# Teaching visual recognition systems

Kristen Grauman
Department of Computer Science
University of Texas at Austin

Work with Sudheendra Vijayanarasimhan, Prateek Jain, Devi Parikh, Adriana Kovashka, and Jeff Donahue
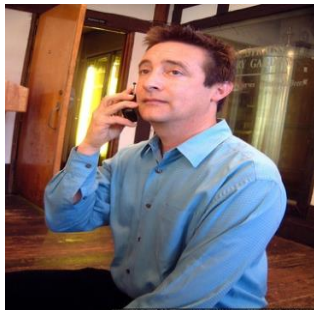
# Visual categories

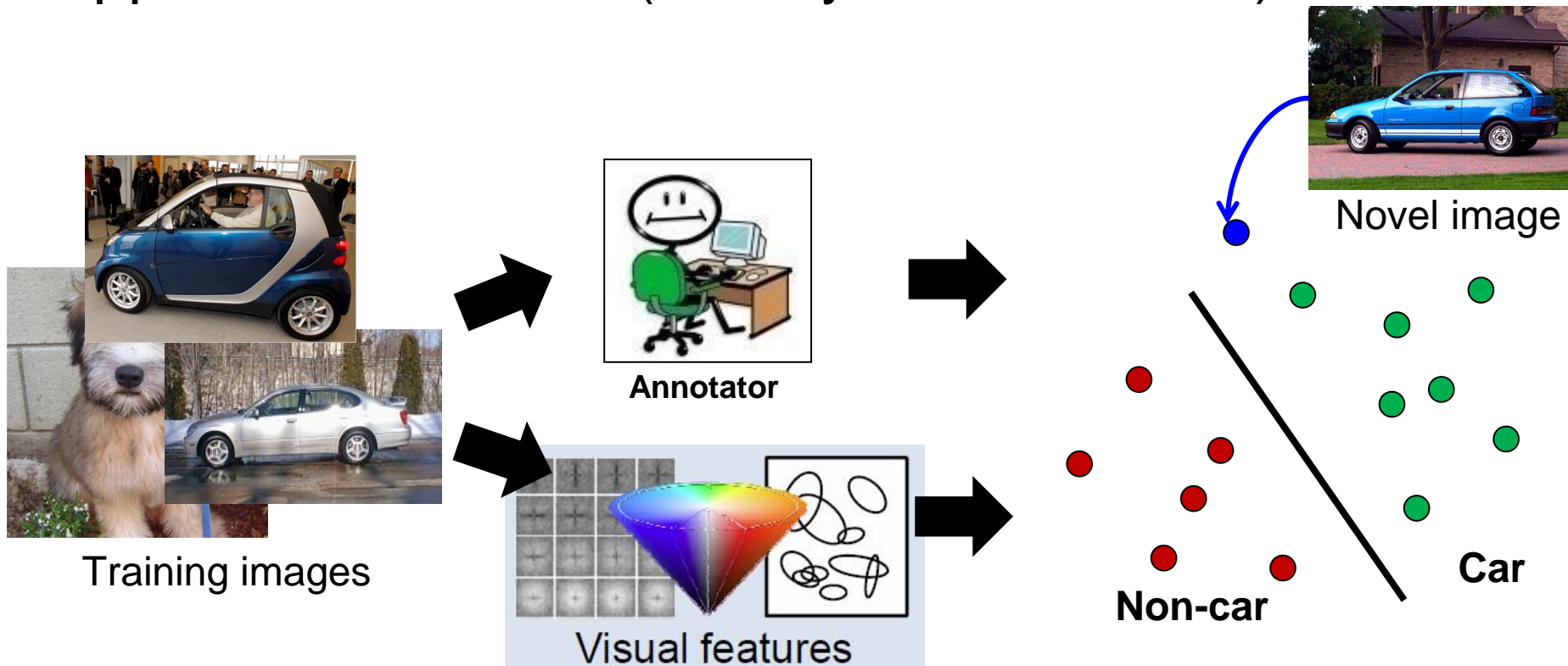Beyond instances, need to recognize and detect *classes* of visually and semantically related…

Objects

Scenes

Activities

# Learning-based methods

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



Training images

Annotator

Visual features

Novel image

Non-car

Car

# Exuberance for image data
# (and their category labels)


**ImageNet**

14M images
1K+ labeled object categories
[Deng et al. 2009-2012]


**80M Tiny Images**

80M images
53K noisily labeled object categories
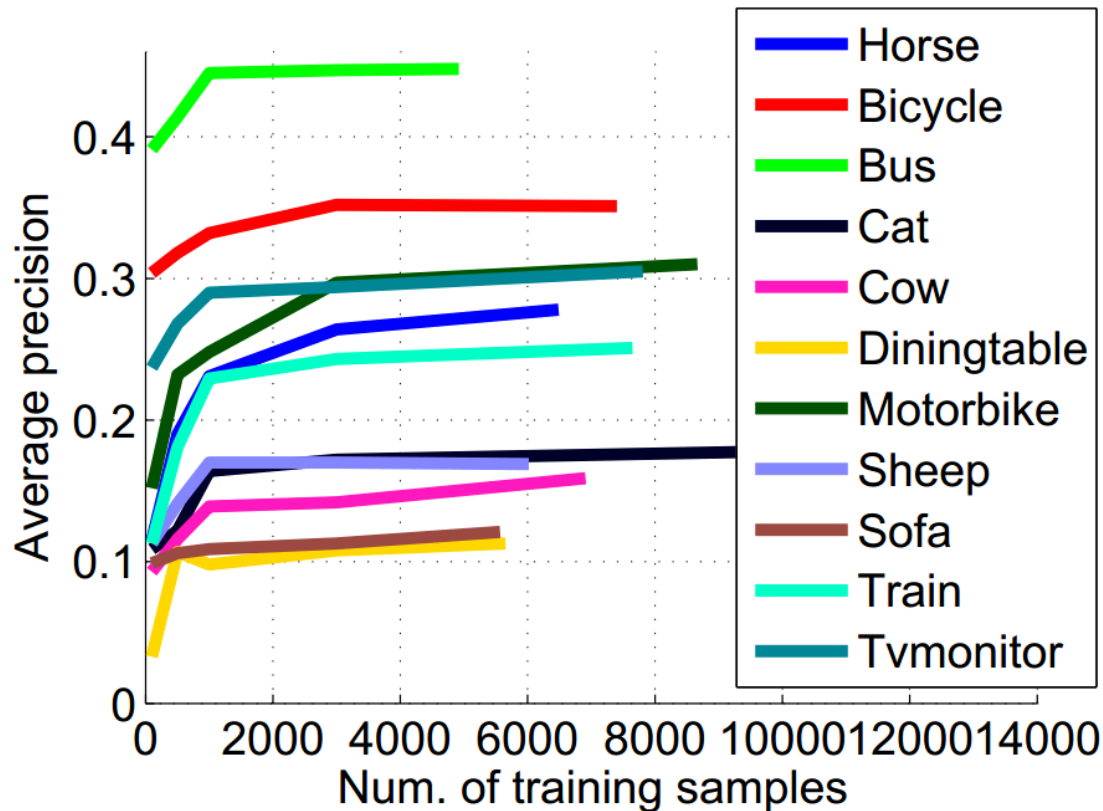[Torralba et al. 2008]


**SUN Database**

131K images
902 labeled scene categories
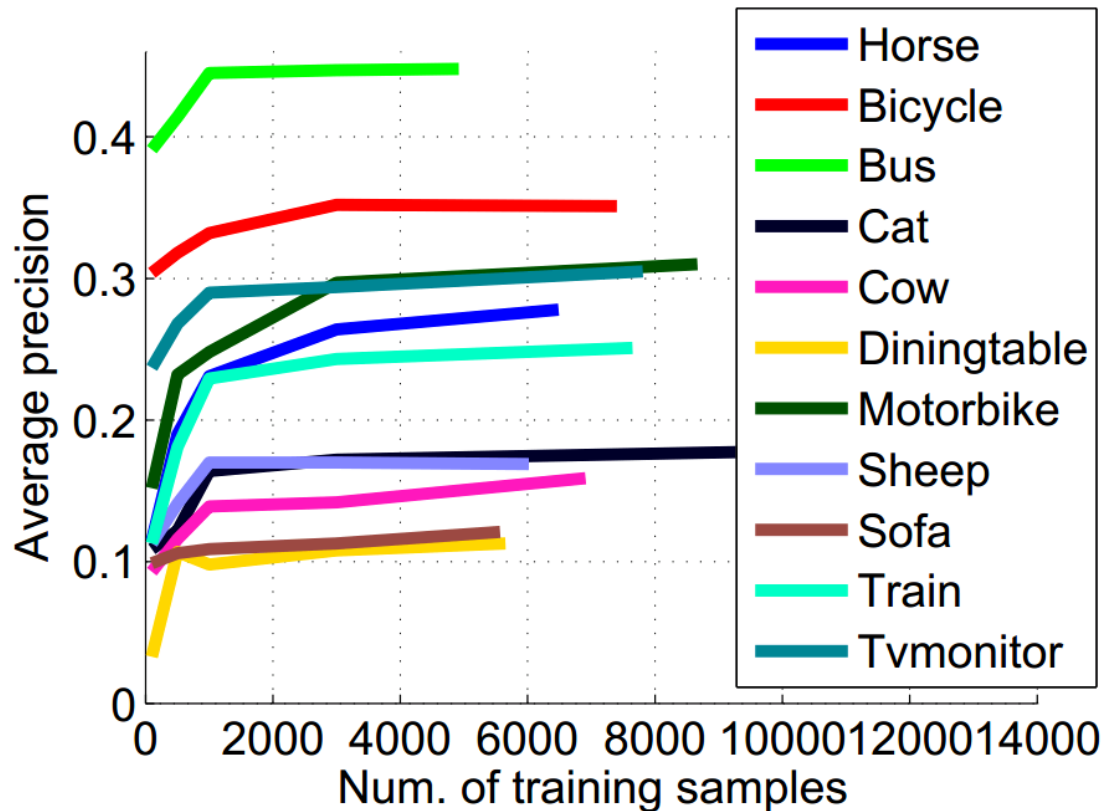4K labeled object categories
[Xiao et al. 2010]

# And yet…

- More data ↔ more accurate visual models?
- Which



*X. Zhu, C. Vondrick, D. Ramanan and C. Fowlkes. Do We Need More Training Data or Better Models for Object Detection?  BMVC 2012.*

# And yet…

- More data ↔ more accurate visual models?



*X. Zhu, C. Vondrick, D. Ramanan and C. Fowlkes. Do We Need More Training Data or Better Models for Object Detection? BMVC 2012.*

# And yet…

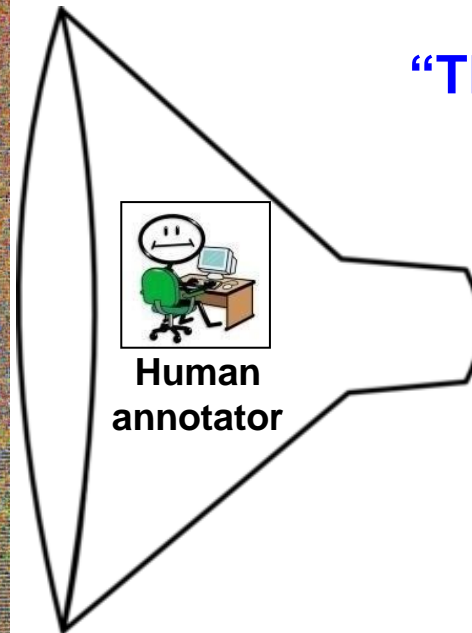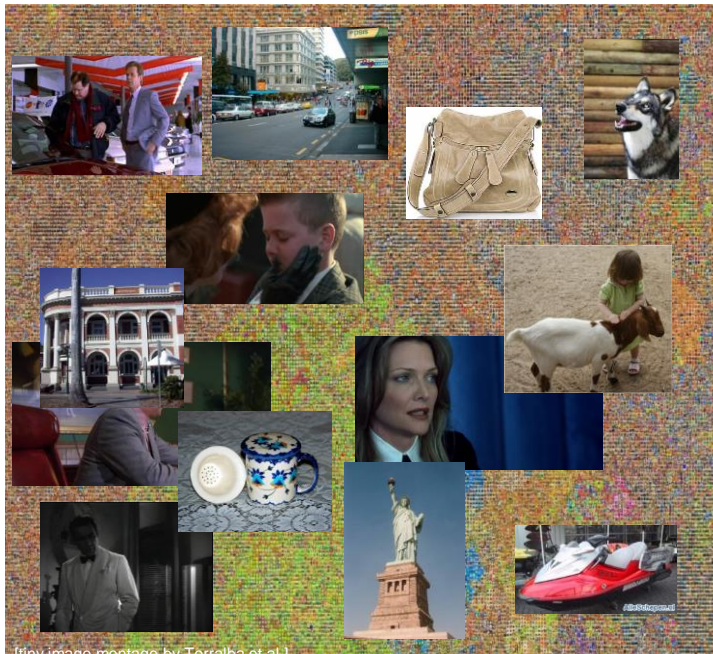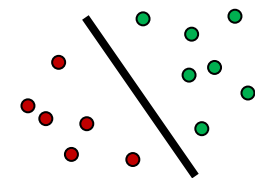- More data ↔ more accurate visual models?
- Which images should be labeled?
- Are labels enough to teach visual concepts?



**Human annotator**

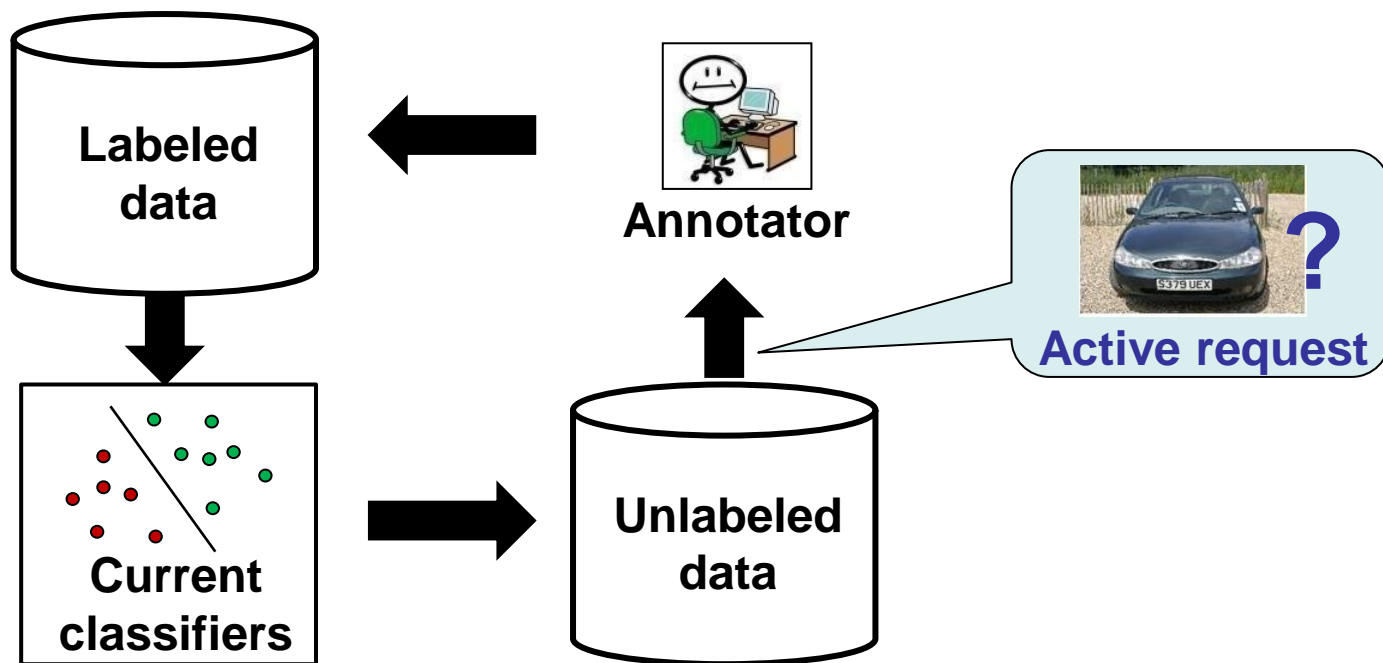**"This image has a cow in it."**

[tiny image montage by Torralba et al.]

# This lecture

Teaching machines visual categories
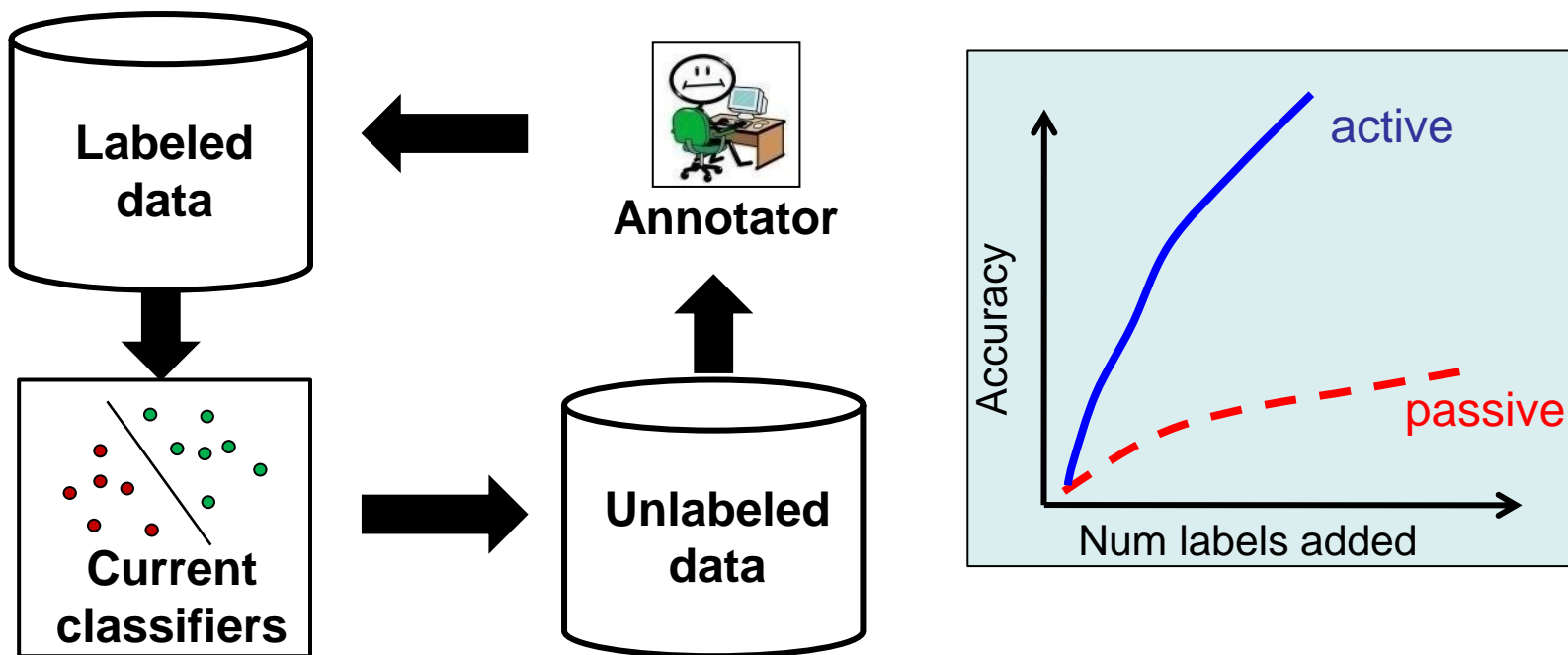
- Active learning to prioritize informative annotations

- Relative attributes to learn from visual comparisons

Kristen Grauman, UT-Austin

# Active learning for image annotation

# Active learning for image annotation

**Labeled data**

**Annotator**

**Current classifiers**

**Unlabeled data**

Accuracy

active

passive

Num labels added

**Intent:** better models, faster/cheaper
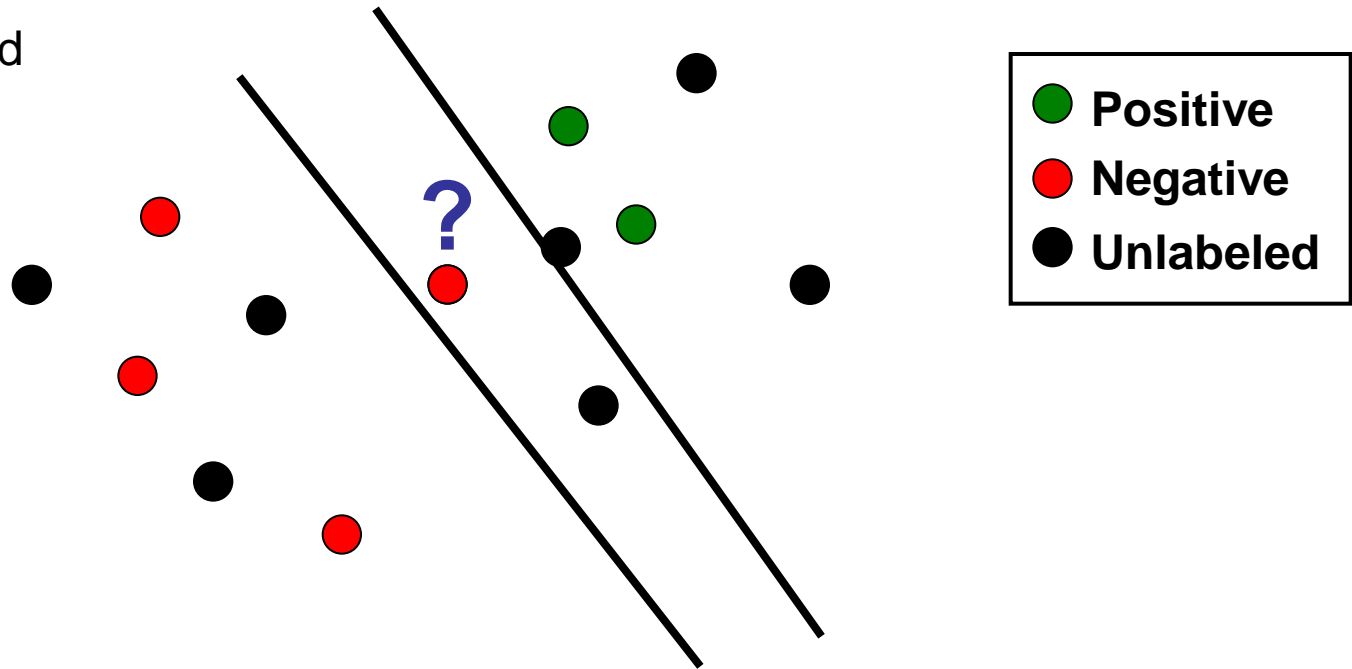
# Active selection

- **Traditional active learning**: obtain most informative labels first.

e.g., margin-based criterion



**Positive**
**Negative**
**Unlabeled**

*[Mackay 1992, Cohn et al. 1996, Freund et al. 1997, Lindenbaum et al. 1999, Tong & Koller 2000, Schohn and Cohn 2000, Campbell et al. 2000, Roy & McCallum 2001, Kapoor et al. 2007,…]*

Kristen Grauman, UT-Austin

# Problem: Active selection and recognition

**Less expensive to obtain**

**More expensive to obtain**

- **Multiple levels** of annotation are possible

- **Variable cost** depending on level *and* example

- **Many annotators** working simultaneously

# Our idea: Cost-sensitive multi-question active learning

- Compute decision-theoretic active selection criterion that weighs both:

    – which *example* to annotate, and

    – what *kind* of annotation to request for it

  as compared to

    – the *predicted effort* the request would require

Kristen Grauman, UT-Austin

# Our idea: Cost-sensitive multi-question active learning



Most regions are understood, but this region is unclear.

This looks expensive to annotate, and it does not seem informative.

This looks expensive to annotate, but it seems very informative.

This looks easy to annotate, but its content is already understood.

Kristen Grauman, UT-Austin

# Multiple-instance learning (MIL)

*negative*

*positive*

*positive bags*

*negative bags*

**Traditional supervised learning**

**Multiple-instance learning**

*[Dietterich et al. 1997]*

# Multiple-instance learning (MIL)



Positive bag
Negative bag

- **Positive instance**: Segment belonging to class
- **Negative instance**: Segment not in class
- **Positive bag**: Image containing class
- **Negative bag**: Image not containing class

*[Dietterich et al.; Maron & Ratan, Yang & Lozano-Perez, Andrews et al.,…]*

Kristen Grauman, UT-Austin

# Multi-question active queries

- Predict which query will be most informative, given the cost of obtaining the annotation.

- Three levels (types) to choose from:

1. Label a region

2. Tag an object in the image

3. Segment the image, name all objects.

Kristen Grauman, UT-Austin

# Decision-theoretic multi-question criterion

$$\mathrm{VALUE}(O, Q) = \mathrm{RISK}(\mathcal{X}_L, \mathcal{X}_U) - \widehat{\mathrm{RISK}}(\mathcal{X}_L \cup O_A, \mathcal{X}_U \setminus O) - \mathrm{COST}(O, Q)$$

Value of asking given question about given data object

Current misclassification risk

Estimated risk if candidate request were answered

Cost of getting the answer

Estimate risk of incorporating the candidate before obtaining true answer $A$ by computing expected value:

$$\widehat{\mathrm{RISK}}(\mathcal{X}_L \cup O_A, \mathcal{X}_U \setminus O) = \sum_{\ell \in \mathbb{L}} \mathrm{RISK}(\mathcal{X}_L \cup O_\ell, \mathcal{X}_U \setminus O) \; p(\ell | O)$$

where $\mathbb{L}$ is set of all possible answers.



For *M* regions $O = \{o_1, \ldots, o_M\}$

$$\approx \frac{1}{S} \sum_{k=1}^{S} \mathrm{RISK}\left(\mathcal{X}_L \cup \{o_1^{(a_1)_k}, \ldots, o_M^{(a_M)_k}\}, \mathcal{X}_U \setminus O\right)$$

# Decision-theoretic multi-question criterion

$$\mathrm{VALUE}(O, Q) = \mathrm{RISK}(\mathcal{X}_L, \mathcal{X}_U) - \widehat{\mathrm{RISK}}(\mathcal{X}_L \cup O_A, \mathcal{X}_U \setminus O) - \mathrm{COST}(O, Q)$$

Current misclassification risk    Estimated risk if candidate request were answered    Cost of getting the answer

Estimate risk of incorporating the candidate before obtaining true answer $A$ by computing expected value:

$$\widehat{\mathrm{RISK}}(\mathcal{X}_L \cup O_A, \mathcal{X}_U \setminus O) = \sum_{\ell \in \mathbb{L}} \mathrm{RISK}(\mathcal{X}_L \cup O_\ell, \mathcal{X}_U \setminus O) \ p(\ell | O)$$

where $\mathbb{L}$ is set of all possible answers.

Cost of the answer: domain knowledge, or directly predict.

Kristen Grauman, UT-Austin

# Predicting effort

- What manual effort cost would we expect to pay for an unlabeled image?



*Which image would you rather annotate?*

# Predicting effort

- What manual effort cost would we expect to pay for an unlabeled image?



*Which image would you rather annotate?*

**Other forms of annotation cost**: expertise required, resolution of data, length of video clips,…

# Learning from annotation examples

Extract cost-indicative image features, train regressor to map features to times.



Interface on Mechanical Turk



…
32 s
24 s
48 s

Collect about 50 responses per training image.

# Predicting effort

# Predicting effort



Predicted Time (Secs) vs Actual Time (Secs)

# Multi-question active learning



**Labeled data**

**Current classifiers**

**Annotator**

**Unlabeled data**

"**Completely segment image #32.**"

"*Does image #7 contain a cow?*"

# Multi-question active learning curves



Region features: texture and color

# Multi-question active learning with objects and attributes

*[Kovashka et al., ICCV 2011]*



Weigh relative impact of an object label or an attribute label, at each iteration.

# Budgeted batch active learning

*[Vijayanarasimhan et al., CVPR 2010]*



$$S^* = \text{argmax } \text{Pred.Gain}(S)$$

$$s.t. \sum_{x \in S} \text{LabelCost}(x) \leq \text{Budget}$$

Select *batch* of examples that together improves classifier objective *and* meets annotation *budget*.

# **Problem**: "Sandbox" active learning

Thus far, tested only in artificial settings:

- Unlabeled data already fixed, small scale, biased



~$10^3$ prepared images

- Computational cost ignored

# **Our idea**: Live active learning

Large-scale active learning of object detectors
with crawled data and crowdsourced labels.

*How to scale active learning to massive unlabeled
  pools of data?*

# SVM margin criterion
# for active selection



Select point nearest to hyperplane decision boundary for labeling.

$$\mathbf{x}^* = \mathrm{argmin}_{\mathbf{x}_i \in \mathcal{U}} \left| \mathbf{w}^T \mathbf{x}_i \right|$$

*[Tong & Koller, 2000; Schohn & Cohn, 2000; Campbell et al. 2000]*

# Sub-linear time active selection

*[Jain, Vijayanarasimhan, Grauman, NIPS 2010]*

We propose a novel hashing approach to identify the most uncertain examples in sub-linear time.



**Current classifier**

$h(w)$

**Unlabeled data**

$h(x)$

| 110 | ■ ■ |
| 101 | ■ ■ ■ |
| 111 | ■ ■ |

**Hash table**

**Actively selected examples**

# Background: Locality-Sensitive Hashing

Probability a *random hyperplane* separates two unit vectors depends on the angle between them:



$x_i$  $r$

$x_j$

Bigger angle:
likely to split

**Corresponding hash function:**

$$h_{\boldsymbol{r}}(\boldsymbol{x}) = \begin{cases} 1, & \text{if } \boldsymbol{r}^T \boldsymbol{x} \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

$$r_i \sim \mathcal{N}(0, 1)$$

**Probability of collision:**

$$\Pr(h_r(x_i) = h_r(x_j)) = 1 - \frac{1}{\pi} \cos^{-1}(x_i^T x_j)$$

*[Goemans and Williamson 1995, Charikar 2004]*

# Hashing a hyperplane query

To retrieve those points for which $\left| \mathbf{w}^T \mathbf{x}_i \right|$ is small, want probable collision for **perpendicular** vectors:



Assuming normalized data.

Should collide

Should not collide

# Hashing a hyperplane query

We generate two independent random vectors **u** and **v:**

- one to constrain angle between **x** and **w**
- one to constrain angle between **x** and **−w**

Collision likely only if <u>neither</u> vector splits

For parallel vectors



Unlikely to split  and  Likely to split

= Likely to split

For perpendicular vectors



Less likely to split  and Less likely to split

= Unlikely to split

# Hashing a hyperplane query

- We define an asymmetric 2-bit hash function:

**H-Hash** family:

$$h_{\mathcal{H}}(\boldsymbol{z}) = \begin{cases} h_{\boldsymbol{u},\boldsymbol{v}}(\boldsymbol{z}, \boldsymbol{z}), & \text{if } \boldsymbol{z} \text{ is a database point vector,} \\ h_{\boldsymbol{u},\boldsymbol{v}}(\boldsymbol{z}, -\boldsymbol{z}), & \text{if } \boldsymbol{z} \text{ is a query hyperplane vector.} \end{cases}$$

where $h_{\boldsymbol{u},\boldsymbol{v}}(\boldsymbol{a}, \boldsymbol{b}) = [h_{\boldsymbol{u}}(\boldsymbol{a}), h_{\boldsymbol{v}}(\boldsymbol{b})] = [\text{sign}(\boldsymbol{u}^T \boldsymbol{a}), \text{sign}(\boldsymbol{v}^T \boldsymbol{b})]$

- We prove necessary conditions for locality sensitivity:

$$\Pr[h_{\mathcal{H}}(\boldsymbol{w}) = h_{\mathcal{H}}(\boldsymbol{x})] = \Pr[h_{\boldsymbol{u}}(\boldsymbol{w}) = h_{\boldsymbol{u}}(\boldsymbol{x})] \Pr[h_{\boldsymbol{v}}(-\boldsymbol{w}) = h_{\boldsymbol{v}}(\boldsymbol{x})]$$

$$= \frac{1}{4} - \frac{1}{\pi^2} \left( \theta_{\boldsymbol{x},\boldsymbol{w}} - \frac{\pi}{2} \right)^2$$

*[Jain, Vijayanarasimhan & Grauman, NIPS 2010]*

# Hashing a hyperplane query

$$h(\mathbf{w}) \rightarrow \{\mathbf{x}_1, \ldots \mathbf{x}_k\}$$



At each iteration of the learning loop, our hash functions map the current hyperplane directly to its nearest unlabeled points.

# Sub-linear time active selection

**Accuracy** improvements as more data labeled

**Time** spent searching for selection

H-Hash result on 1M Tiny Images

By minimizing **both** selection and labeling time, obtain the best accuracy per unit time.

# PASCAL Visual Object Categorization

- Closely studied object detection benchmark
- Original image data from Flickr



http://pascallin.ecs.soton.ac.uk/challenges/VOC/

# Live active learning



**Consensus (Mean shift)**

amazon mechanical turk
Artificial Artificial Intelligence

**Annotated data**

For 4.5 million unlabeled instances,
10 minutes machine time per iter,
vs. 60 hours for a linear scan.

"bicyc

flickr

**Jumping window candidates**

$h(\varphi(O_i))$

**Hash table of image windows**

1111

**Actively selected examples**

**Unlabeled images**

**Unlabeled windows**

*[Vijayanarasimhan & Grauman CVPR 2011]*

# Live active learning results

PASCAL VOC objects - Flickr test set



**Outperforms status quo data collection approach**

# Live active learning results

What does the live learning system ask first?

**Live active learning (ours)**



**Keyword+image baseline**



First selections made when learning "boat"

# PASCAL Live active learning results

Live learning improves some of most difficult
PASCAL VOC categories:

|  | bird | boat | dog | potted plant | sheep | chair |
|---|---|---|---|---|---|---|
| Ours | **15.8**\* | **18.9**\* | **25.3**\* | 11.6\* | **28.4**\* | 9.1\* |
| Previous best | 15.3 | 16.8 | 21.5 | **14.6** | 23.9 | **17.9** |

Our approach's efficiency makes live learning feasible

|  | Active selection | Training | Detection per image |
|---|---|---|---|
| Ours + active | 10 mins | 5 mins | 150 secs |
| LSVM [Felzenszwalb et al. 2009] | 3 hours | 4 hours | 2 secs |
| SP+MKL [Vedaldi et al. 2009] | 93 hours | > 2 days | 67 secs |

Previous best : [Vedaldi et al. ICCV 2009] or [Felzenszwalb et al. PAMI 2009]

# Summary so far

Actively eliciting human insight for visual recognition algorithms.

- **Multi-question active learning** to formulate annotation requests that specify the example *and* the task.

- **Budgeted batch selection** for effective joint selection of multiple requests suited for online annotators.

- **Live active learning** shows large-scale practical impact.

# Ongoing challenges in active visual learning

- Crowdsourcing: reliability, expertise, economics

- Utility tied to specific classifier or model

- Joint batch selection ("non-myopic") expensive, remains challenging

- Active annotations for objects/activity in video

# This lecture

Teaching machines visual categories

- Active learning to prioritize informative annotations

- Relative attributes to learn from visual comparisons

# 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, …]*

Kristen Grauman, UT-Austin

Mule

# Attributes

*A mule…*

Is furry                    Has four-legs

Legs shorter                Tail longer
than horses'                than donkeys'

Has tail

# Binary attributes

*A mule…*

### Is furry

Legs shorter
than horses'

### Has four-legs

Tail longer
than donkeys'

### Has tail

*[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, …]*

# Relative attributes

*A mule…*

Is furry        Has four-legs

**Legs shorter** than horses'       **Tail longer** than donkeys'

Has tail

# Relative attributes

*[Parikh & Grauman, ICCV 2011]*

- Represent ***visual comparisons*** between classes, images, and their properties.

How should relative attributes
be learned?

What do we need to capture
from human annotators?

Kristen Grauman, UT-Austin

# Learning relative attributes

- Learn a ranking function for each attribute, e.g. "brightness".

- Supervision consists of:



$O_m$: ordered pairs

$E_m$: Similar pairs

Kristen Grauman, UT-Austin

# Learning relative attributes

Learn a ranking function

$$a_m(\boldsymbol{x_i}) = \boldsymbol{w}_m^T \boldsymbol{x_i}$$

*Image features*

*Learned parameters*

that best satisfies the constraints:

$$\forall (i,j) \in O_m : \boldsymbol{w}_m^T \boldsymbol{x_i} > \boldsymbol{w}_m^T \boldsymbol{x_j}$$

$$\forall (i,j) \in E_m : \boldsymbol{w}_m^T \boldsymbol{x_i} = \boldsymbol{w}_m^T \boldsymbol{x_j}$$

Kristen Grauman, UT-Austin

# Learning relative attributes

## Max-margin learning to rank formulation

$$\min \quad \left( \frac{1}{2} \|\boldsymbol{w}_{\boldsymbol{m}}^{T}\|_{2}^{2} + C \left( \sum \xi_{ij}^{2} + \sum \gamma_{ij}^{2} \right) \right)$$

$$s.t. \quad \boldsymbol{w}_{\boldsymbol{m}}^{T}(\boldsymbol{x_i} - \boldsymbol{x_j}) \geq 1 - \xi_{ij}$$

$$|\boldsymbol{w}_{\boldsymbol{m}}^{T}(\boldsymbol{x_i} - \boldsymbol{x_j})| \leq \gamma_{ij}$$

$$\xi_{ij} \geq 0; \gamma_{ij} \geq 0$$

Rank margin

$\boldsymbol{W}_m$

Image $\rightarrow$ Relative attribute score

*Joachims, KDD 2002; Parikh and Grauman, ICCV 2011*

Kristen Grauman, UT-Austin

# Relating images

**bright** →



**formal** →



**natural** →



• We can rank images according to attribute strength

# Relating images



Density

Novel image

Conventional binary description: *not dense*

# Relating images

Density

Novel image



*more dense than*          *less dense than*

# Relating images



Density

Novel
image

C C H H *H* C F *H* *H* M *F* *F* I F

*more dense* than **Highways**,
*less dense* than **Forests**

# Relating images

Multi-attribute descriptions offer greater precision when they are relative

**Binary (existing):**

Not Young

BushyEyebrows

RoundFace

(Viggo)

**Relative (ours):**

More Young than CliveOwen
Less Young than ScarlettJohansson

More BushyEyebrows than ZacEfron
Less BushyEyebrows than AlexRodriguez

More RoundFace than CliveOwen
Less RoundFace than ZacEfron

# Applications of relative attributes

Enable new modes of human-system communication

- **Training category models through descriptions**:

  "Rabbits are furrier than dogs."

- **Rationales to explain image labels:**

  "It's not a coastal scene because it's too cluttered."

- **Semantic relative feedback for image search:**

  "I want shoes like these, but shinier."

# Relative zero-shot learning

**Training**: Images from **S seen** categories and

Descriptions of **U unseen** categories



Age:    **Hugh**⤞**Clive**⤞**Scarlett**        **Jared**⤞**Miley**
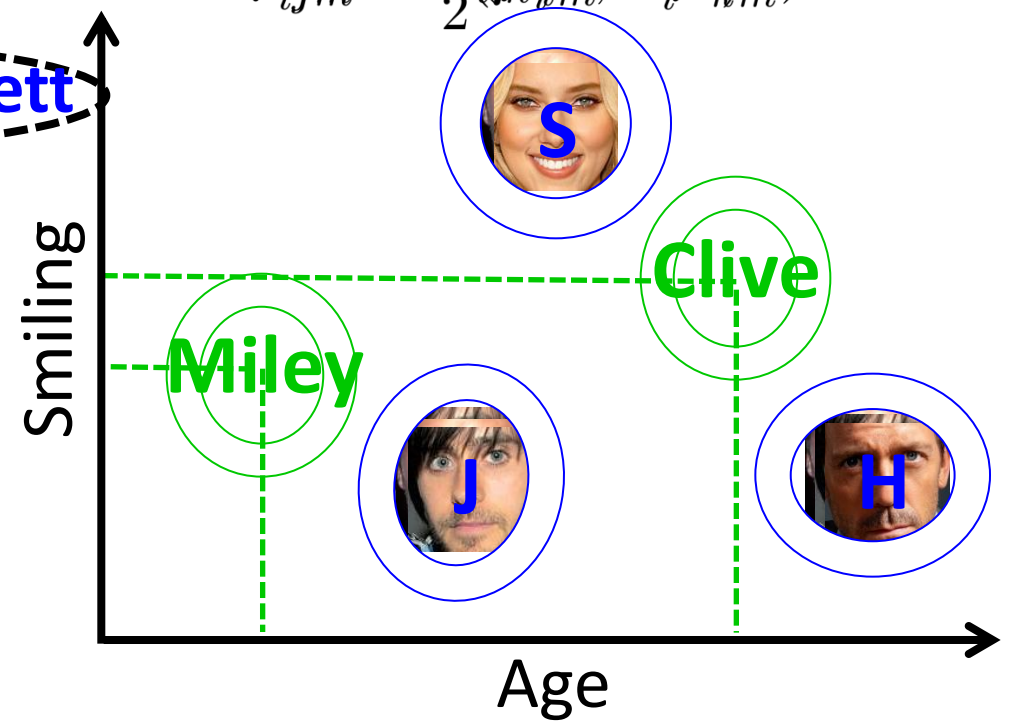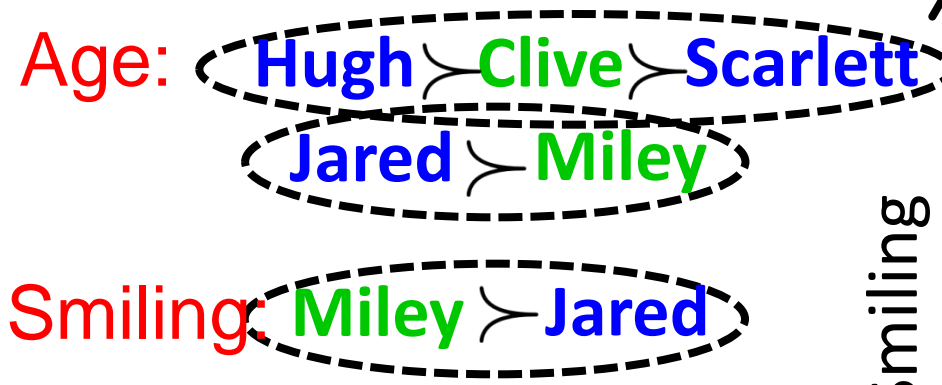
Smiling:

**Miley**⤞**Jared**

Need not use all attributes, nor all seen categories

**Testing**: Categorize image into one of **S**+**U** classes

# Relative zero-shot learning

We can predict new classes based on their **relationships** to existing classes – even without training images.

$$\mu_{ijm}^{(s)} \sim N(\mu_{ijm}^{(s)}, \Sigma_{km}^{(s)})$$

Age: **Hugh** ⪰ **Clive** ⪰ **Scarlett**

**Jared** ⪰ **Miley**

Smiling: **Miley** ⪰ **Jared**



Infer image category using max-likelihood

# Datasets

Outdoor Scene Recognition
(OSR) [Oliva 2001]



8 classes, ~2700 images, Gist
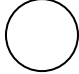6 attributes: open, natural, etc.

Public Figures Faces
(PubFig) [Kumar 2009]



8 classes, ~800 images,
Gist+color

11 attributes: white, chubby, etc.

# Baselines

- **Binary** attributes:
  Direct Attribute Prediction
  [Lampert et al. 2009]

bear  turtle  rabbit

furry

big

- Relative attributes via
  classifier scores

# Relative zero-shot learning



An attribute is more discriminative when used relatively

# Bootstrapped scene learning
# with relative attribute constraints

*[Gupta et al. ECCV 2012]*

Semantic supervision:  **Is More Open**

Amphitheatre  **>**  Barn

Amphitheatre  **>**  Conference Room

Desert  **>**  Barn

**Has Taller Structures**

Church (Outdoor)  **>**  Cemetery

Barn  **>**  Cemetery

# Bootstrapped scene learning

**Labeled Seed Examples**

**Bootstrapping**



[Gupta et al. ECCV 2012]

# Bootstrapped scene learning

**Labeled Seed Examples**

**Bootstrapping**

**Constrained Bootstrapping**

**Amphitheatre**



**Amphitheatre** **Amphitheatre**



Attributes

Indoor

Has Seat Rows

**Auditorium**



**Auditorium** **Auditorium**



Attributes

Has Larger Circular Structures

[Gupta et al. ECCV 2012]

# Applications of relative attributes

Enable new modes of human-system communication

- **Training category models through descriptions**:

    "Rabbits are furrier than dogs."

- **Rationales to explain image labels:**

    "It's not a coastal scene because it's too cluttered."

- **Semantic relative feedback for image search:**

    "I want shoes like these, but shinier."

# Complex visual recognition tasks

*[Donahue and Grauman, ICCV 2011]*



Is the team winning?
**How can you tell?**

Is it a safe route?
**How can you tell?**

Is her form good?
**How can you tell?**

## Main idea:

- Solicit a visual rationale for the label.

- Ask the annotator not just *what*, but also *why.*

# Soliciting visual rationales

**Annotation task**: Is her form good? How can you tell?



☑ `pointed toes`
☑ `balanced`
☐ `falling`
☐ `knee angled`

## Spatial rationale



Synthetic contrast example

## Attribute rationale



Synthetic contrast example

*[Annotator Rationales for Visual Recognition. J. Donahue and K. Grauman, ICCV 2011]*

Kristen Grauman, UT-Austin

# Rationales' influence on the classifier



Decision boundary *refined* in order to satisfy "secondary" margin

$$\text{minimize} \quad \left( \frac{1}{2}||\boldsymbol{w}||^2 + C \sum_i \xi_i + C_c \sum_i \gamma_i \right)$$

$$s.t. \quad y_i \boldsymbol{w}^T \boldsymbol{x_i} \geq 1 - \xi_i; \quad \forall i \in \mathcal{L}$$

$$\boxed{y_i(\boldsymbol{w}^T \boldsymbol{x_i} - \boldsymbol{w}^T \boldsymbol{v_i}) \geq \mu(1 - \gamma_i); \quad \forall i \in \mathcal{C}}$$

$$\xi_i \geq 0; \gamma_i \geq 0,$$

[Zaidan et al. Using Annotator Rationales to Improve Machine Learning for Text Categorization, NAACL HLT 2007]

# Rationale results

- **Scene Categories**: How can you tell the scene category?



- **Hot or Not**: What makes them hot (or not)?



- **Public Figures**: What attributes make them (un)attractive?



Collect rationales from hundreds of MTurk workers.

*[Annotator Rationales for Visual Recognition. J. Donahue and K. Grauman, ICCV 2011]*

Kristen Grauman, UT-Austin

# Example rationales from MTurk

**Scene categories**



Typical    Tight    "Artistic"

**Hot or Not**



Hot, Male  Not, Male  Hot, Female  Not, Female

**PubFig Attractiveness**



*Youth*
*Smiling*
*Straight Hair*
*Narrow Eyes*

*Youth*
*Black Hair*
*Goatee*
*Square Face*
*Shiny Skin*
*High Cheekbones*

# Rationale results

Mean AP



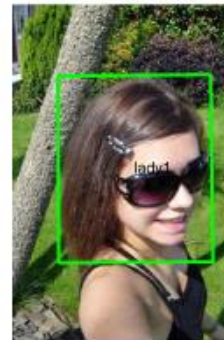| Scenes | Originals | +Rationales |
|---|---|---|
| Kitchen | 0.1196 | **0.1395** |
| Living Rm | 0.1142 | **0.1238** |
| Inside City | 0.1299 | **0.1487** |
| Coast | 0.4243 | **0.4513** |
| Highway | 0.2240 | **0.2379** |
| Bedroom | 0.3011 | **0.3167** |
| Street | 0.0778 | **0.0790** |
| Country | 0.0926 | **0.0950** |
| Mountain | 0.1154 | **0.1158** |
| Office | 0.1051 | **0.1052** |
| Tall Building | 0.0688 | **0.0689** |
| Store | 0.0866 | **0.0867** |
| Forest | 0.3956 | **0.4006** |



| Hot or Not | Originals | +Rationales |
|---|---|---|
| Male | 54.86% | **60.01%** |
| Female | 55.99% | **57.07%** |



| PubFig | Originals | +Rationales |
|---|---|---|
| Male | 64.60% | **68.14%** |
| Female | 51.74% | **55.65%** |

*[Donahue & Grauman, ICCV 2011]*

# Rationale results

Why not just use discriminative feature selection?

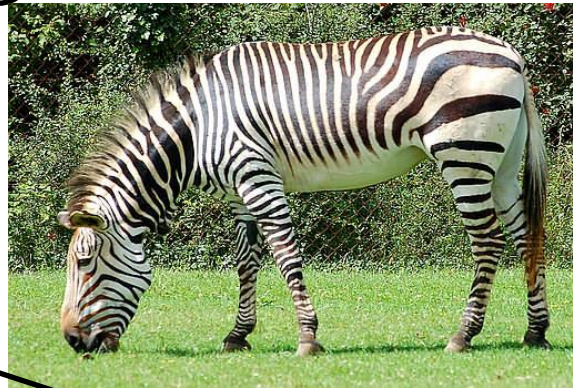| Scenes | Originals | +Rationales | Mutual information |
|---|---|---|---|
| Kitchen | 0.1196 | **0.1395** | 0.1202 |
| Living Rm | 0.1142 | **0.1238** | 0.1159 |
| Inside City | 0.1299 | **0.1487** | 0.1245 |
| Coast | 0.4243 | **0.4513** | 0.4129 |
| Highway | 0.2240 | **0.2379** | 0.2112 |
| Bedroom | 0.3011 | **0.3167** | 0.2927 |
| Street | 0.0778 | **0.0790** | 0.0775 |
| Country | 0.0926 | **0.0950** | 0.0941 |
| Mountain | 0.1154 | **0.1158** | 0.1154 |
| Office | 0.1051 | **0.1052** | 0.1048 |
| Tall Building | 0.0688 | **0.0689** | 0.0686 |
| Store | 0.0866 | **0.0867** | 0.0866 |
| Forest | 0.3956 | **0.4006** | 0.3897 |

Mean AP        *[Donahue & Grauman, ICCV 2011]*

# Relative feedback for object learning

[Parkash & Parikh, ECCV 2012]

Current belief

I think this is a giraffe. What do you think?

No, its neck is too short for it to be a giraffe.

Knowledge of the world

Ah! These must not be giraffes either then.

[Animals with even shorter necks]

......

Feedback on one, transferred to many

Biswas & Parikh, CVPR 2013; Parkash & Parikh, ECCV 2012]

Slide credit: Devi Parikh

# Applications of relative attributes

Enable new modes of human-system communication

- **Training category models through descriptions**:

    "Rabbits are furrier than dogs."

- **Rationales to explain image labels:**

    "It's not a coastal scene because it's too cluttered."

- **Semantic relative feedback for image search:**

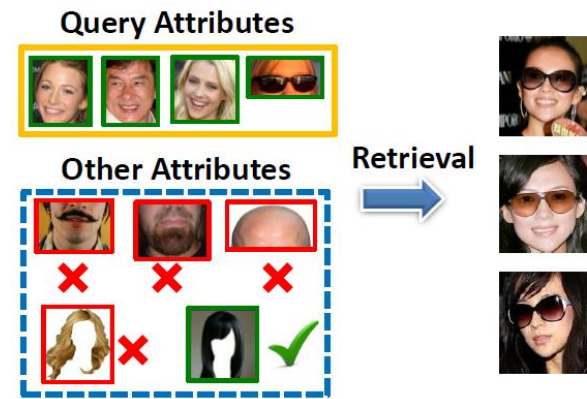    "I want shoes like these, but shinier."

# Attributes for search

Previously, attributes serve as keywords for one-shot search


Kumar et al. 2008
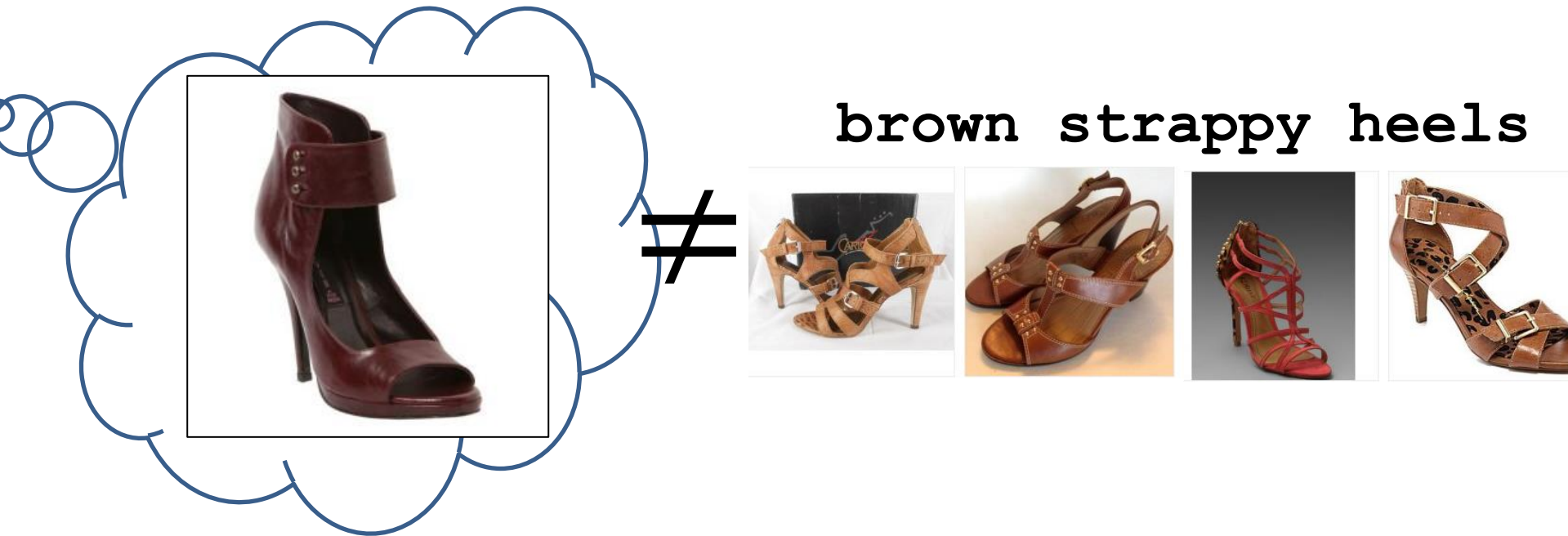

Siddiquie et al. 2011


Vaquero et al. 2009

Kristen Grauman, UT-Austin

# Problem with one-shot visual search

- But keywords (including attributes) can be insufficient to capture target in one shot.



**brown strappy heels**

# Interactive visual search



- Interactive search can help iteratively refine

- …but traditional **binary relevance feedback** offers only coarse communication between user and system

# WhittleSearch: Relative attribute feedback

*[Kovashka et al. CVPR 2012]*

**Query: "white high-heeled shoes"**



*Initial top search results*

**Feedback:**
**"more formal than these"**

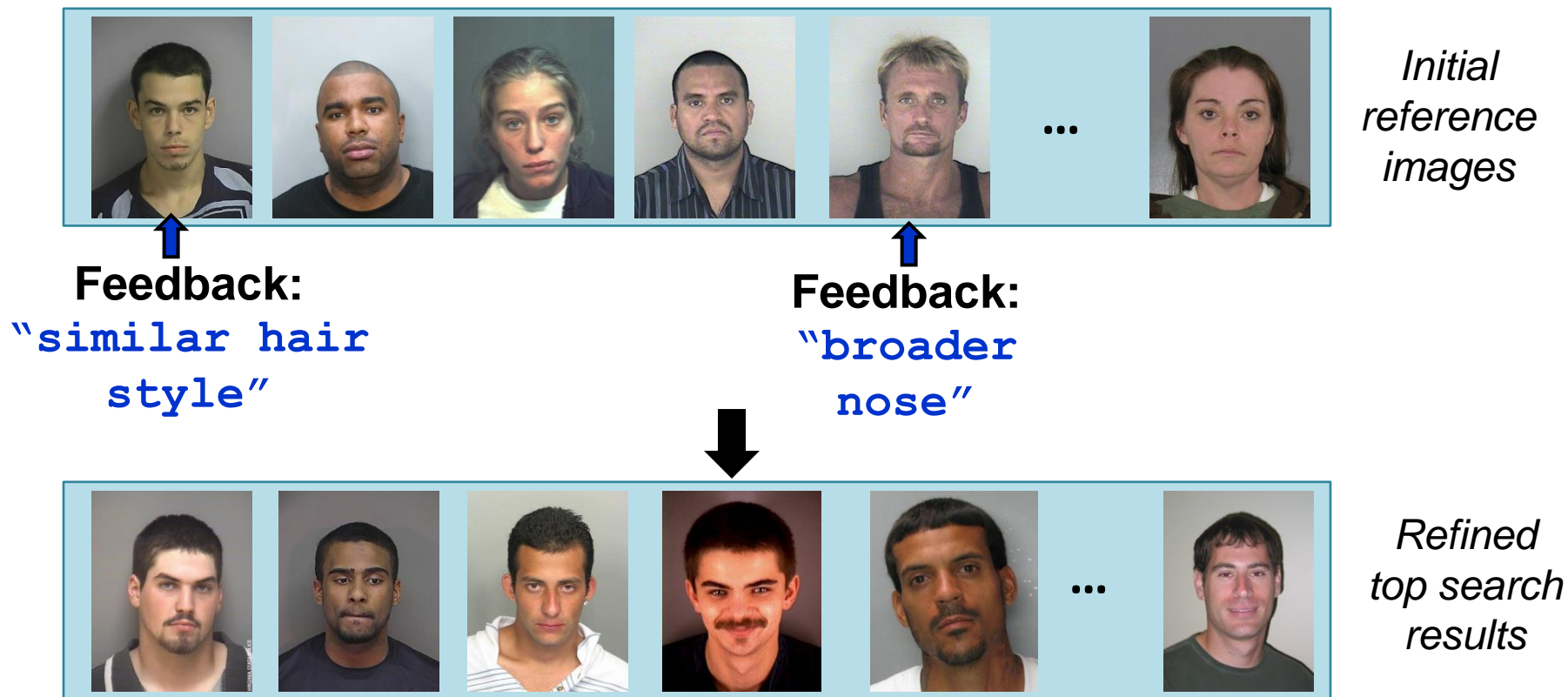**Feedback:**
**"shinier than these"**

*Refined top search results*

Whittle away irrelevant images via precise semantic feedback

# WhittleSearch: Relative attribute feedback

*[Kovashka et al. CVPR 2012]*



*Initial reference images*

**Feedback:**
"similar hair style"

**Feedback:**
"broader nose"

*Refined top search results*

Whittle away irrelevant images via precise semantic feedback

*Kovashka, Parikh, and Grauman, CVPR 2012*

# WhittleSearch with relative attribute feedback

*natural* →



**scores = scores + 1**

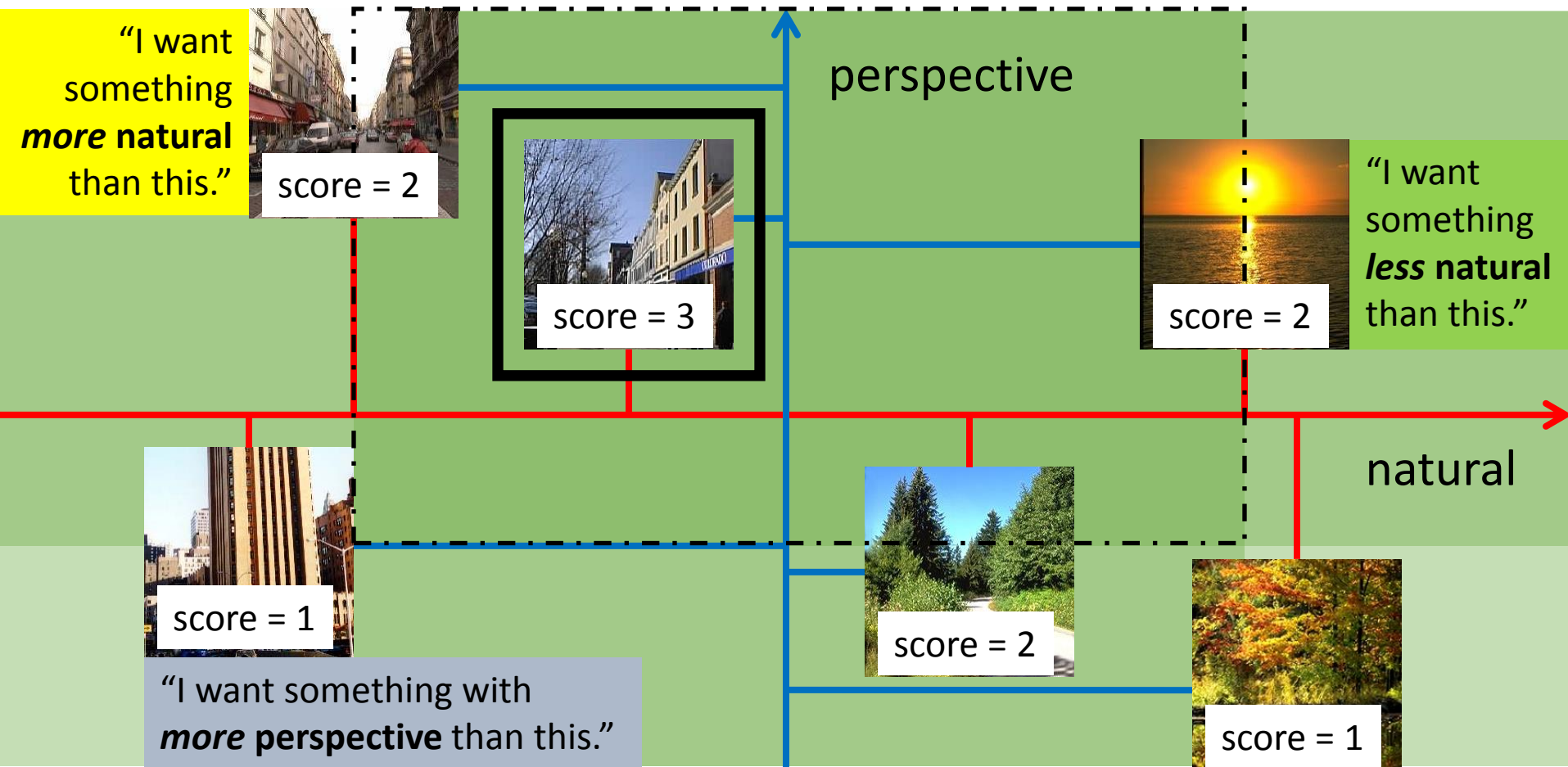"I want something *less* **natural** than this."

**scores = scores + 0**

Offline:

> We learn a spectrum for each attribute

During search:

1. User selects some reference images and marks *how they differ from the* desired target

2. We update the scores for each database image
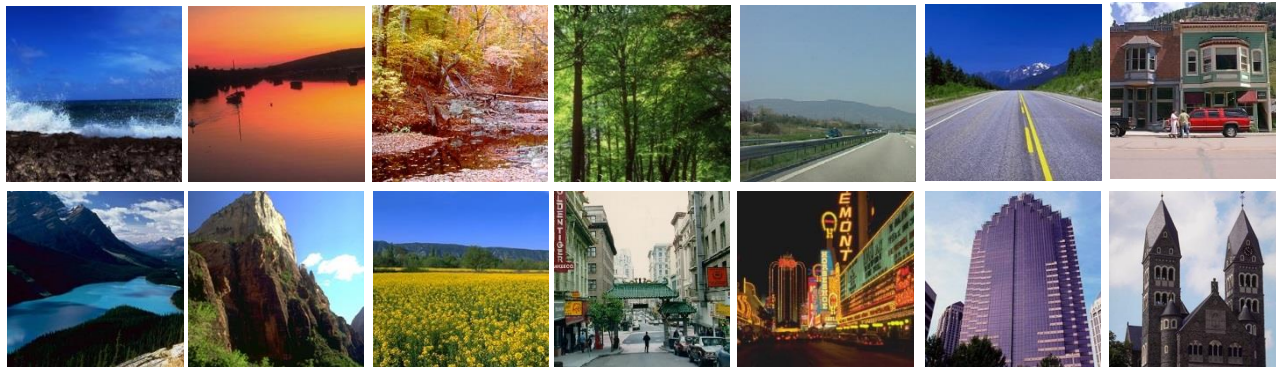
# WhittleSearch with relative attribute feedback



"I want something **more natural** than this."

perspective

score = 2

score = 3

"I want something **less natural** than this."

score = 2

natural

score = 1

"I want something with **more perspective** than this."

score = 2

score = 1

# Datasets



**Shoes: [Berg; Kovashka]**
14,658 shoe images;
10 attributes:
"pointy", "bright", "high-heeled", "feminine" etc.

**OSR: [Oliva & Torralba]**
2,688 scene images;
6 attributes:
"natural", "perspective", "open-air", "close-depth" etc.
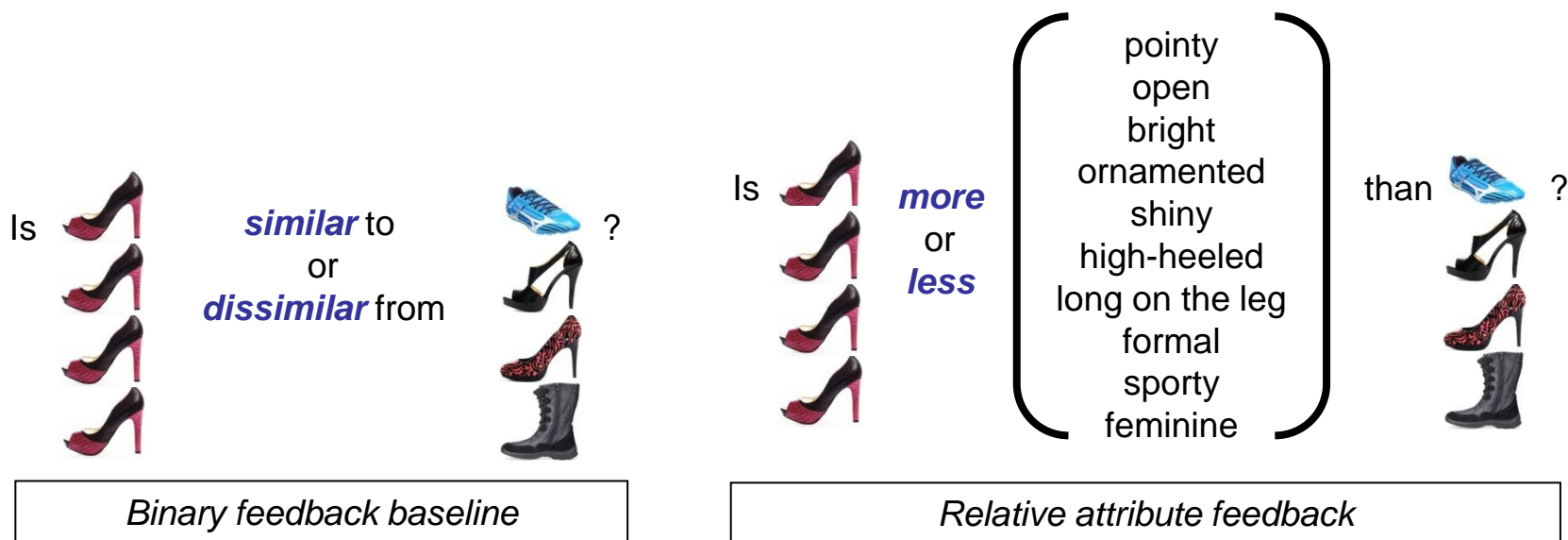
**PubFig: [Kumar et al.]**
772 face images;
11 attributes:
"masculine", "young", "smiling", "round-face", etc.

# Experimental setup

- Give the user the target image to look for

- Pair each target image with 16 reference images
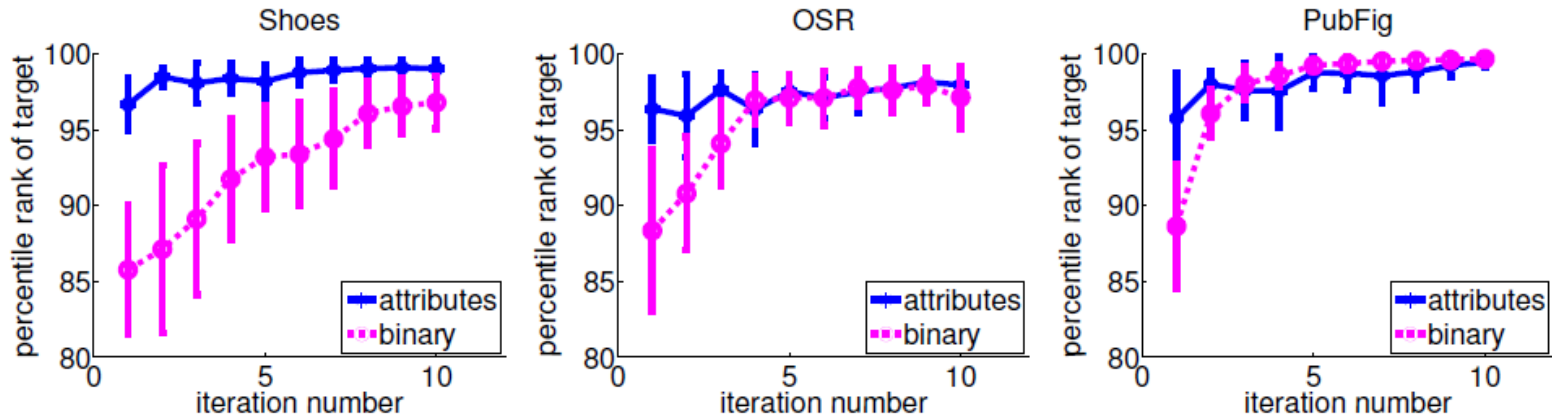
- Get judgments on pairs from users on MTurk

Is *similar* to or *dissimilar* from ?

*Binary feedback baseline*

Is *more* or *less* [ pointy open bright ornamented shiny high-heeled long on the leg formal sporty feminine ] than ?

*Relative attribute feedback*

# WhittleSearch Results



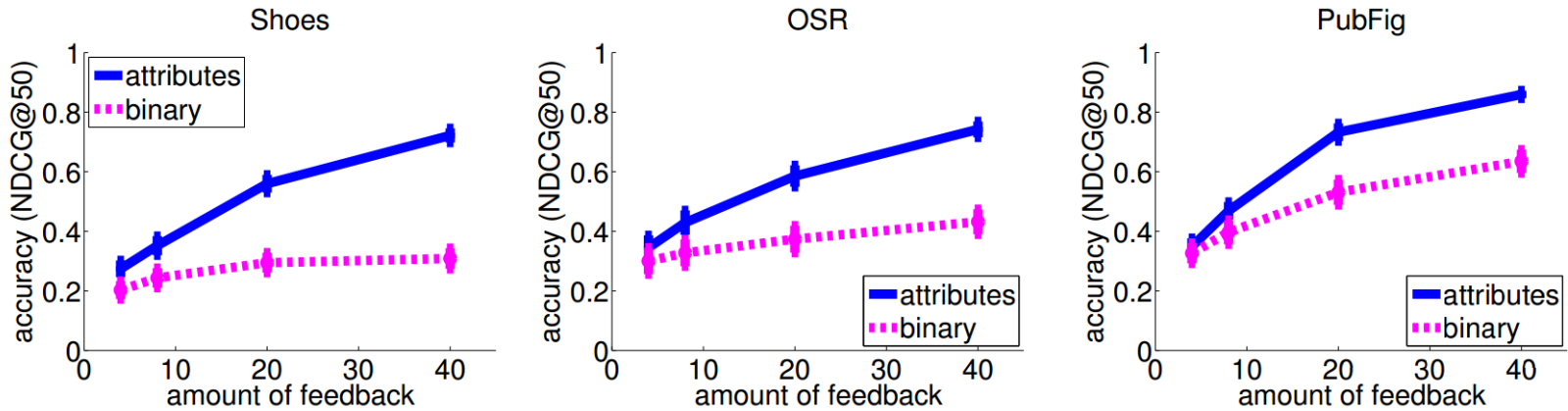Binary relevance feedback

Relative attribute feedback

Kristen Grauman, UT-Austin

# WhittleSearch Results



**We more rapidly converge on the envisioned visual content.**



**Richer feedback → faster gains per unit of user effort.**

Kristen Grauman, UT-Austin

# Example WhittleSearch

**Query:** "I want a bright, open shoe that is short on the leg."

Round 1

More open than

More bright in color than

*Selected feedback*

Less high at the heel than

Less ornaments than

Round 2

Round 3
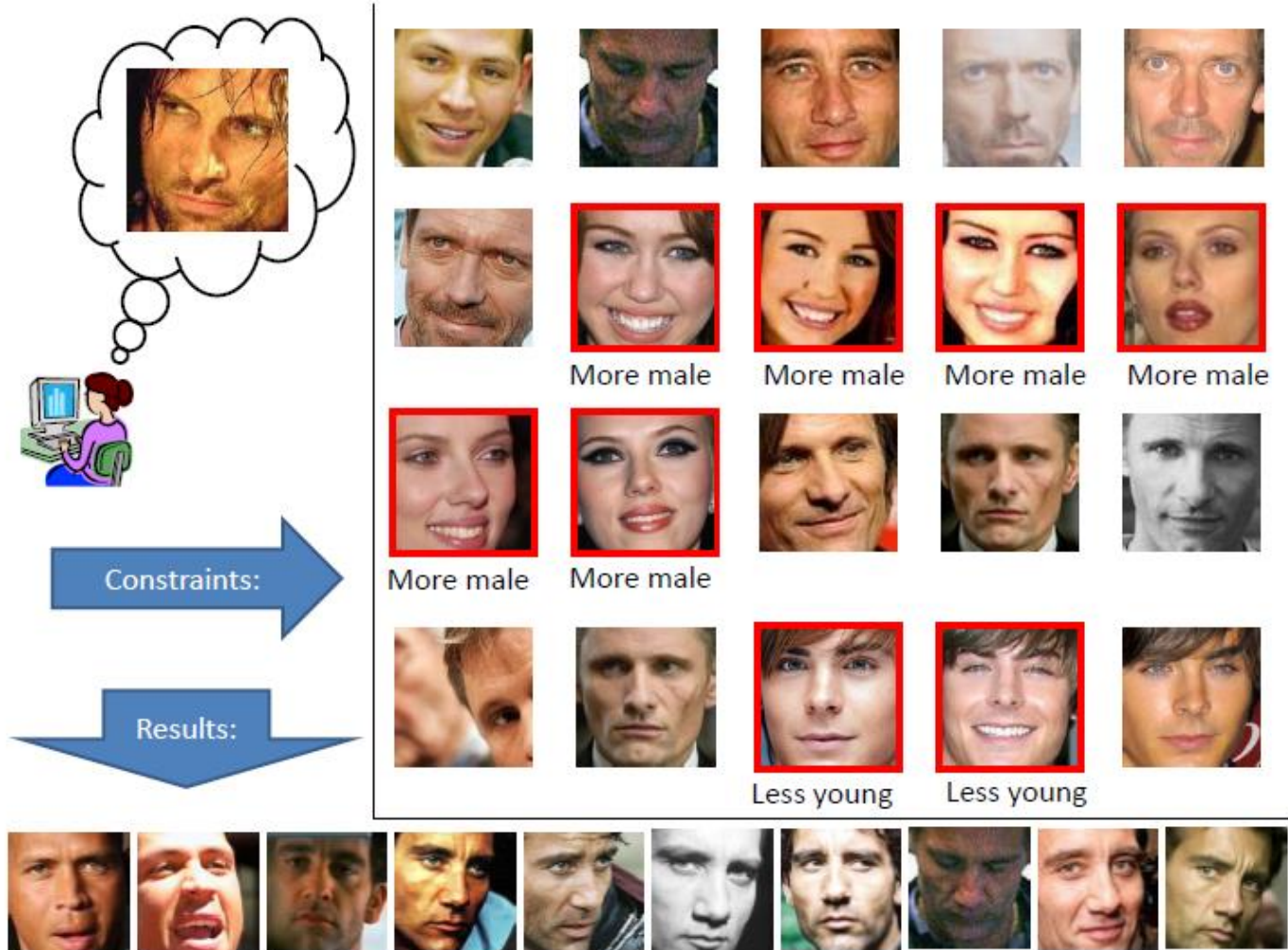
Match

More formal than

More bright in color than

Higher at the heel than

More open than

# Failure case (?)



Is the user searching for a specific person (identity), or an image similar to the specific target image?

# WhittleSearch Demo

## http://godel.ece.vt.edu/whittle/

# Problem: Where is feedback most useful?



Page 1

"*More* **open** than this."

"*Less* **shiny** than this."

"*Less* **sporty** than this."

- The most *relevant* images might not be most *informative*

- Existing active methods focus on binary relevance, expensive selection procedures

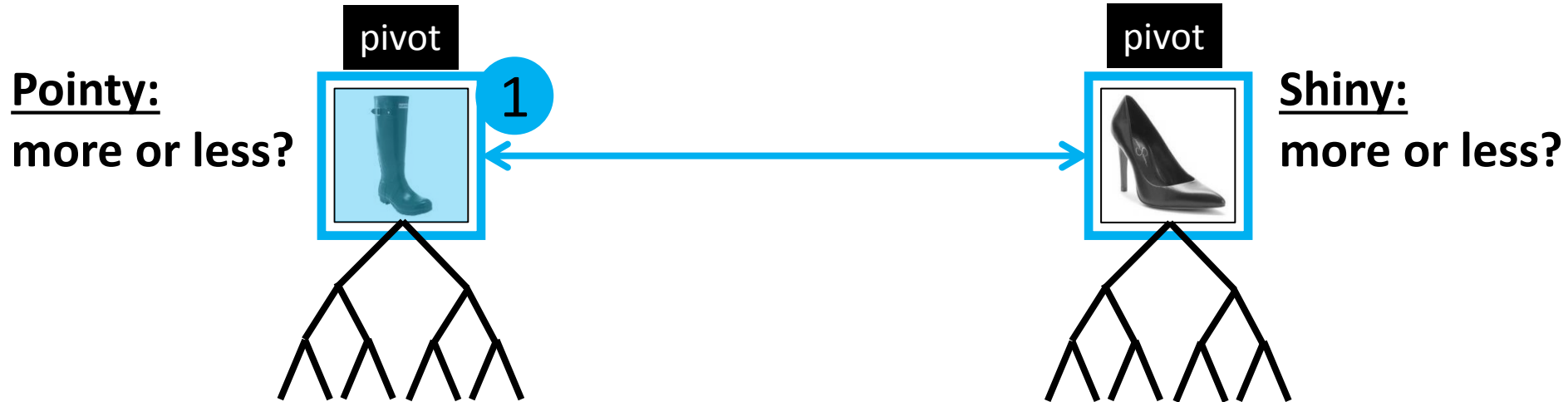    *[Tong & Chang 2001, Li et al. 2001, Cox et al. 2000, Ferecatu & Geman 2007, …]*

# Idea: Attribute Pivots for Guiding Feedback

*[Kovashka and Grauman, 2013]*



**Are the shoes you seek *more* or *less* <u>feminine</u> than** ?

*More*
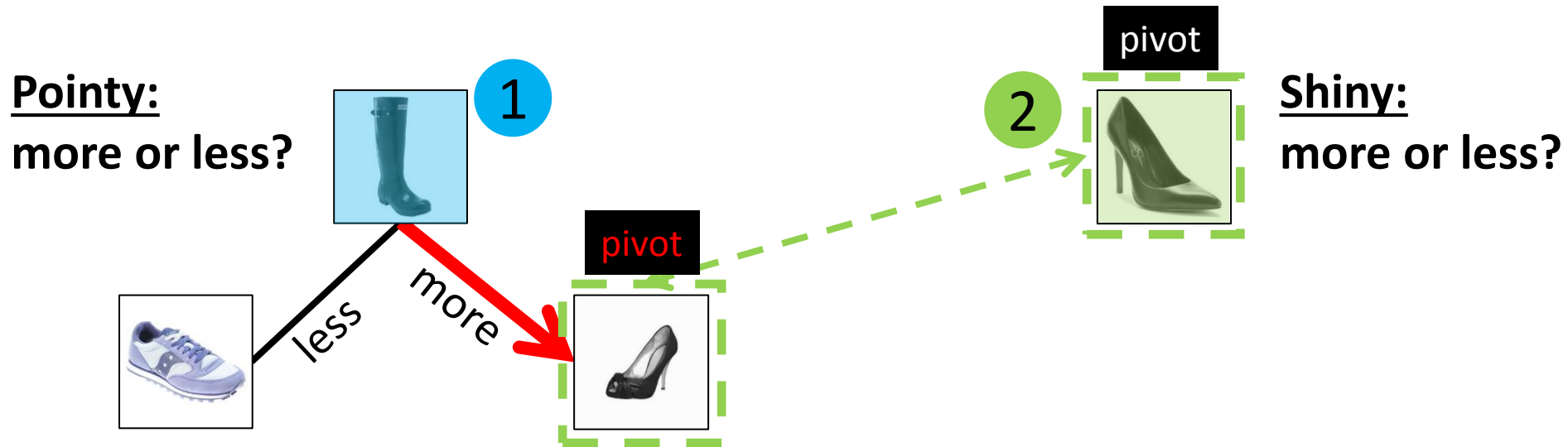
*... more* or *less* <u>bright</u> than ?

*Less*

- Select series of most informative *visual comparisons* that user should make to help deduce target
- Use binary search trees in attribute space for rapid selection

# Selecting a Series of Informative Comparisons



**Pointy:** more or less?

**Shiny:** more or less?

# Selecting a Series of Informative Comparisons



**Pointy:**
**more or less?**

**Shiny:**
**more or less?**

pivot

pivot

less

more

1

2

# Selecting a Series of Informative Comparisons

# Selecting a Series of Informative Comparisons



Kristen Grauman, UT-Austin
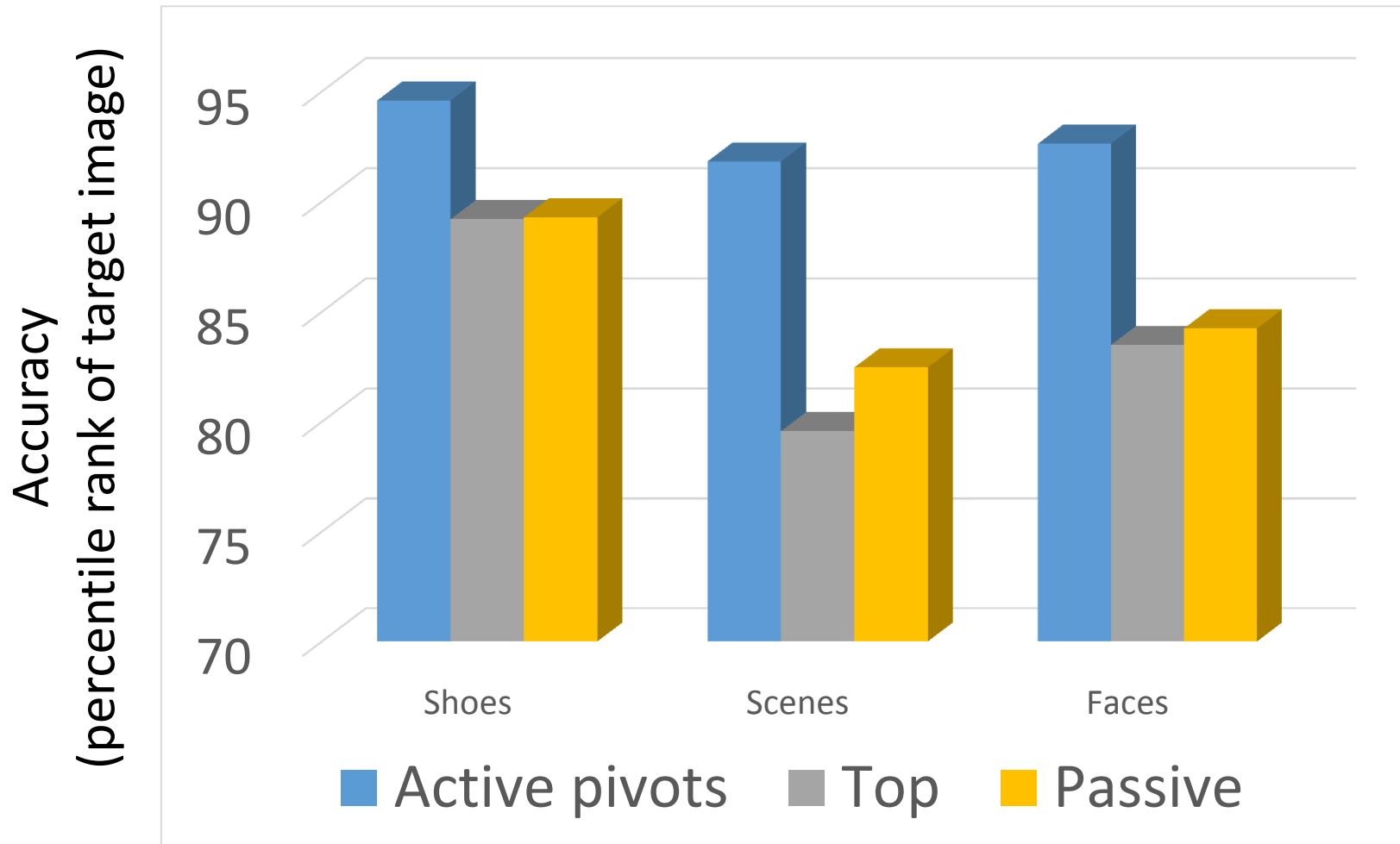
# Attribute Pivots for Active WhittleSearch

## Active feedback requests zero in on target more quickly

# Ongoing issues with attributes

- What attributes should be in the vocabulary?

- How to align user's attribute language with the visual attribute models?

- Joint learning of multiple attributes?

- Category-based vs. image-based comparative constraints?

- Class-specific attributes?

- How do we make sure we're learning the "right" thing?

# Summary

- Humans are not simply "label machines"

- More data need not mean better learning

- Active learning focuses annotator effort

- Widen access to visual knowledge by modeling visual comparisons

- Relative attributes enable new applications for recognition and visual search

# References

- WhittleSearch: Image Search with Relative Attribute Feedback. A. Kovashka, D. Parikh, and K. Grauman. CVPR 2012.

- Hashing Hyperplane Queries to Near Points with Applications to Large-Scale Active Learning. P. Jain, S. Vijayanarasimhan, and K. Grauman. NIPS 2010.

- Annotator Rationales for Visual Recognition. J. Donahue and K. Grauman. ICCV 2011.

- Actively Selecting Annotations Among Objects and Attributes. A. Kovashka, S. Vijayanarasimhan, and K. Grauman. ICCV 2011.

- Large-Scale Live Active Learning: Training Object Detectors with Crawled Data and Crowds. S. Vijayanarasimhan and K. Grauman. CVPR 2011.

- Cost-Sensitive Active Visual Category Learning. S. Vijayanarasimhan and K. Grauman. *International Journal of Computer Vision* (IJCV), Vol. 91, Issue 1 (2011), p. 24.

- What's It Going to Cost You?: Predicting Effort vs. Informativeness for Multi-Label Image Annotations. S. Vijayanarasimhan and K. Grauman. CVPR 2009.

- Multi-Level Active Prediction of Useful Image Annotations for Recognition. S. Vijayanarasimhan and K. Grauman. NIPS 2008.

- Far-Sighted Active Learning on a Budget for Image and Video Recognition. S. Vijayanarasimhan, P. Jain, and K. Grauman. CVPR 2010.

- Relative Attributes. D. Parikh and K. Grauman. ICCV 2011.