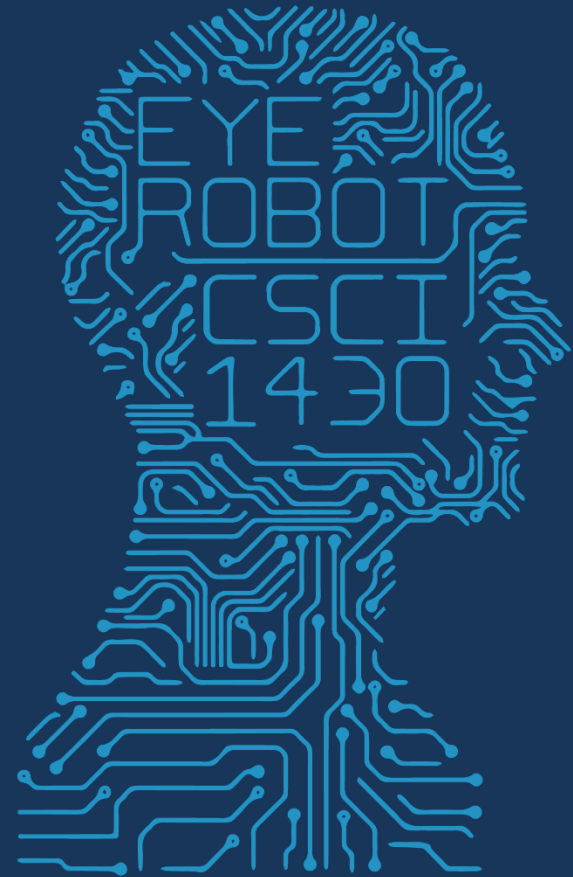




1950

FUTURE VISION



2017 MWF 1PM 368

COMPUTER VISION



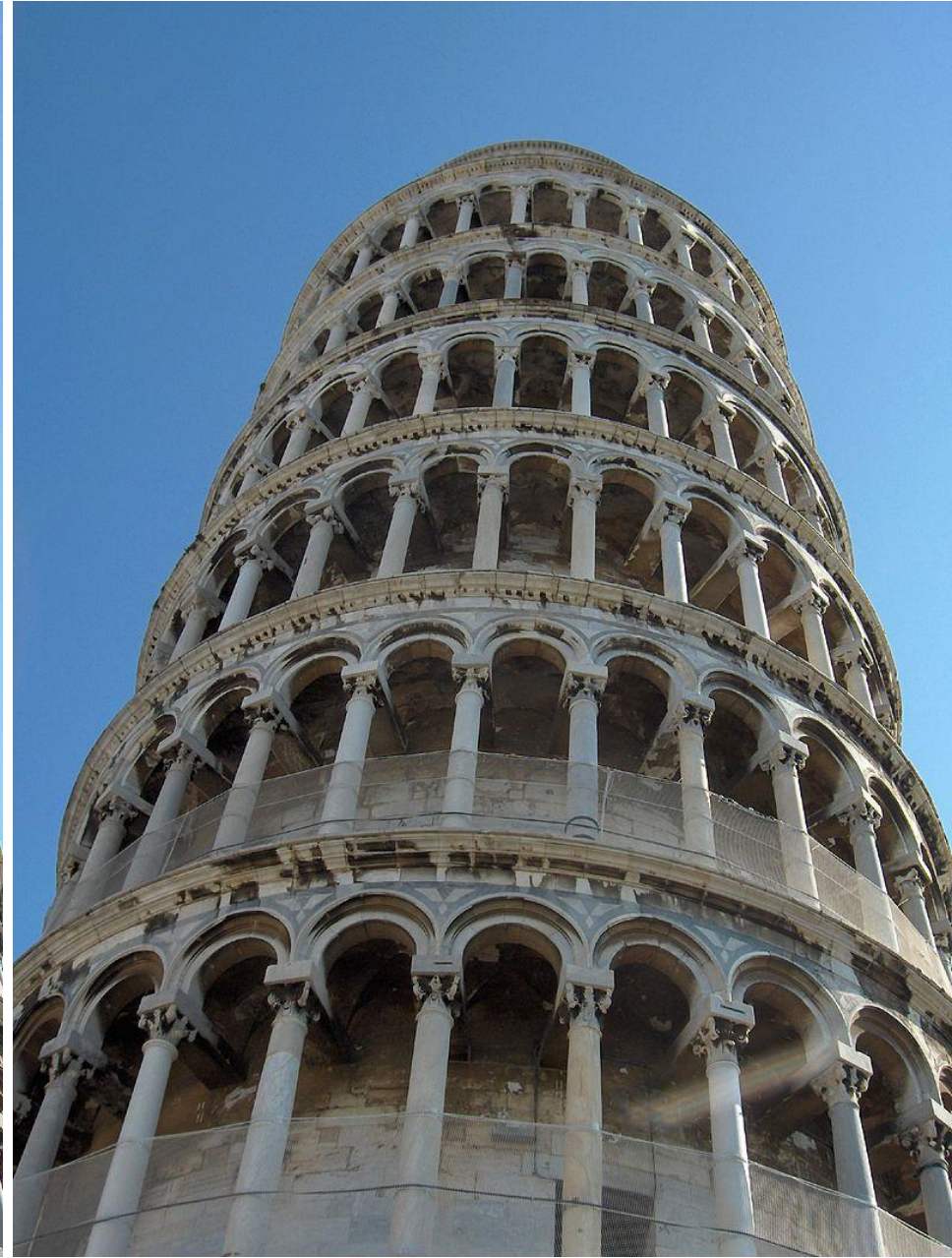




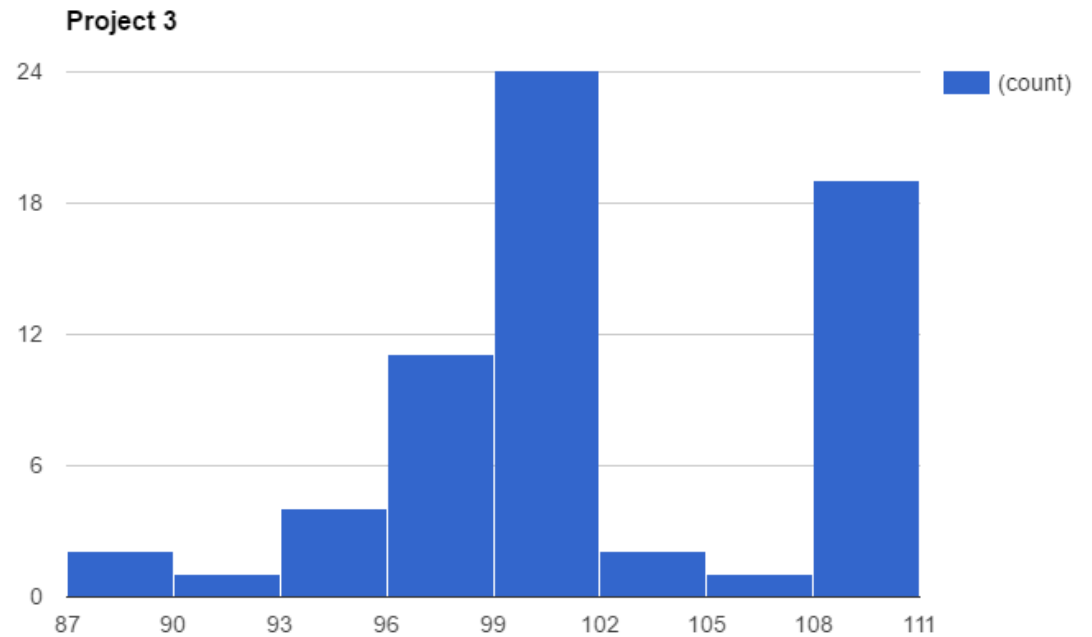




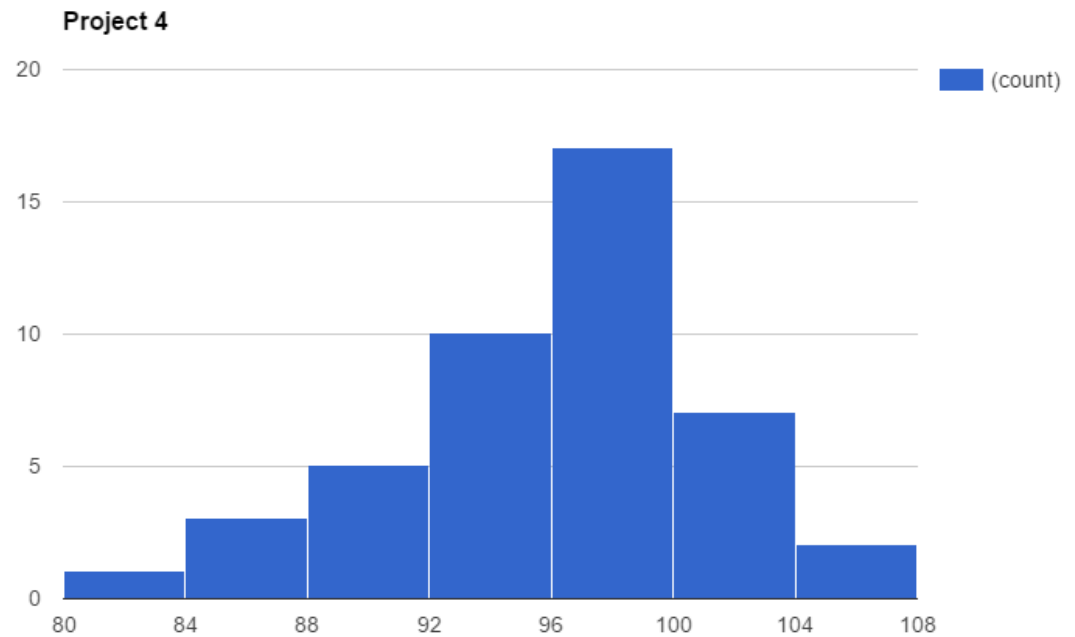
Photo: Georges Jansoon. Illusion: Frederick Kingdom, Ali Yoonessi and Elena Gheorghiu



Project 3: Camera calibration



Project 4: Scene Recognition



Normalization – why?

Required by some underlying property of the learning mechanism.

- E.G., removing hyperplane bias in SVM to aid fitting.

Also called ‘feature scaling’.

- Many methods, e.g.,

$$x' = \frac{x - \bar{x}}{\sigma} \qquad x' = \frac{x}{||x||}$$

Normalization can be implemented *wrt. other data points*, and sometimes *wrt. other features*.

Wrt. other data points – human weight

Feature Scaling Formula

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

Annotations:

- X' : new (rescaled) feature
- X : info taken from old feature(s)
- X_{\min} : 115
- X_{\max} : 175

Quiz

old weights: [115, 140, 175]

what is X'_{140} ?



Tiny Image as a feature vector

- Each pixel is treated as a different feature.
- Feature vector is matrix reshaped into an array.

Normalization – how?

```
function image_feats = get_tiny_images( image_paths )
```

```
size = 16;
```

```
N = size(image_paths, 1);
```

```
image_feats = zeros(N, size * size);
```

```
for i = 1:N
```

```
    img = im2double( imread(image_paths{i, 1}) );
```

```
    img = imresize( img, [size, size] );
```

```
    fv = reshape( img, [1, size * size] );
```

```
    fv = fv - mean( fv );
```

```
    image_feats(i, :) = fv ./ norm( fv );
```

```
end
```

Normalization – how?

```
function image_feats = get_tiny_images( image_paths )
```

```
size = 16;
```

```
N = size(image_paths, 1);
```

```
image_feats = zeros(N, size * size);
```

```
for i = 1:N
```

```
    img = im2double( imread(image_paths{i, 1}) );
```

```
    img = imresize( img, [size, size] );
```

```
    fv = reshape( img, [1, size * size] );
```

```
    fv = fv - mean( fv );
```

```
    image_feats(i, :) = fv ./ norm( fv );
```

```
end
```

Mean across *features*

```
function image_feats = get_tiny_images( image_paths )
```

```
size = 16;
```

```
N = size(image_paths, 1);
```

```
image_feats = zeros(N, size * size);
```

```
for i = 1:N
```

```
    img = im2double( imread(image_paths{i, 1}) );
```

```
    img = imresize( img, [size, size] );
```

```
    fv = reshape( img, [1, size * size] );
```

```
    fv = fv - mean( fv );
```

% Mean across features (pixels) per data point

```
    image_feats(i, :) = fv ./ norm( fv );
```

```
end
```

Mean across *data points*

```
function image_feats = get_tiny_images( image_paths )
```

```
size = 16;
```

```
N = size(image_paths, 1);
```

```
image_feats = zeros(N, size * size);
```

```
for i = 1:N
```

```
    img = im2double( imread(image_paths{i, 1}) );
```

```
    img = imresize( img, [size, size] );
```

```
    fv = reshape( img, [1, size * size] );
```

```
end
```

```
mean_img = mean(image_feats); % Mean of each feature (pixel) across data points
```

```
var_img = std(image_feats);
```

```
for i = 1:N
```

```
    image_feats(i,:) = ( image_feats(i,:) - mean_img ) ./ var_img;
```

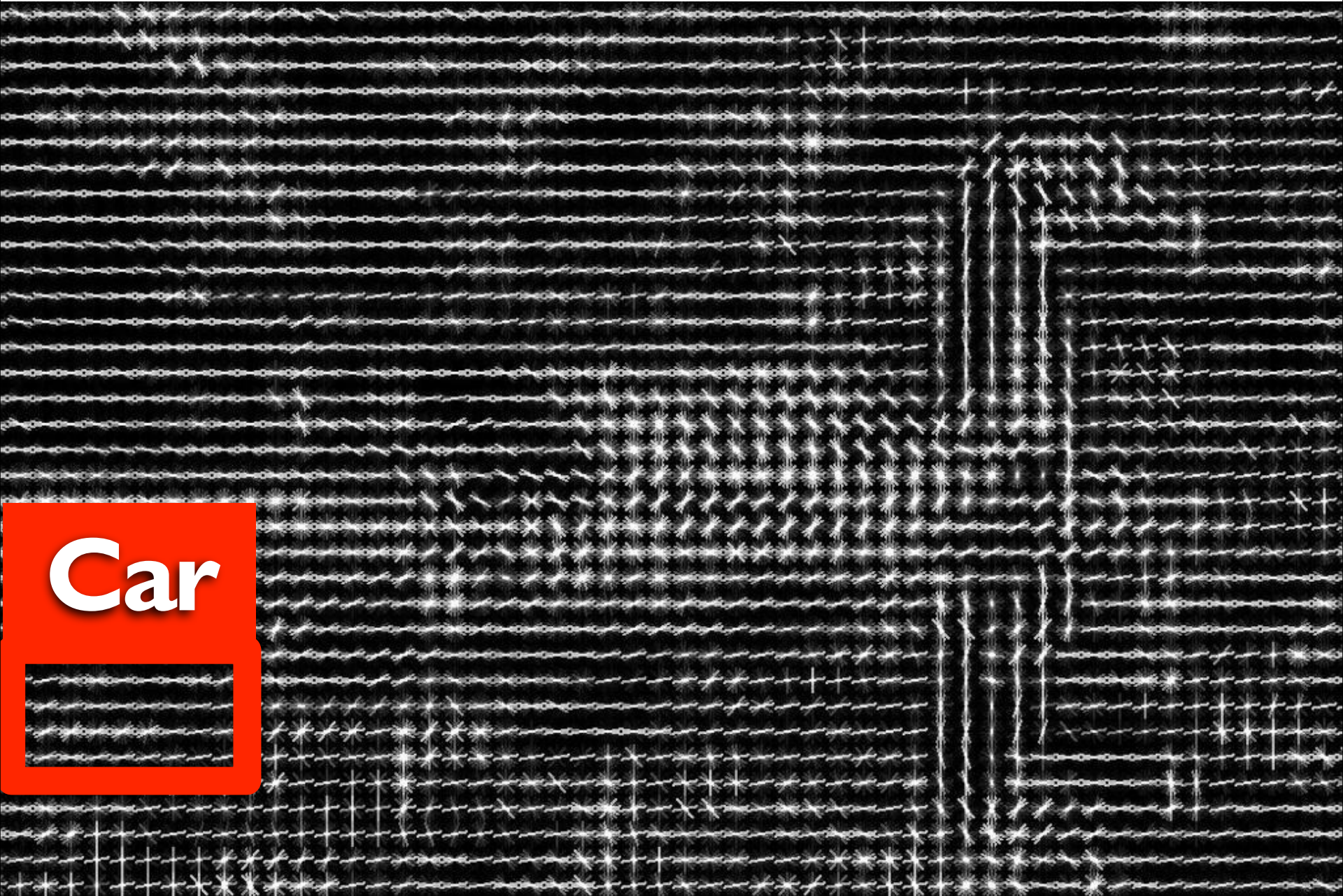
```
end
```

Friday: CV for Social ~~Good~~ Bad

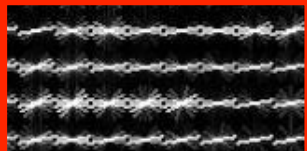
- We saw how dataset bias can introduce error.
- But what about the underlying feature representation itself?
- How might I discover why my CV system is bad?
- How might I explain the failure?



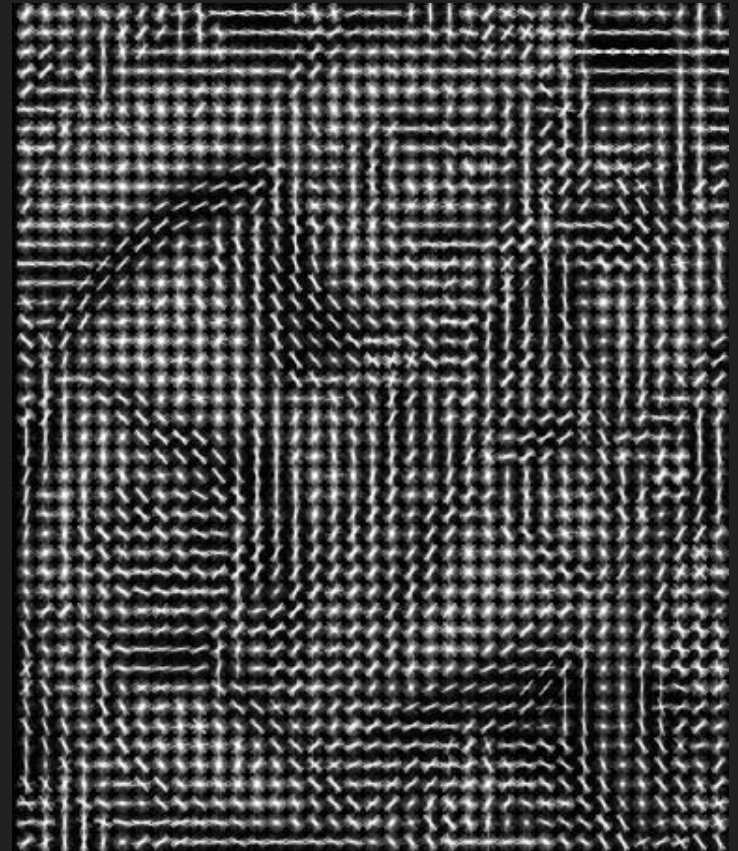




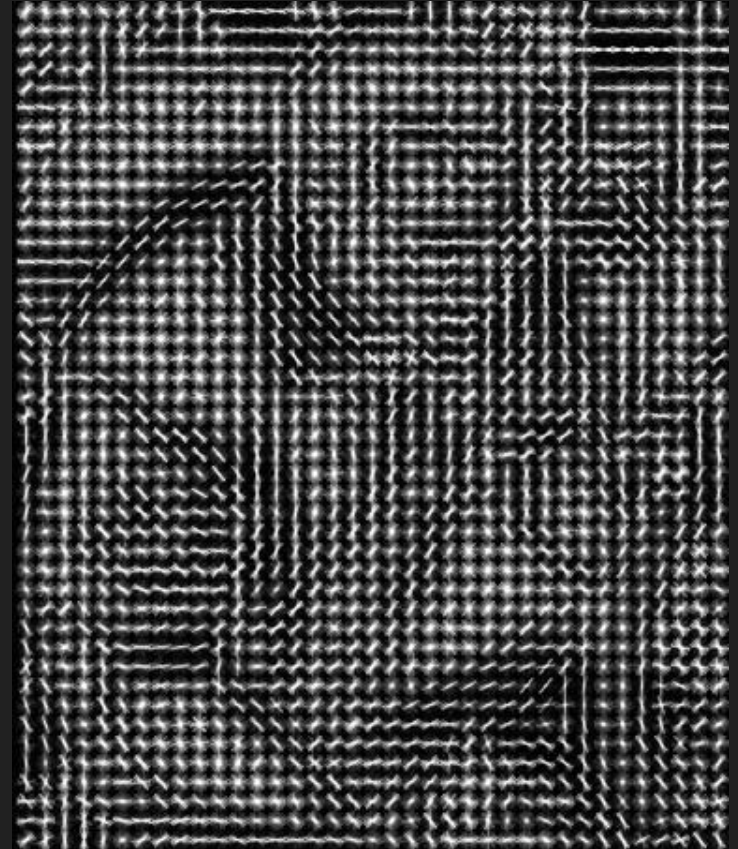
Car



What information is lost?



What information is lost?



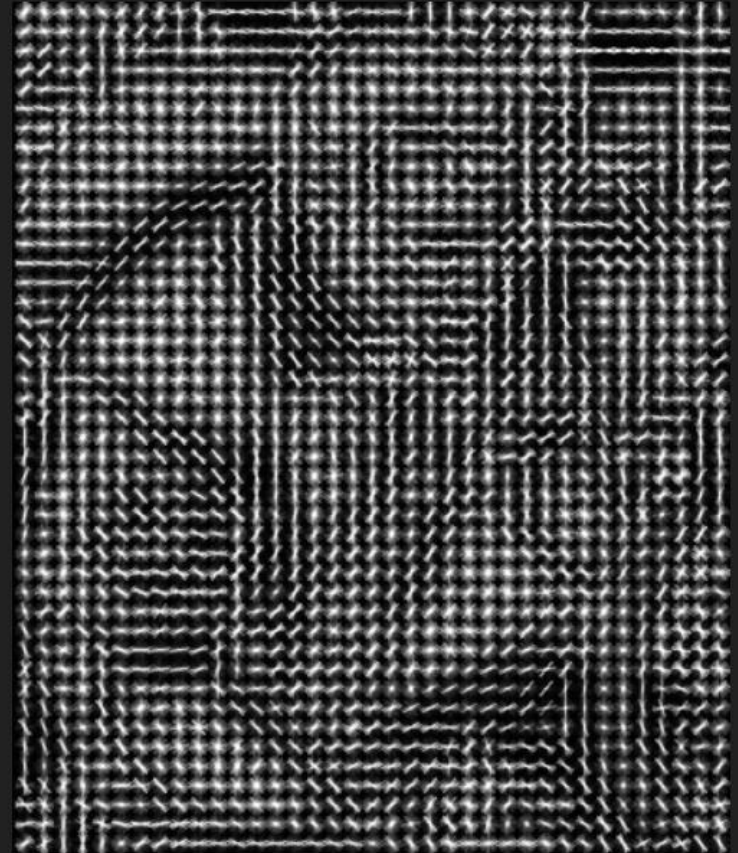
What information is lost?

$$\min_{x \in \mathbb{R}^d} ||\phi(x) - y||_2^2$$

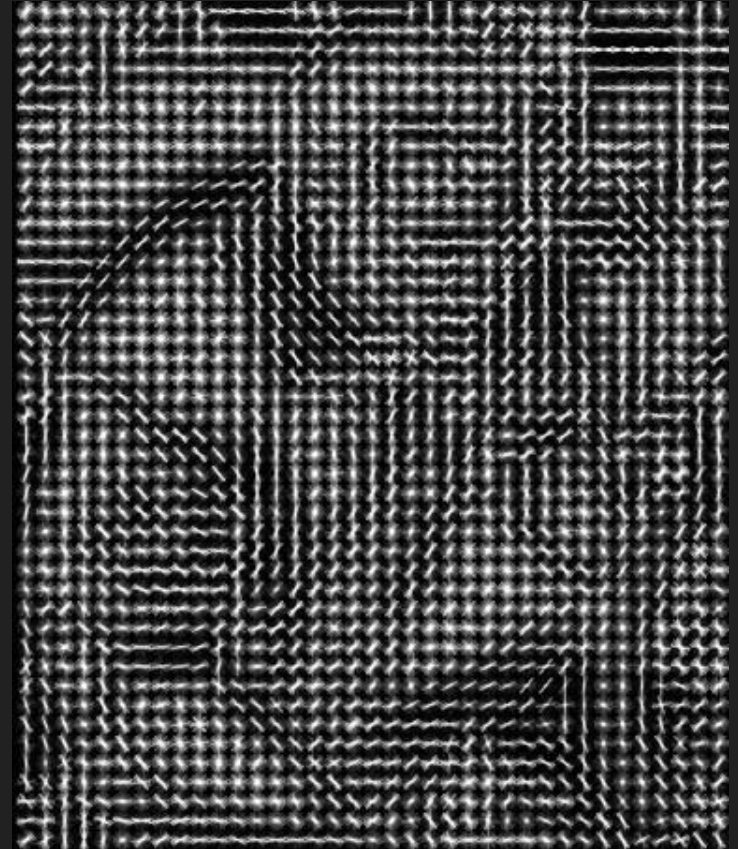
HOG = ϕ

Many-to-one function

No inverse

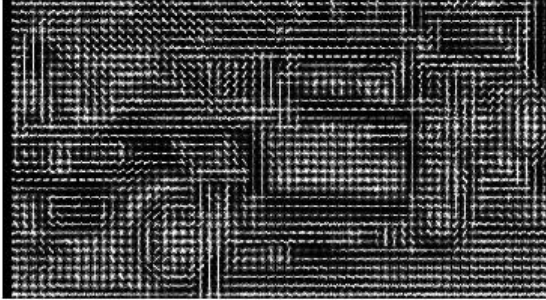


What information is lost?



HOGgles (Vondrick et al. ICCV 2013)

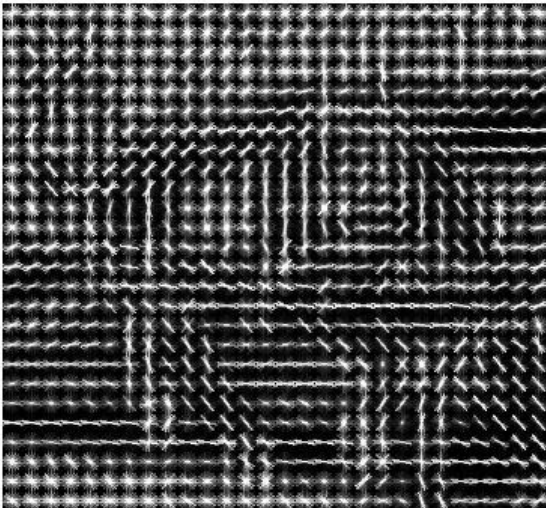
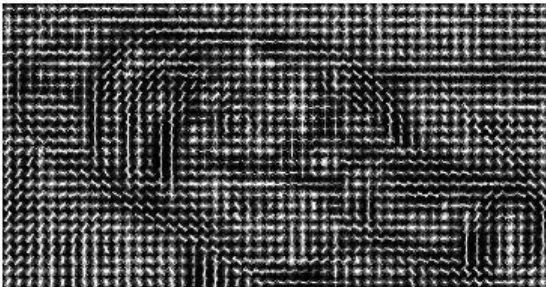
HOG [1]



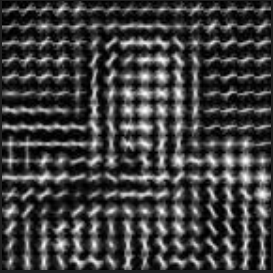
Inverse (Us)



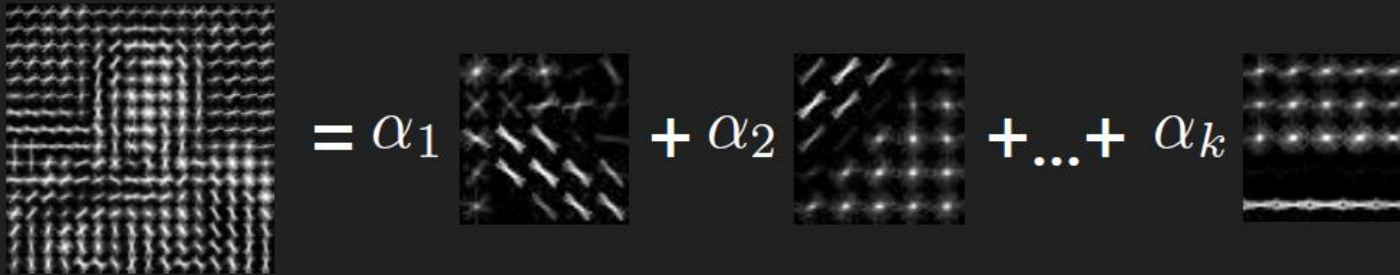
Original



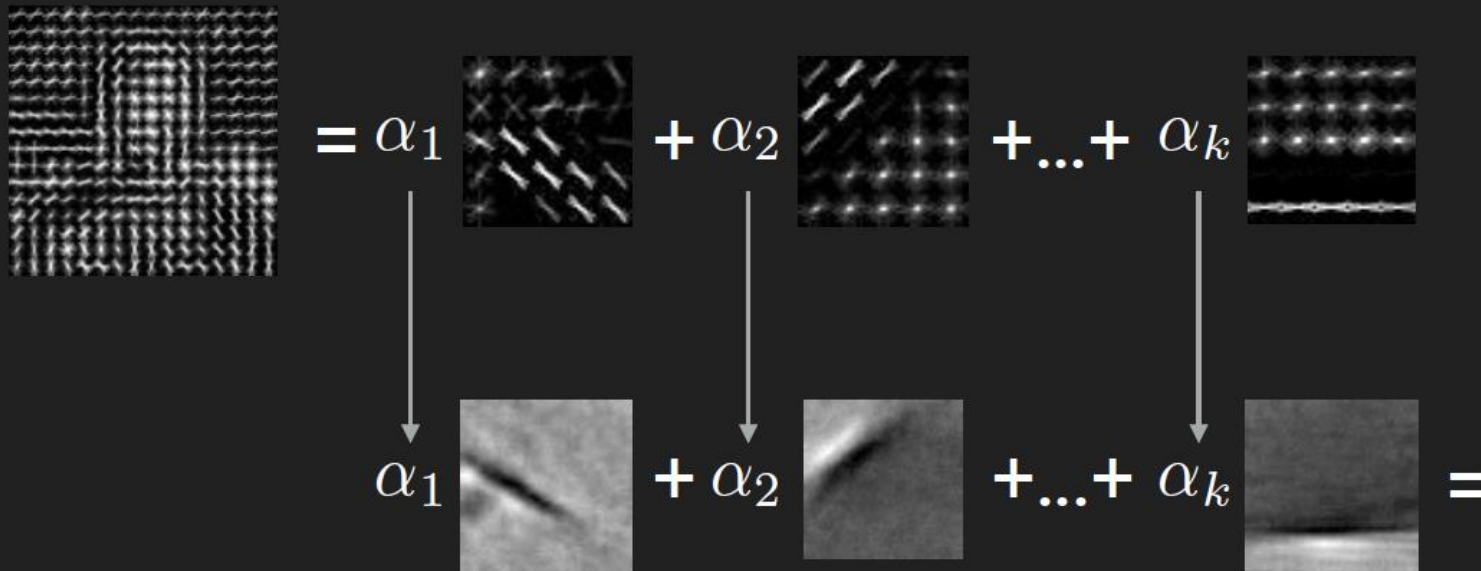
Method: Paired Dictionary



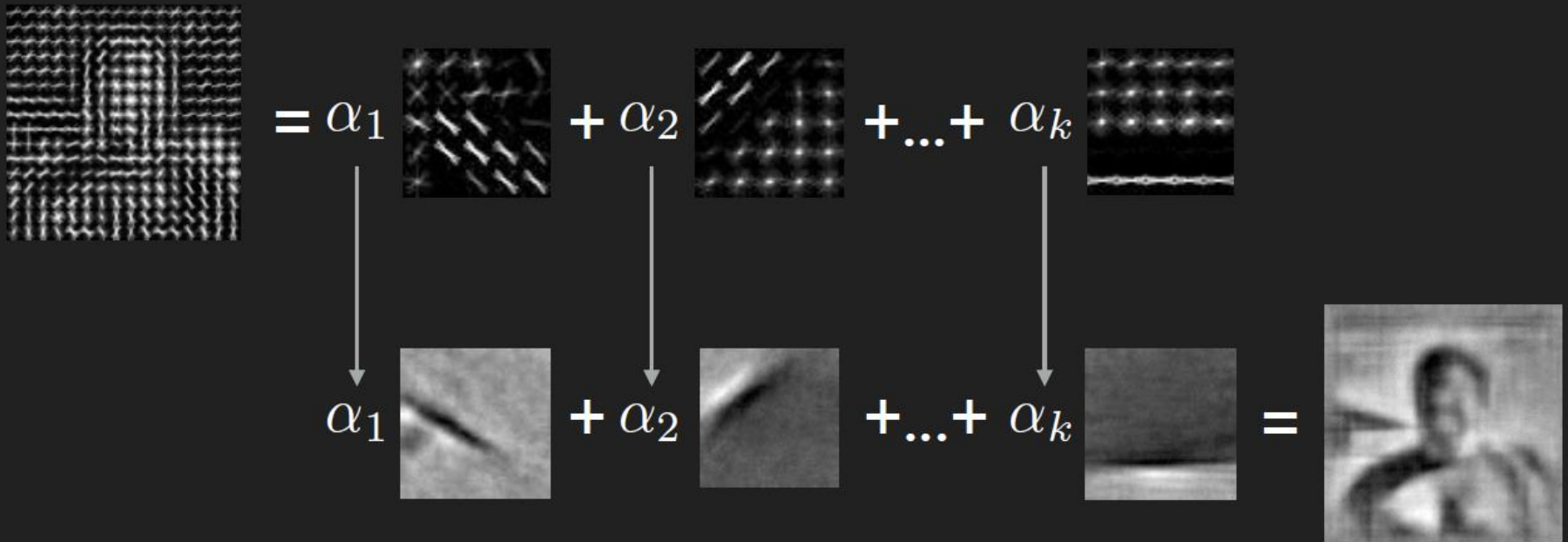
Method: Paired Dictionary

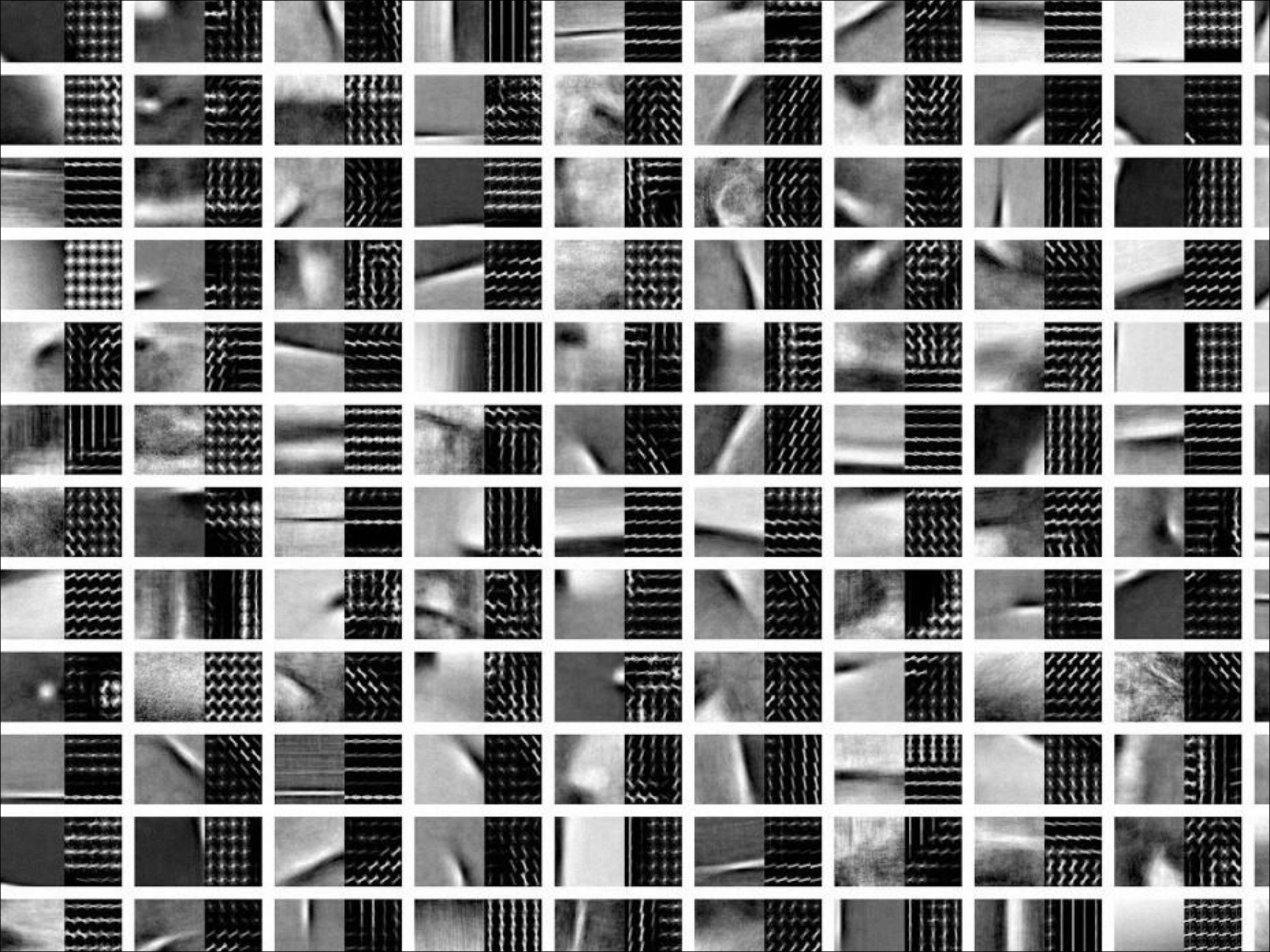

$$\text{Target Image} = \alpha_1 \text{Dict}_1 + \alpha_2 \text{Dict}_2 + \dots + \alpha_k \text{Dict}_k$$

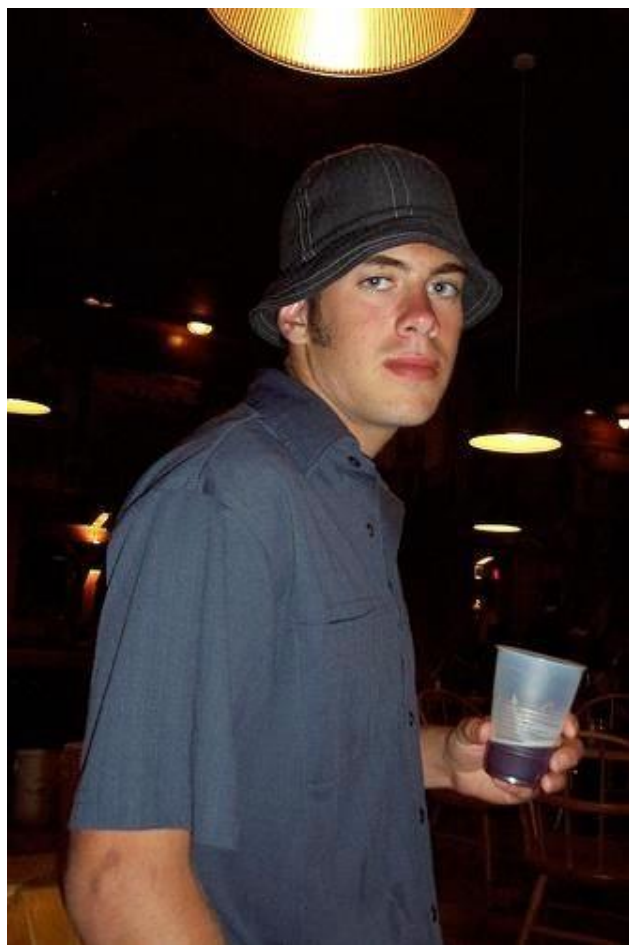
Method: Paired Dictionary



Method: Paired Dictionary







HumanVision



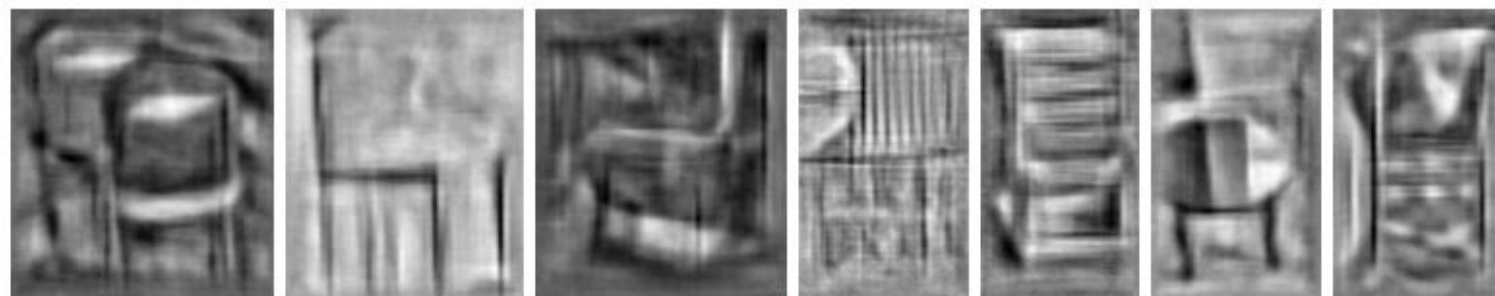
HOGVision

Visualizing Top Detections

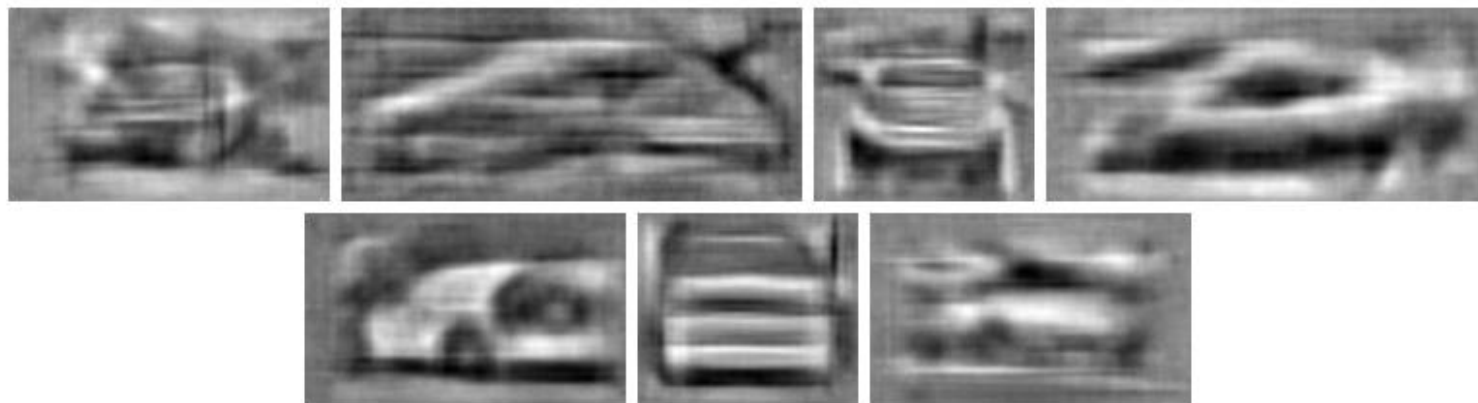
We have visualized some high scoring detections from the deformable parts model. Can you guess which are false alarms? Click on the images below to reveal the corresponding RGB patch. You might be surprised!



Person

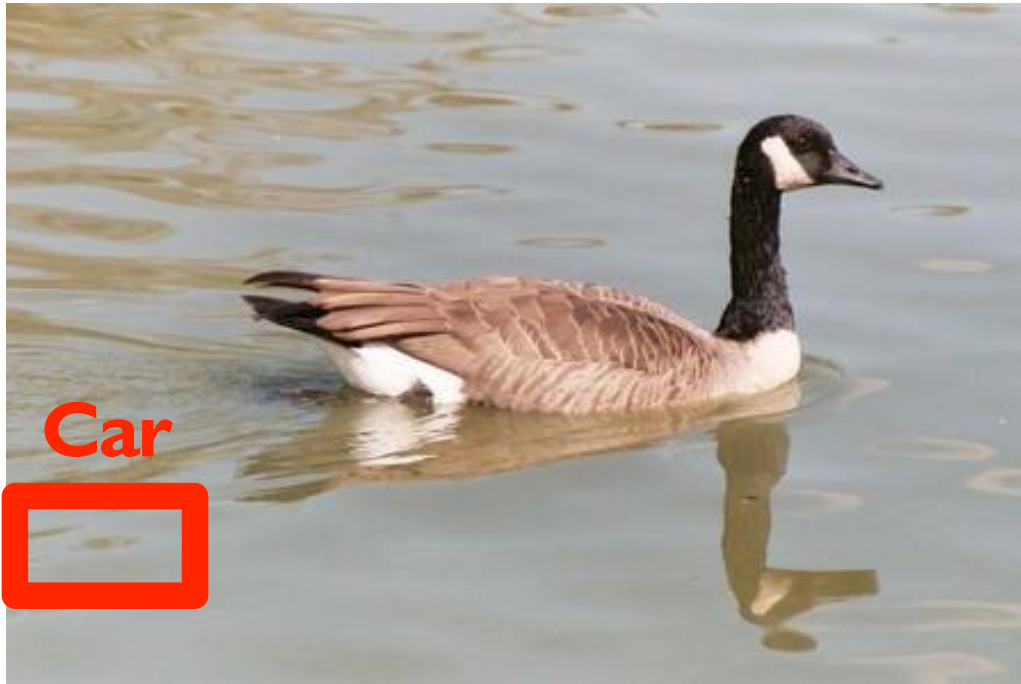


Chair

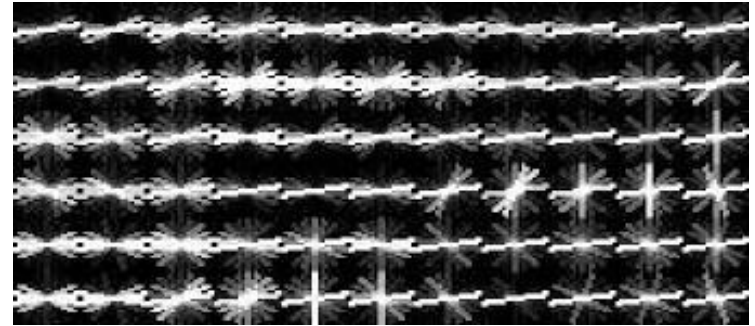
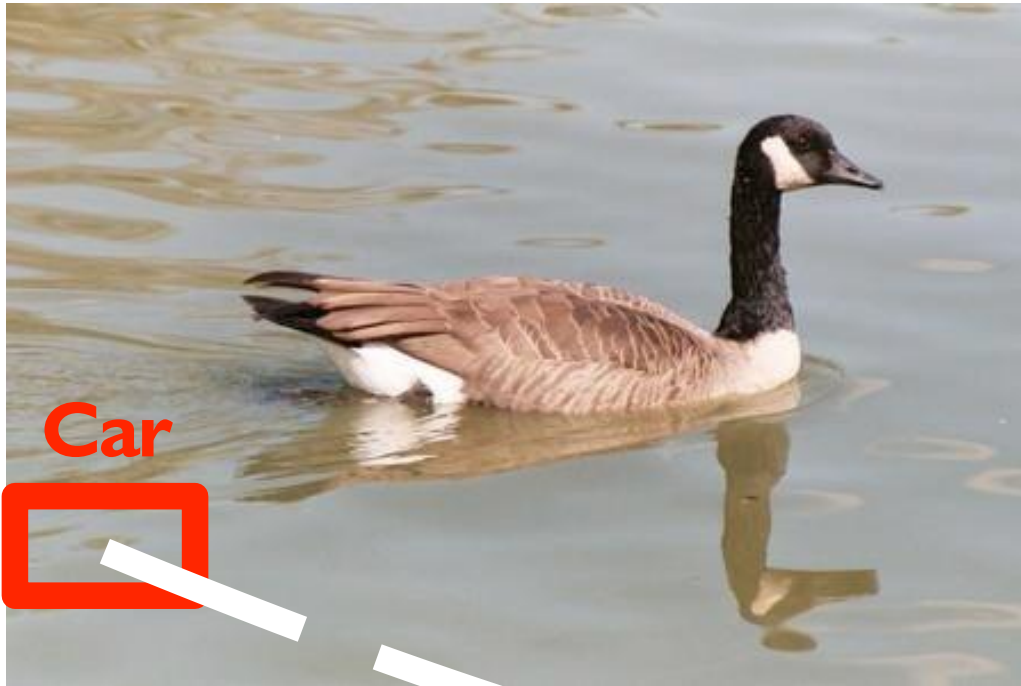


Car

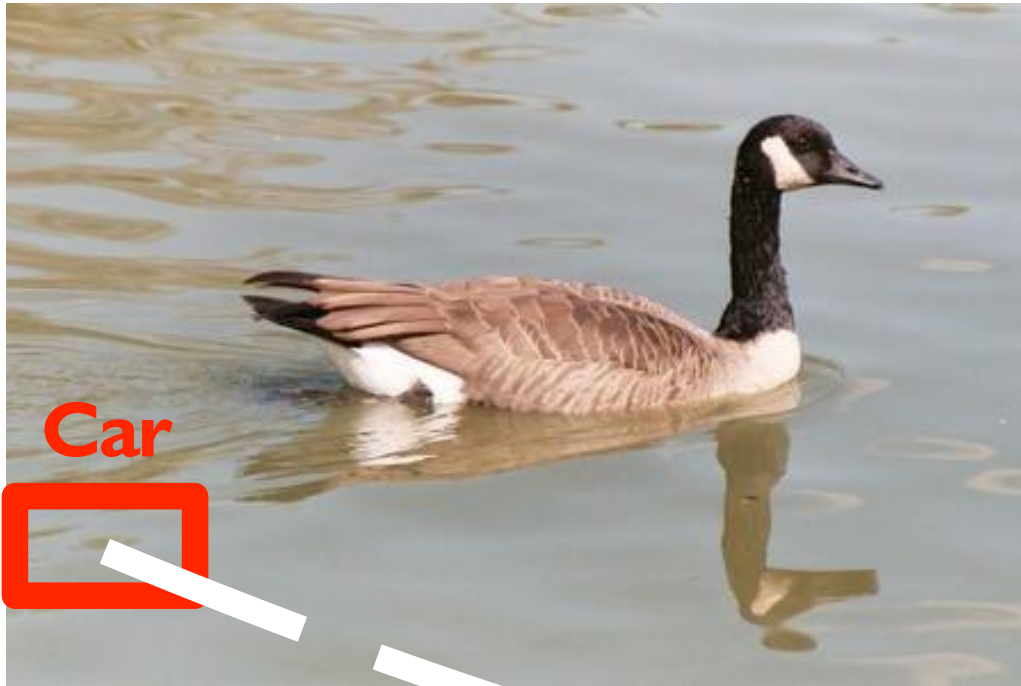
Why did the detector fail?



Why did the detector fail?



Why did the detector fail?



Code Available

Try it on your project 5!

<http://web.mit.edu/vondrick/ihog/>

```
ihog = invertHOG(feats);
```



Friday: CV for Social ~~Good~~ Bad

- We saw how dataset bias can introduce error.
- ...and how features can be ambiguous.
- What about label bias errors?
- How might I move towards a more flexible label system?

Describing Objects by their Attributes

Ali Farhadi, Ian Endres,
Derek Hoiem, David Forsyth

CVPR 2009





What do we want to know about this object?



What do we want to know about this object?

Object recognition expert:
“Dog”



What do we want to know about this object?

Object recognition expert:
“Dog”

Person in the Scene:
“Big pointy teeth”, “Can move fast”, “Looks angry”

Goal: Infer Object Properties



Can I **draw with it**?

Can I **put stuff in it**?

What **shape** is it?

Is it **alive**?

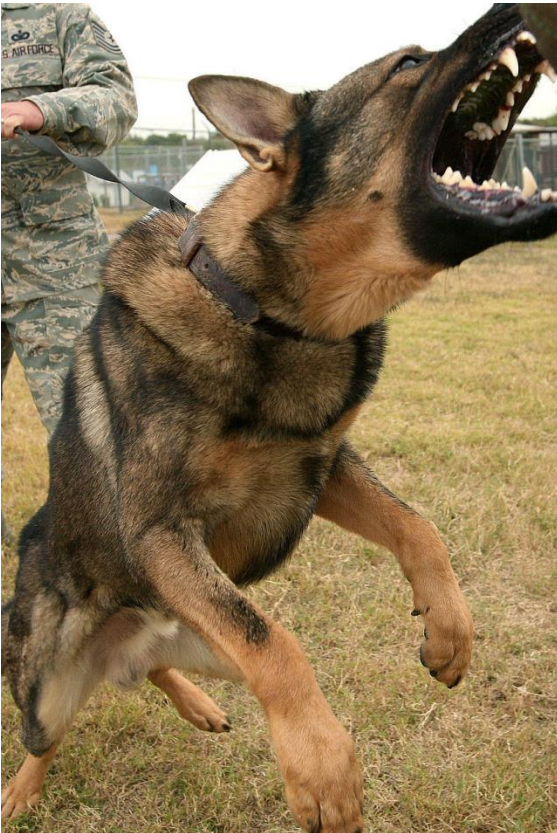
Is it **soft**?

Does it have a **tail**?

Will it **blend**?

Why Infer Properties

1. We want detailed information about objects



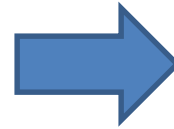
“Dog”
vs.

“Large, angry animal with pointy teeth”

Why Infer Properties

2. We want to be able to infer something about unfamiliar objects – “zero shot learning”

Familiar Objects



New Object



Why Infer Properties

2. We want to be able to infer something about unfamiliar objects – “zero shot learning”

If we can infer category names...

Familiar Objects



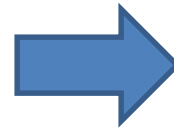
Cat



Horse



Dog



New Object



???

Why Infer Properties

2. We want to be able to infer something about unfamiliar objects – “zero shot learning”

If we can infer properties...

Familiar Objects



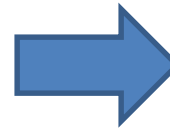
Has Stripes
Has Ears
Has Eyes
....



Has Four Legs
Has Mane
Has Tail
Has Snout
....



Brown
Muscular
Has Snout
....



New Object



Has Stripes (like cat)
Has Mane and Tail (like horse)
Has Snout (like horse and dog)

Why Infer Properties

3. We want to make comparisons between objects or categories

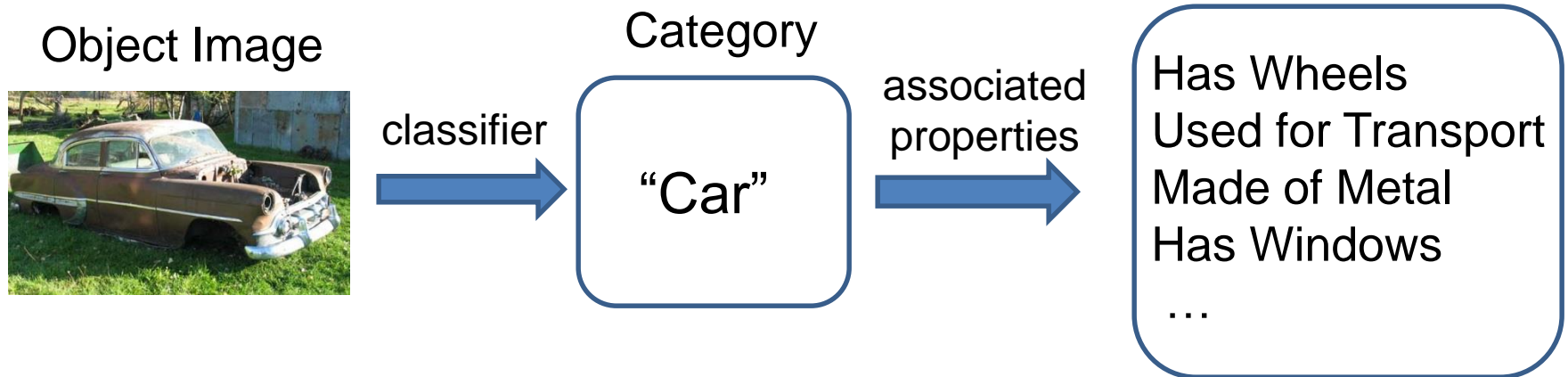


What is unusual about this dog?

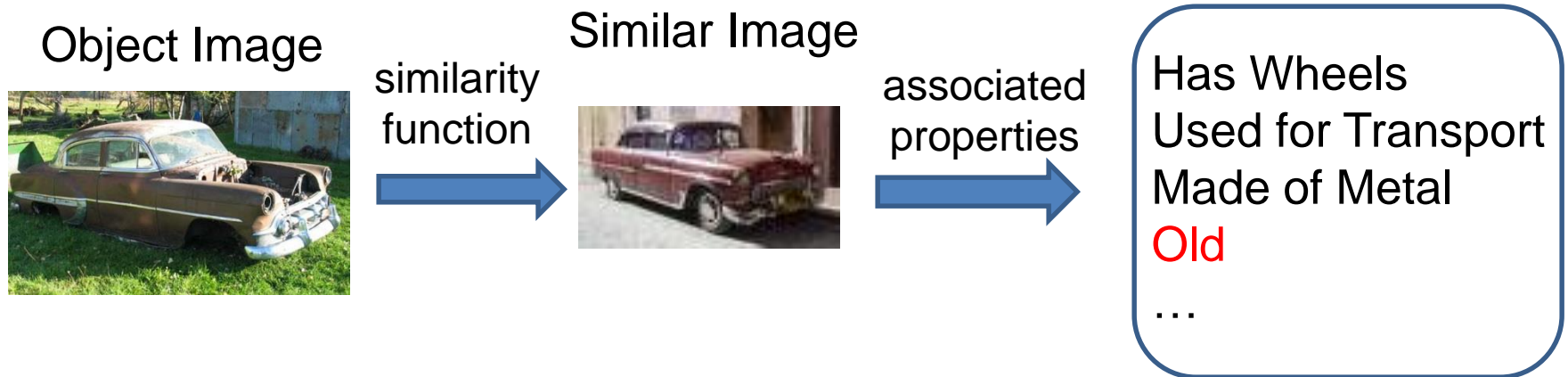


What is the difference between horses and zebras?

Strategy 1: Category Recognition



Strategy 2: Exemplar Matching



Malisiewicz Efros 2008

Hays Efros 2008

Efros et al. 2003

Strategy 3: Infer Properties Directly

Object Image



classifier for each attribute



No Wheels

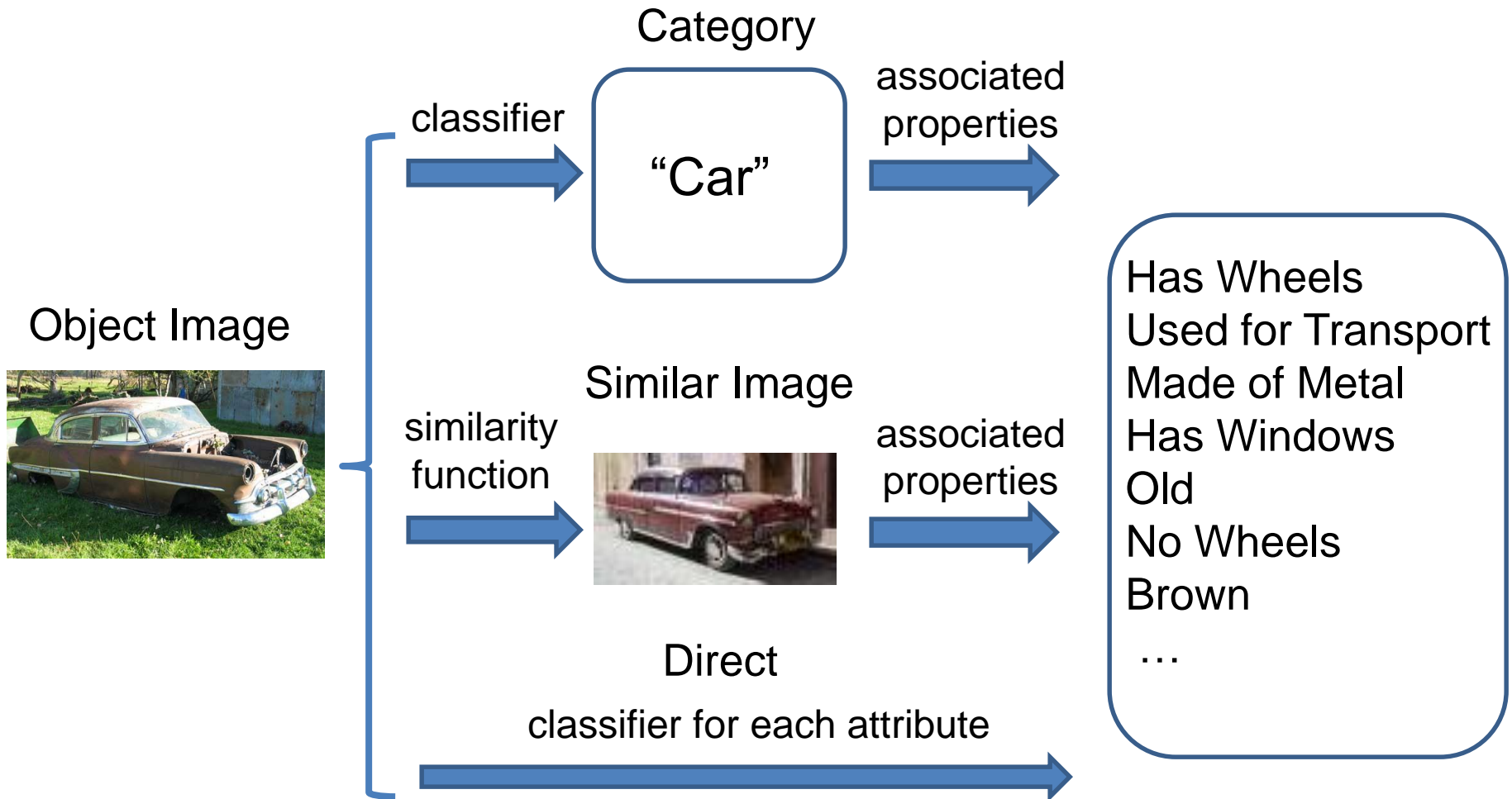
Old

Brown

Made of Metal

...

The Three Strategies



Candidate attributes

- Visible parts: “wheels”, “snout”, “eyes”
- Visible materials or material properties: “made of metal”, “shiny”, “clear”, “made of plastic”
- Shape: “3D boxy”, “round”

Attribute Examples



Shape: Horizontal Cylinder

Part: Wing, Propeller, Window, *Wheel*

Material: *Metal*, Glass



Shape:

Part: Window, *Wheel*, Door, Headlight, Side Mirror

Material: *Metal*, Shiny

Attribute Examples



Shape:

Part: Head, Ear, Nose,
Mouth, Hair, Face,
Torso, Hand, Arm

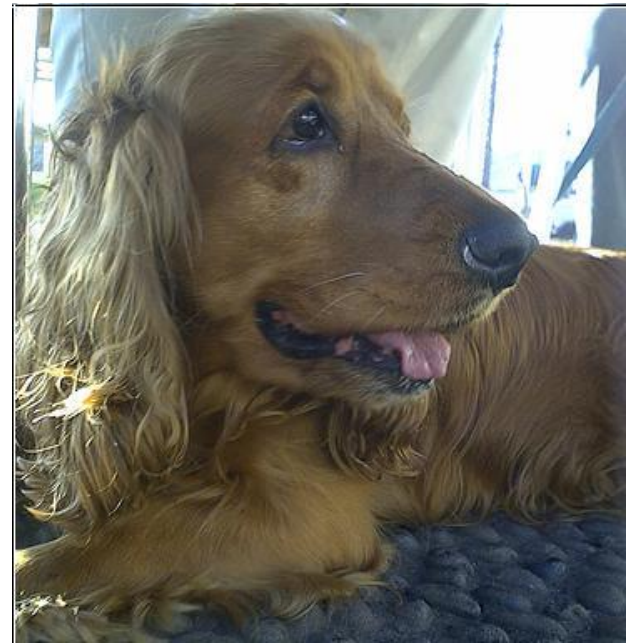
Material: Skin, Cloth



Shape:

Part: Head, Ear, Snout,
Eye

Material: Furry



Shape:


Part: Head, Ear, Snout,
Eye, Torso, Leg

Material: Furry

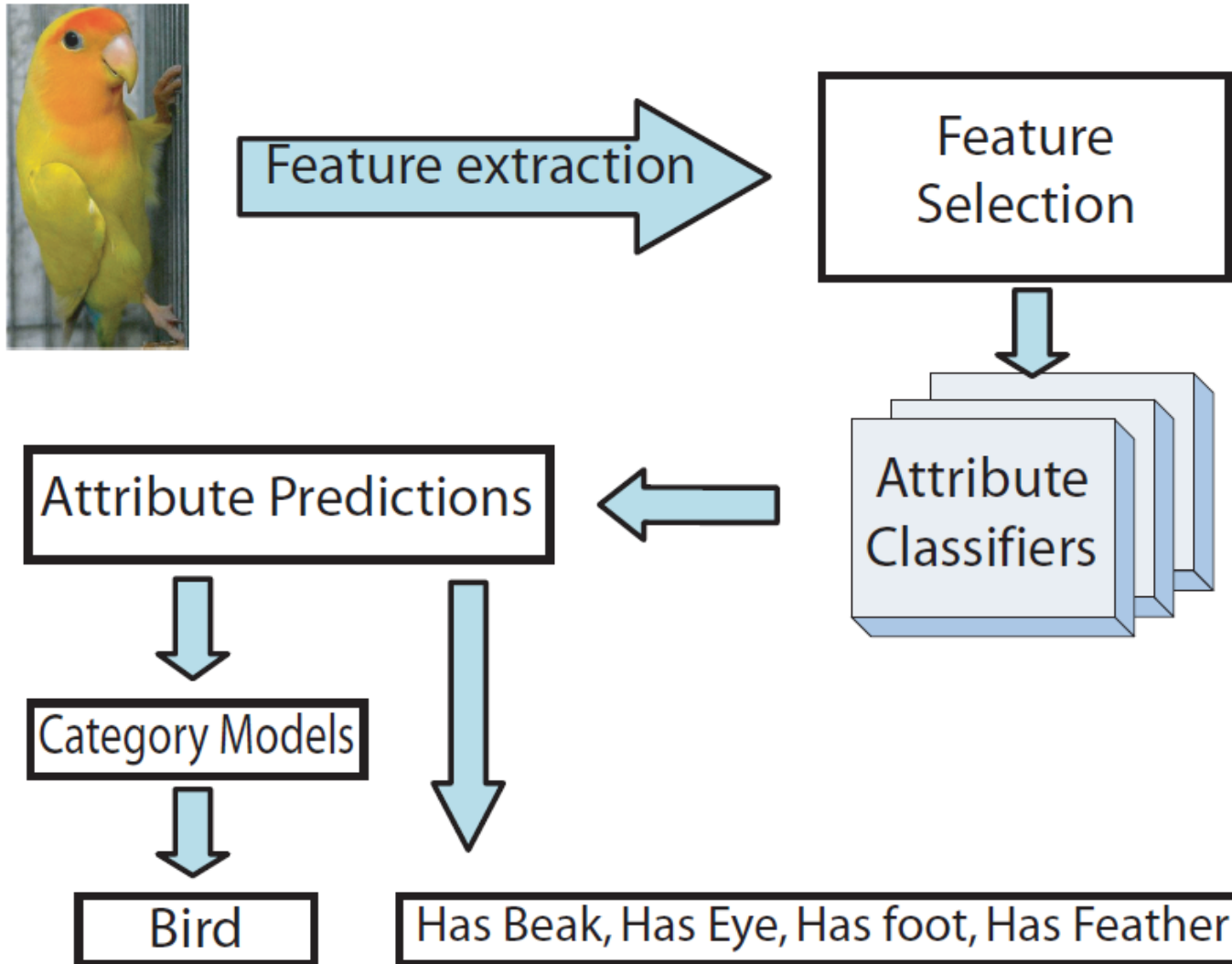
Datasets

- a-Pascal
 - 20 categories from PASCAL 2008 trainval dataset (10K object images)
 - airplane, bicycle, bird, boat, bottle, bus, car, cat, chair, cow, dining table, dog, horse, motorbike, person, potted plant, sheep, sofa, train, tv monitor
 - ‘Ground truth’ for 64 attributes
 - Annotation via Amazon’s Mechanical Turk
- a-Yahoo
 - 12 new categories from Yahoo image search
 - bag, building, carriage, centaur, donkey, goat, jet ski, mug, monkey, statue of person, wolf, zebra
 - Categories chosen to share attributes with those in Pascal
- Attribute labels are somewhat ambiguous
 - Agreement among “experts” 84.3
 - Between experts and Turk labelers 81.4
 - Among Turk labelers 84.1

Annotation on Amazon Turk

cow																																																																	
	<table><tr><td>-Viewpoint-</td><td>-Viewpoint-</td></tr><tr><td><input type="checkbox"/> facing me</td><td><input type="checkbox"/> facing me</td></tr><tr><td><input checked="" type="checkbox"/> * away from me</td><td><input type="checkbox"/> * away from me</td></tr><tr><td><input checked="" type="checkbox"/> facing left</td><td><input checked="" type="checkbox"/> facing left</td></tr><tr><td><input type="checkbox"/> facing right</td><td><input type="checkbox"/> facing right</td></tr><tr><td><input type="checkbox"/> from above</td><td><input type="checkbox"/> from above</td></tr><tr><td><input type="checkbox"/> from below</td><td><input type="checkbox"/> from below</td></tr><tr><td>-Context-</td><td>-Context-</td></tr><tr><td><input checked="" type="checkbox"/> Grass/Field</td><td><input checked="" type="checkbox"/> Grass/Field</td></tr><tr><td><input type="checkbox"/> Street/road</td><td><input type="checkbox"/> Street/road</td></tr><tr><td>--Shape--</td><td>--Shape--</td></tr><tr><td><input type="checkbox"/> Occluded</td><td><input type="checkbox"/> Occluded</td></tr><tr><td>--Part--</td><td>--Part--</td></tr><tr><td><input type="checkbox"/> Tail</td><td><input type="checkbox"/> Tail</td></tr><tr><td><input checked="" type="checkbox"/> Head</td><td><input checked="" type="checkbox"/> Head</td></tr><tr><td><input checked="" type="checkbox"/> Ear</td><td><input checked="" type="checkbox"/> Ear</td></tr><tr><td><input checked="" type="checkbox"/> Snout</td><td><input checked="" type="checkbox"/> Snout</td></tr><tr><td><input checked="" type="checkbox"/> Eye</td><td><input checked="" type="checkbox"/> Eye</td></tr><tr><td><input checked="" type="checkbox"/> Torso</td><td><input checked="" type="checkbox"/> Torso</td></tr><tr><td><input checked="" type="checkbox"/> Leg</td><td><input checked="" type="checkbox"/> Leg</td></tr><tr><td><input checked="" type="checkbox"/> Foot/Shoe</td><td><input type="checkbox"/> Foot/Shoe</td></tr><tr><td><input checked="" type="checkbox"/> Horn</td><td><input type="checkbox"/> Horn</td></tr><tr><td><input type="checkbox"/> Rein</td><td><input type="checkbox"/> Rein</td></tr><tr><td>--Material--</td><td>--Material--</td></tr><tr><td><input checked="" type="checkbox"/> Furry</td><td><input checked="" type="checkbox"/> Furry</td></tr><tr><td>--Pose--</td><td>--Pose--</td></tr><tr><td><input checked="" type="checkbox"/> Standing</td><td><input checked="" type="checkbox"/> Standing</td></tr><tr><td><input type="checkbox"/> Sitting</td><td><input type="checkbox"/> Sitting</td></tr><tr><td><input type="checkbox"/> Walking</td><td><input type="checkbox"/> Walking</td></tr><tr><td><input type="checkbox"/> Lying Straight</td><td><input type="checkbox"/> Lying Straight</td></tr><tr><td><input type="checkbox"/> Lying Curled</td><td><input type="checkbox"/> Lying Curled</td></tr><tr><td><input type="checkbox"/> Open Mouth</td><td><input type="checkbox"/> Open Mouth</td></tr></table>	-Viewpoint-	-Viewpoint-	<input type="checkbox"/> facing me	<input type="checkbox"/> facing me	<input checked="" type="checkbox"/> * away from me	<input type="checkbox"/> * away from me	<input checked="" type="checkbox"/> facing left	<input checked="" type="checkbox"/> facing left	<input type="checkbox"/> facing right	<input type="checkbox"/> facing right	<input type="checkbox"/> from above	<input type="checkbox"/> from above	<input type="checkbox"/> from below	<input type="checkbox"/> from below	-Context-	-Context-	<input checked="" type="checkbox"/> Grass/Field	<input checked="" type="checkbox"/> Grass/Field	<input type="checkbox"/> Street/road	<input type="checkbox"/> Street/road	--Shape--	--Shape--	<input type="checkbox"/> Occluded	<input type="checkbox"/> Occluded	--Part--	--Part--	<input type="checkbox"/> Tail	<input type="checkbox"/> Tail	<input checked="" type="checkbox"/> Head	<input checked="" type="checkbox"/> Head	<input checked="" type="checkbox"/> Ear	<input checked="" type="checkbox"/> Ear	<input checked="" type="checkbox"/> Snout	<input checked="" type="checkbox"/> Snout	<input checked="" type="checkbox"/> Eye	<input checked="" type="checkbox"/> Eye	<input checked="" type="checkbox"/> Torso	<input checked="" type="checkbox"/> Torso	<input checked="" type="checkbox"/> Leg	<input checked="" type="checkbox"/> Leg	<input checked="" type="checkbox"/> Foot/Shoe	<input type="checkbox"/> Foot/Shoe	<input checked="" type="checkbox"/> Horn	<input type="checkbox"/> Horn	<input type="checkbox"/> Rein	<input type="checkbox"/> Rein	--Material--	--Material--	<input checked="" type="checkbox"/> Furry	<input checked="" type="checkbox"/> Furry	--Pose--	--Pose--	<input checked="" type="checkbox"/> Standing	<input checked="" type="checkbox"/> Standing	<input type="checkbox"/> Sitting	<input type="checkbox"/> Sitting	<input type="checkbox"/> Walking	<input type="checkbox"/> Walking	<input type="checkbox"/> Lying Straight	<input type="checkbox"/> Lying Straight	<input type="checkbox"/> Lying Curled	<input type="checkbox"/> Lying Curled	<input type="checkbox"/> Open Mouth	<input type="checkbox"/> Open Mouth
-Viewpoint-	-Viewpoint-																																																																
<input type="checkbox"/> facing me	<input type="checkbox"/> facing me																																																																
<input checked="" type="checkbox"/> * away from me	<input type="checkbox"/> * away from me																																																																
<input checked="" type="checkbox"/> facing left	<input checked="" type="checkbox"/> facing left																																																																
<input type="checkbox"/> facing right	<input type="checkbox"/> facing right																																																																
<input type="checkbox"/> from above	<input type="checkbox"/> from above																																																																
<input type="checkbox"/> from below	<input type="checkbox"/> from below																																																																
-Context-	-Context-																																																																
<input checked="" type="checkbox"/> Grass/Field	<input checked="" type="checkbox"/> Grass/Field																																																																
<input type="checkbox"/> Street/road	<input type="checkbox"/> Street/road																																																																
--Shape--	--Shape--																																																																
<input type="checkbox"/> Occluded	<input type="checkbox"/> Occluded																																																																
--Part--	--Part--																																																																
<input type="checkbox"/> Tail	<input type="checkbox"/> Tail																																																																
<input checked="" type="checkbox"/> Head	<input checked="" type="checkbox"/> Head																																																																
<input checked="" type="checkbox"/> Ear	<input checked="" type="checkbox"/> Ear																																																																
<input checked="" type="checkbox"/> Snout	<input checked="" type="checkbox"/> Snout																																																																
<input checked="" type="checkbox"/> Eye	<input checked="" type="checkbox"/> Eye																																																																
<input checked="" type="checkbox"/> Torso	<input checked="" type="checkbox"/> Torso																																																																
<input checked="" type="checkbox"/> Leg	<input checked="" type="checkbox"/> Leg																																																																
<input checked="" type="checkbox"/> Foot/Shoe	<input type="checkbox"/> Foot/Shoe																																																																
<input checked="" type="checkbox"/> Horn	<input type="checkbox"/> Horn																																																																
<input type="checkbox"/> Rein	<input type="checkbox"/> Rein																																																																
--Material--	--Material--																																																																
<input checked="" type="checkbox"/> Furry	<input checked="" type="checkbox"/> Furry																																																																
--Pose--	--Pose--																																																																
<input checked="" type="checkbox"/> Standing	<input checked="" type="checkbox"/> Standing																																																																
<input type="checkbox"/> Sitting	<input type="checkbox"/> Sitting																																																																
<input type="checkbox"/> Walking	<input type="checkbox"/> Walking																																																																
<input type="checkbox"/> Lying Straight	<input type="checkbox"/> Lying Straight																																																																
<input type="checkbox"/> Lying Curled	<input type="checkbox"/> Lying Curled																																																																
<input type="checkbox"/> Open Mouth	<input type="checkbox"/> Open Mouth																																																																

Approach



Features + classifiers

Spatial pyramid histograms of quantized

- Color and texture for **materials**
- Histograms of gradients (HOG) for **parts**
- Canny edges for **shape**

Learn presence / absence of attribute.

- Train one classifier (linear SVM) per attribute

Average ROC Area

Trained on a-PASCAL objects

Test Objects	Parts	Materials	Shape
a-PASCAL	0.794	0.739	0.739
a-Yahoo	0.726	0.645	0.677

Describing Objects by their Attributes



'is 3D Boxy'
'is Vert Cylinder'
'has Window'
'has Row Wind'
X'has Headlight'



'has Hand'
'has Arm'
X'has Screen'
'has Plastic'
'is Shiny'



'has Head'
'has Hair'
'has Face'
X'has Saddle'
'has Skin'

No examples from these object categories were seen during training

Describing Objects by their Attributes



'is 3D Boxy'
'has Wheel'
'has Window'
'is Round'
'has Torso'



'has Tail'
'has Snout'
'has Leg'
X 'has Text'
X 'has Plastic'

No examples from these object categories were seen during training

Category Recognition

- Semantic attributes not enough
 - 74% accuracy even with ground truth attributes
- Introduce discriminative attributes
 - Trained by selecting subset of classes and features
 - Dogs vs. sheep using color
 - Cars and buses vs. motorbikes and bicycles using edges
 - Train 10,000 and select 1,000 most reliable, according to a validation set

Attributes not big help when sufficient data

- Use attribute predictions as features
- Train linear SVM to categorize objects

PASCAL 2008	Base Features	Semantic Attributes	All Attributes
Classification Accuracy	58.5%	54.6%	59.4%
Class-normalized Accuracy	35.5%	28.4%	37.7%

Absence of typical attributes



Aeroplane
No "wing"



Car
No "window"



Boat
No "sail"



Aeroplane
No "jet engine"



Motorbike
No "side mirror"



Car
No "door"



Sheep
No "wool"

752 reports

68% are correct

Presence of atypical attributes



Motorbike
"cloth"



People
"label"



Bird
"Leaf"



Bus
"face"



Aeroplane
"beak"



Sofa
"wheel"



Bike
"Horn"

951 reports
47% are correct

Visual Recognition with Humans in the Loop

**Steve Branson, Catherine Wah, Florian Schroff,
Boris Babenko, Peter Welinder, Pietro Perona,
Serge Belongie**

Part of the [Visipedia project](#)

Introduction:

(A) Easy for Humans



Chair? Airplane? ...

Computers starting
to get good at this.

(B) Hard for Humans



Finch? Bunting?...

If it's hard for humans,
it's probably too hard
for computers.

(C) Easy for Humans

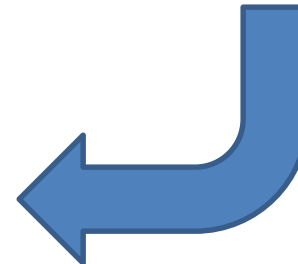


Yellow Belly? Blue Belly? ...

Semantic feature
extraction difficult for
computers.



Combine strengths
to solve this
problem.

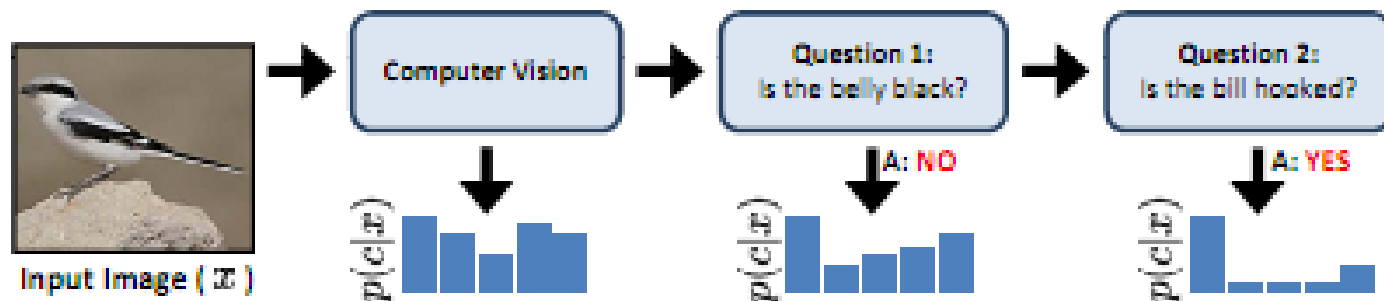


The Approach: What is progress?

- Supplement visual recognition with the human capacity for visual feature extraction to tackle difficult (fine-grained) recognition problems.
- Typical progress is viewed as increasing data difficulty while maintaining full autonomy
- Reduction in human effort on difficult data.

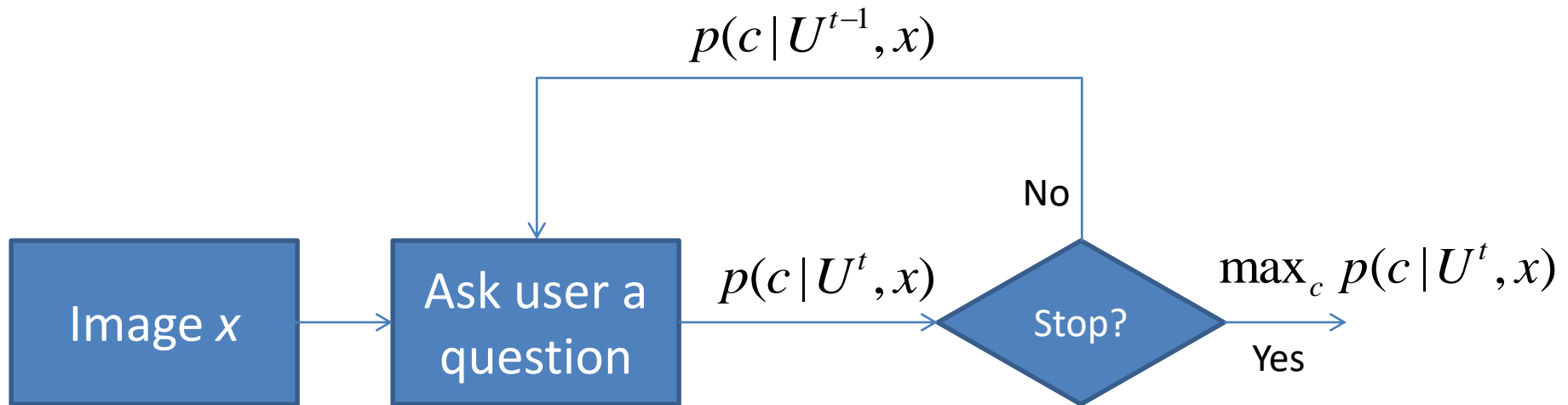
The Approach: 20 Questions

- Ask the user a series of discriminative visual questions to make the classification.



Which 20 questions?

- At each step, exploit the image itself and the user response history to select the most informative question to ask next.



Which question to ask?

- The question that will reduce entropy the most, taking into consideration the computer vision classifier confidences for each category.

The Dataset: Birds-200

- 6033 images of 200 species



Implementation



- Assembled 25 visual questions encompassing 288 visual attributes extracted from www.whatbird.com
- Mechanical Turk users asked to answer questions and provide confidence scores.

User Responses.

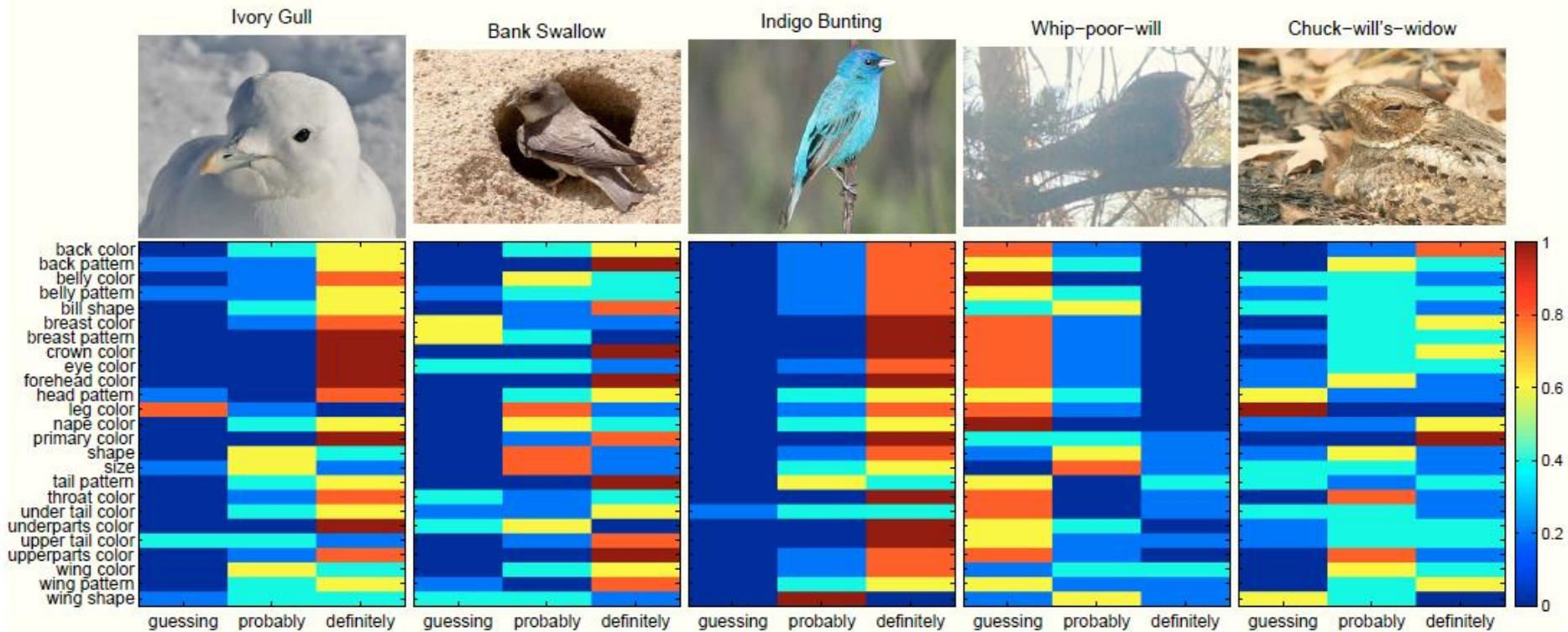
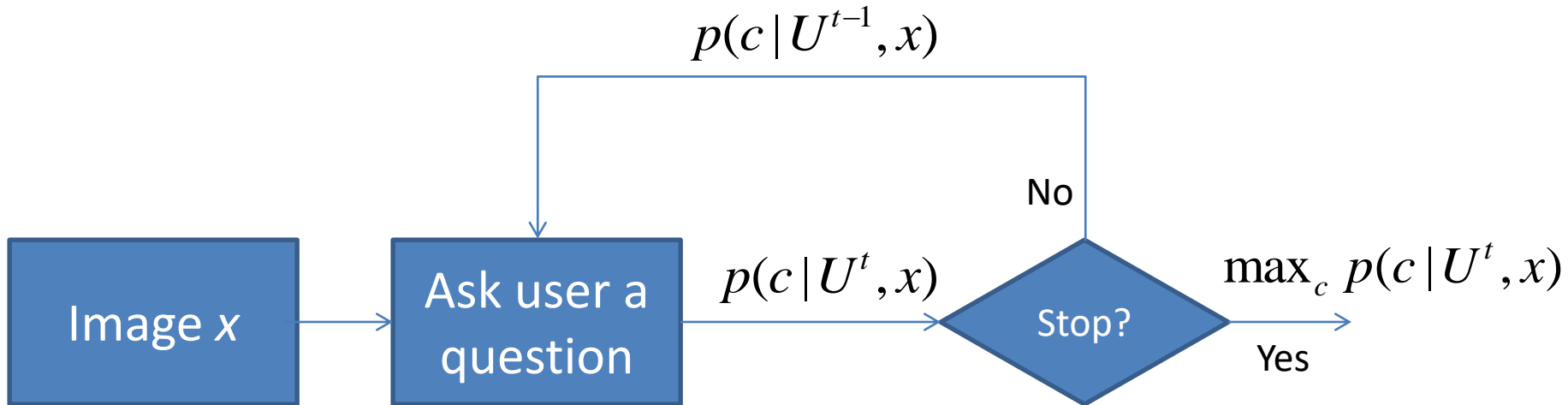


Fig. 4. Examples of user responses for each of the 25 attributes. The distribution over $\{Guessing, Probably, Definitely\}$ is color coded with blue denoting 0% and red denoting 100% of the five answers per image attribute pair.

Visual recognition

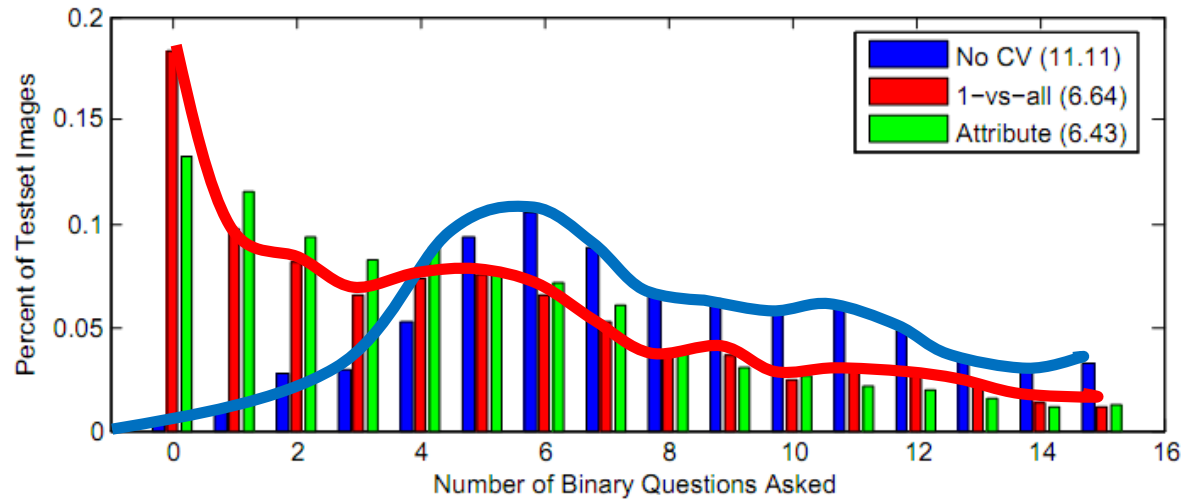
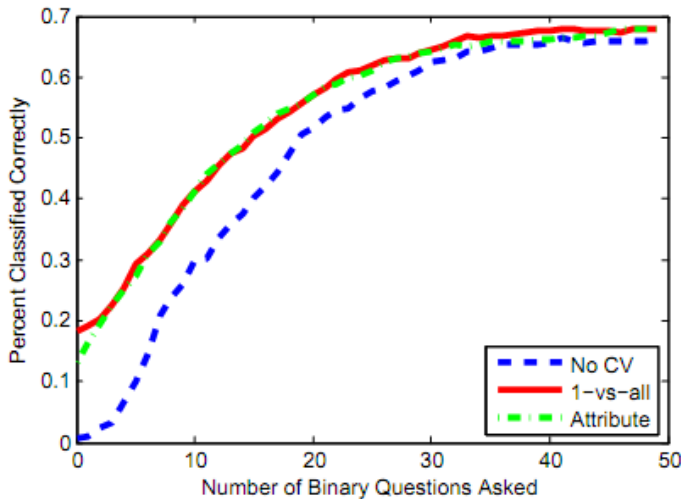
- Any vision system that can output a probability distribution across classes will work.
- Authors used Andrea Vedaldi's code.
 - Color/gray SIFT
 - VQ geometric blur
 - 1 v All SVM
- Authors added full image color histograms and VQ color histograms

Experiments



- 2 Stop criteria:
 - Fixed number of questions – evaluate accuracy
 - User stops when bird identified – measure number of questions required.

Results



- Average number of questions to make ID reduced from 11.11 to 6.43
- Method allows CV to handle the easy cases, consulting with users only on the more difficult cases.

Key Observations

- Visual recognition reduces labor over a pure “20 Q” approach.
- Visual recognition improves performance over a pure “20 Q” approach. (69% vs 66%)
- User input dramatically improves recognition results. (66% vs 19%)

Strengths and weaknesses

- Handles very difficult data and yields excellent results.
- Plug-and-play with many recognition algorithms.
- Requires significant user assistance
- Reported results assume humans are perfect verifiers
- Is the reduction from 11 questions to 6 really that significant?