
Relative Attributes

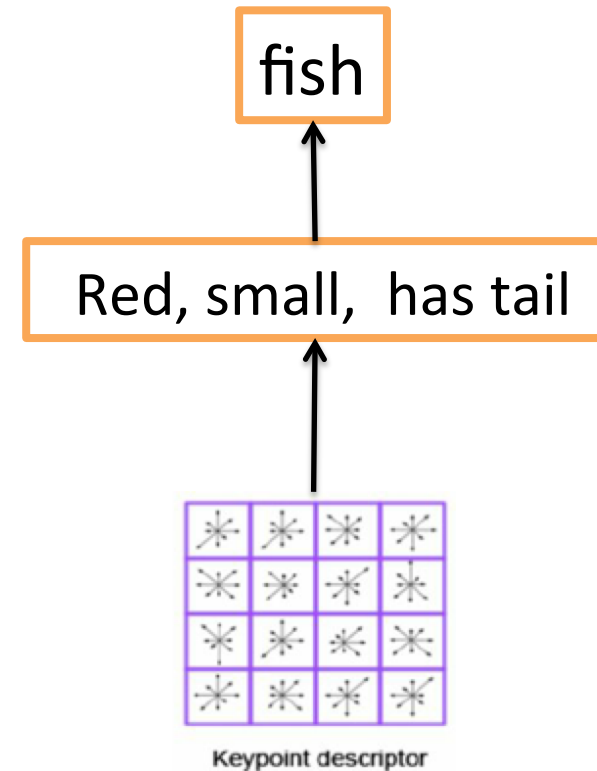
— **Shuangfei Fan** —

Electrical & Computer Engineering
Virginia Tech

- **Introduction**
- Learning Relative Attributes
- Relative Zero-shot Learning
- Automatic Relative Image Description
- Datasets
- Experiments
- Conclusion

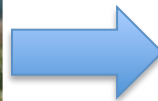
What are Attributes

- ❖ Low-level concepts: features
- ❖ High-level concepts: labels, categories
- ❖ Mid-level concepts: attributes
 - ❖ Shared across categories
 - ❖ Have semantic meanings
 - ❖ Visual concepts (machine detectable)



Why Attributes?

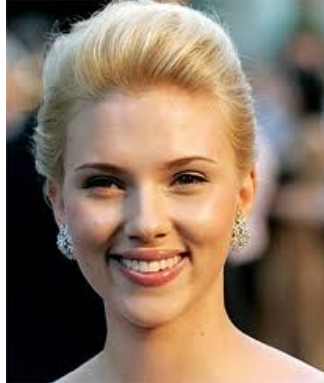
- ❖ How humans naturally describe natural concepts
 - ❖ Image search
 - ❖ Describe unknown objects



Has Horn
Has leg
Has Head
Has Wool



Relative Attributes



Smiling



???



Not smiling



Natural



???



Not natural

Figure Credit: Devi Parikh

Relative Attributes

Smiling



>



>



Natural



>



>



OVERVIEW

- Introduction
- **Learning Relative Attributes**
- Relative Zero-shot Learning
- Automatic Relative Image Description
- Datasets
- Experiments
- Conclusion

Learning Relative Attributes

For each attribute a_m , **open**

Supervision is

$$O_m: \left\{ \left(\left[\text{Cathedral} \right] \succ \left[\text{City} \right] \right), \dots \right\},$$

$$S_m: \left\{ \left(\left[\text{Beach} \right] \sim \left[\text{Field} \right] \right), \dots \right\}$$

Slide Credit: Devi Parikh

Learning Relative Attributes

Learn a scoring function $r_m(\mathbf{x}_i) = \mathbf{w}_m^T \mathbf{x}_i$

that best satisfies constraints:

$$\forall (i, j) \in O_m : \mathbf{w}_m^T \mathbf{x}_i > \mathbf{w}_m^T \mathbf{x}_j$$

$$\forall (i, j) \in S_m : \mathbf{w}_m^T \mathbf{x}_i = \mathbf{w}_m^T \mathbf{x}_j$$

Learning Relative Attributes

Max-margin learning to rank formulation

$$\begin{aligned} \min \quad & \left(\frac{1}{2} \|\mathbf{w}_m^T\|_2^2 + C \left(\sum \xi_{ij}^2 + \sum \gamma_{ij}^2 \right) \right) \\ \text{s.t} \quad & \mathbf{w}_m^T (\mathbf{x}_i - \mathbf{x}_j) \geq 1 - \xi_{ij}, \forall (i, j) \in O_m \\ & |\mathbf{w}_m^T (\mathbf{x}_i - \mathbf{x}_j)| \leq \gamma_{ij}, \forall (i, j) \in S_m \\ & \xi_{ij} \geq 0; \gamma_{ij} \geq 0 \end{aligned}$$

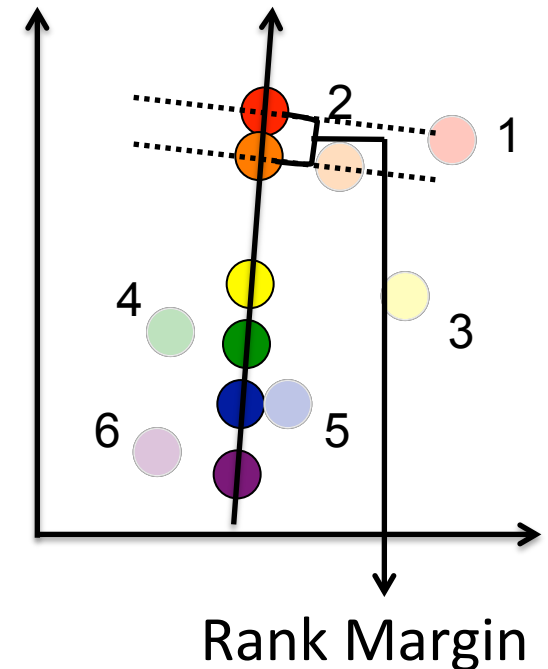


Image \rightarrow Relative Attribute Score

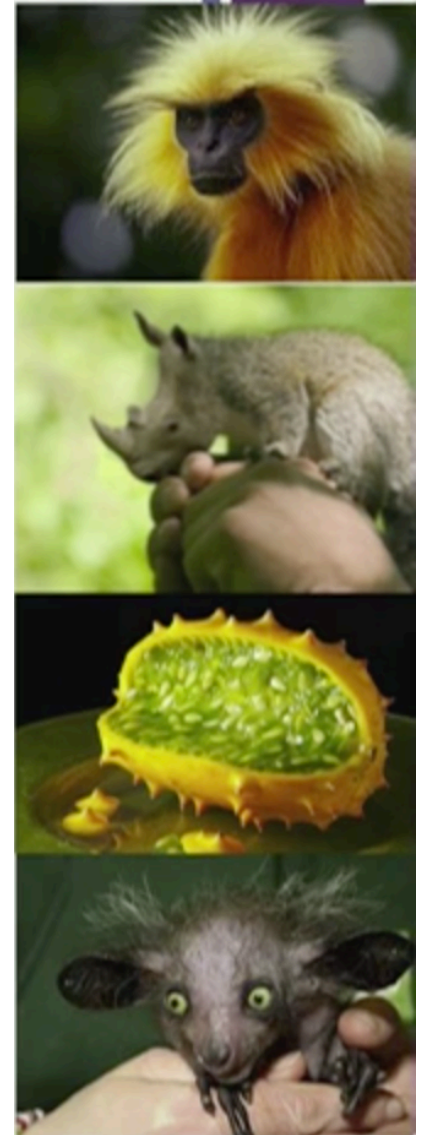
Slide Credit: Devi Parikh

OVERVIEW

- Introduction
- Learning Relative Attributes
- **Relative Zero-shot Learning**
- Automatic Relative Image Description
- Datasets
- Experiments
- Conclusion

Zero-shot Learning

- ❖ Recognize the Wampimuk
 - ❖ Impossible?
- ❖ Solution: semantic transfer
 - ❖ Wampimuk: small, horn, furry, cute
- ❖ Zero-Shot:
 - ❖ Pattern recognition with no training examples
 - ❖ Solved by semantic transfer



Slide Credit: Timothy Hospedales

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, or all seen categories

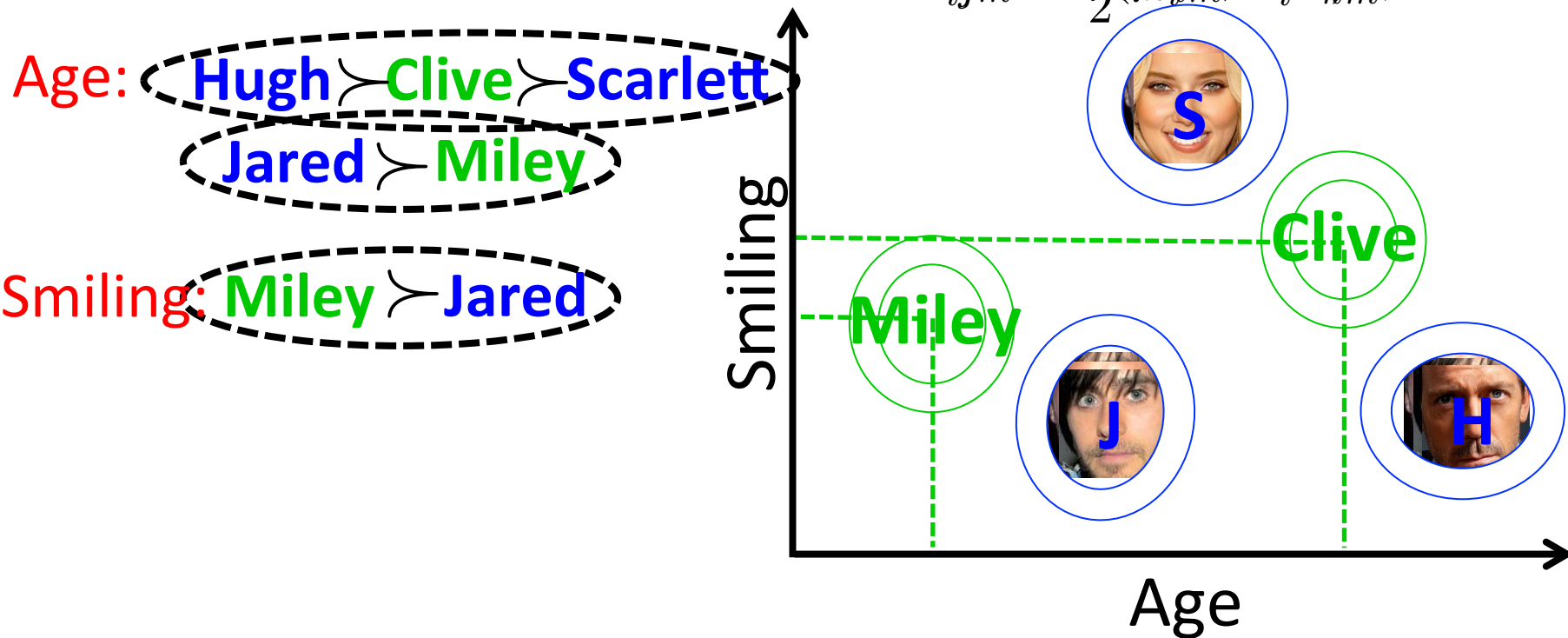
Testing: Categorize image into one of **S+U** categories

Slide Credit: Devi Parikh

Relative Zero-shot Learning

Can predict new classes based on their relationships to existing classes – without training images

$$\mu_{ijm}^{(s)} = \mathcal{N}(\mu_{ijm}^{(s)}, \Sigma_{ijm}^{(s)})$$



Infer image category using max-likelihood

$$c^* = \operatorname{argmax}_{j \in \{1, \dots, N\}} P(\tilde{x}_i | \mu_j, \Sigma_j)$$

Slide Credit: Devi Parikh

OVERVIEW

- Introduction
- Learning Relative Attributes
- Relative Zero-shot Learning
- **Automatic Relative Image Description**
- Datasets
- Experiments
- Conclusion

Automatic Relative Image Description

Density



Novel image

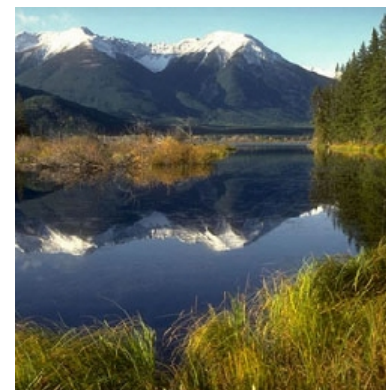


Conventional binary description: *not dense*

Dense:



Not dense:



Slide Credit: Devi Parikh

Automatic Relative Image Description

Density

Novel image



more dense than

less dense than



Slide Credit: Devi Parikh

Automatic Relative Image Description

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***

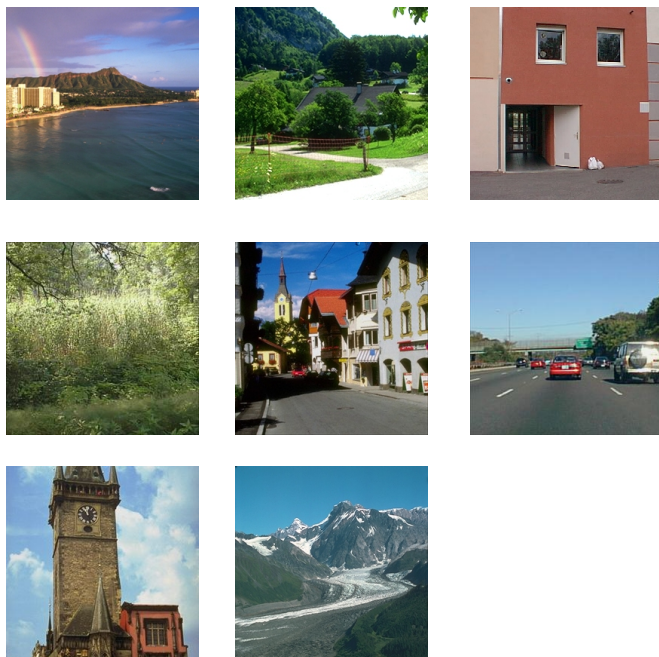
Slide Credit: Devi Parikh

OVERVIEW

- Introduction
- Learning Relative Attributes
- Relative Zero-shot Learning
- Automatic Relative Image Description
- **Datasets**
- Experiments
- Conclusion

Datasets

Outdoor Scene Recognition (OSR) [Oliva 2001]



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

Public Figures Face (PubFig) [Kumar 2009]



8 classes, ~800 images, Gist+color
11 attributes: white, chubby, etc.

Attributes labeled at category level

Slide Credit: Devi Parikh

Category level annotation

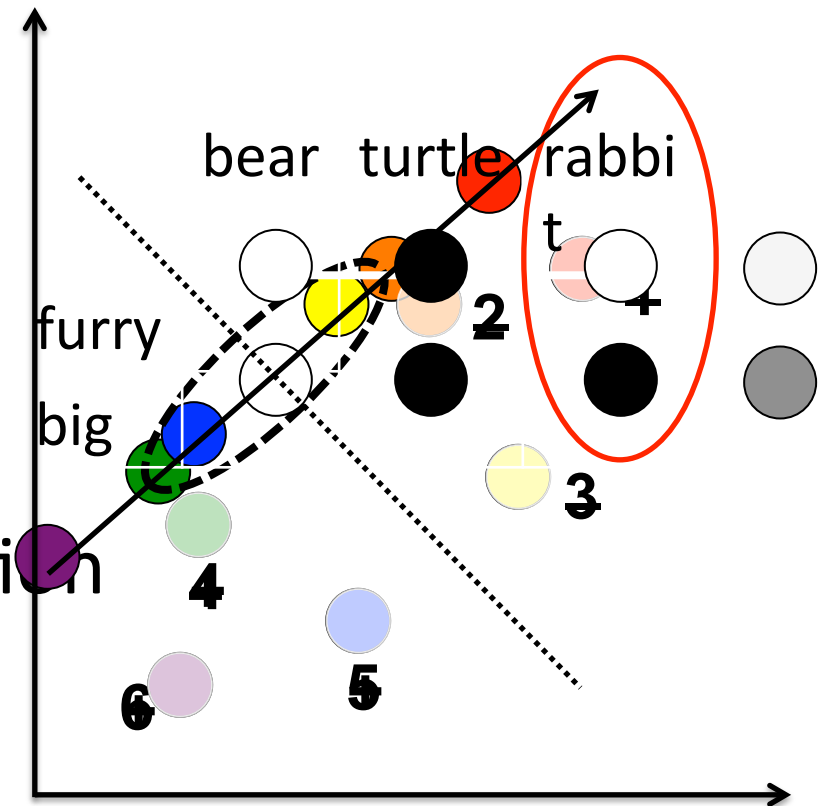
	Binary	Relative
OSR	TI S HC OMF	
natural	00001111	T<I~S<H<C~O~M~F
open	00011110	T~F<I~S<M<H~C~O
perspective	11110000	O<C<M~F<H<I<S<T
large-objects	11100000	F<O~M<I~S<H~C<T
diagonal-plane	11110000	F<O~M<C<I~S<H<T
close-depth	11110001	C<M<O<T~I~S~H~F
PubFig	ACHJ MS VZ	
Masculine-looking	11110011	S<M<Z<V<J<A<H<C
White	01111111	A<C<H<Z<J<S<M<V
Young	00001101	V<H<C<J<A<S<Z<M
Smiling	11101101	J<V<H<A~C<S~Z<M
Chubby	10000000	V<J<H<C<Z<M<S<A
Visible-forehead	11101110	J<Z<M<S<A~C~H~V
Bushy-eyebrows	01010000	M<S<Z<V<H<A<C<J
Narrow-eyes	01100011	M<J<S<A<H<C<V<Z
Pointy-nose	00100001	A<C<J~M~V<S<Z<H
Big-lips	10001100	H<J<V<Z<C<M<A<S
Round-face	10001100	H<V<J<C<Z<A<S<M

OVERVIEW

- Introduction
- Learning Relative Attributes
- Relative Zero-shot Learning
- Automatic Relative Image Description
- Datasets
- **Experiments**
- Conclusion

Experiments: Baselines

- Zero-shot learning
 - Binary attributes:
Direct Attribute Prediction
 - Relative attributes via
classifier scores
- Automatic image-description
 - Binary attributes



Slide Credit: Devi Parikh

Experiments: Zero-shot learning

- Robustness:
 - Fewer comparisons to train relative attributes
 - More unseen (fewer seen) categories
- Flexibility in supervision:
 - ‘Looseness’ in description of unseen
 - Fewer attributes used to describe unseen

Slide Credit: Devi Parikh

Experiments: Zero-shot learning

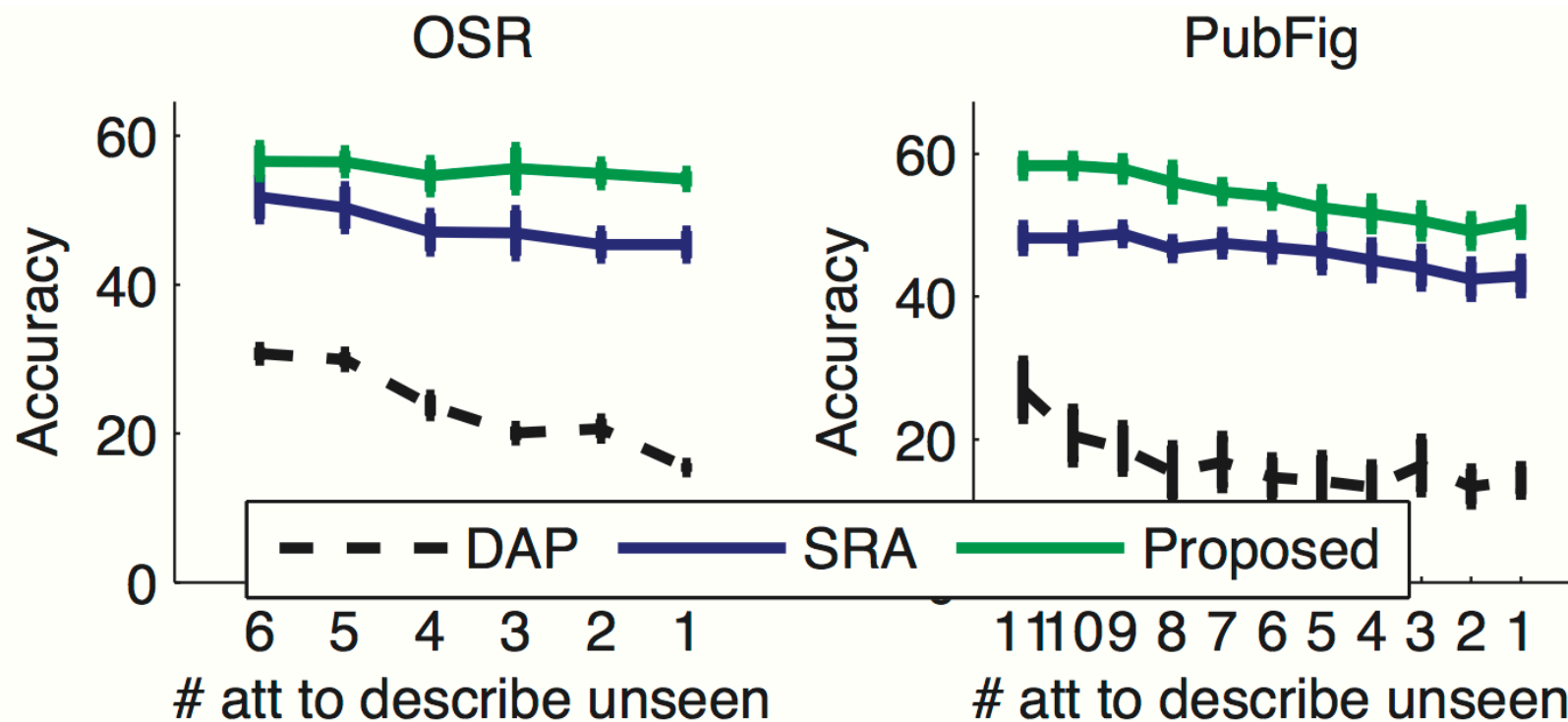


Figure 5. Zero-shot learning performance as fewer attributes are used to describe the unseen categories.

Experiments: Describe images

Binary attribute:

Not natural

Not open

Has perspective

Relative attribute:

More natural than insidicity

Less natural than highway

More open than street

Less open than coast

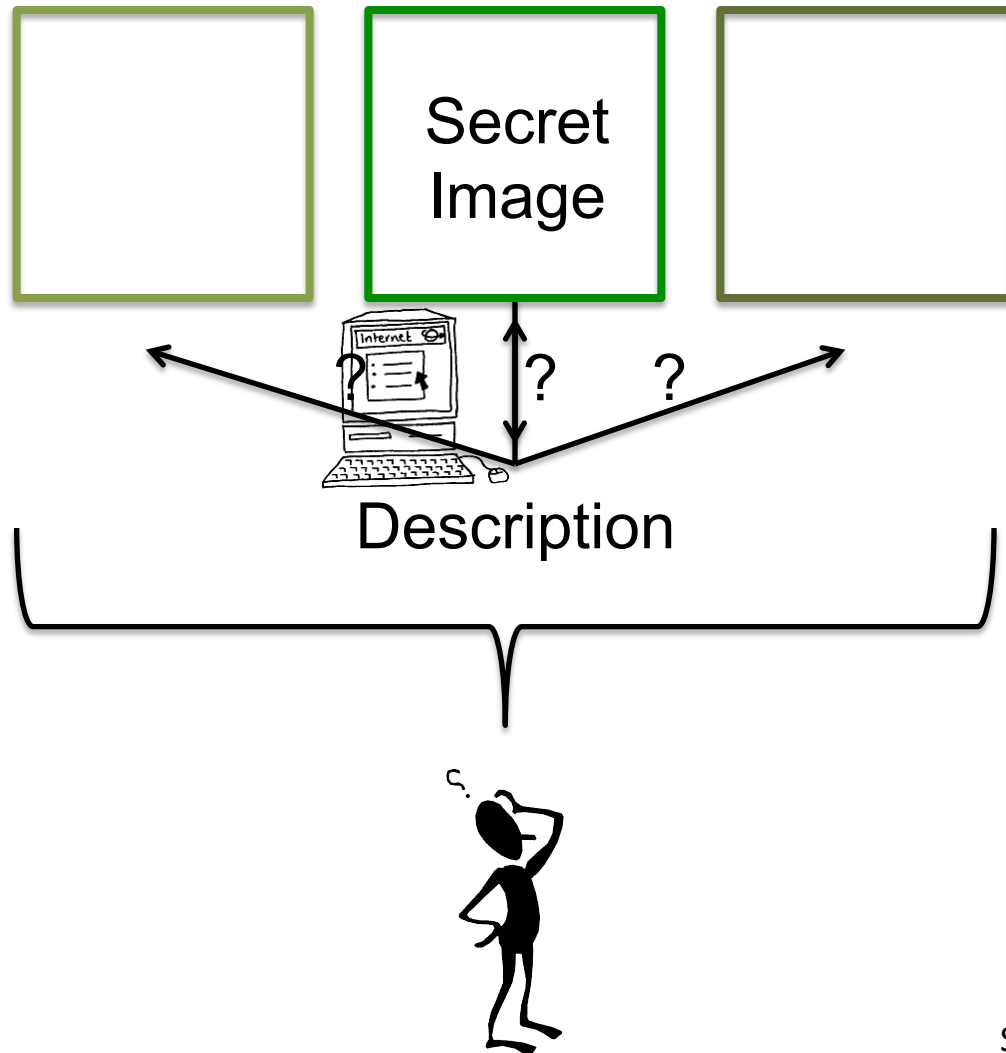
Has more perspective than highway

Has less perspective than insidicity



Experiments: Describe images

Human Studies: Which Image is Being Described?



Slide Credit: Devi Parikh

Experiments: Describe images

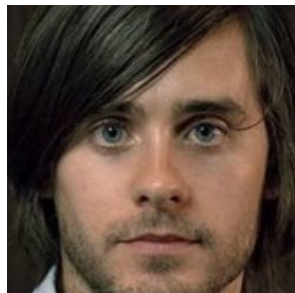


Binary: **Smiling, Young**

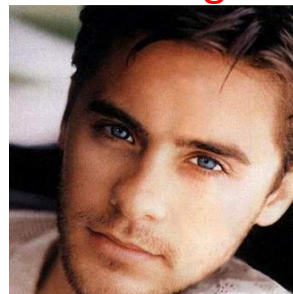
Smiling



Not Smiling



Young



Not Young

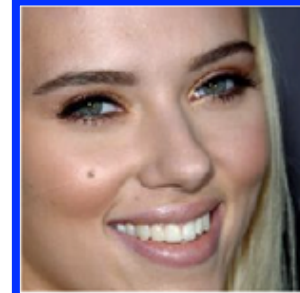


Relative

More Smiling than



Less Smiling than



Younger than



Older than



Slide Credit: Devi Parikh

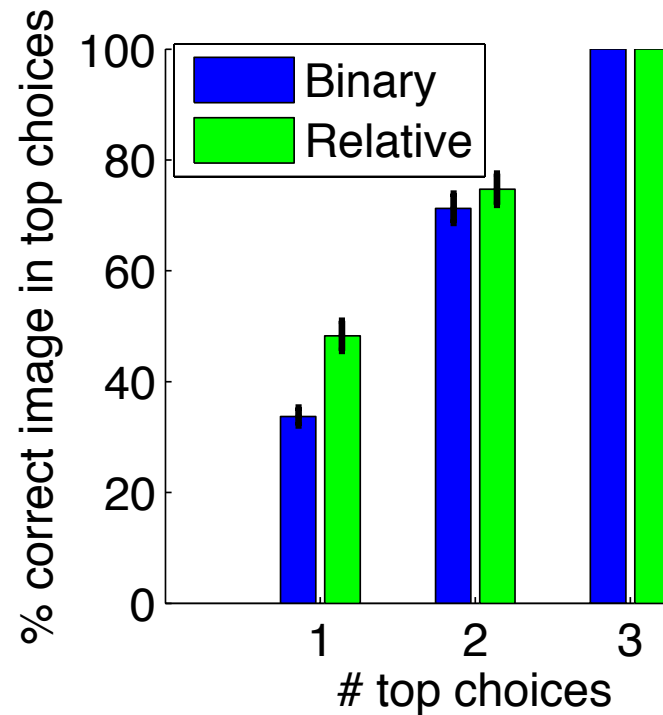
Experiments: Describe images

Human Studies: Which Image is Being Described?

18 subjects

Test cases:

10 OSR, 20 PubFig



Slide Credit: Devi Parikh

OVERVIEW

- Introduction
- Learning Relative Attributes
- Relative Zero-shot Learning
- Automatic Relative Image Description
- Datasets
- Experiments
- **Conclusion**

Conclusion

- Relative attributes
 - Allow relating images and categories to each other
 - Learn ranking function for each attribute
- Novel applications
 - Natural and accurate zero-shot learning from attribute comparisons
 - Automatically generating precise relative image descriptions for human interpretation

Questions?

BACKUP

Experiments: Zero-shot learning

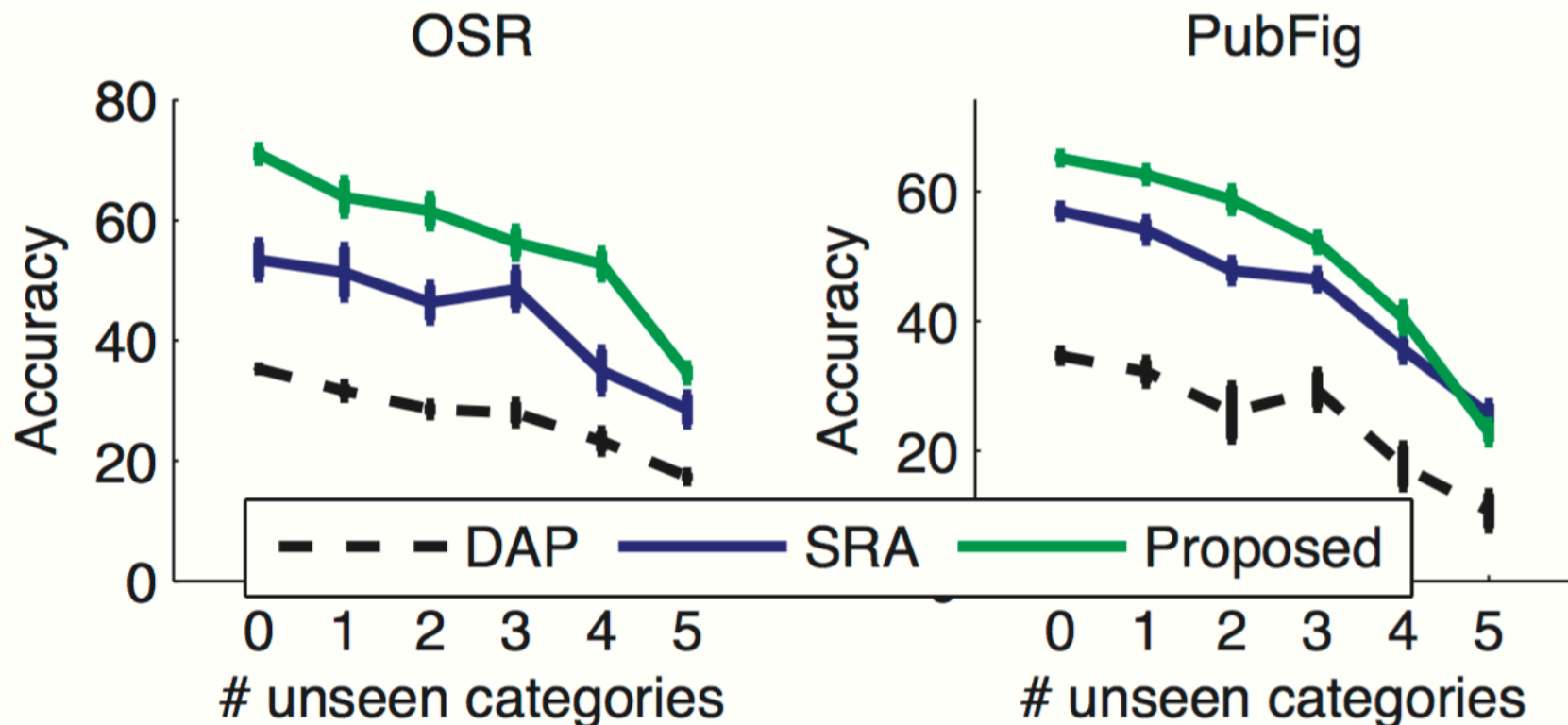


Figure 3. Zero-shot learning performance as the proportion of unseen categories increases. Total number of classes N remains constant at 8.

Experiments: Zero-shot learning

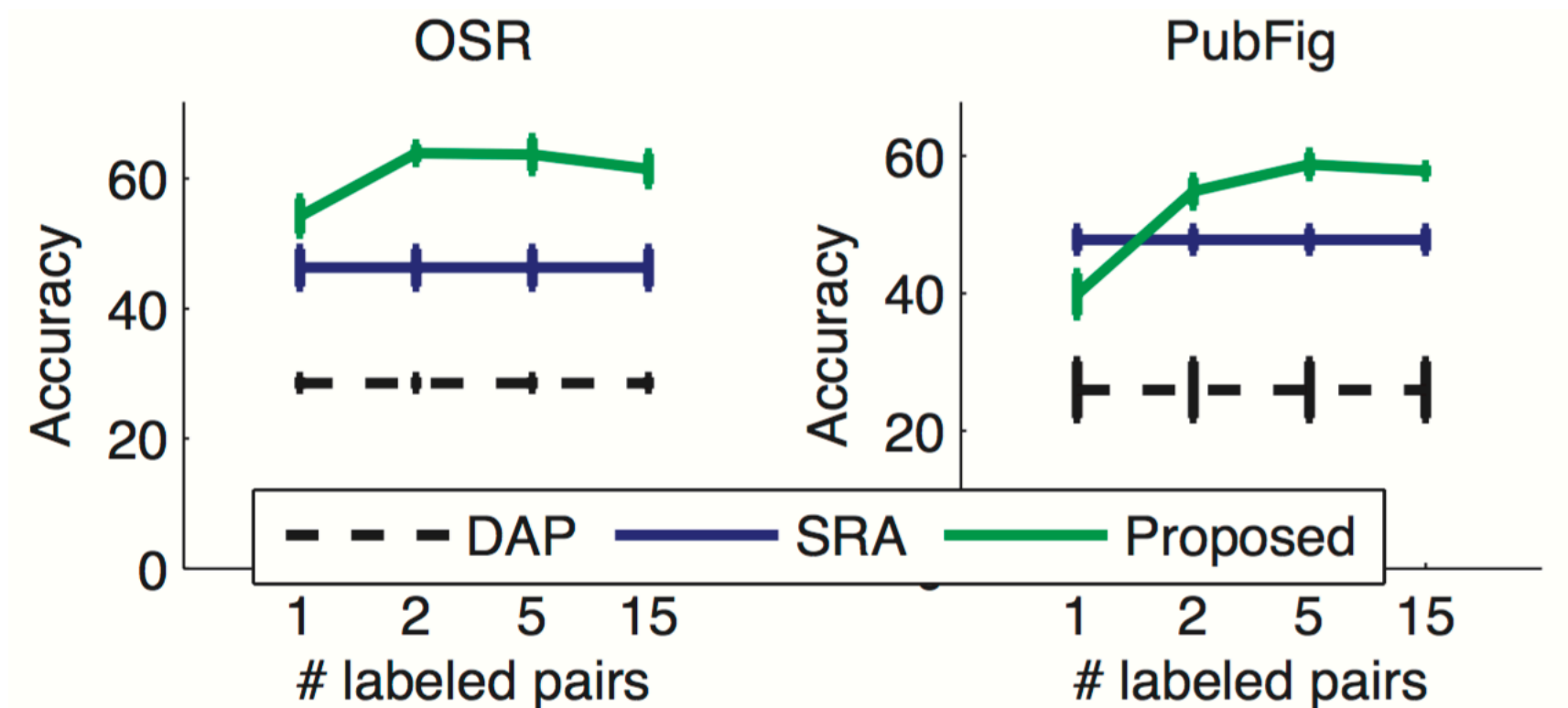


Figure 4. Zero-shot learning performance as more pairs of seen categories are related (*i.e.* labeled) during training.

Experiments: Zero-shot learning

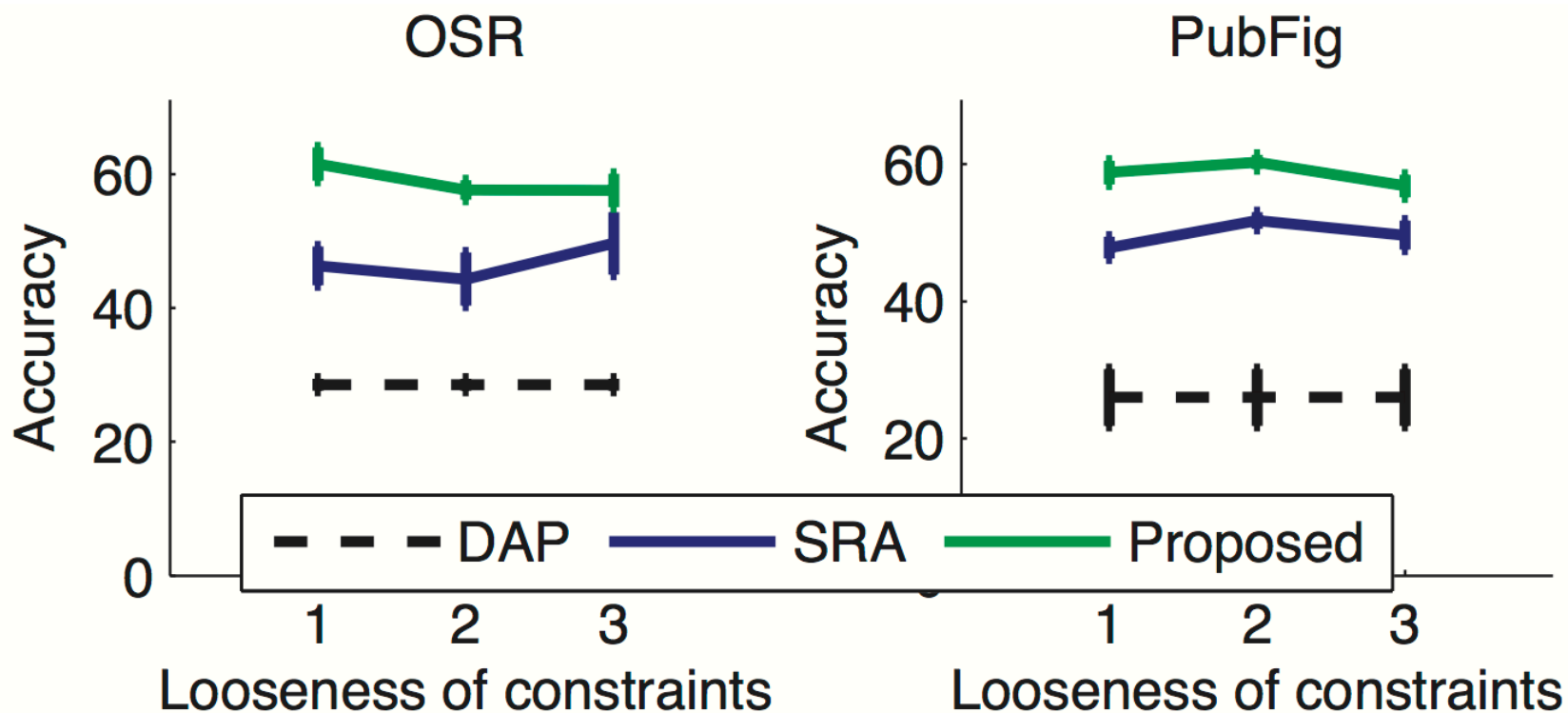


Figure 6. Zero-shot learning performance as the unseen categories are described via looser relationships.

- ❖ GIST is a steerable filter (Gabor filter) response of an image.
- ❖ Any image has 1 GIST descriptor of 512 dimensions.
- ❖ GIST was developed to provide a holistic descriptor that provides a simpler representation.
- ❖ Compared to SIFT features:
 - ❖ SIFT is a localized image patch descriptor. A typical image has a few thousand SIFT descriptors, each of 128 dimensions.
 - ❖ SIFT was designed for scale and affine invariance in wide baseline image matching tasks, which were part of stereo vision.