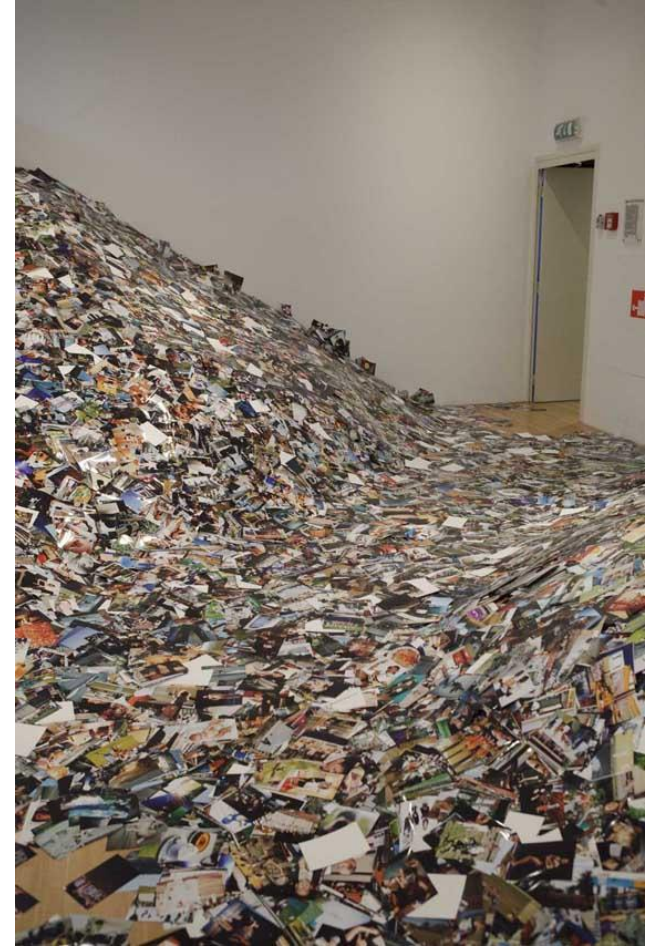# Human Computation and Computer Vision

CS143 Computer Vision
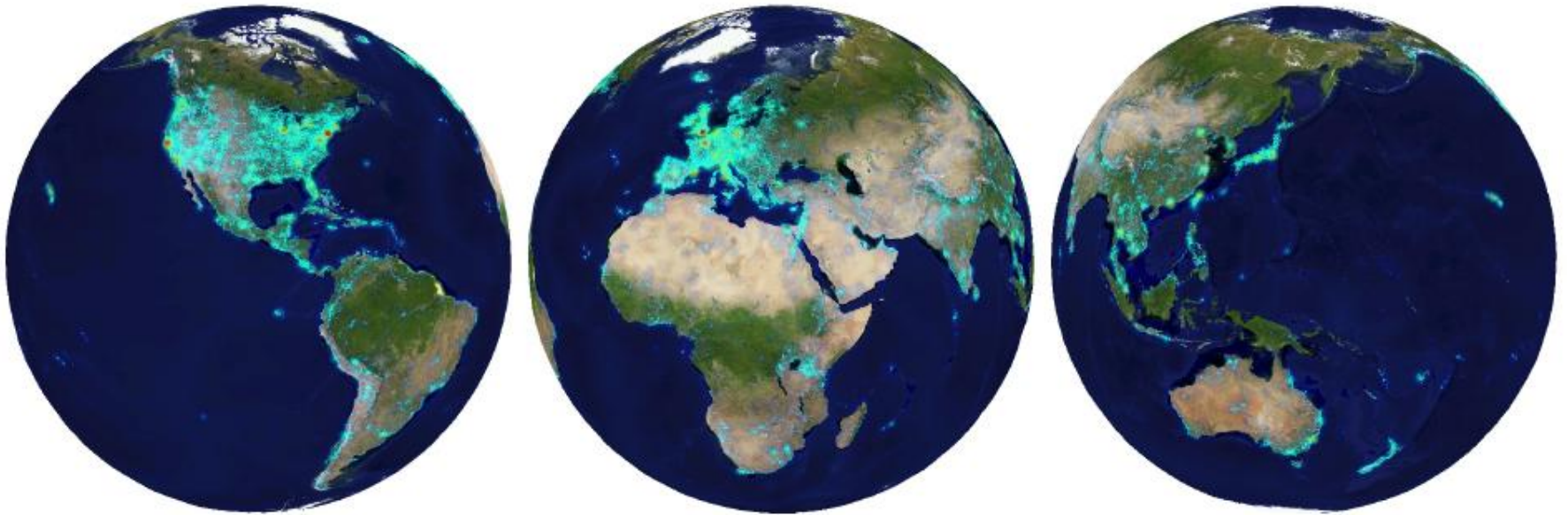
James Hays, Brown University

# 24 hours of Photo Sharing



installation by Erik Kessels

# And sometimes Internet photos have useful labels



Im2gps. Hays and Efros. CVPR 2008

# But what if we want more?

# Image Categorization

Training



**Training Images**

Training Labels → Classifier Training

Image Features → Classifier Training → Trained Classifier

# Image Categorization

## Training

Training Images

Training Labels → 

Image Features → Classifier Training → Trained Classifier

## Testing

Test Image → Image Features → Trained Classifier → Prediction **Outdoor**

# Human Computation for Annotation

Unlabeled Images

Show images, Collect and filter labels

Training Images

Training Labels

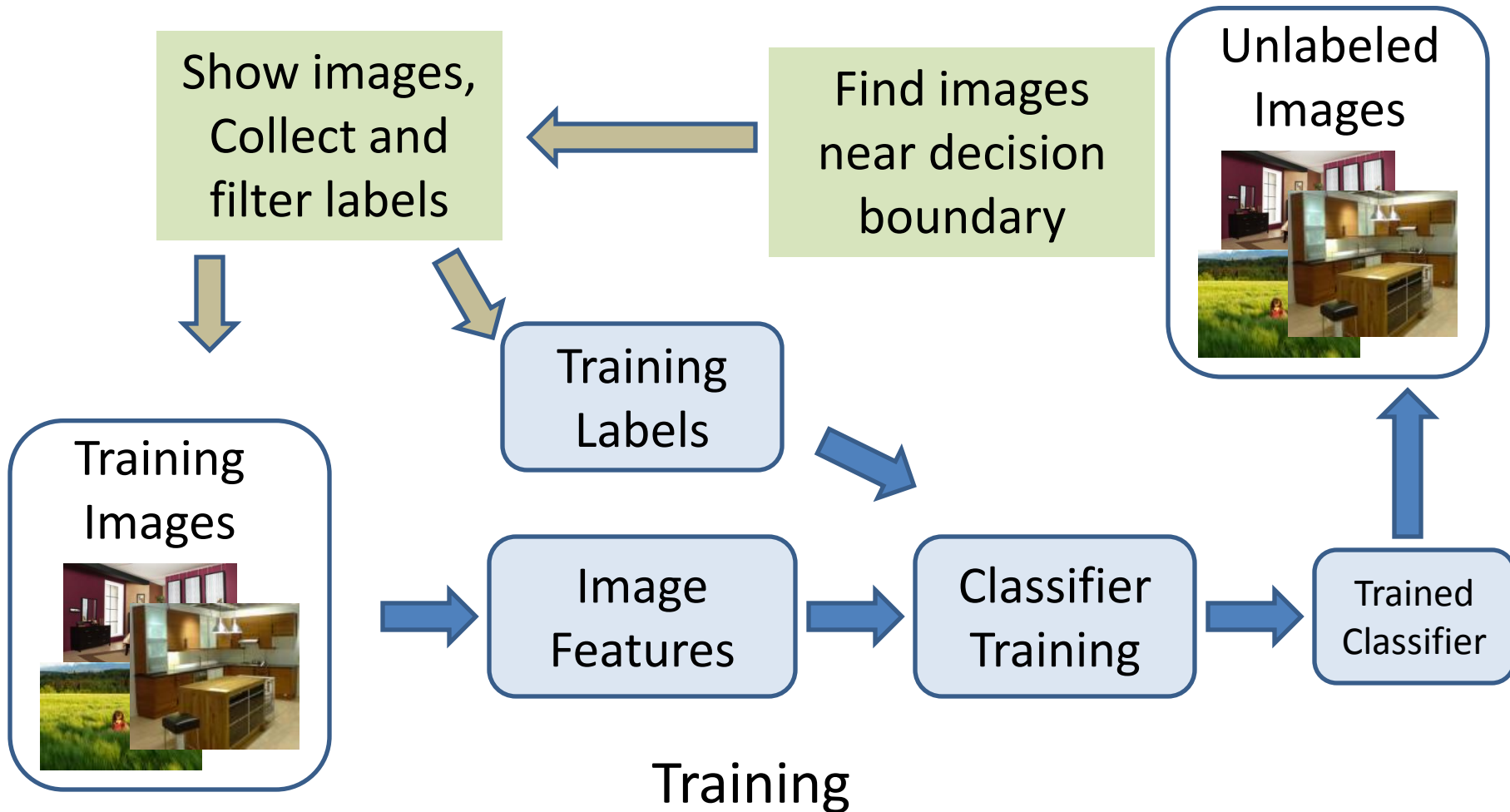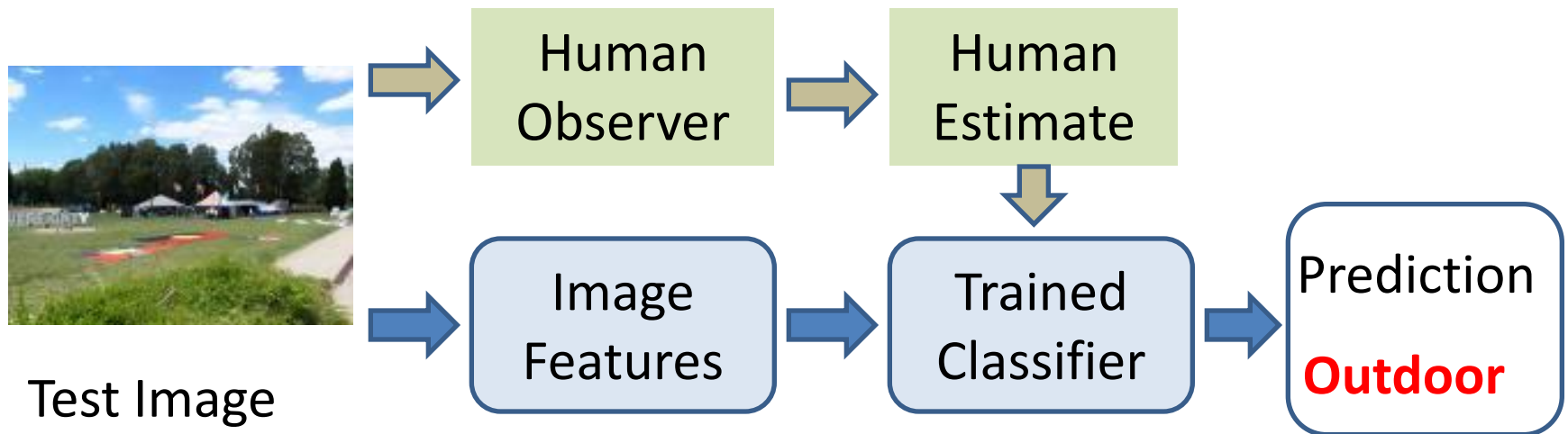Image Features

Classifier Training

Trained Classifier

Training

# Active Learning

# Human-in-the-loop Recognition



Test Image → Human Observer → Human Estimate

Test Image → Image Features → Trained Classifier → Prediction **Outdoor**

Testing

# Outline

- Human Computation for Annotation
  - ESP Game
  - Mechanical Turk
- Human-in-the-loop Recognition
  - Visipedia

Luis von Ahn and Laura Dabbish. Labeling Images with a Computer Game.
ACM Conf. on Human Factors in Computing Systems, CHI 2004

Building datasets

6000 images from flickr.com

Annotators

amazonmechanical turk
Artificial Artificial Intelligence
beta

Is there an Indigo bunting in the image?

100s of training images

Slide credit: Welinder et al

# Task: Find the Indigo Bunting



hit rate (correct detection) vs rate of correct rejection

# Task: Find the Indigo Bunting



6% error

15% error

31% error

50% error

hit rate (correct detection)

rate of correct rejection

Task: Find the Indigo Bunting

hit rate (correct detection)

rate of correct rejection

6% error
15% error
31% error
50% error

competent

Slide credit: Welinder et al

Task: Find the Indigo Bunting

bots

6% error

15% error

31% error

50% error

hit rate (correct detection)

rate of correct rejection

Slide credit: Welinder et al

Task: Find the Indigo Bunting

optimists

6% error

15% error

31% error

50% error

hit rate (correct detection)

rate of correct rejection

Slide credit: Welinder et al

Task: Find the Indigo Bunting

hit rate (correct detection) vs. rate of correct rejection

6% error
15% error
31% error
50% error

pessimists

Slide credit: Welinder et al

Task: Find the Indigo Bunting

hit rate (correct detection) vs rate of correct rejection

6% error
15% error
31% error
50% error

adversaries

Slide credit: Welinder et al

# Utility data annotation via Amazon Mechanical Turk

 X   100 000   =   $5000

Alexander Sorokin
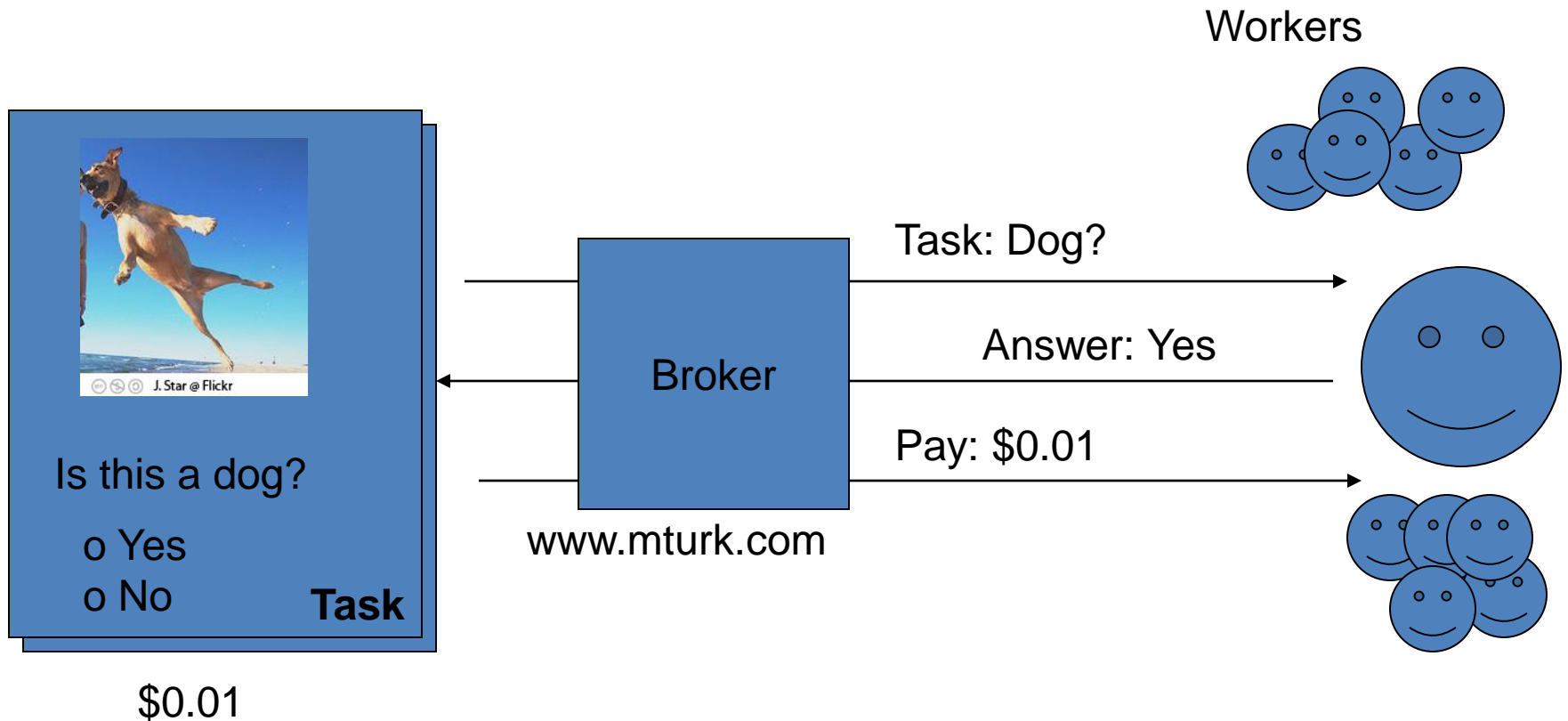
David Forsyth

University of Illinois at Urbana-Champaign

Slides by Alexander Sorokin

# Amazon Mechanical Turk

Workers

Is this a dog?

o Yes
o No

**Task**

$0.01

J. Star @ Flickr

Broker

www.mturk.com

Task: Dog?

Answer: Yes

Pay: $0.01

# Annotation protocols

- Type keywords
- Select relevant images
- Click on landmarks
- Outline something
- Detect features

……….. anything else ………

# Type keywords



$0.01

# Select examples



Joint work with Tamara and Alex Berg

http://visionpc.cs.uiuc.edu/~largescale/data/simpleevaluation/html/horse.html

# Select examples

# Click on landmarks

http://vision-app1.cs.uiuc.edu/mt/results/people14-batch11/p7/

# Outline something



$0.01

Data from Ramanan NIPS06

# Motivation



X   100 000   =   $5000

Custom
annotations

Large scale

Low price

# Issues

- Quality?
  - How good is it?
  - How to be sure?
- Price?
  - How to price it?

# Annotation quality



Agree within 5-10 pixels on 500x500 screen

There are bad ones.

A          C          E          G

# How do we get quality annotations?

# Ensuring Annotation Quality

- Consensus / Multiple Annotation / "Wisdom of the Crowds"

- Gold Standard / Sentinel
  - Special case: qualification exam

- Grading Tasks
  - A second tier of workers who grade others

# Pricing

- Trade off between throughput and cost
- Higher pay can actually attract scammers

# 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](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

- Here, the authors view progress as 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.

# Some definitions:

$Q = \{q_1 ... q_n\}$ • Set of possible questions

$a_i \in A_i$ • Possible answers to question *i*

$r_i \in V$ • Possible confidence in answer *i* (Guessing, Probably, Definitely)

$u_i = (a_i, r_i)$ • User response

$U^t$ • History of user responses at time *t*

# Question selection

- Seek the question that gives the maximum information gain (entropy reduction) given the image and the set of previous user responses.

$$I\left(c; u_i \mid x, U^{t-1}\right) = \sum_{u_i \in A_i \times V} p\left(u_i \mid x, U^{t-1}\right) \quad H\left(c \mid x, u_i \cup U^{t-1}\right) - H\left(c \mid x, U^{t-1}\right)$$

Probability of obtaining Response $u_i$ given the image And response history

Entropy when response is Added to history

Entropy before response is added.

where $\quad H\left(c \mid x, U^{t-1}\right) = -\sum_{c=1}^{C} p\left(c \mid x, U^{t-1}\right) \log p\left(c \mid x, U^{t-1}\right)$

# Incorporating vision

- Bayes Rule
- A visual recognition algorithm outputs a probability distribution across all classes that is used as the prior.
- A posterior probability is then computed based on the probability of obtaining a particular response history given each class.

$$p\left(c \mid x, U\right) = \eta\, p\left(U \mid c, x\right)\, p\left(c \mid x\right) = \eta\, p\left(U \mid c\right)\, p\left(c \mid x\right)$$

# Modeling user responses

- Assume that the questions are answered independently.

$$p\left(U^{t-1} \mid c\right) = \prod_{i}^{t-1} p\left(u_i \mid c\right)$$    Required for posterior computation

$$p\left(u_i \mid x, U^{t-1}\right) = \sum_{c=1}^{C} p\left(u_i \mid c\right) p\left(c \mid x, U^{t-1}\right)$$    Required for information gain computation

# The Dataset: Birds-200

- 6033 images of 200 species

# Implementation

**amazon**mechanical turk

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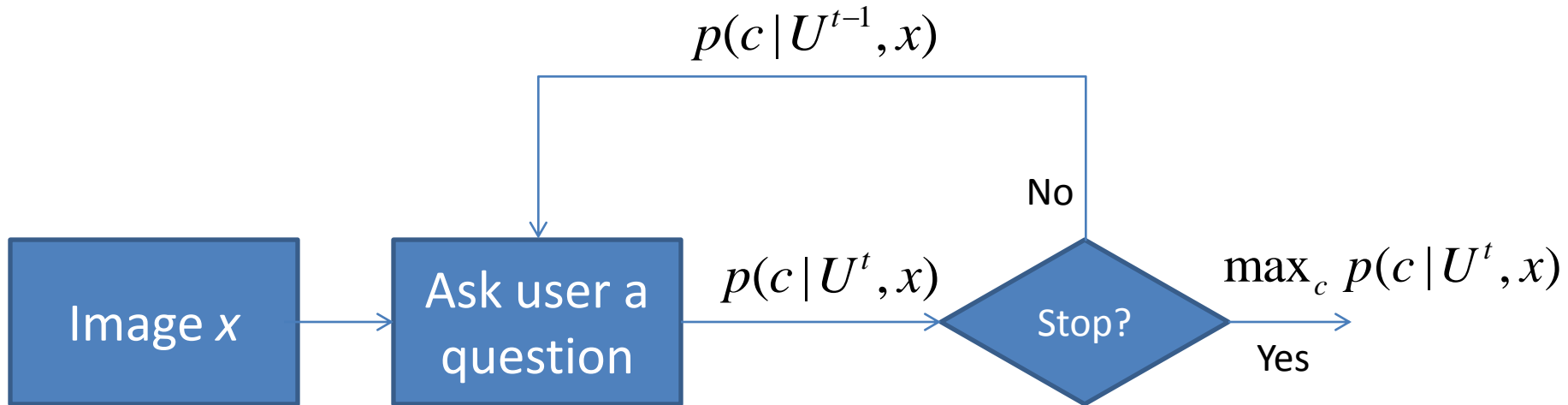- Mechanical Turk users asked to answer questions and provide confidence scores.

# 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.
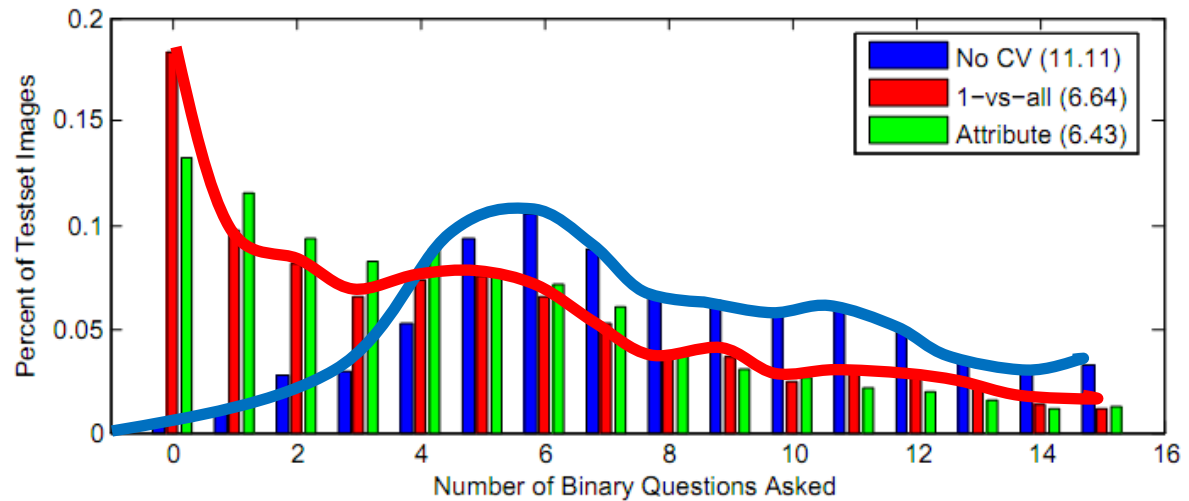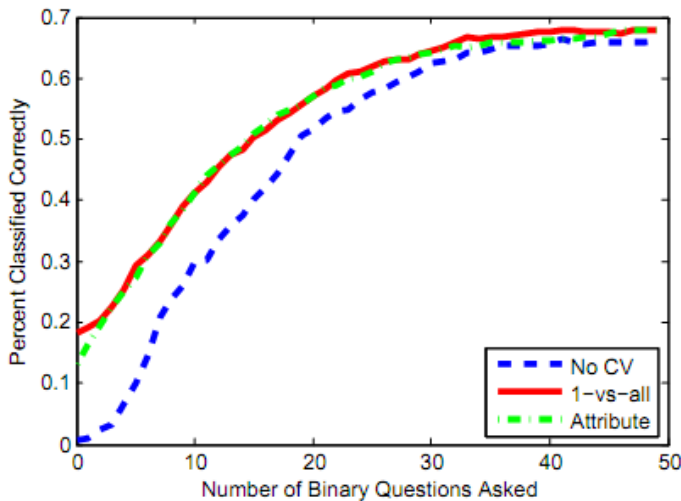
# Visual recognition

- Any vision system that can output a probability distribution across classes will work.

- Authors used Andrea Vedaldis's code.
  - Color/gray SIFT
  - VQ geometric blur
  - 1 v All SVM

- Authors added full image color histograms and VQ color histograms

# Experiments



$$p(c \mid U^{t-1}, x)$$

Image $x$ → Ask user a question → $p(c \mid U^t, x)$ → Stop?

No

Yes → $\max_c \ p(c \mid U^t, x)$

- 2 Stop criteria:
  - Fixed number of questions – evaluate accuacy
  - 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?