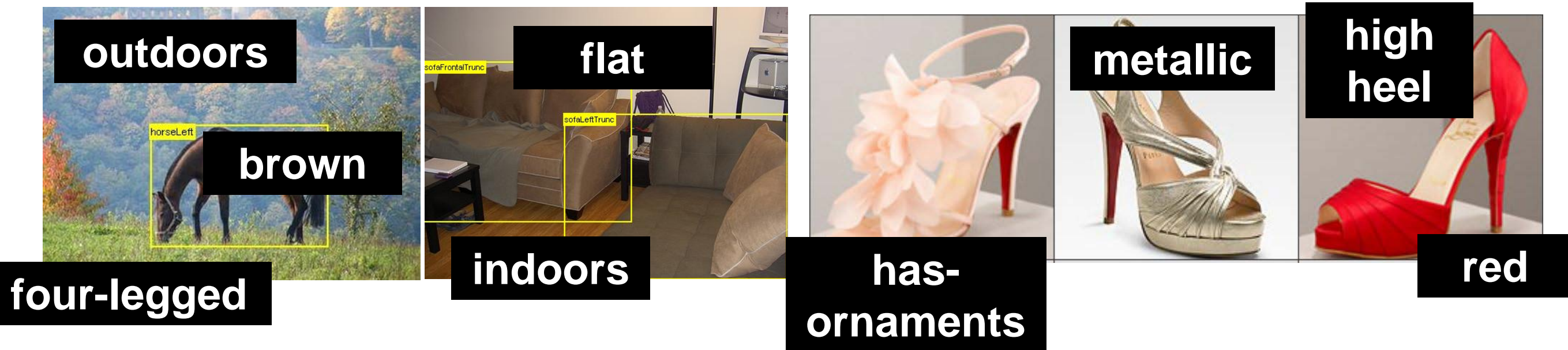
maximally distinct & learnable

Typical assumptions

1. Test set will look like the training set.
2. Human labelers “see” the same thing.

Visual attributes

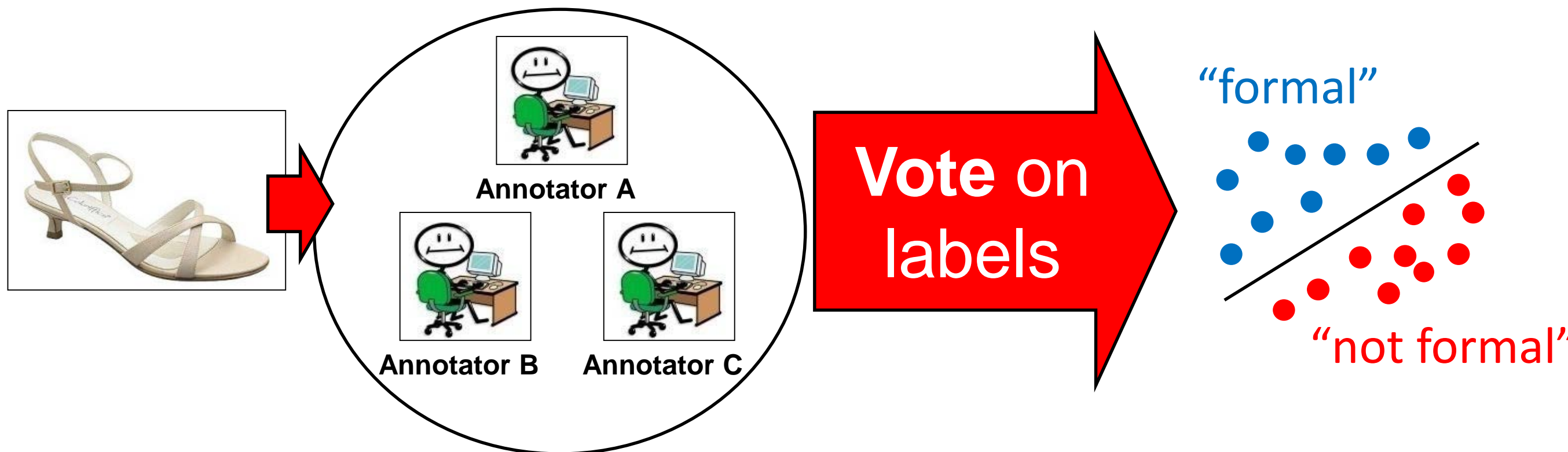
- High-level semantic properties shared by objects
- Human-understandable and machine-detectable



[Oliva et al. 2001, Ferrari & Zisserman 2007, Kumar et al. 2008, Farhadi et al. 2009, Lampert et al. 2009, Endres et al. 2010, Wang & Mori 2010, Berg et al. 2010, Branson et al. 2010, Parikh & Grauman 2011, ...]

Standard approach

Learn one monolithic model per attribute



Problem

There may be valid perceptual differences within an attribute.

Formal? User labels:
50% "yes"
50% "no"



Binary attribute

More ornamented? User labels:
50% "first"
20% "second"
30% "equally"

or



Relative attribute

Imprecision of attributes

Fine-grained meaning



Overweight?

or just

Chubby?

Imprecision of attributes

Context

Is  *formal*?

= *formal* wear for a **conference**? OR

= *formal* wear for a **wedding**?

Imprecision of attributes

Cultural

Is  *blue or green?*

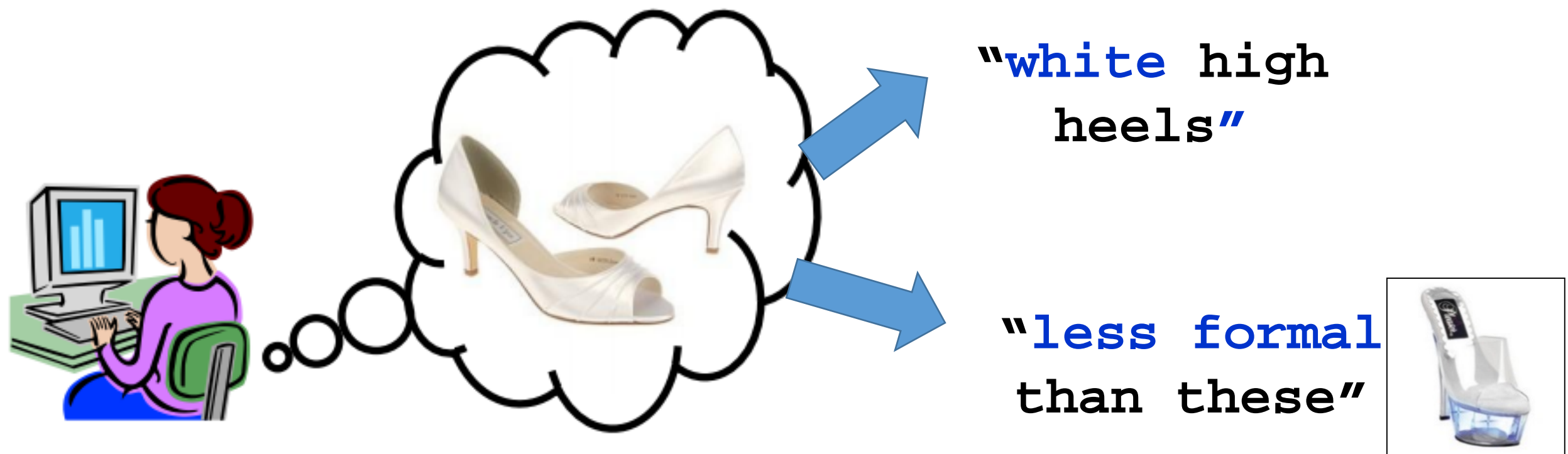
English: “blue”

Russian: “neither”
 (“голубой” vs. “синий”)

Japanese: “both”
 (“青” = blue *and* green)

But do we need to be that precise?

Yes. Applications like image search require that **user's perception** matches **system's predictions**.



Our idea

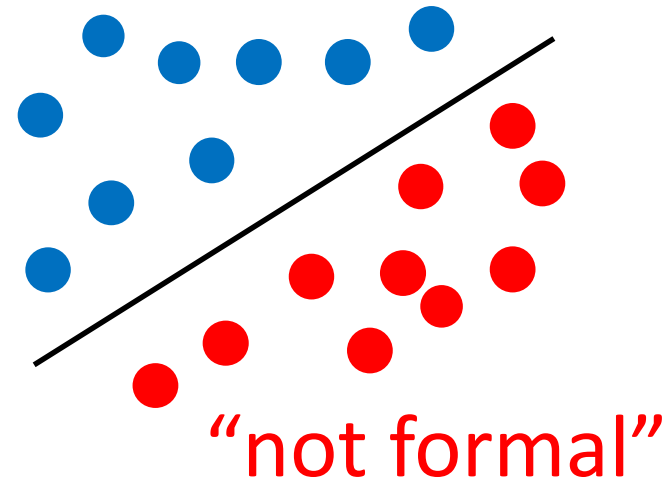
- Treat learning **perceived attributes** as an adaptation problem.
- **Adapt generic attribute model** with minimal user-specific labeled examples.
- Obtain **implicit user-specific labels** from user's search history

Our idea



Vote on labels

“formal”



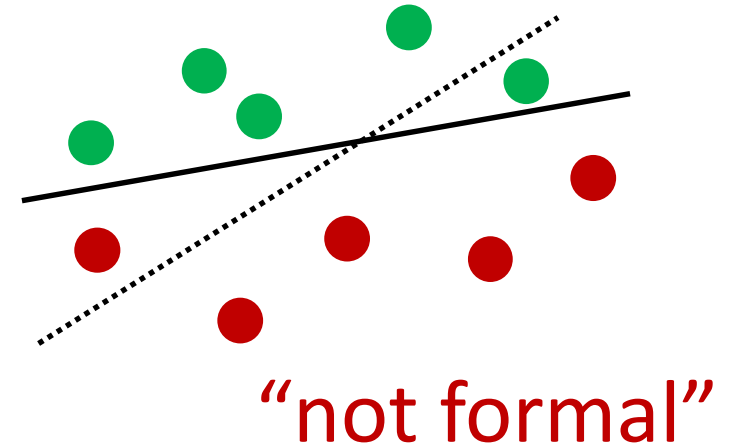
Adapt



“formal”

“not formal”

“formal”



Learning adapted attributes

- Adapting **binary attribute** classifiers:

Given user-labeled data $D_b = \{\mathbf{x}_i, y_i\}_{i=1}^N$

and generic model \mathbf{w}'_b ,

$$\min_{\mathbf{w}_b} \frac{1}{2} \|\mathbf{w}_b - \mathbf{w}'_b\|^2 + C \sum_{i=1}^N \xi_i,$$

subject to $y_i \mathbf{x}_i^T \mathbf{w}_b \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad \forall i$

Learning adapted attributes

- Adapting **relative attribute** rankers:

Given user-labeled data $D_r = \{(x_{i_1} \succ x_{j_1})\}_{i=1}^N$

and generic model w'_r ,

$$\min_{w_r} \frac{1 - \delta}{2} \|w_r\|^2 + \frac{\delta}{2} \|w_r - w'_r\|^2 + C \sum_{i=1}^N \xi_i$$

subject to $w_r^T x_{i_1} - w_r^T x_{i_2} \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad \forall i,$

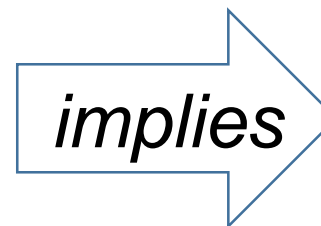
Collecting user-specific labels

- **Explicitly** from actively requested labels
 - Seek labels on uncertain and diverse images
- **Implicitly** from search history

- Transitivity

“My target is...

less formal than



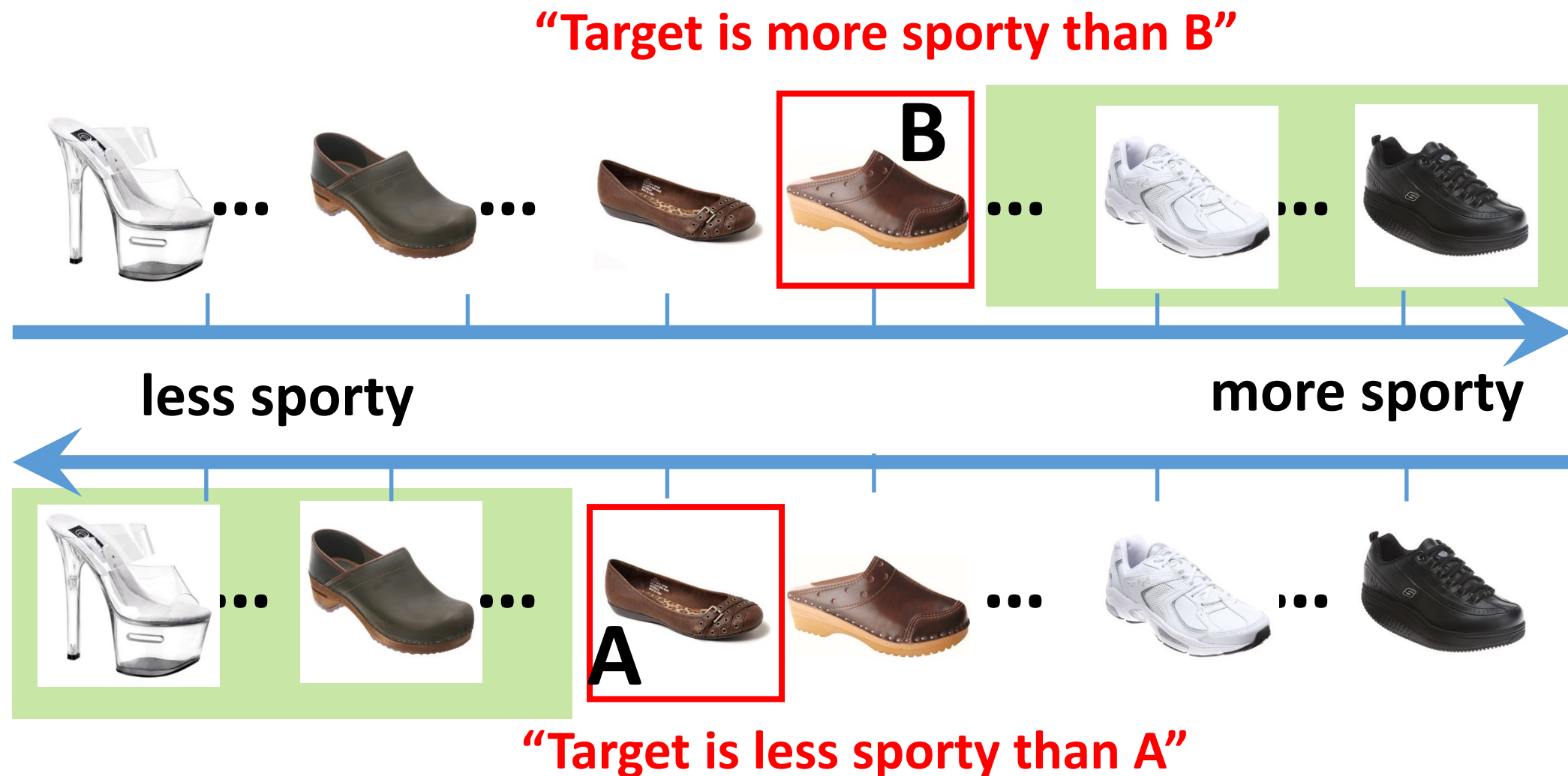
more formal than



“

- Contradictions

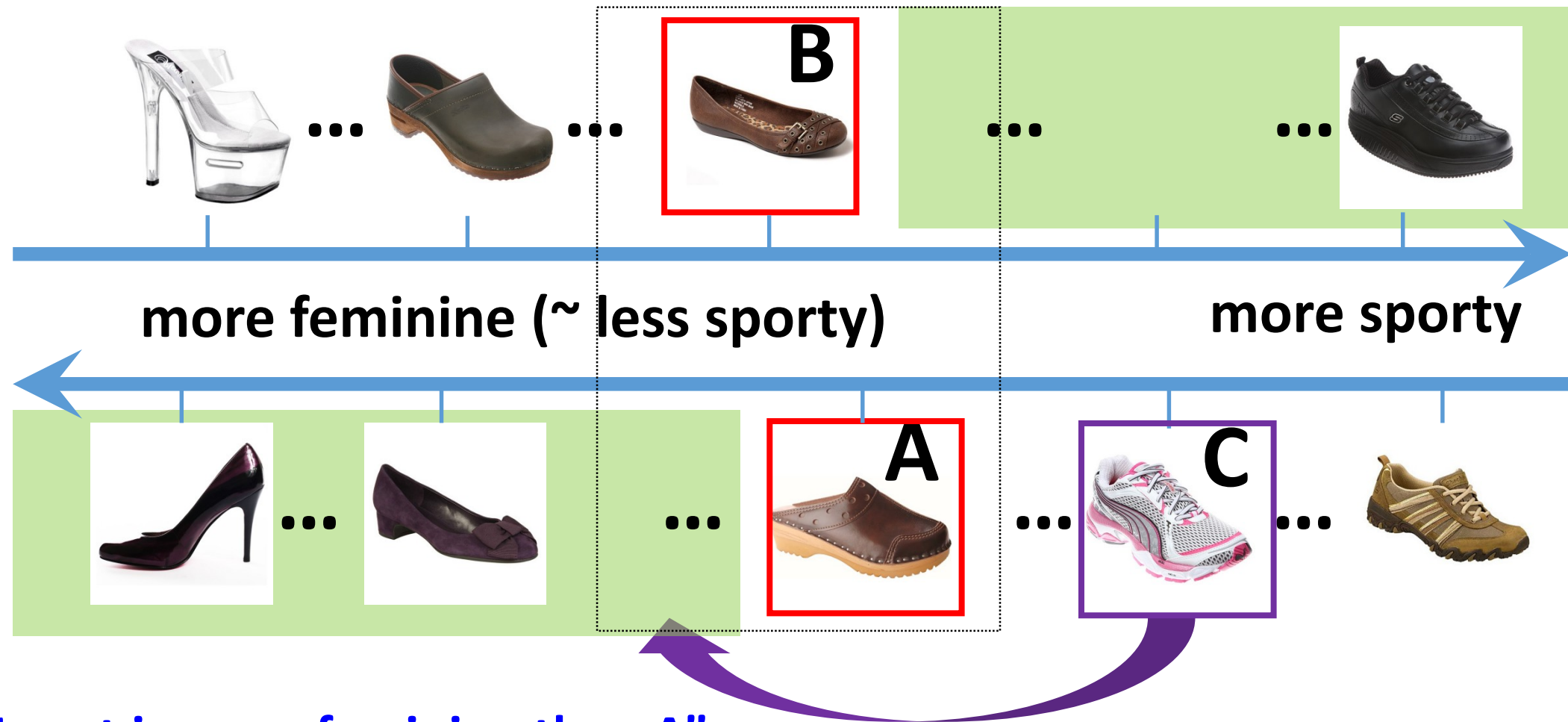
Inferring implicit labels



User's feedback history can reveal mismatch in perceived and predicted attributes

Inferring implicit labels

“Target is more sporty than B”



User’s feedback history can reveal mismatch in perceived and predicted attributes

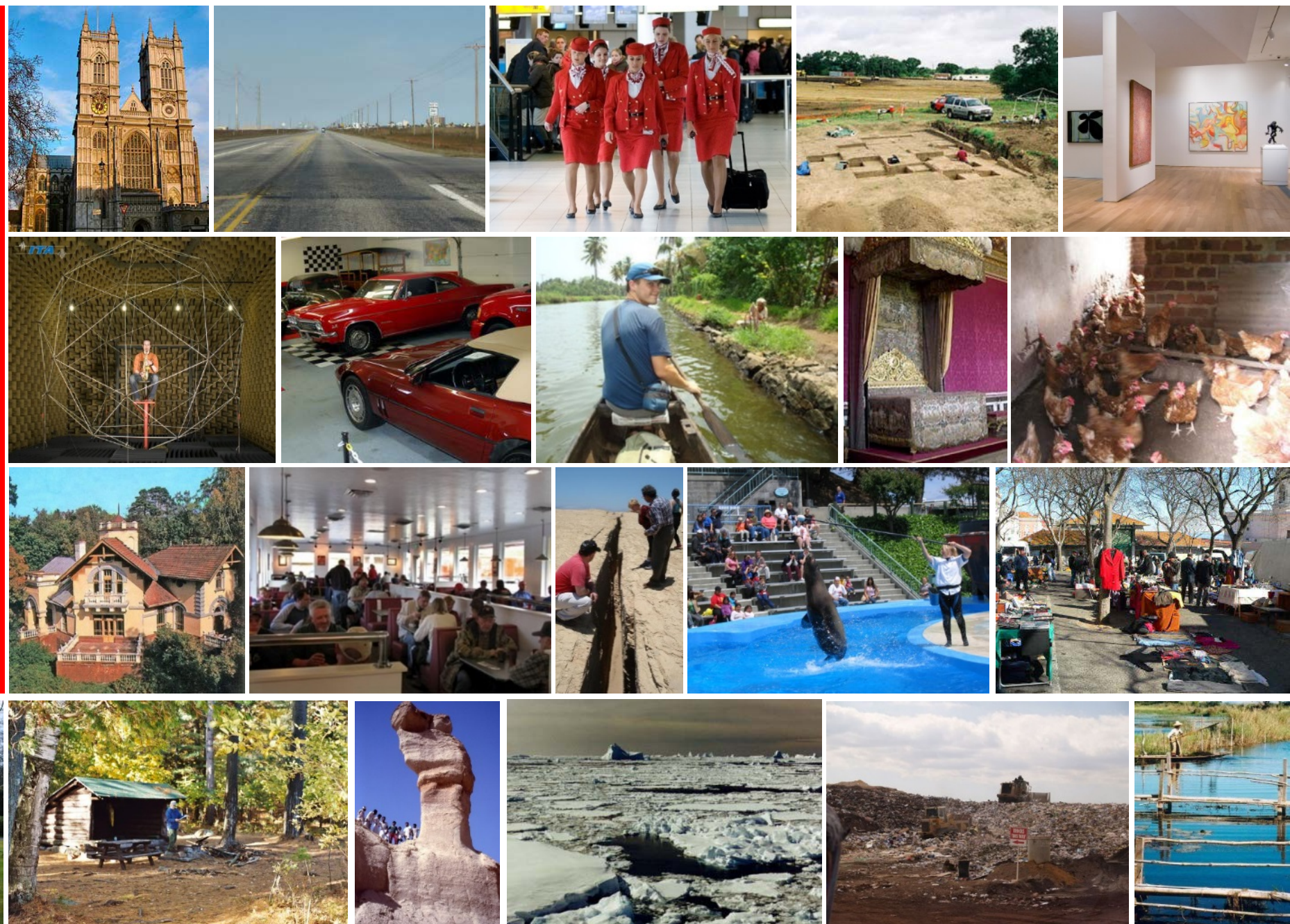
Datasets

SUN Attributes:

14,340 scene images

12 attributes:

“sailing”, “hiking”,
“vacationing”, “open area”,
“vegetation”, etc.



Shoes:

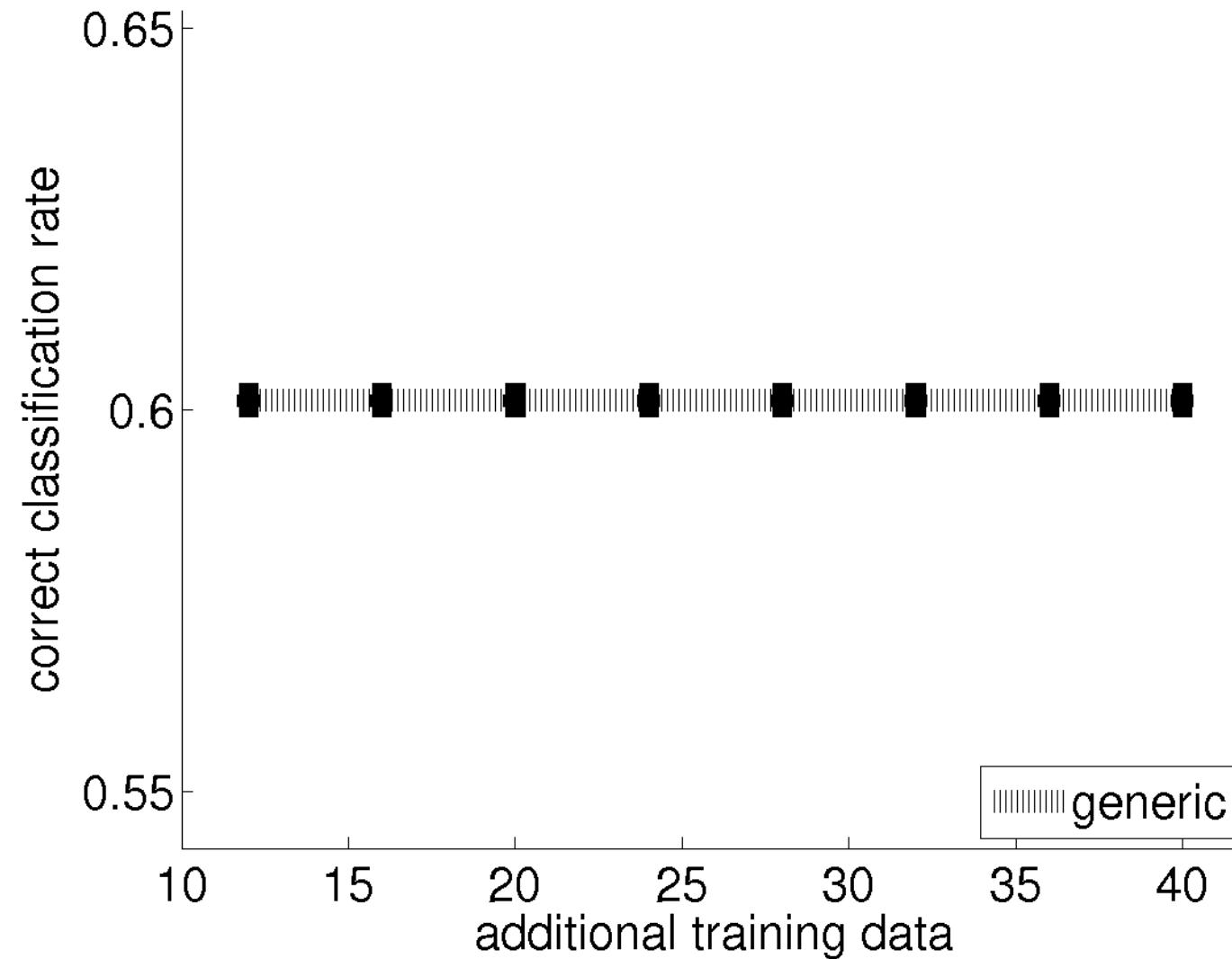
14,658 shoe images;

10 attributes:

“pointy”, “bright”, “high-
heeled”, “feminine” etc.

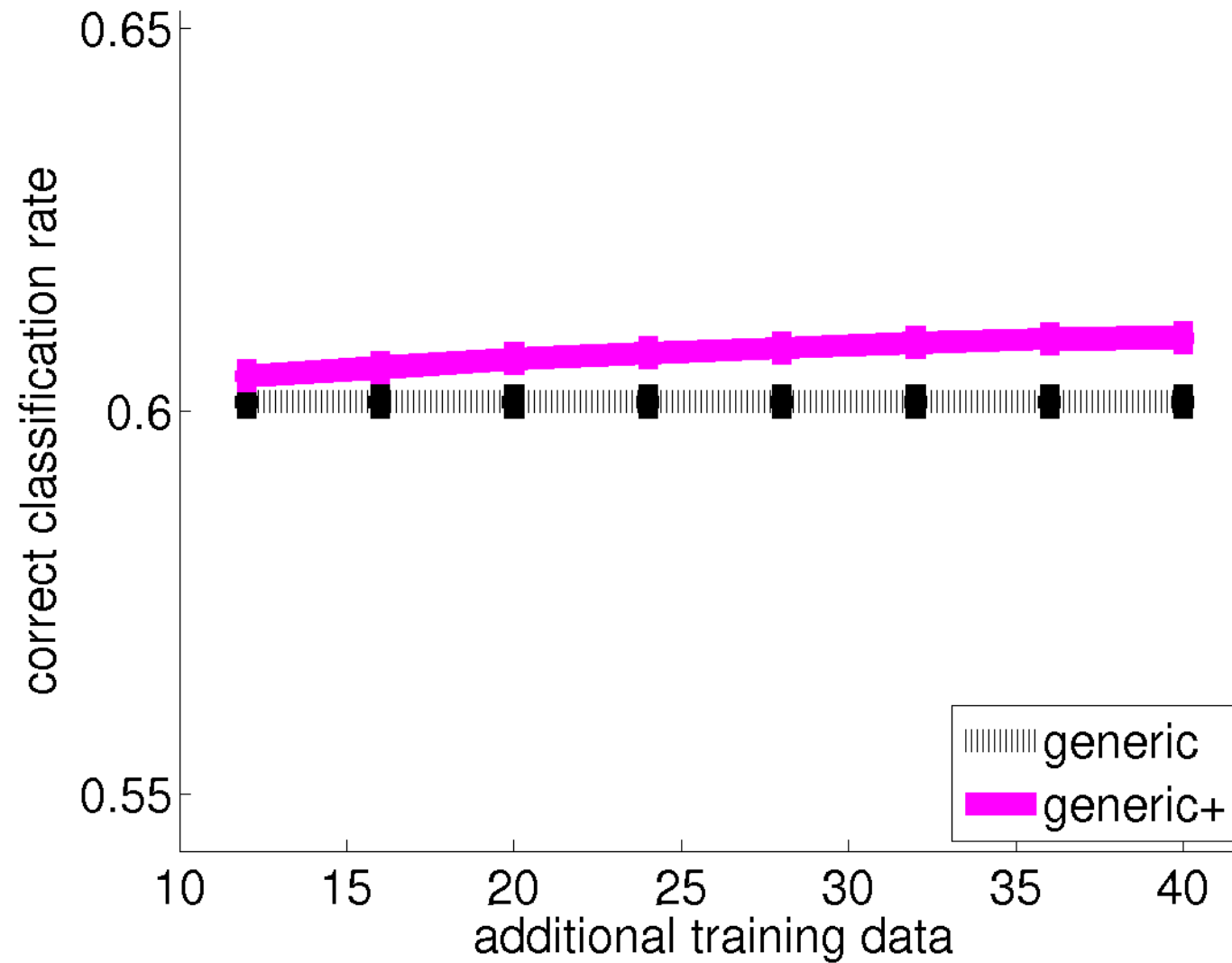


Adapted attribute accuracy



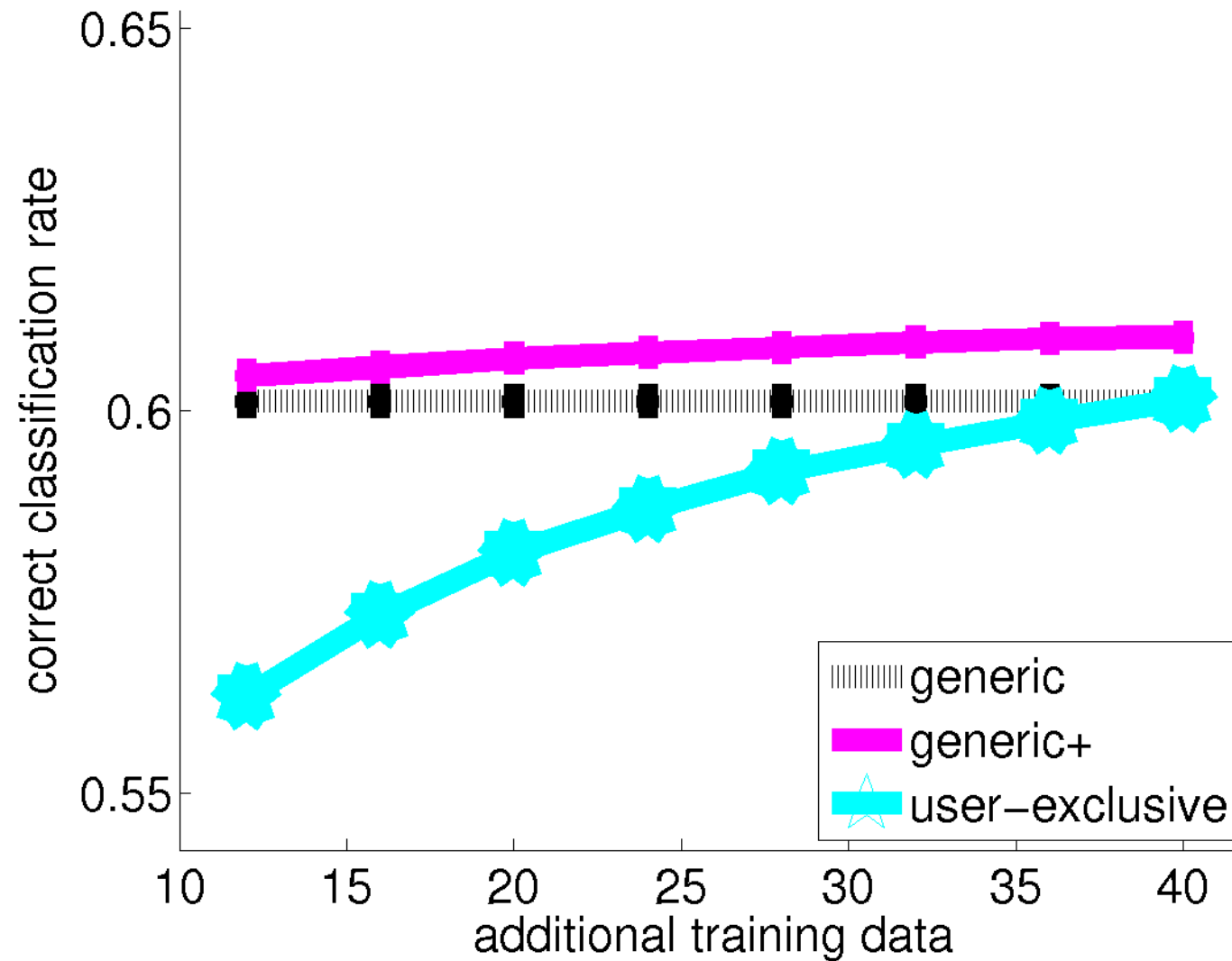
- 3 datasets
- 22 attributes
- 75 total users

Adapted attribute accuracy



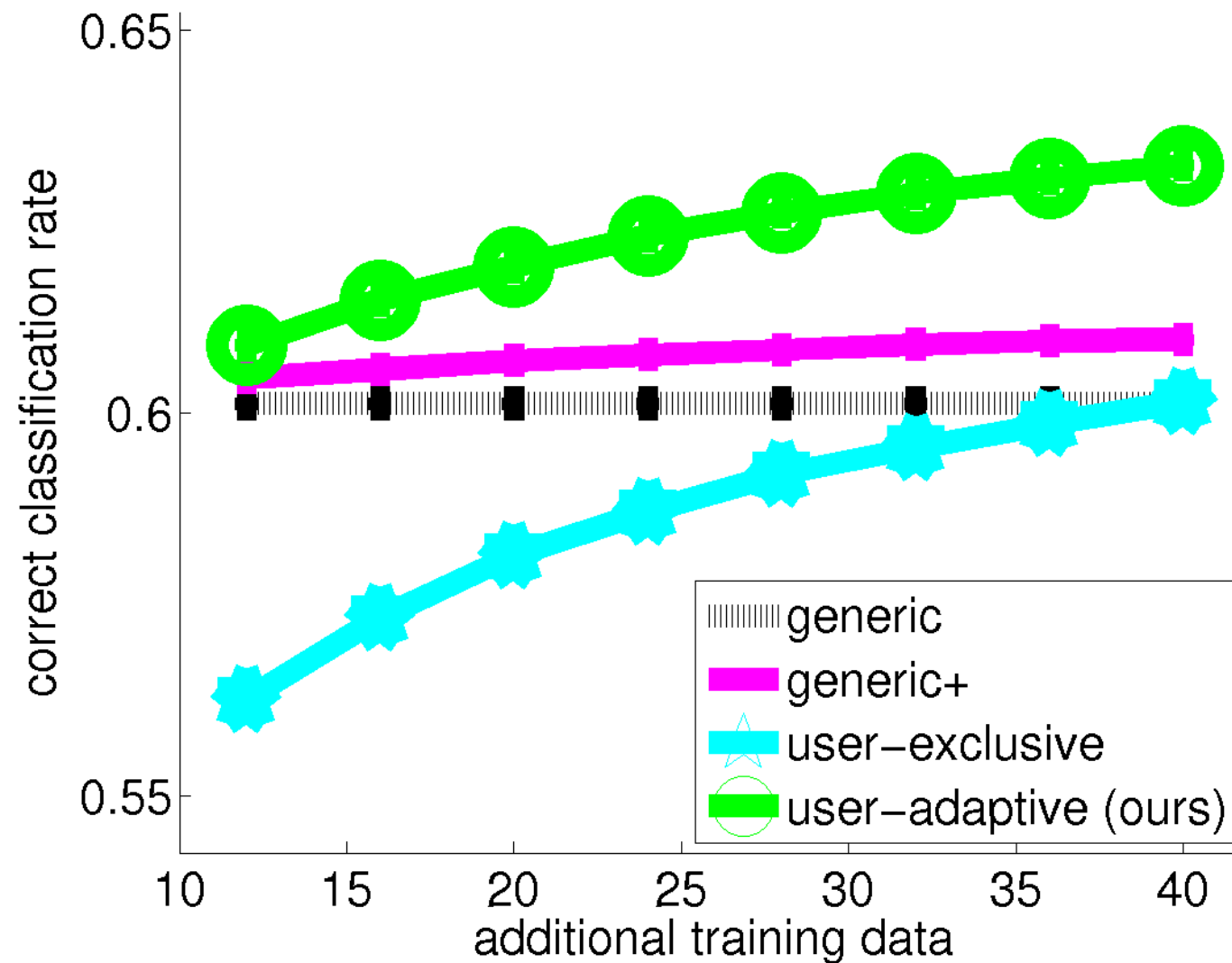
- 3 datasets
- 22 attributes
- 75 total users

Adapted attribute accuracy



- 3 datasets
- 22 attributes
- 75 total users

Adapted attribute accuracy



Adaptation approach most accurately captures perceived attributes

Which images most influence adaptation?



pointy



open



bright



ornamented



shiny



high-heeled



long



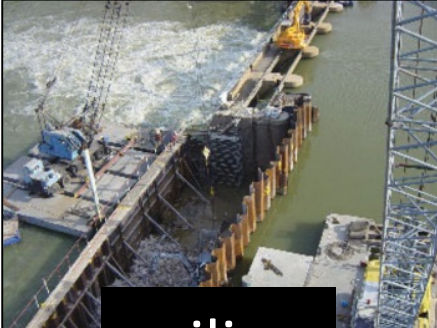
formal



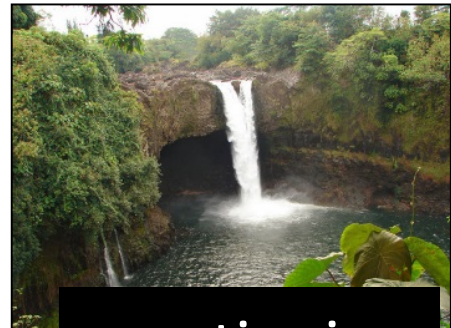
sporty



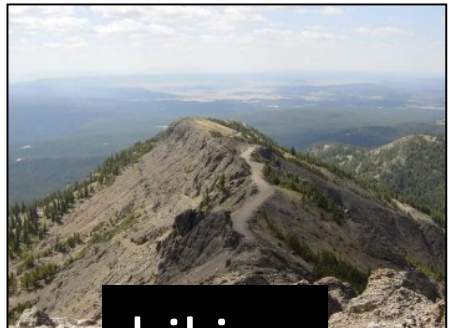
feminine



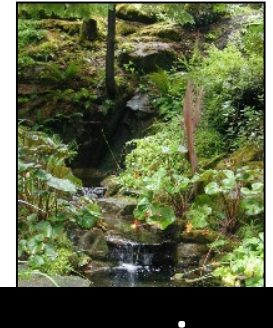
sailing



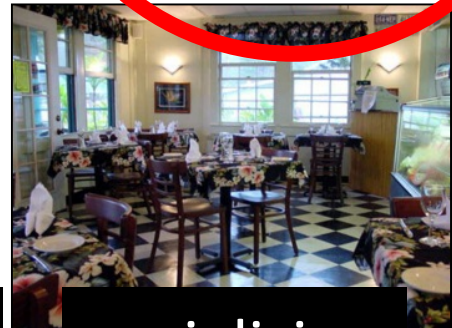
vacationing



hiking



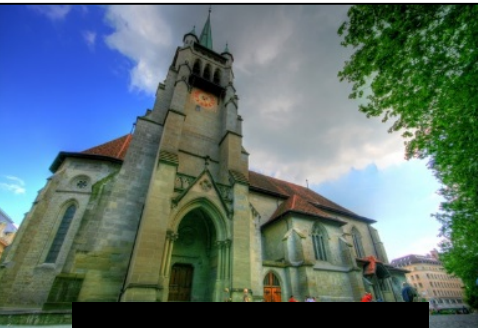
camping



socializing



shopping



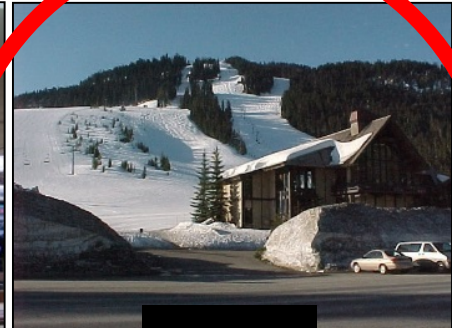
vegetation



clouds



natural light



cold



open area



horizon far

Visualizing adapted attributes

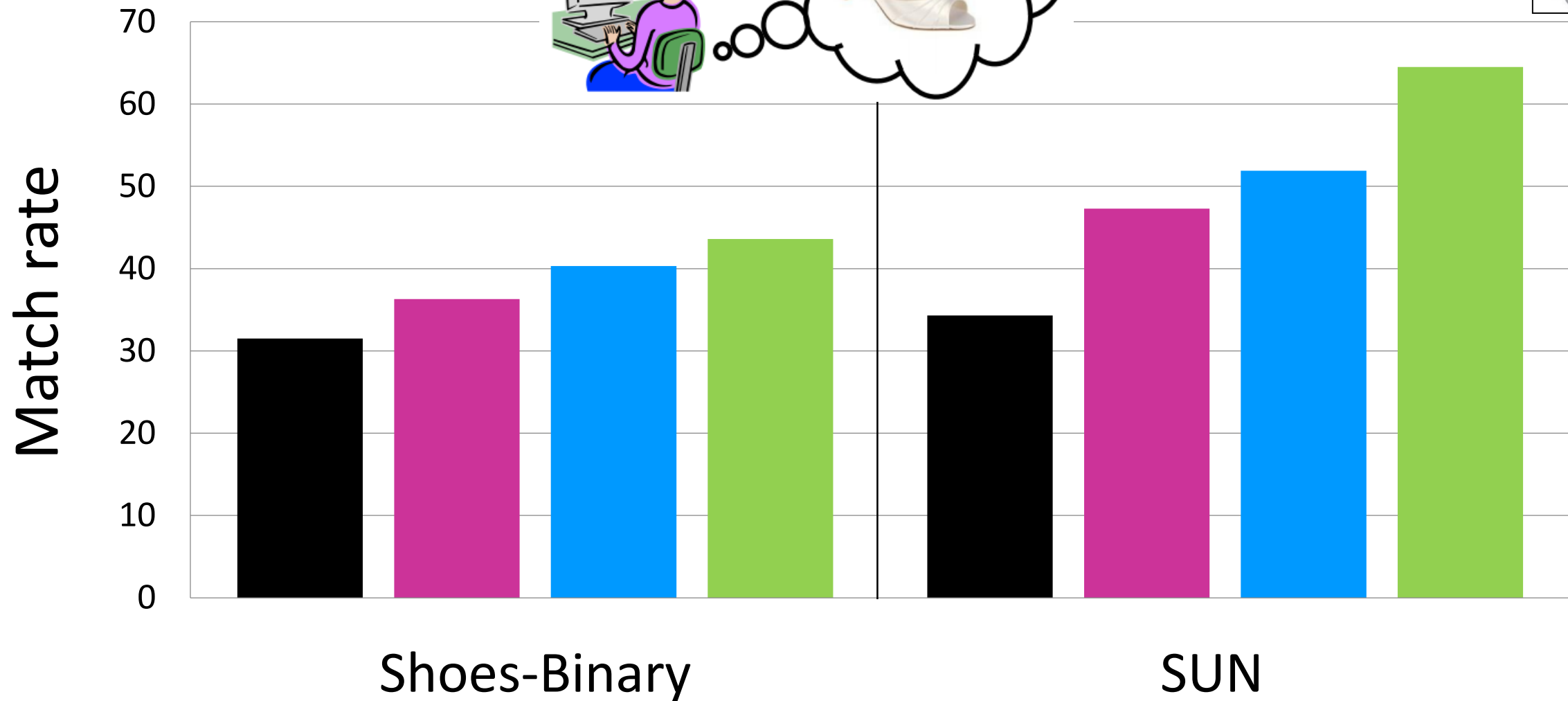
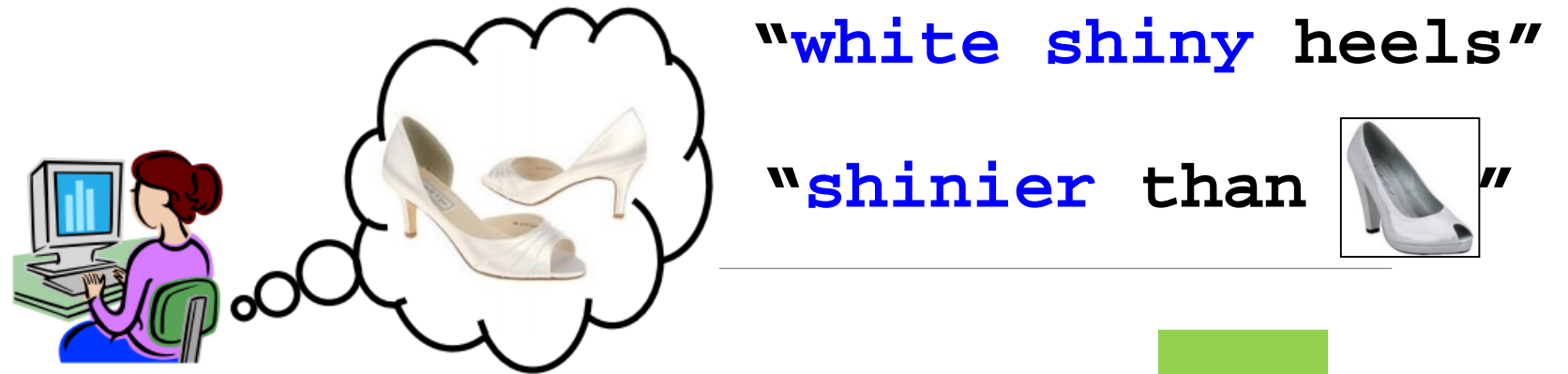
Shoes – Relative Attributes – “Formal”



SUN – Binary Attributes – “Vacationing”

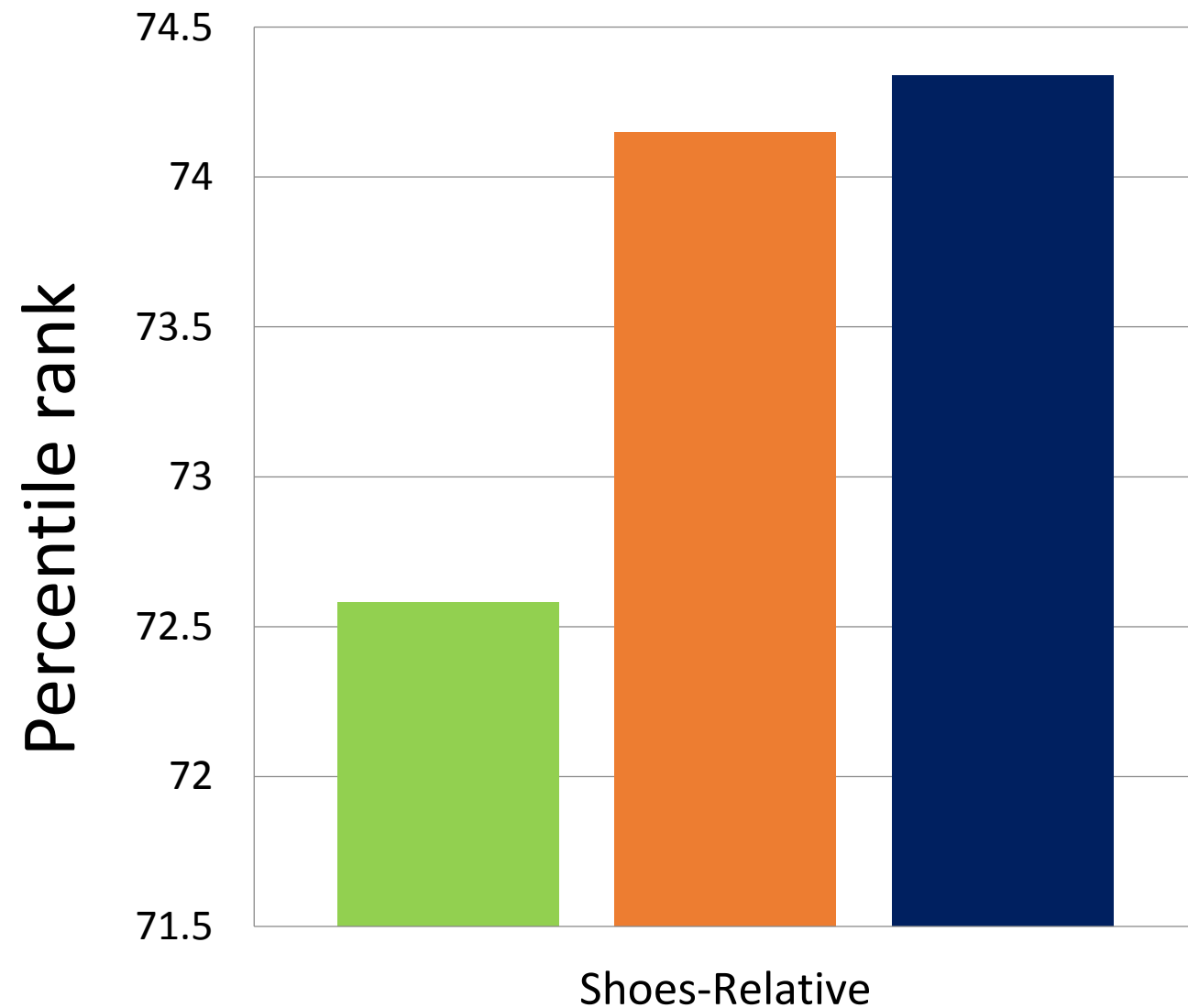
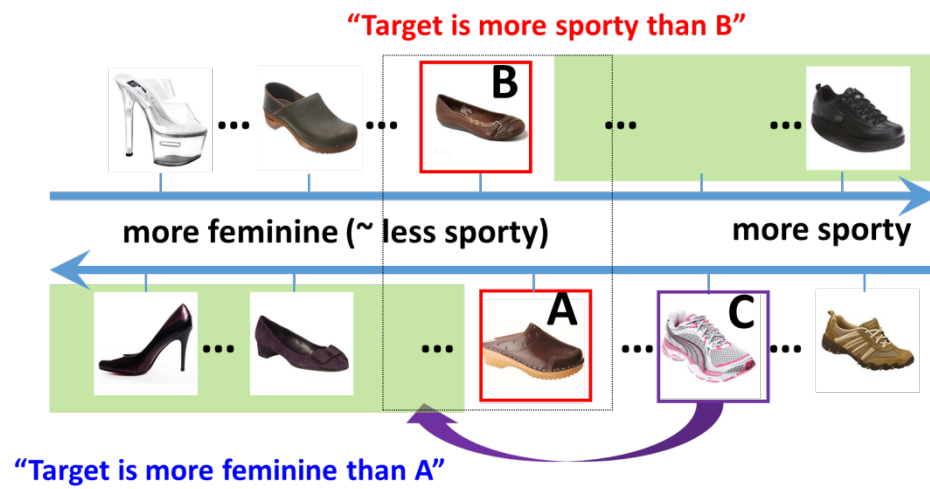


Personalizing image search with adapted attributes



■ generic ■ generic+ ■ user-exclusive ■ user-adaptive

Impact of implicit labels



- explicit labels only
- +contradictions
- +transitivity

Summary

- **Practical concerns if learning visual categories:**
 - Test images can look different from training images!
 - People do not perceive image labels universally!
- **Domain adaptation methods help address them**
 - Landmark-based unsupervised adaptation
 - Reshaping datasets into latent domains
 - Adapt generic models to account for user-specific perception of attributes

References

- [Attribute Adaptation for Personalized Image Search](#). A. Kovashka and K. Grauman. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Sydney, Australia, December 2013.
- [Reshaping Visual Datasets for Domain Adaptation](#). B. Gong, K. Grauman, and F. Sha. In Proceedings of Advances in Neural Information Processing Systems (NIPS), Tahoe, Nevada, December 2013.
- [Connecting the Dots with Landmarks: Discriminatively Learning Domain-Invariant Features for Unsupervised Domain Adaptation](#). B. Gong, K. Grauman, and F. Sha. In International Conference on Machine Learning (ICML), Atlanta, GA, June 2013.
- [Geodesic Flow Kernel for Unsupervised Domain Adaptation](#). B. Gong, Y. Shi, F. Sha, and K. Grauman. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Providence, RI, June 2012.