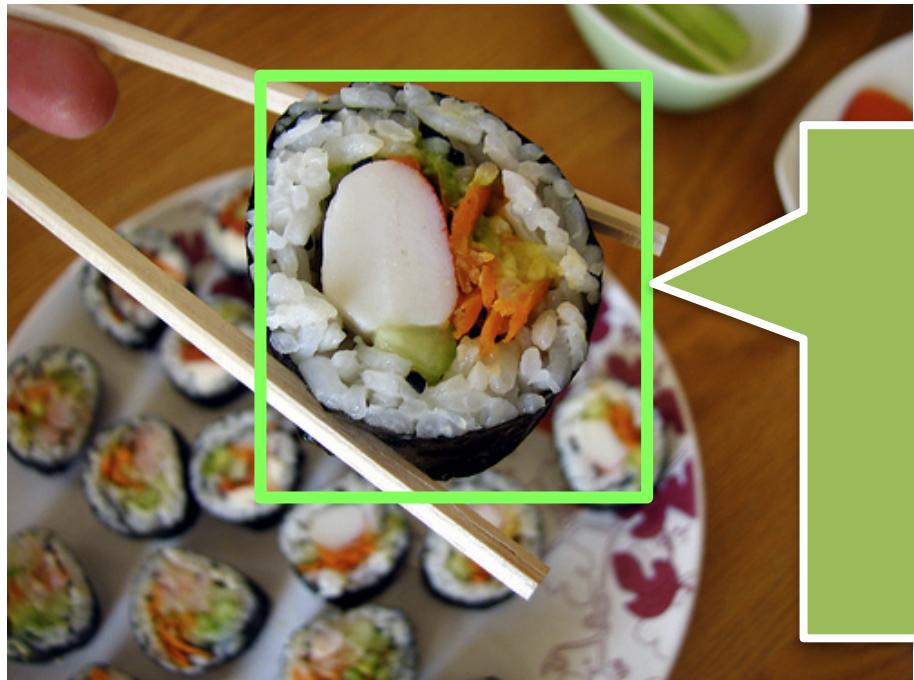


# Lecture: Large Scale Recognition

Juan Carlos Niebles and Ranjay Krishna  
Stanford Vision and Learning Lab



## California Roll

Ingredients: Rice, Seaweed, Crab, Cucumber, Avocado

Calories: 40

Fat: 7g

Carb: 40g

Protein: 5g

**Gluten Free**

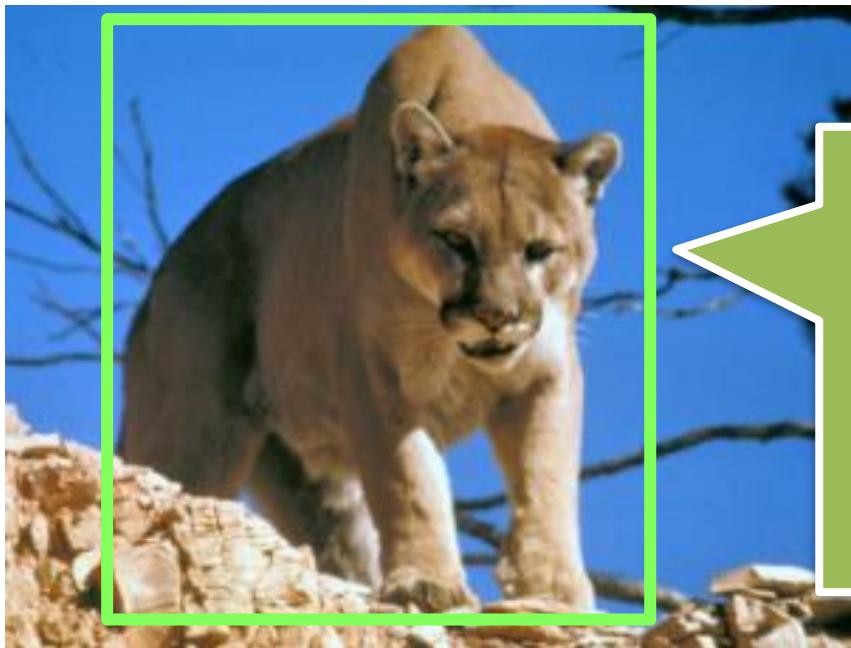


## Amanita phalloides

[http://en.wikipedia.org/wiki/  
Amanita\\_phalloides](http://en.wikipedia.org/wiki/Amanita_phalloides)

**TOXIC. DO NOT EAT**





## Mountain Lion

**DO NOT RUN**

Raise arms to appear larger.  
Show your teeth



**IKEA POANG Chair  
ON SALE  
\$29.00 at [ikea.com](http://ikea.com)**



## Mornonga (Japanese flying squirrel)

Inhabits sub-alpine forests in Japan.  
Nocturnal. Eats seeds, fruit, tree leaves  
(Wikipedia)

I wish my computer could recognize  
EVERYTHING





Surveillance



Robotics



Assistive tools



Wearable devices



Smart photo album

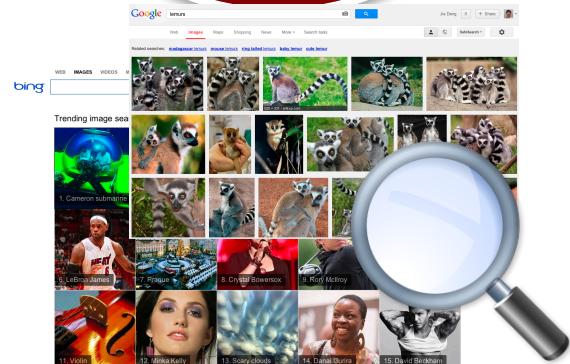


Image search



Driverless cars



Mining social media

# **What can computers already recognize?**



[]  
AUTO



Nikon

The Nikon S60. Detects up to 12 faces.



# Google Goggles

Use pictures to search the web.

New!



[Text](#)



[Landmarks](#)



[Books](#)



[Contact Info](#)



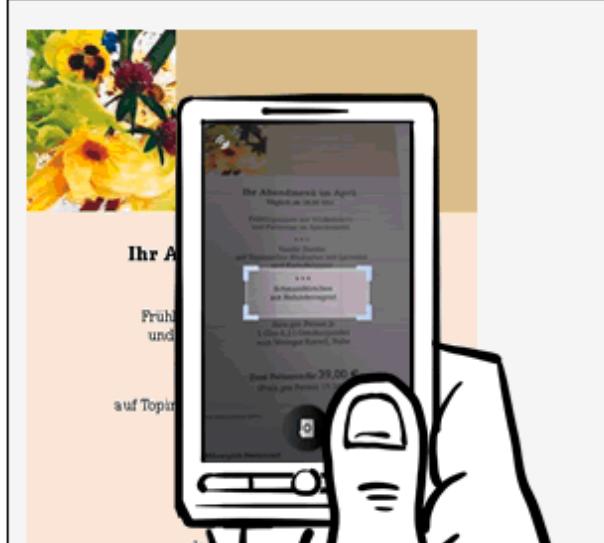
[Artwork](#)



[Wine](#)

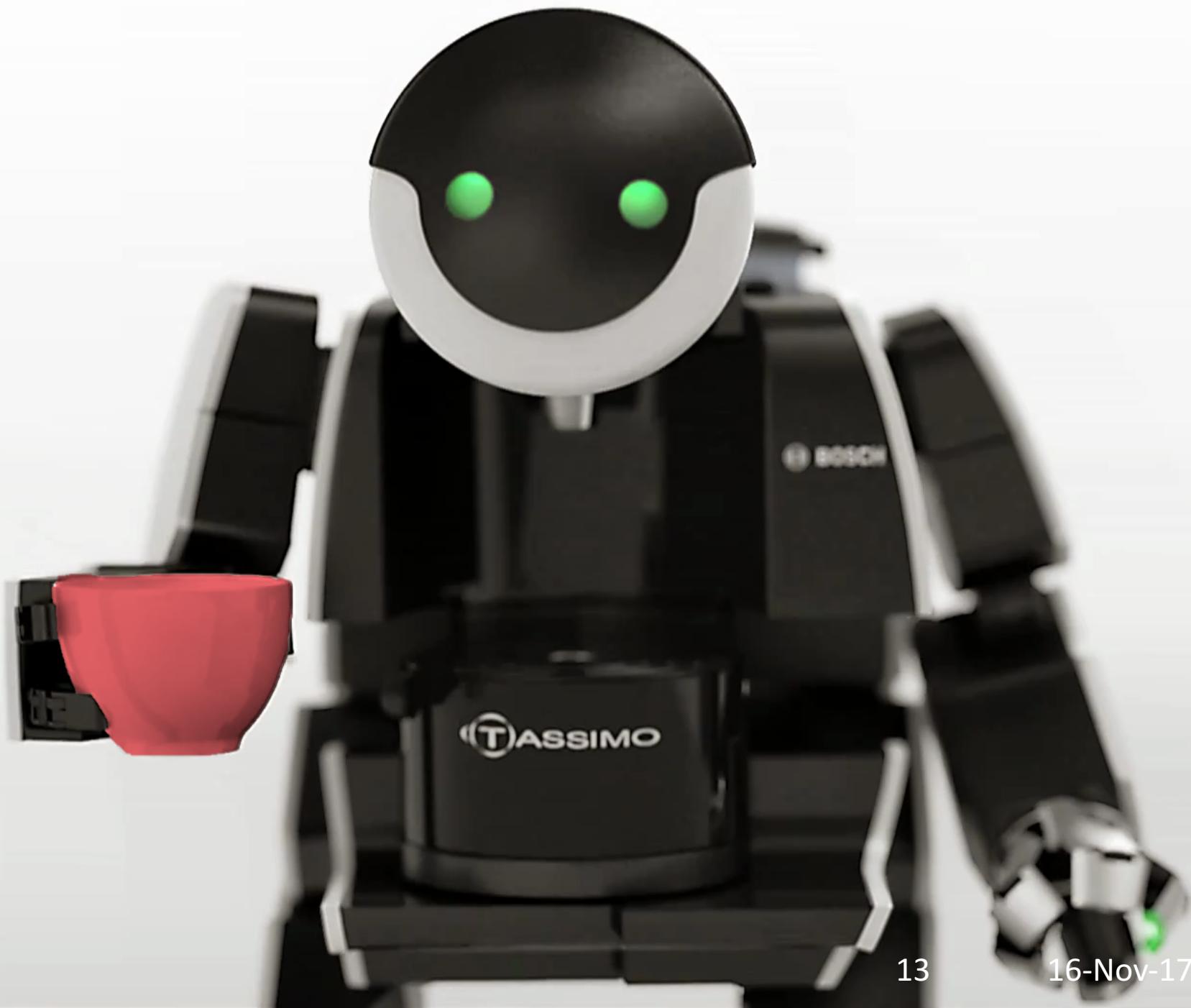


[Logos](#)



**What's the next to work on?**

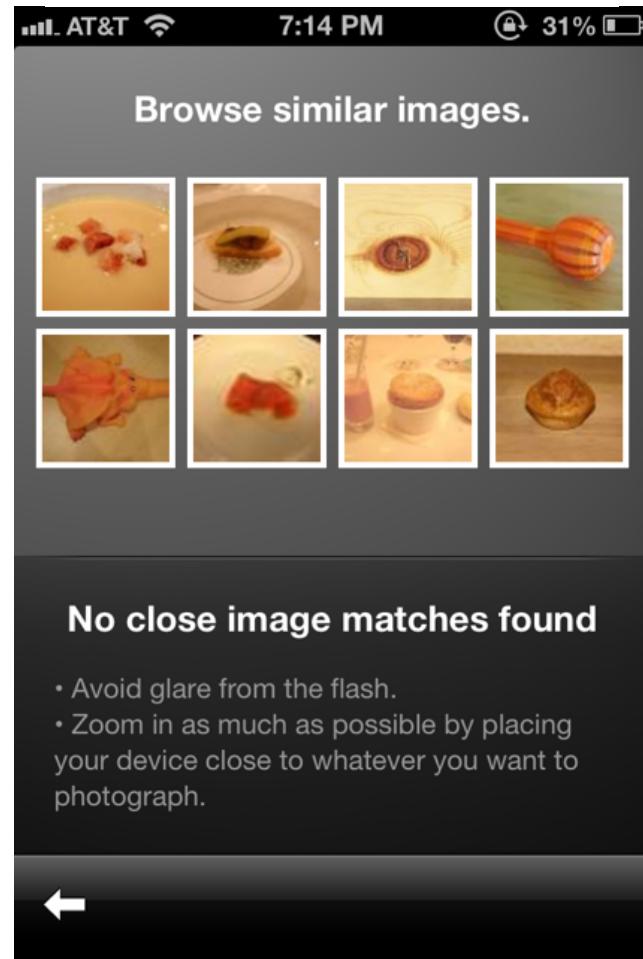
**Coffee Mugs!**





## Google Goggles

Use pictures to search the web.



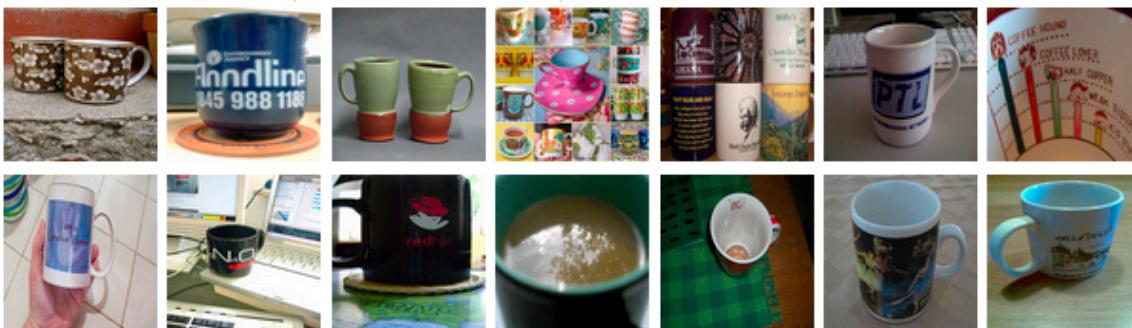
# PASCAL VOC [Everingham et al. 2006-2012]



Airplane	Dining table
Bird	Dog
Boat	Horse
Bike	Motorbike
Bottle	Person
Bus	Potted plant
Car	Sheep
Cat	Sofa
Chair	Train
Cow	TV monitor

No Coffee Mugs!

# The rest of the talk will be about Coffee Mugs!



# What about Gas Pumps!



Image size:  
401 × 604

No other sizes of this image found

Google  
images

[Visually similar images](#) - Report images



The rest of the talk will be about **Coffee Mugs**

**And Gas Pumps**

**And Solar arrays**

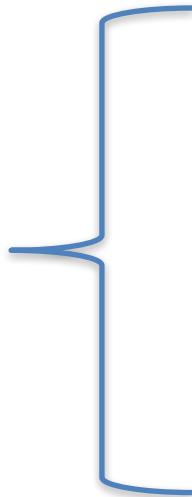
**Radio**

**First aid kit**

**Spacesuit**

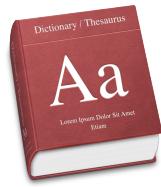
**Oxygen Cylinder**

What do they have in common?



**Let's work on recognizing EVERYTHING**

# How many things are there?



**10K+**

[Biederman '87]

**60K+**

product  
categories

**80K+**

English nouns  
[Miller '95; Fellbaum  
'98]

**3.5M+**

unique tags  
[Sigurbjörnsson &  
Zwol '08]

**4.1M+**

articles

# From 20 classes to Millions?



4 September 2008 | www.nature.com/nature | £10

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

# nature

## THE BITER BIT

Viral infections for viruses

## TROPICAL CYCLONES

The strong get stronger

## BLACK HOLE PHYSICS

A new window on the  
Galactic Centre

# BIG DATA

NATUREJOBS  
Minnesota musings

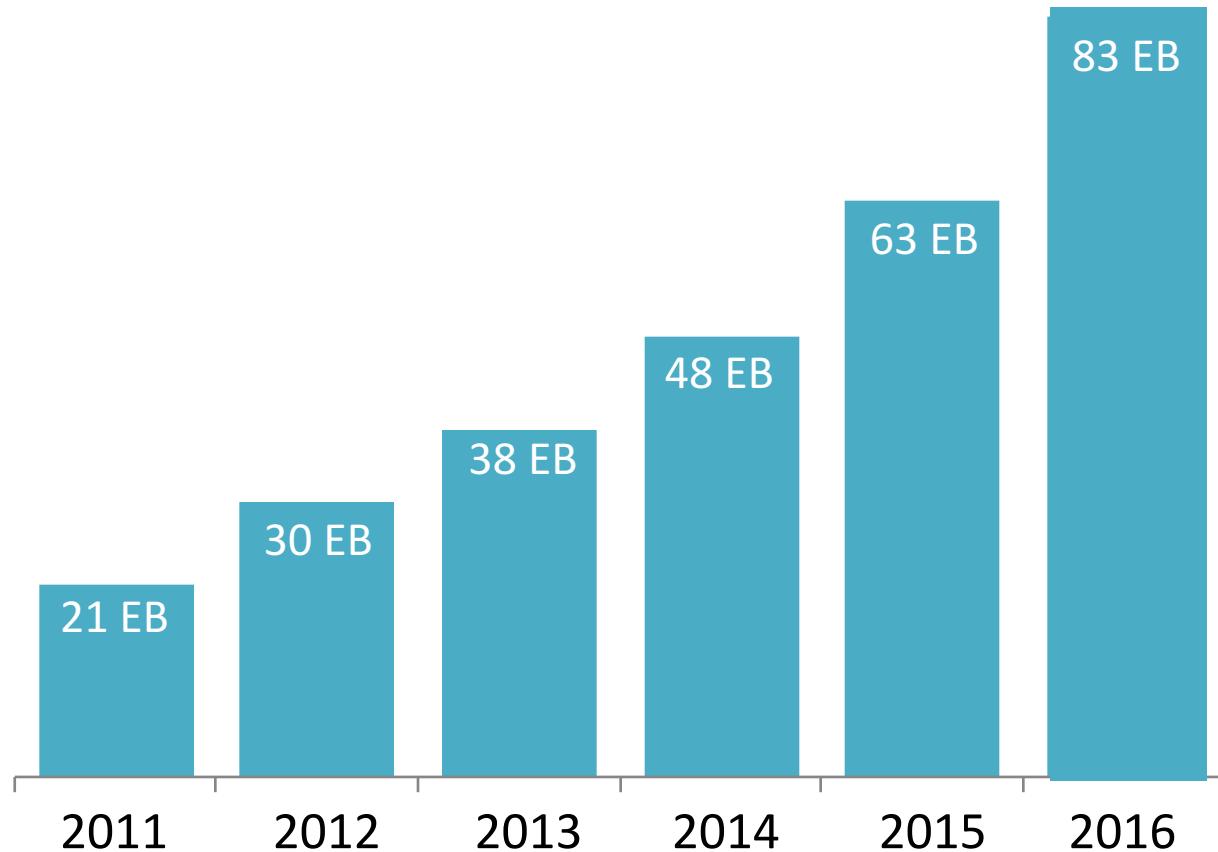
## SCIENCE IN THE PETABYTE ERA



9 770028 085095

# **Big Data from the Internet**

# Global Consumer Internet Traffic Per Month

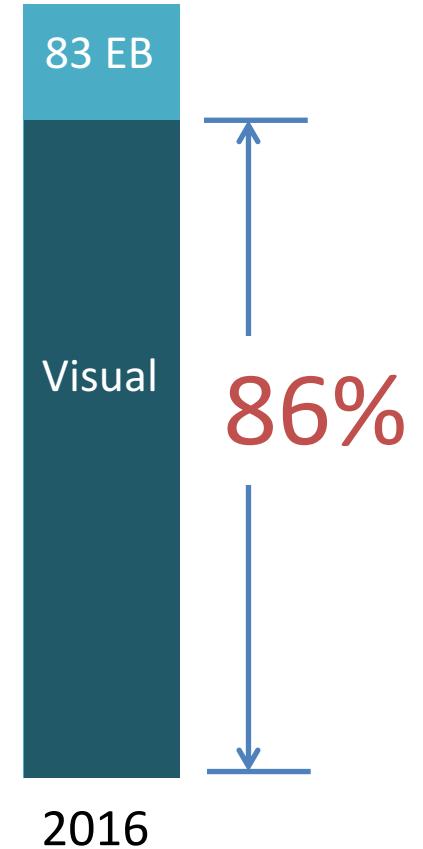




72 hours of videos / min



300 million images / day



## **Big Data from the Internet**

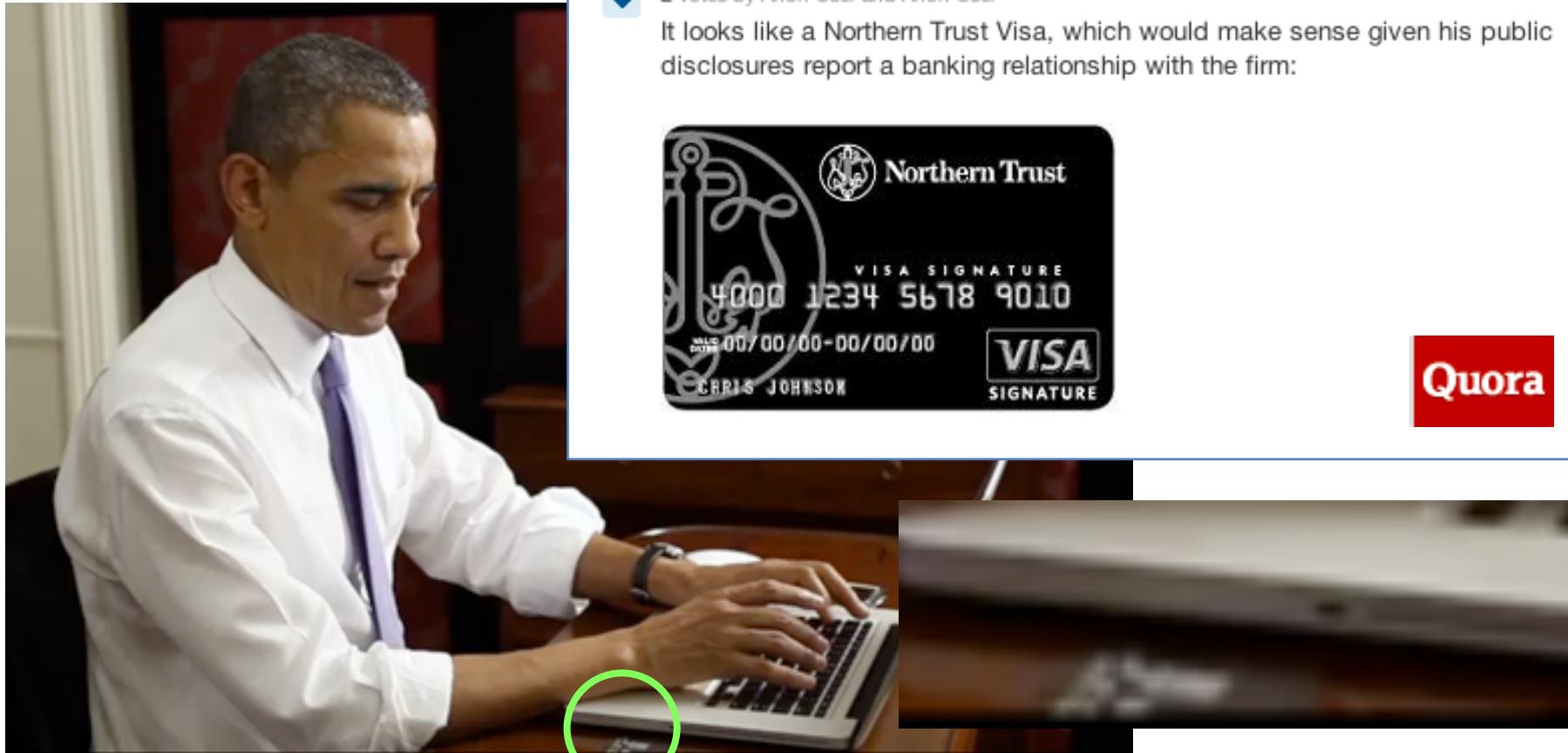
→ The Internet can teach **EVERYTHING**



**Evolution Gone Wild**

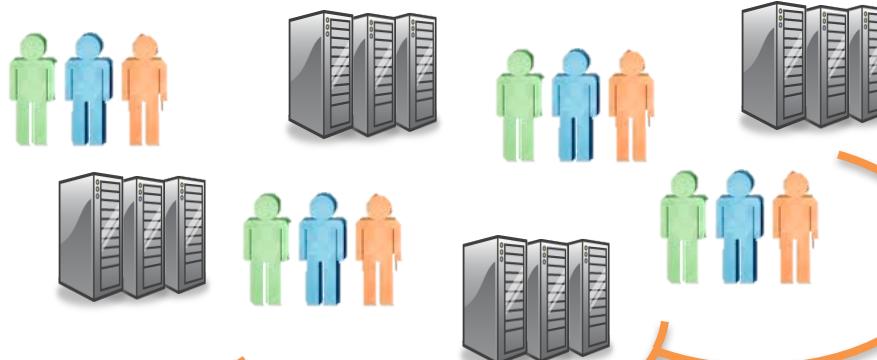
*Future plants and animals*

<http://www.worth1000.com/contests/12705/contest>



**What kind of credit card is President Obama using in this video of him donating to his campaign?**

## The Internet: Machines + Crowd



Big Data

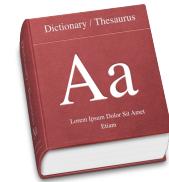
Teach machines to recognize **EVERYTHING**

## PASCAL VOC



20

[Everingham et al.'06-'12]



10K+

[Biederman '87]

Goal: Build a recognition engine on ~~EVERYTHING~~  
**10K classes**



[Deng et al. 2009]

[www.image-net.org](http://www.image-net.org)

**22K** categories and **14M** images

- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
- Food
- Materials
- Structures
  - Artifact
  - Tools
  - Appliances
  - Structures
- Person
- Scenes
  - Indoor
  - Geological Formations
  - Sport Activities

# Number of Labeled Images

SUN, **131K**

[Xiao et al. '10]

LabelMe, **37K**

[Russell et al. '07]

PASCAL VOC, **30K**

[Everingham et al. '06-'12]

Caltech101, **9K**

[Fei-Fei, Fergus, Perona, '03]

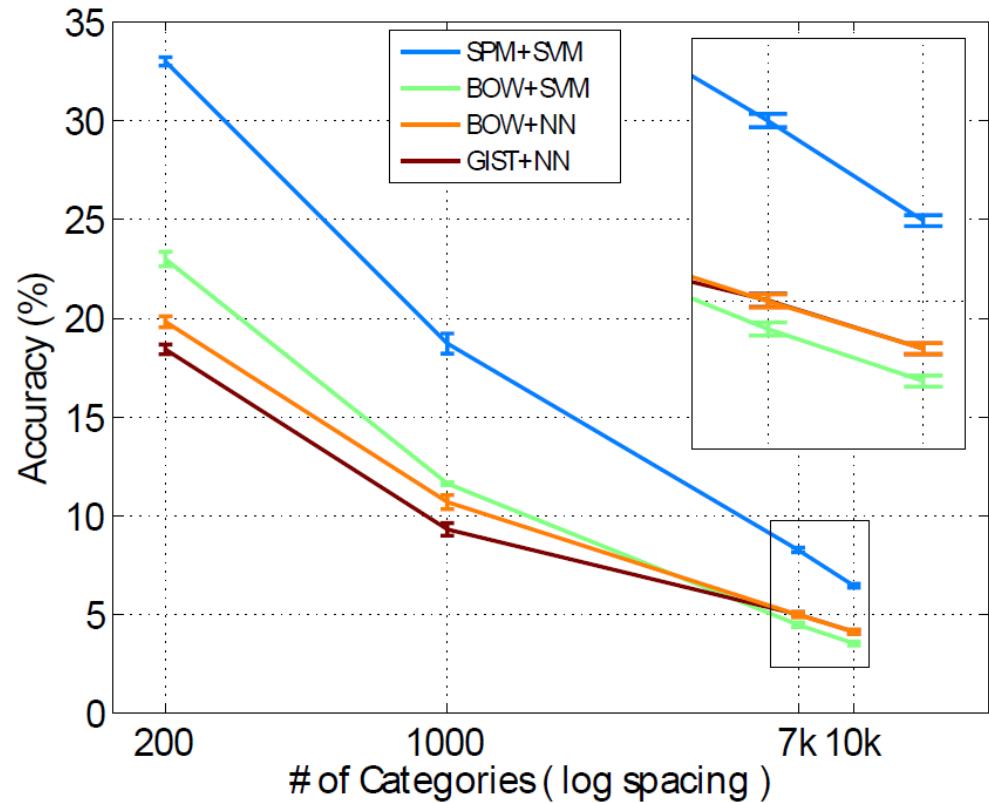
**ImageNet, 14M**

[Deng et al. '09]

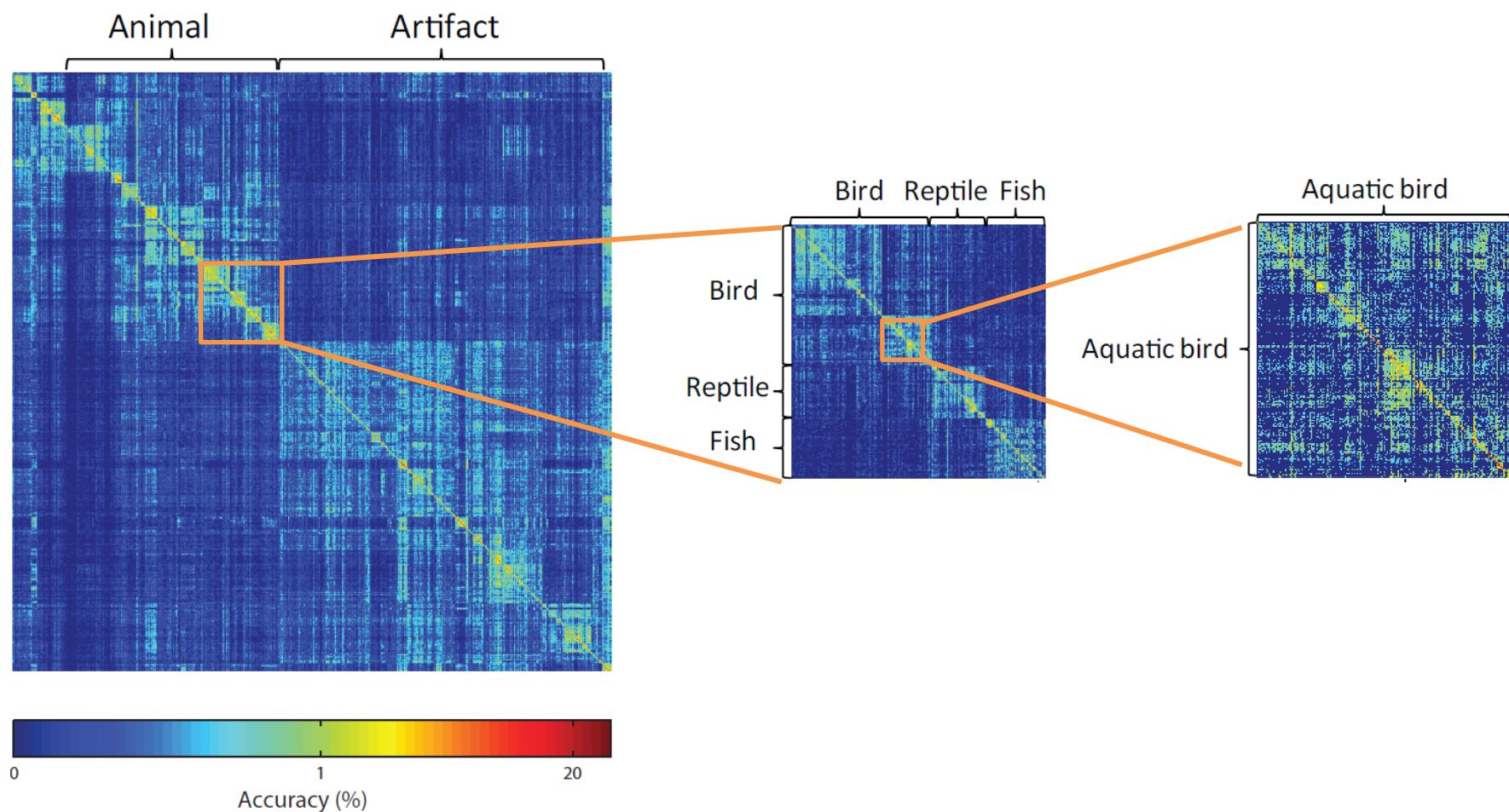


# Learn to Classify 10K Classes

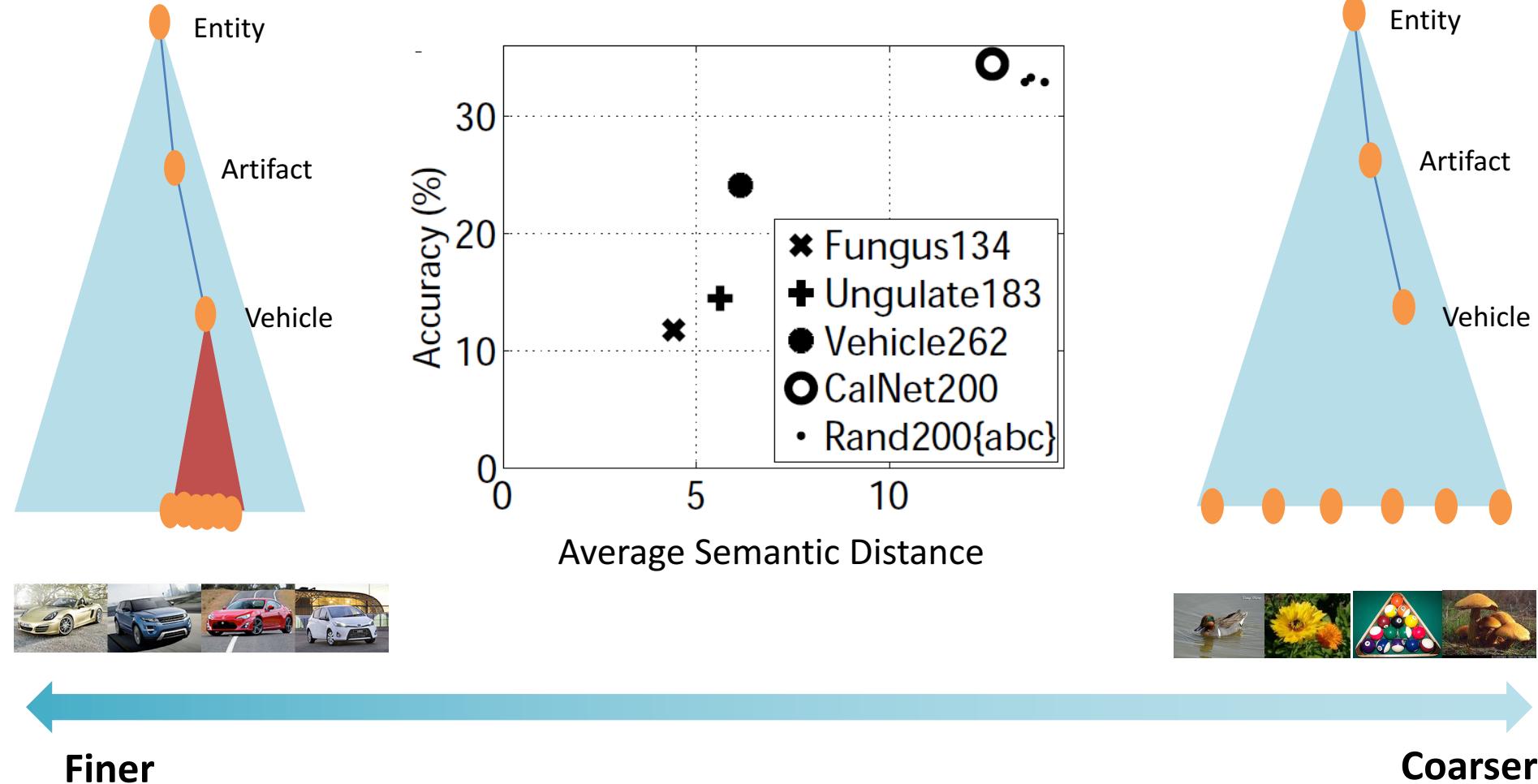
- 9 Million images
- 4 methods
  - SPM+SVM [Lazebnik et al. '06]
  - BOW+SVM [Csurka et al. '04]
  - BOW+NN
  - GIST+NN [Oliva et al. '01]
- 6.4% for 10K categories



# Learn to Classify 10K Classes

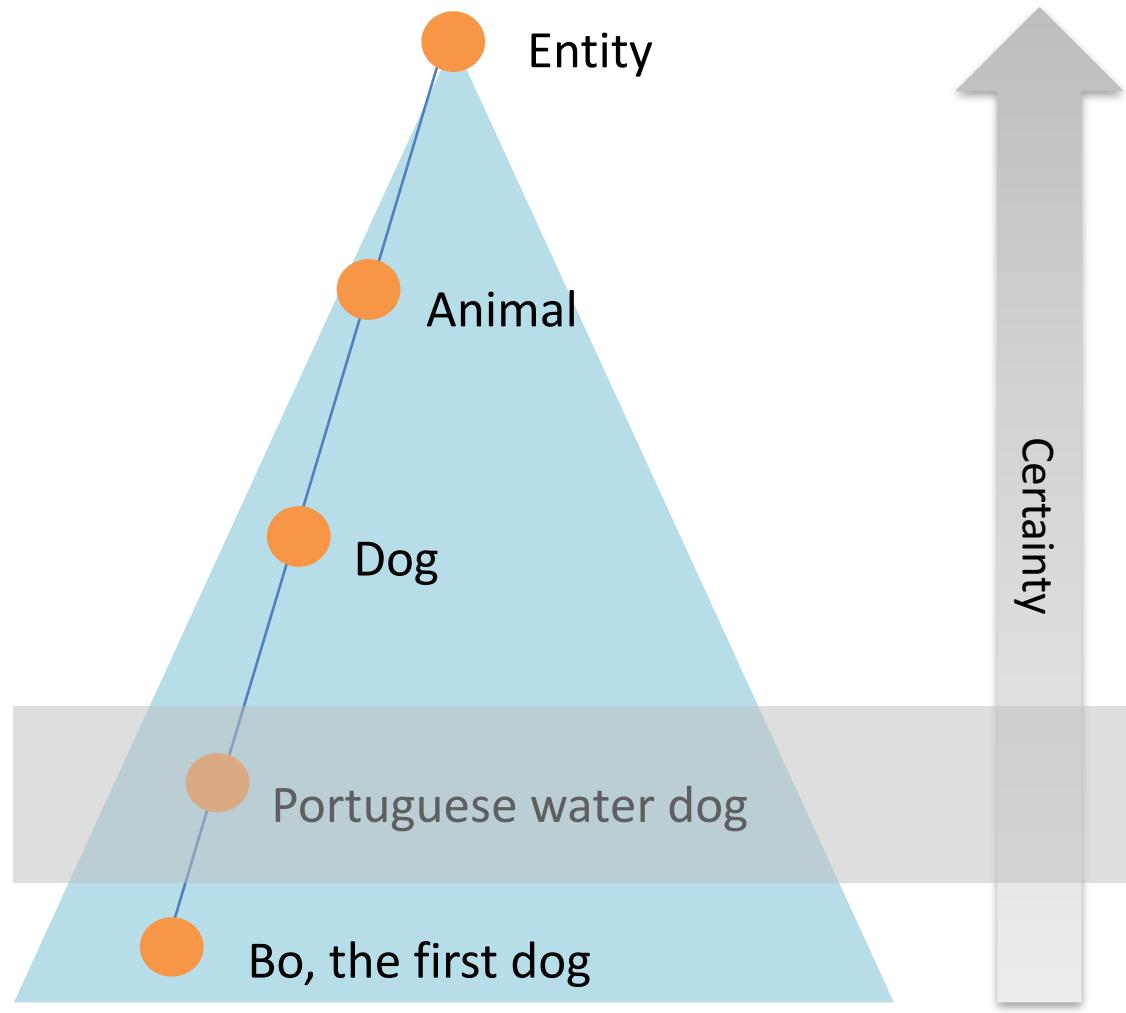
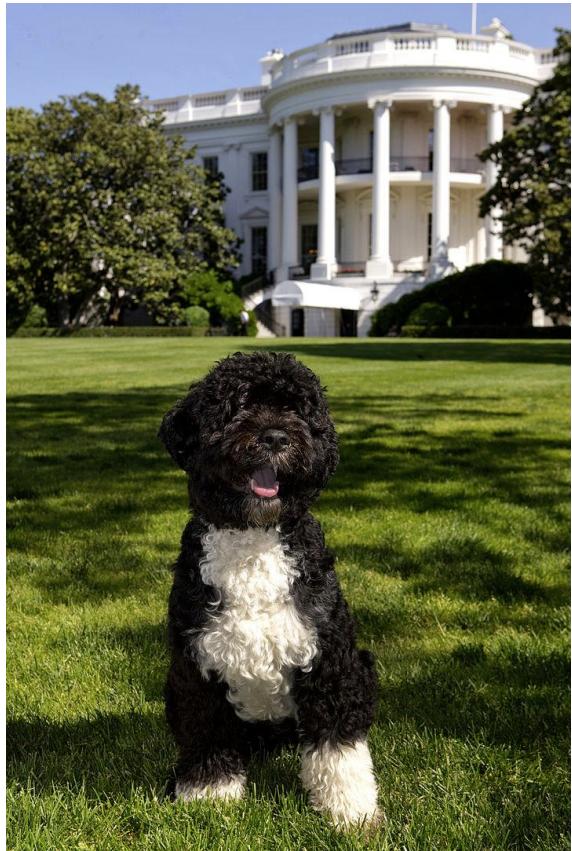


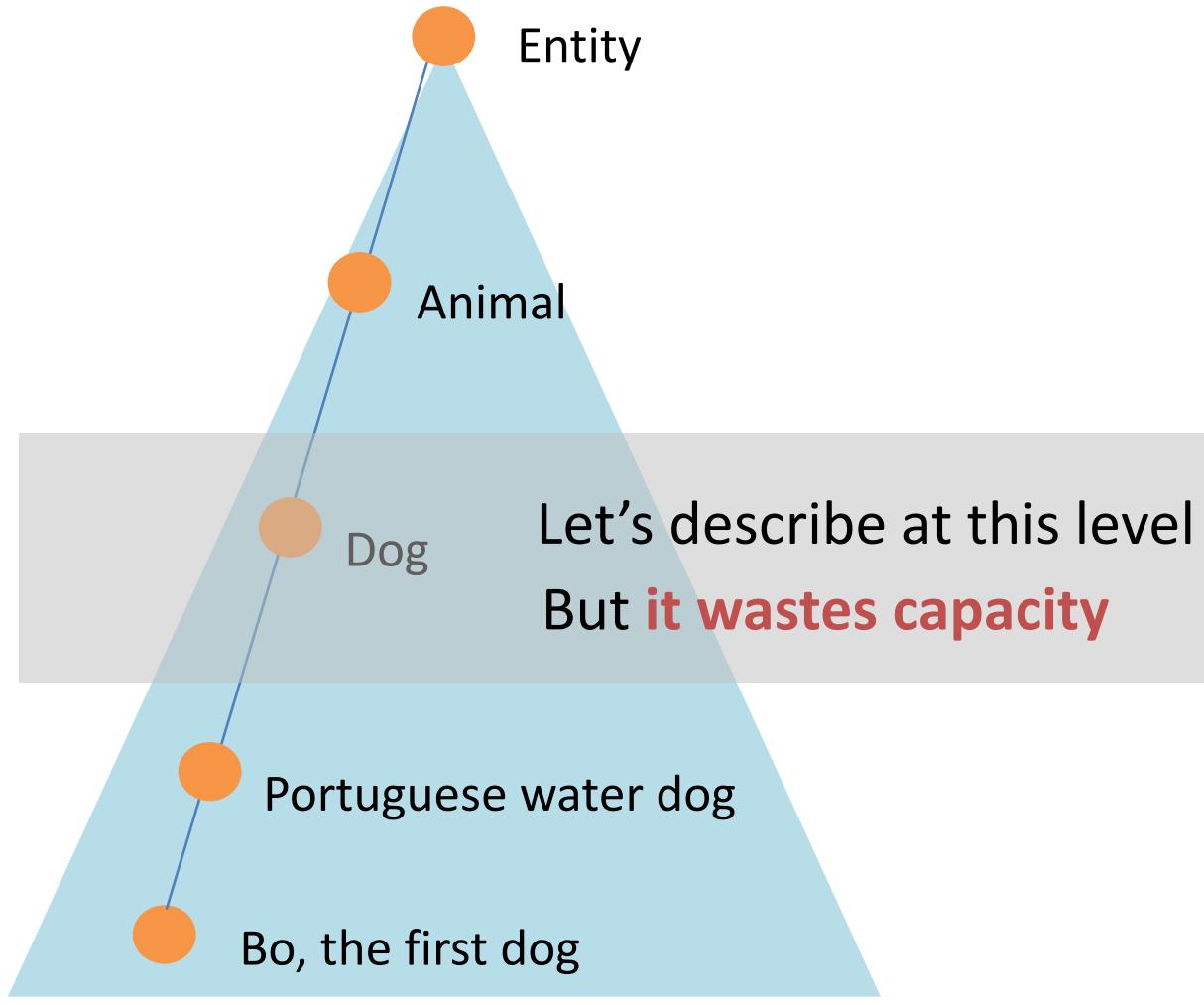
# Fine-grained categories are a lot harder



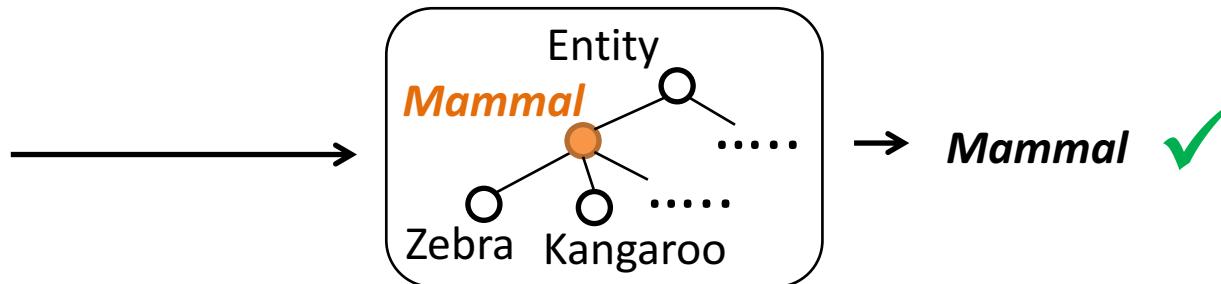
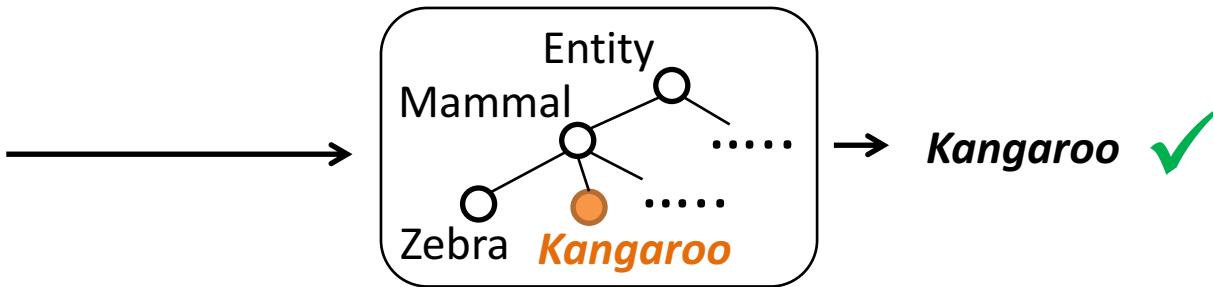
# Challenges

- Semantic hierarchy
- Fine-grained classes
- Large-scale Learning

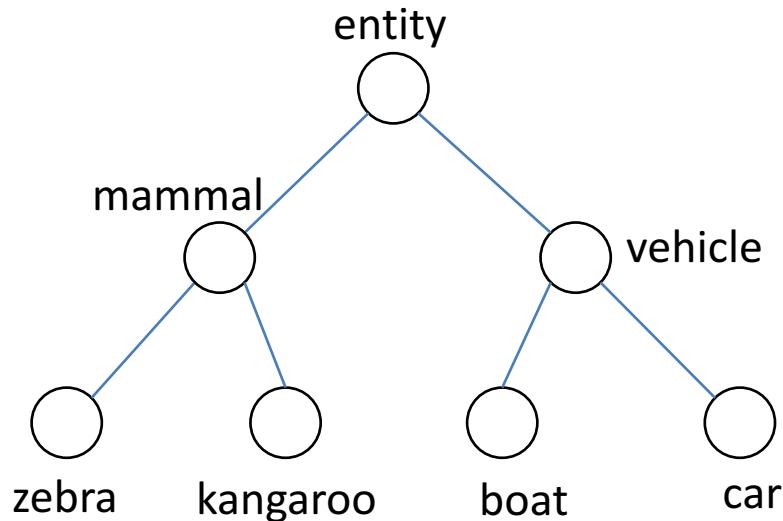




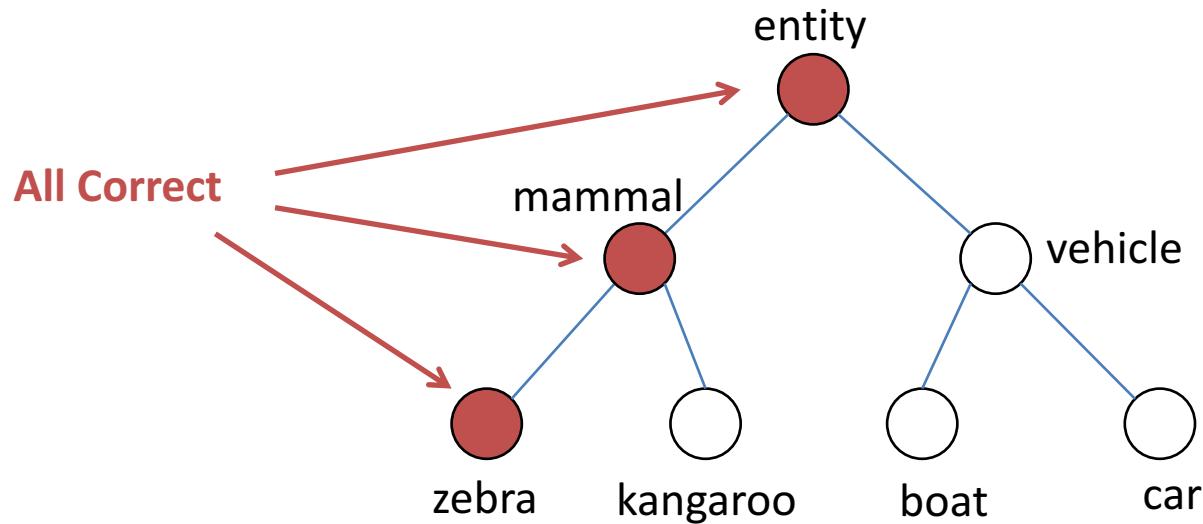
# Hedging: Be as informative as possible with few mistakes



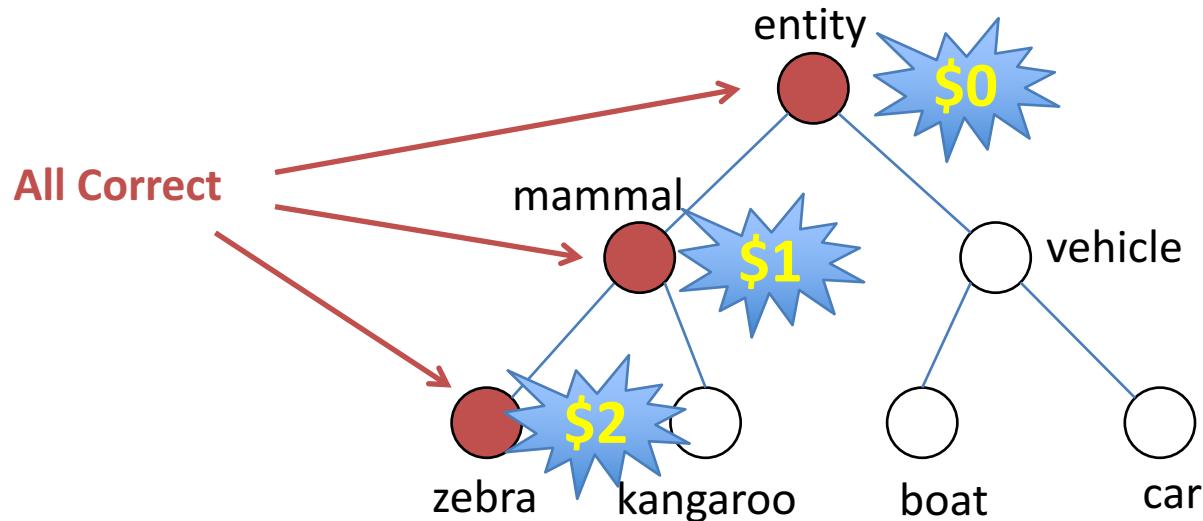
# Formal Problem Statement



# Formal Problem Statement



# Formal Problem Statement



# Formal Problem Statement

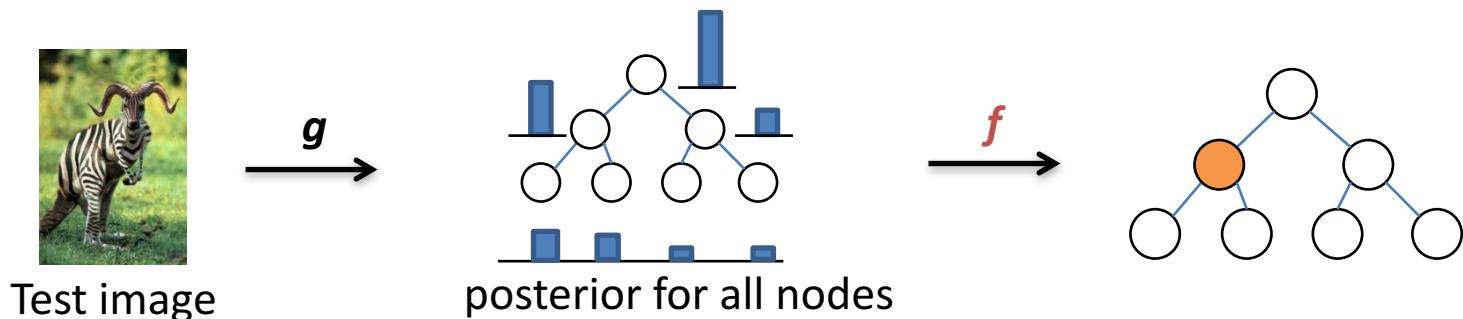
## Assumptions

- Same distribution for training and test.
- A base classifier  $g$  that gives posterior probability on the hierarchy.

## Goal

- Find a **decision rule**  $f$ 
  - Expected accuracy  $A(f)$  is at least  $1-\varepsilon$
  - Maximize expected reward  $R(f)$

$$\begin{aligned} & \underset{f}{\text{Maximize}} \quad R(f) \\ & \text{Subject to } A(f) \geq 1 - \varepsilon \end{aligned}$$

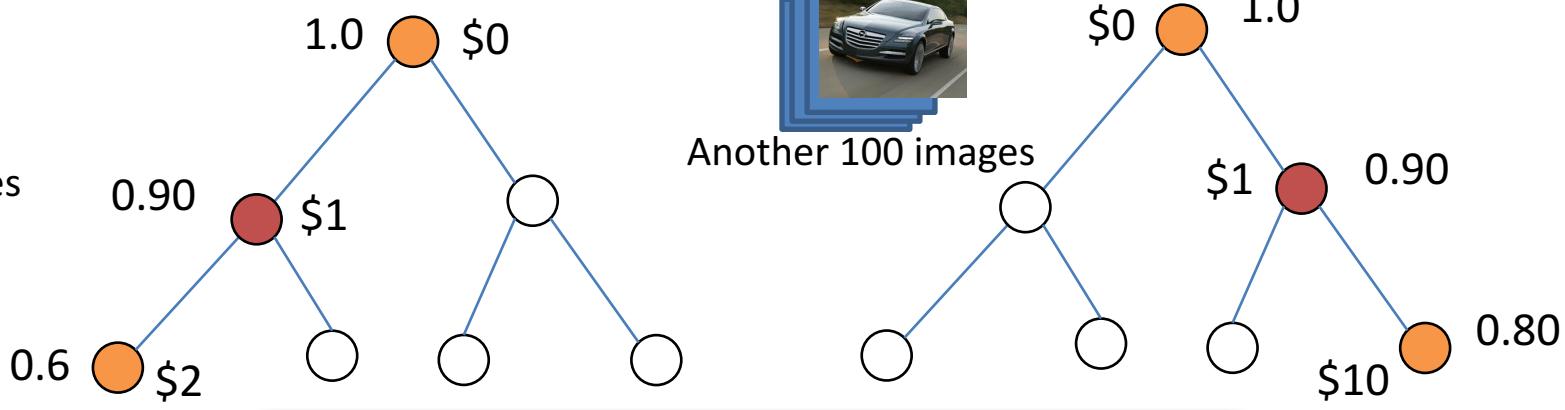


Deng, Krause, Berg, Fei-Fei, CVPR2012

**Pick a global confidence threshold  $T=0.9$**  [Vailaya et al. '99]



100 images



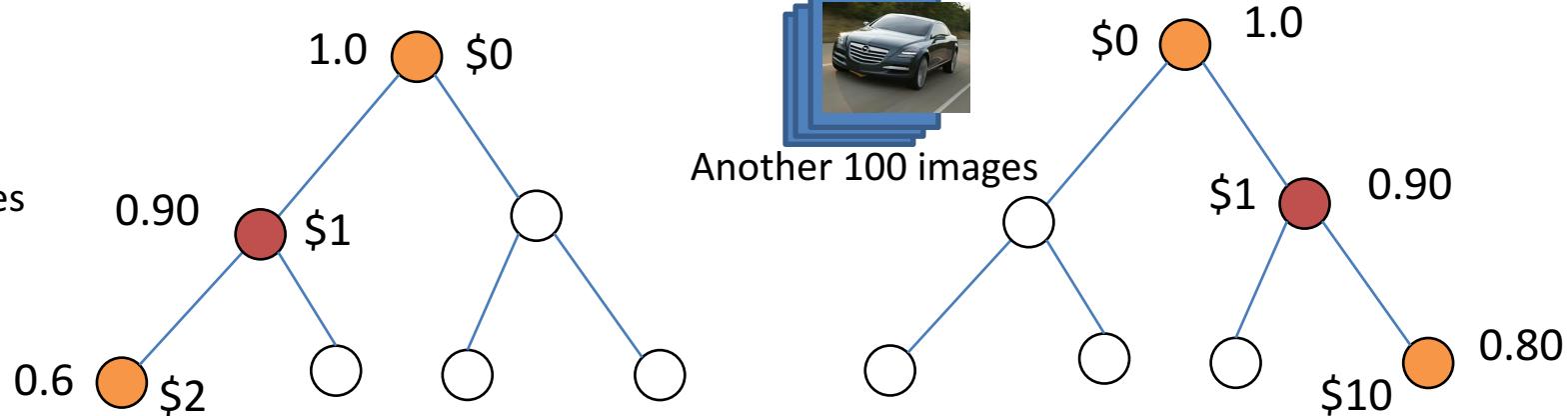
$$\text{Reward} = (\$1 * 0.90 + \$1 * 0.90) / 2 = \$0.90$$

$$\text{Accuracy} = (0.90 + 0.90) / 2 = 0.90$$

## Pick a global confidence threshold $T=0.9$ [Vailaya et al. '99]



100 images

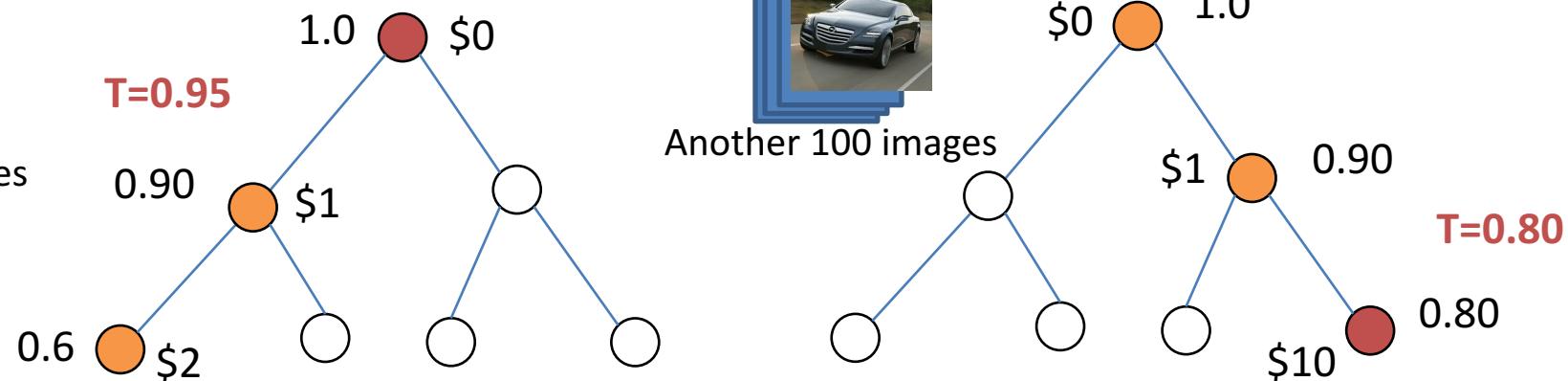


$$\text{Reward} = (\$1 * 0.90 + \$1 * 0.90) / 2 = \$0.90$$

$$\text{Accuracy} = (0.90 + 0.90) / 2 = 0.90$$



100 images

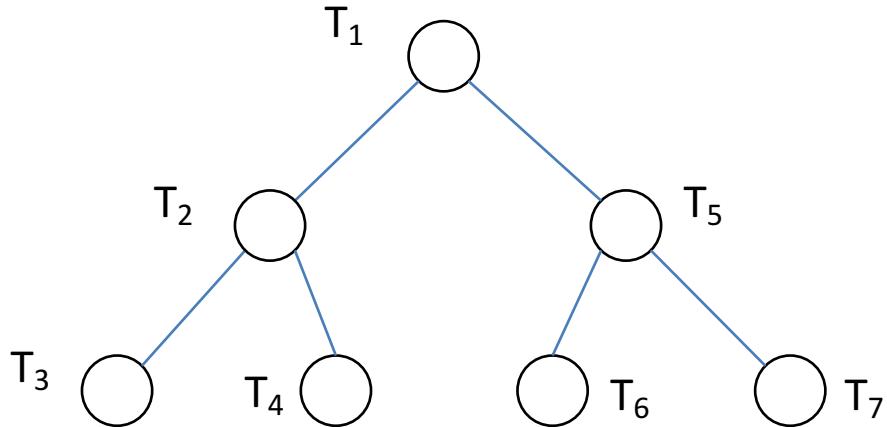


$T=0.95$

$T=0.80$

$$\text{Reward} = (\$0 * 1.0 + \$10 * 0.80) / 2 = \$4$$

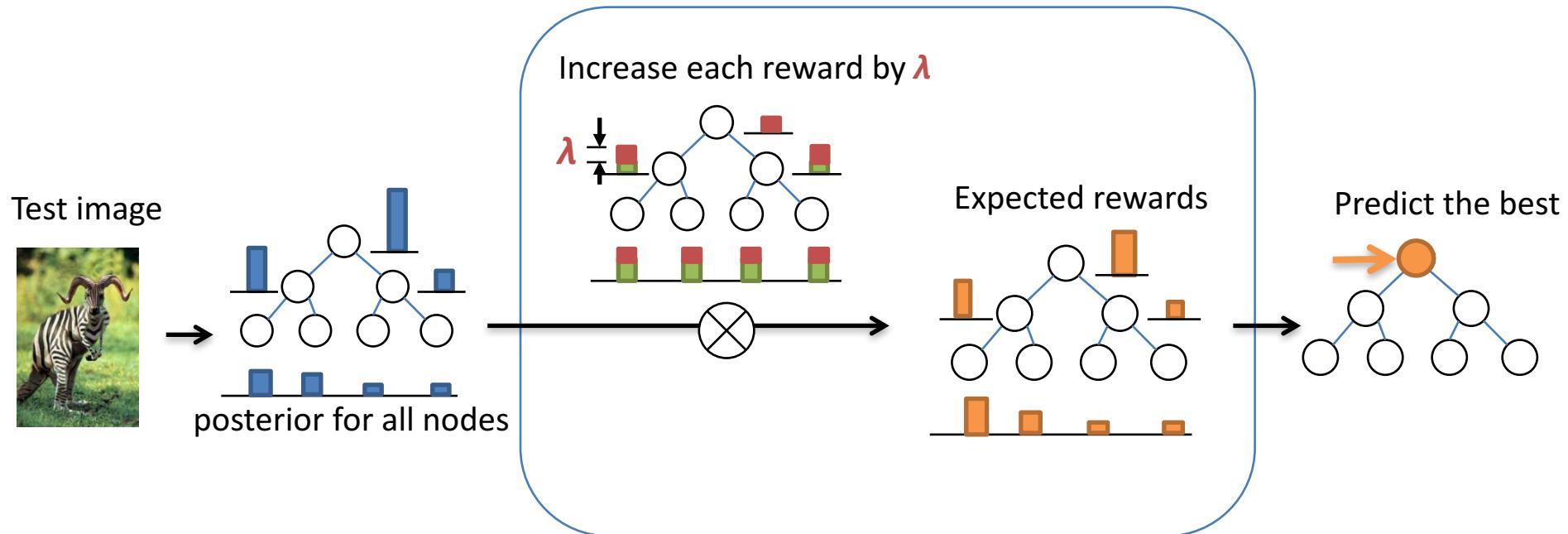
$$\text{Accuracy} = (1.0 + 0.80) / 2 = 0.90$$

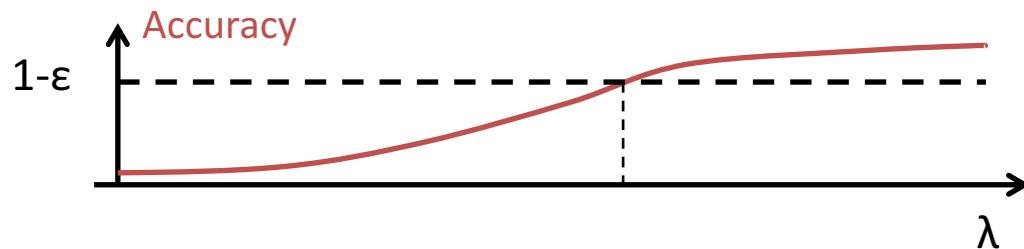
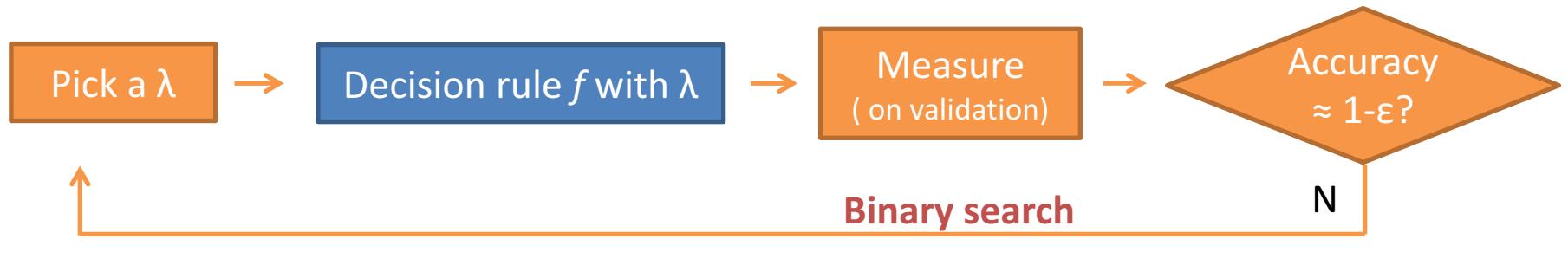


We can optimize individual thresholds...

But actually we don't need to.  
There is a simpler and provably optimal solution

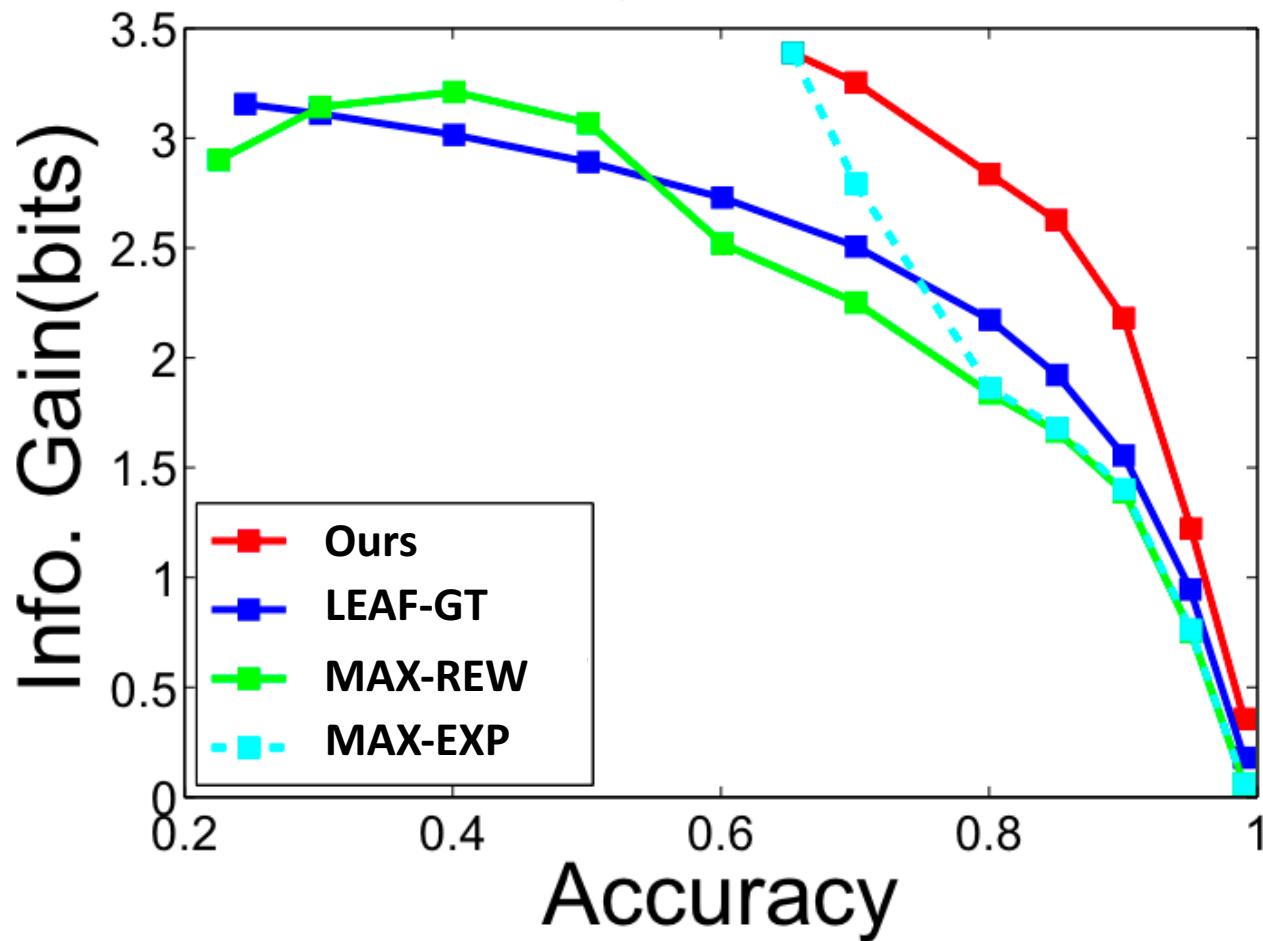
# A global, fixed scalar parameter $\lambda \geq 0$





**Theorem:** Under very mild conditions, this is optimal.

# ImageNet10K



Deng, Krause, Berg, Fei-Fei, CVPR2012

[www.image-net.org/eva](http://www.image-net.org/eva)



The **EVA system**, powered by **ImageNet**, can annotate images with guaranteed accuracies. It currently recognizes over **10,000** visual categories. See the **project** page to find out more.

Paste a URL | Upload an image

ANNOTATE



AT&T 7:14 PM 31%

Google Goggles  
Use pictures to search the web.

Browse similar images.

No close image matches found

- Avoid glare from the flash.
- Zoom in as much as possible by placing your device close to whatever you want to photograph.

←



0.95 coffee mug  
0.97 mug  
0.99 drinking vessel



Image size:  
401 × 604

No other sizes of this image found.

#### [Visually similar images](#) - Report images



0.87 face , gas pump, person  
0.90 face , gas pump



0.75 artifact, crater, matter, vertebrate

0.77 crater, matter, vertebrate

0.78 chordate, crater, matter

0.86 animal, matter

0.87 animal



0.78 person, instrument

0.84 person

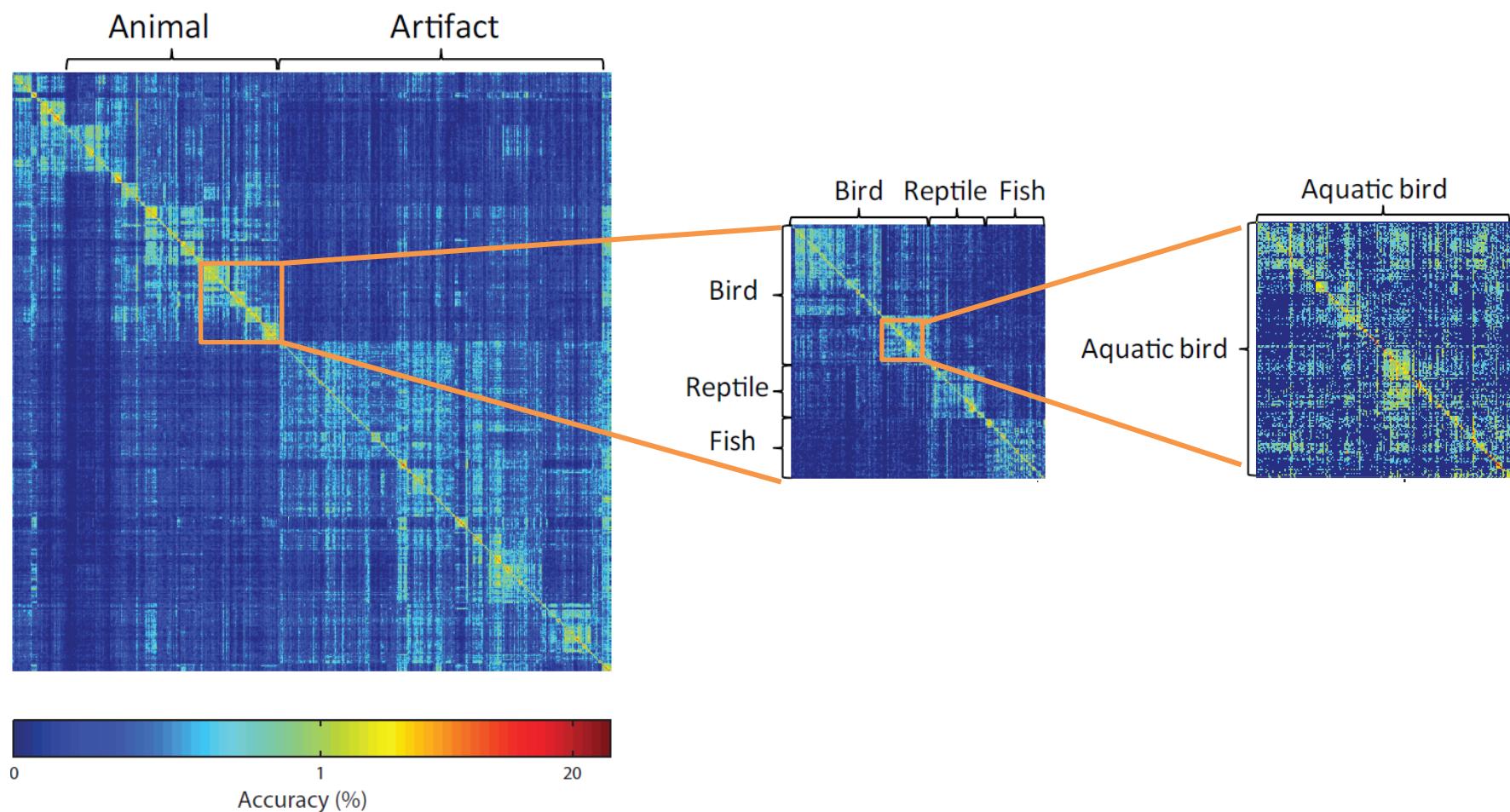
# Challenges

- Semantic hierarchy

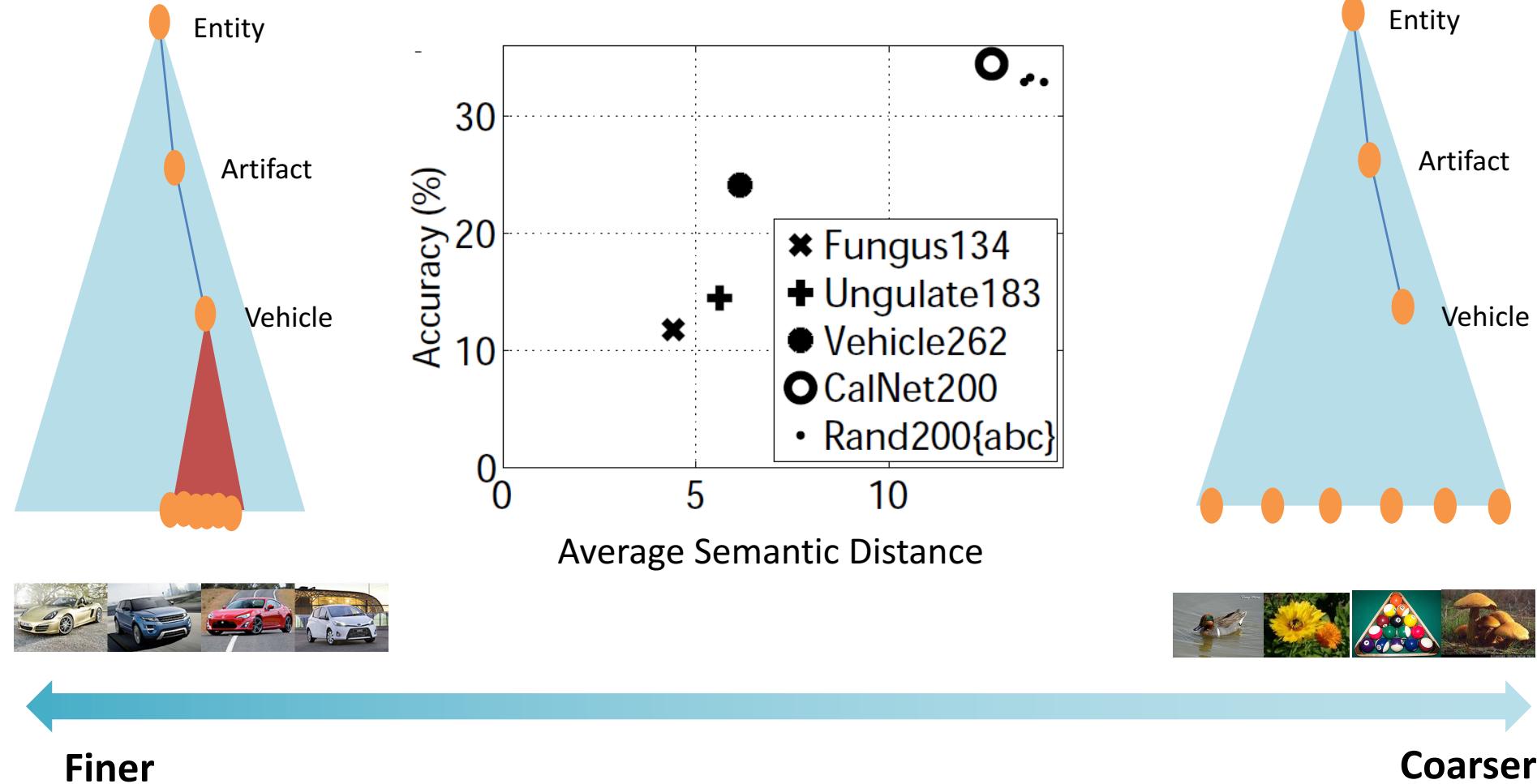
- Fine-grained classes

- Large-scale Learning

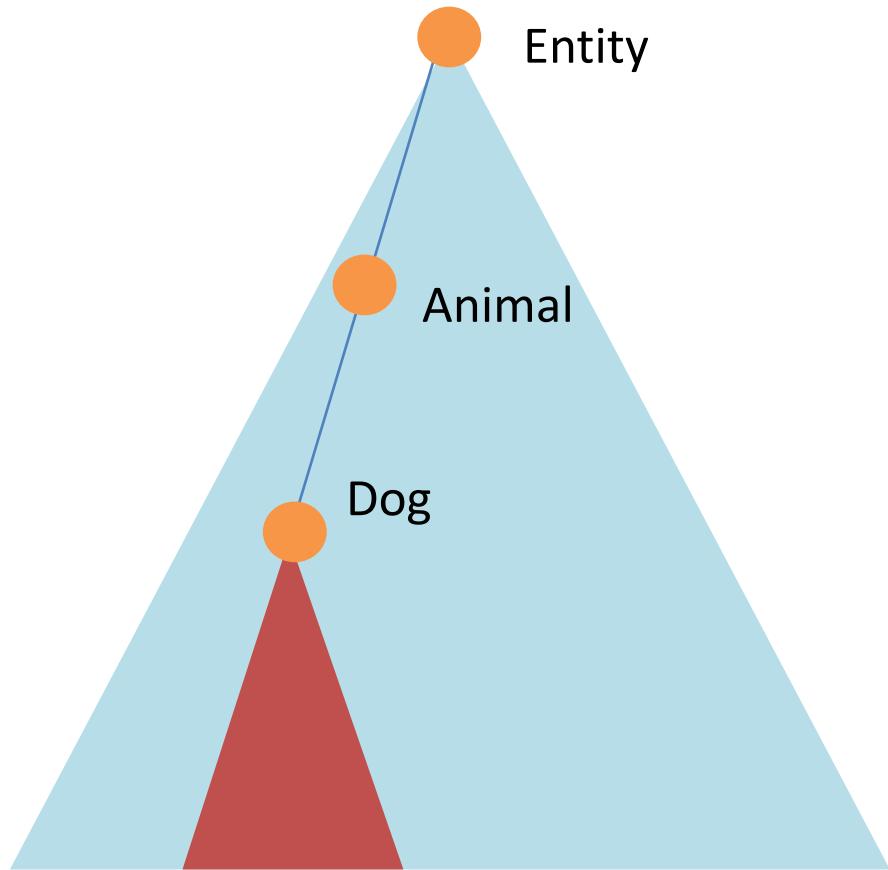
# Learn to Classify 10K Classes



# Fine-grained categories are a lot harder



# Fine-Grained Recognition



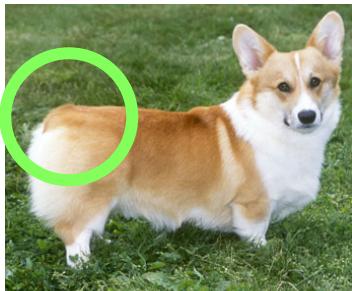
**What breed is this dog?**

- Basic category known.
- Object already localized.

# Fine-Grained Recognition



...



...



?

**What breed is this dog?**

# Fine-Grained Recognition



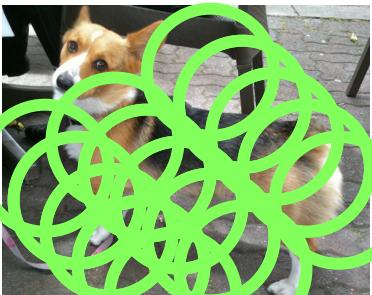
Cardigan Welsh Corgi



...

## Existing Work

Feature Selection  
from all possible  
locations



Pembroke Welsh Corgi



...

- [Branson et al. '10]
- [Bo et al. '10]
- [Farrell et al. '11]
- [Yao et al. '11]
- [Yao et al. '12]

# Fine-Grained Recognition



Cardigan Welsh Corgi



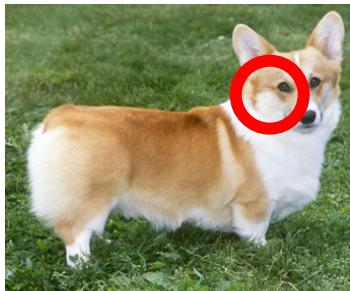
...

## Existing Work

Feature Selection  
from all possible  
locations



Pembroke Welsh Corgi

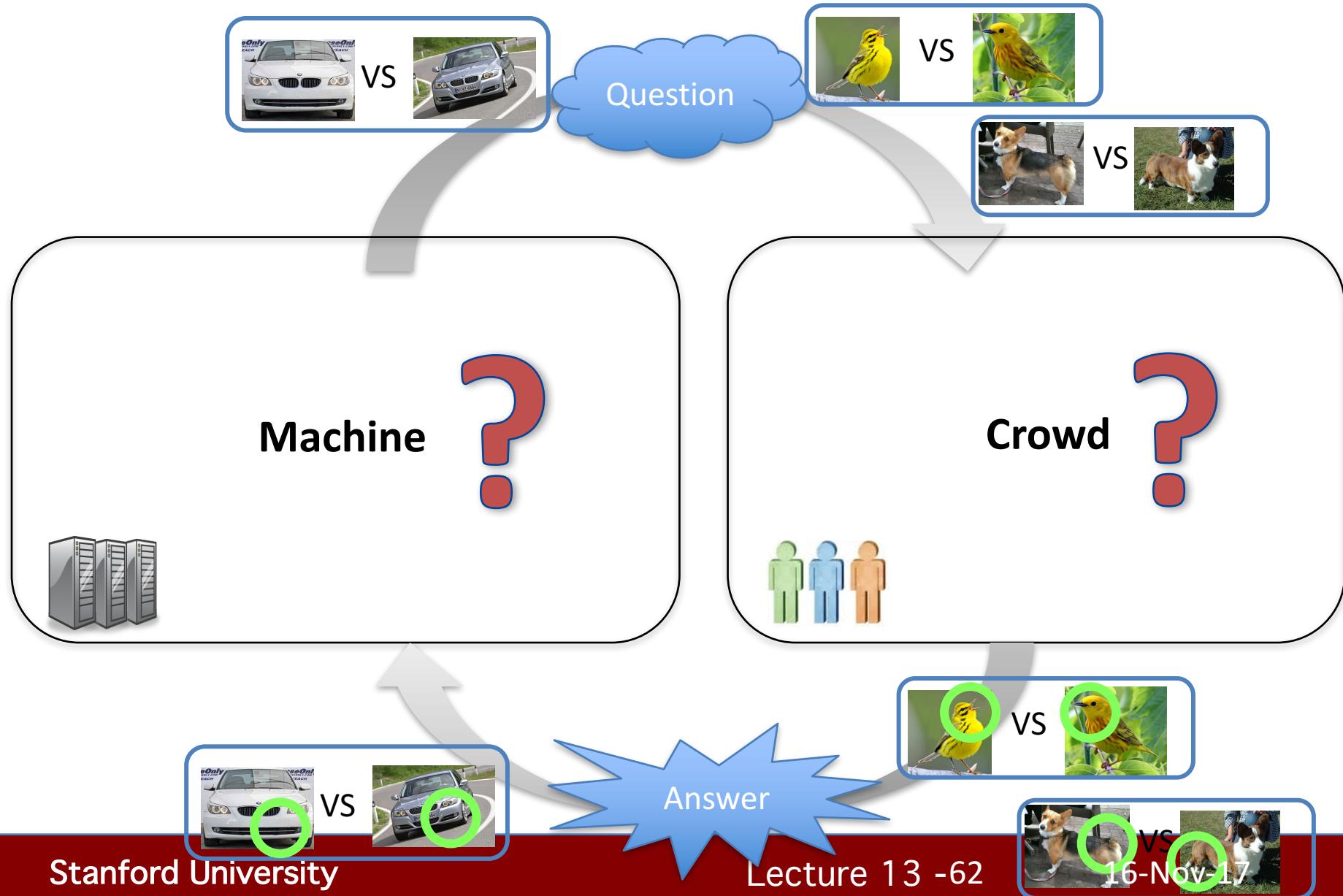


...

Can fail to find the right features.

- [Branson et al. '10]
- [Bo et al. '10]
- [Farrell et al. '11]
- [Yao et al. '11]
- [Yao et al. '12]

# Crowd-Machine Collaboration



# Crowd-Machine Collaboration



VS



VS



Crowd



VS

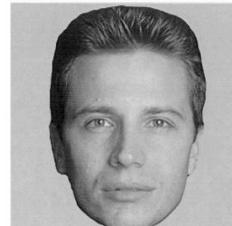


# Bubbles [Gosselin & Schyns, '01]

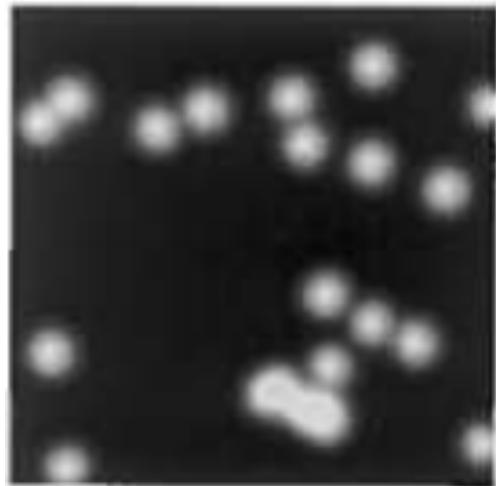


Smiling

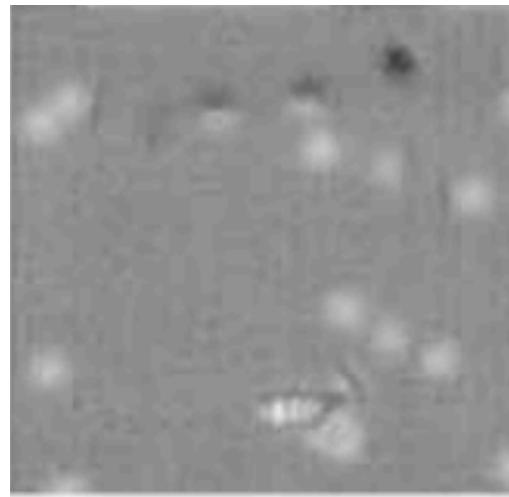
VS



Neutral



Random Bubble Mask



?

# Bubbles [Gosselin & Schyns, '01]



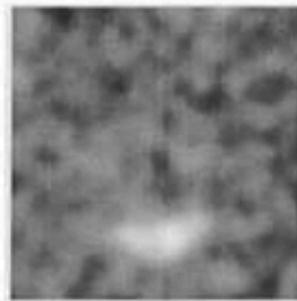
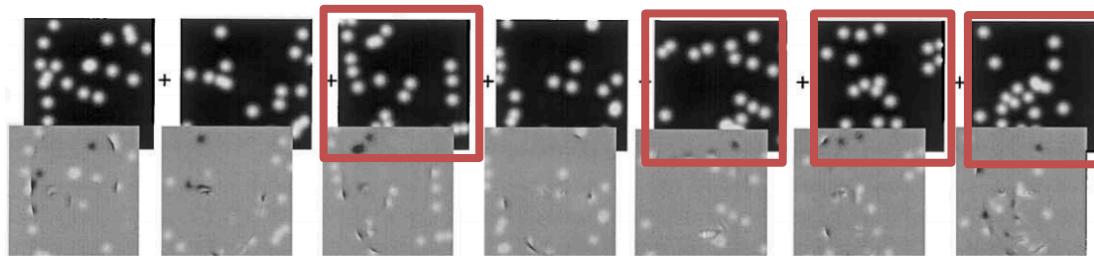
Smiling

VS



Neutral

Too costly



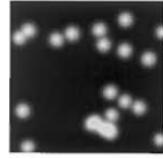
# Annotation Rationale

[Donahue & Grauman '11]



No easy quality assurance

What makes the form of the skater good?



Gosselin &  
Schyns '01



Donahue &  
Grauman '11

**Cost Effective**



**Quality Assurance**



**\* YOU AND A RANDOM PARTNER TAKE TURNS PEEKING AND BOOMING \***

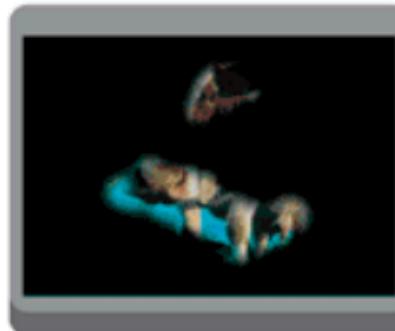
**BOOM** : REVEAL PARTS OF THE IMAGE TO YOUR PARTNER

TIME LEFT  
2:23



**PEEK** : GUESS WHAT YOUR PARTNER IS REVEALING

TIME LEFT  
2:23



GIVE HINTS  
IF NECESSARY

TELL YOUR PARTNER IF  
A GUESS IS **HOT** OR **COLD**

HINTS HELP  
YOU GUESS

PASS FOR  
DIFFICULT IMAGES

## Peekaboom [Ahn, Liu, Blum '06]

Does not work for fine-grained classes



Gosselin &  
Schyns '01



Donahue &  
Grauman '11



Ahn, Liu,  
Blum '06



### Cost Effective

✗

✓

✓

✓

### Quality Assurance

✓

✗

✓

✓

### Fine Grained

✓

✓

✗

✓

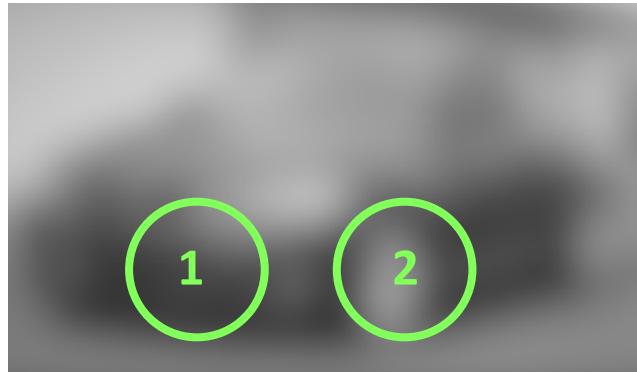
# The Bubbles Game

J. Deng, J. Krause, L. Fei-Fei. **Fine-Grained Crowdsourcing for Fine-Grained Recognition.** CVPR 2013.

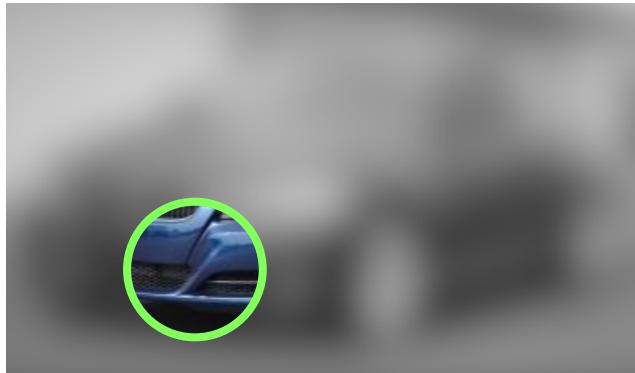
# Let's play



# Let's play



# Let's play



# Let's play





[Click Me or Press 1](#)

[Prairie Warbler \(wikipedia\)](#)



Draw and identify the category.



[Click Me or Press 2](#)

[Yellow Warbler \(wikipedia\)](#)



Bubble: [Smaller  
\(Press - or 's'\)](#)

[Bigger  
\(Press + or 'w'\)](#)

Bubble cost: 25

Points for correct identification: 100

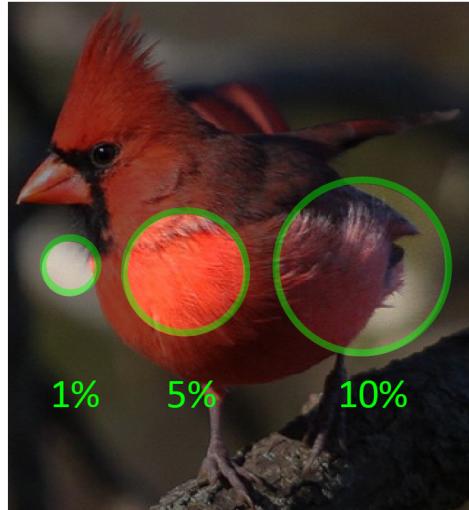
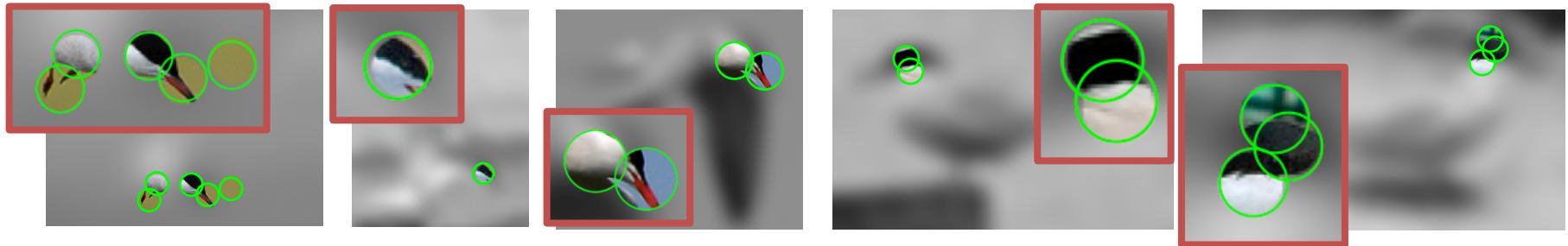
**Total score: 0**

**Total score needed to submit: 1000**

[Pass this image](#)

[Change the pair of categories](#)

# Crowd Picked Bubbles (AMT for now)



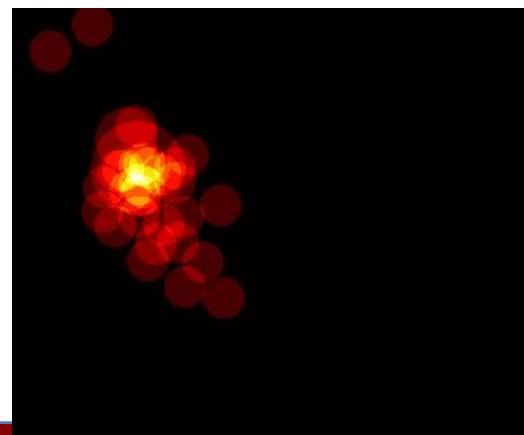
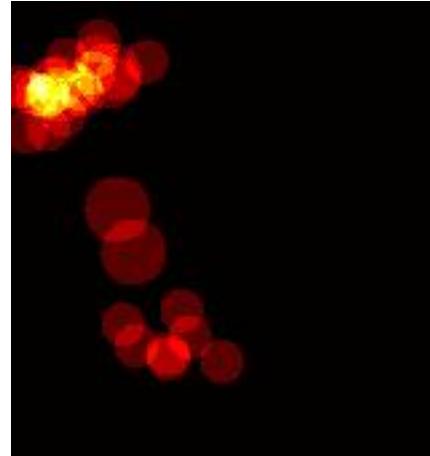
70% of games are successful

>90% of successful games use <10% of area

Bubble sizes as proportions of image

Deng, Krause, & Fei-Fei, CVPR2013

# Bubble Heatmaps



# Crowd-Machine Collaboration



VS



VS



## Bubbles Game



VS



# Crowd-Machine Collaboration



VS

Question



VS



VS

Machine



Bubbles Game



VS

Answer



VS



16-Nov-17

# BubbleBank Representation



A test image

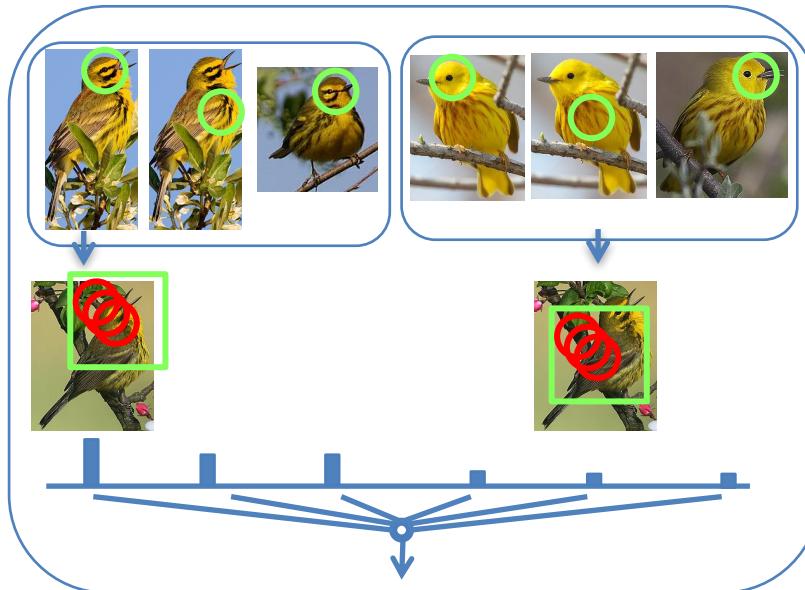
Crowd-picked  
Bubbles  
(on training)



Max-pooling



## V1-like models



## BubbleBank

Crowd picked

Pool over a single  
region (spatial prior)

[Deng et al. '13]

## Prior Work

Clustering random  
patches

Pool over multiple  
uniformly sampled  
regions (e.g. SPM)

[Lecun et al. '98]

[Csurka et al. '04]

[Lazebnik et al. '06]

[Wang et al. '09]

[Lee et al. '09]

[Pinto et al. '09]

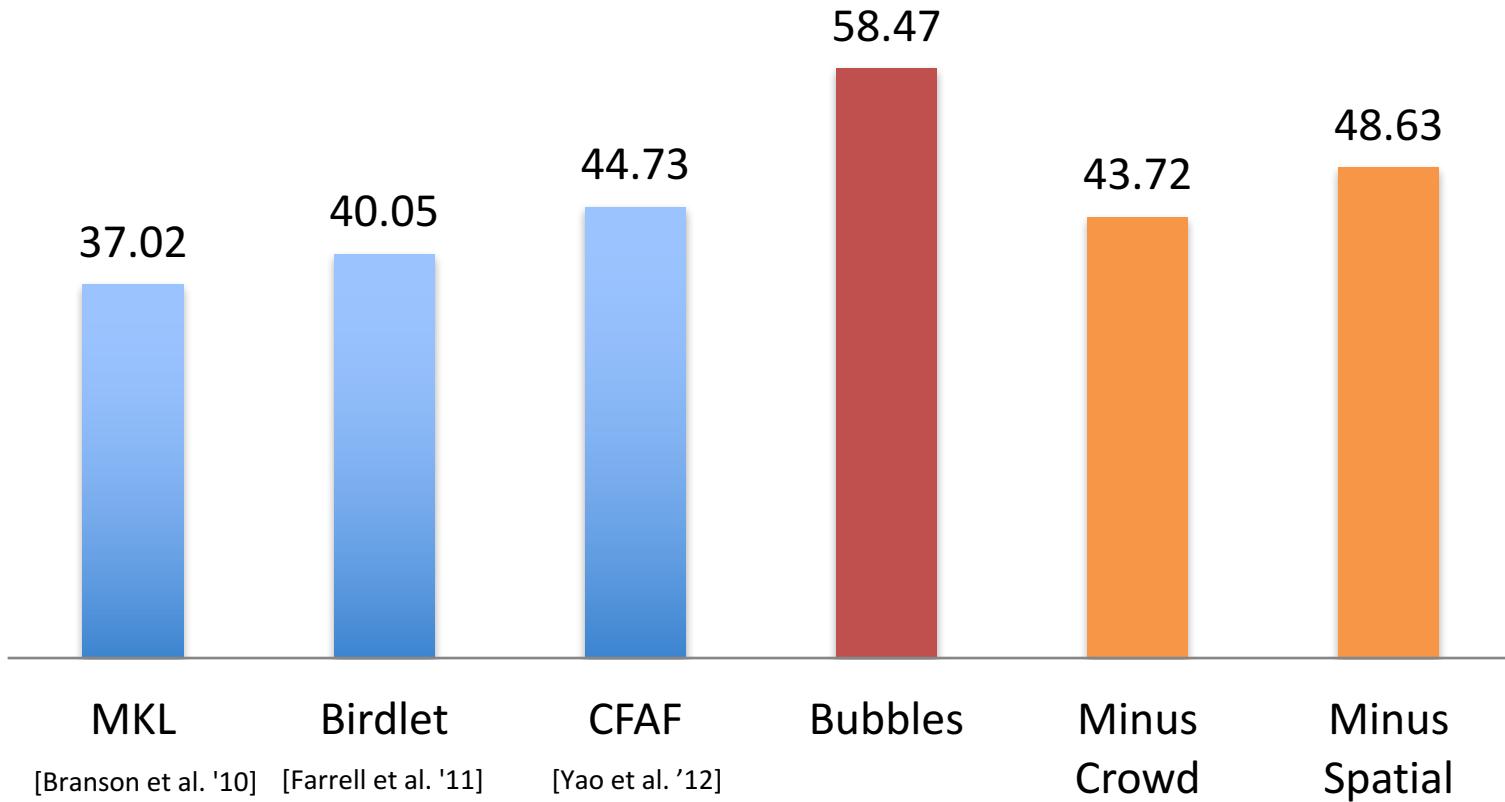
[Perronnin et al. '10]

[Li et al. '10]

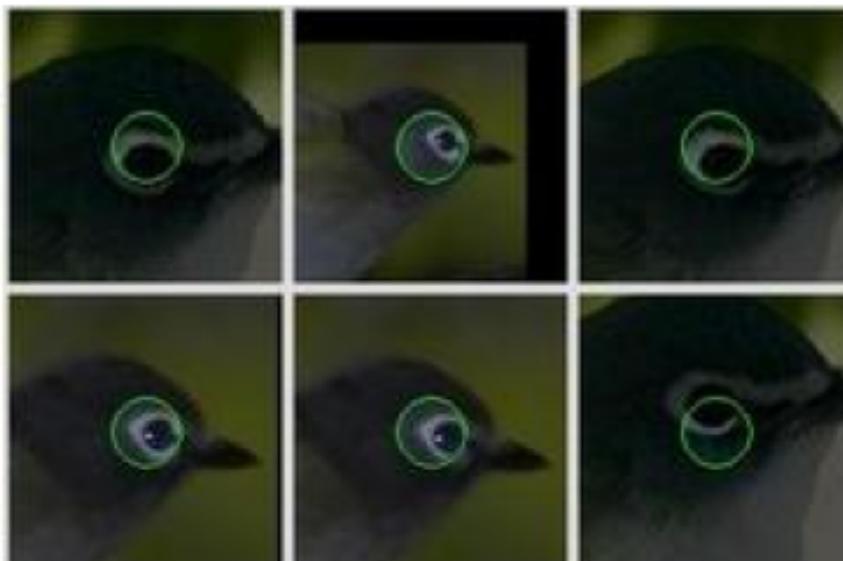
[Coates & Ng '11].

[Yao et al. '12]

## Mean Average Precision (Caltech-USCD-Bird)



# Top Activated Bubbles



# Crowd-Machine Collaboration



VS

Question



VS



VS

## BubbleBank



## Bubbles Game



VS

Answer



VS



16-Nov-17

# Challenges

- Semantic hierarchy
- Fine-grained classes
- Large-scale Learning