

The Multidimensional Wisdom of Crowds

Welinder P., Branson S., Belongie S., Perona, P

Experiment Presentation
[CS395T] Visual Recognition Fall 2012

Presented by: Niveda Krishnamoorthy

Problem Overview

Indigo Bunting



Search

Photos Groups People

Everyone's Uploads

indigo bunting

SEARCH

Full Text | Tags Only
Advanced Search

Sort: Relevant Recent Interesting

View: Small Medium Detail Slideshow



From I Bird 2



From The Nature...



From Momba...



From davidcreebir...



From naturelover2...



From K_Alanka



From davidcreebir...



From William ...



From violetfm



From jbobbe



From RitaK.



From mayalu



From ff151



From prairiedog



From redow



From ajnaturephot...



From Ken...



From R Hanson



From [Christine]



From davidcreebir...



From reemac640

5,926 results

6000 images
from flickr.com

Building datasets

100s of
training images



Annotators



amazon **mechanical turk**
beta Artificial Intelligence

Is there an Indigo bunting in the image?

Find the Indigo Bunting



A



A



A



A



A



A

Find the Indigo Bunting



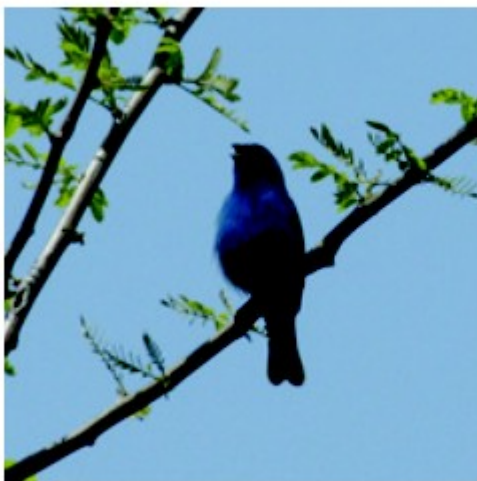
A B C D



A B C D



A B C D



A B C D



A B C D

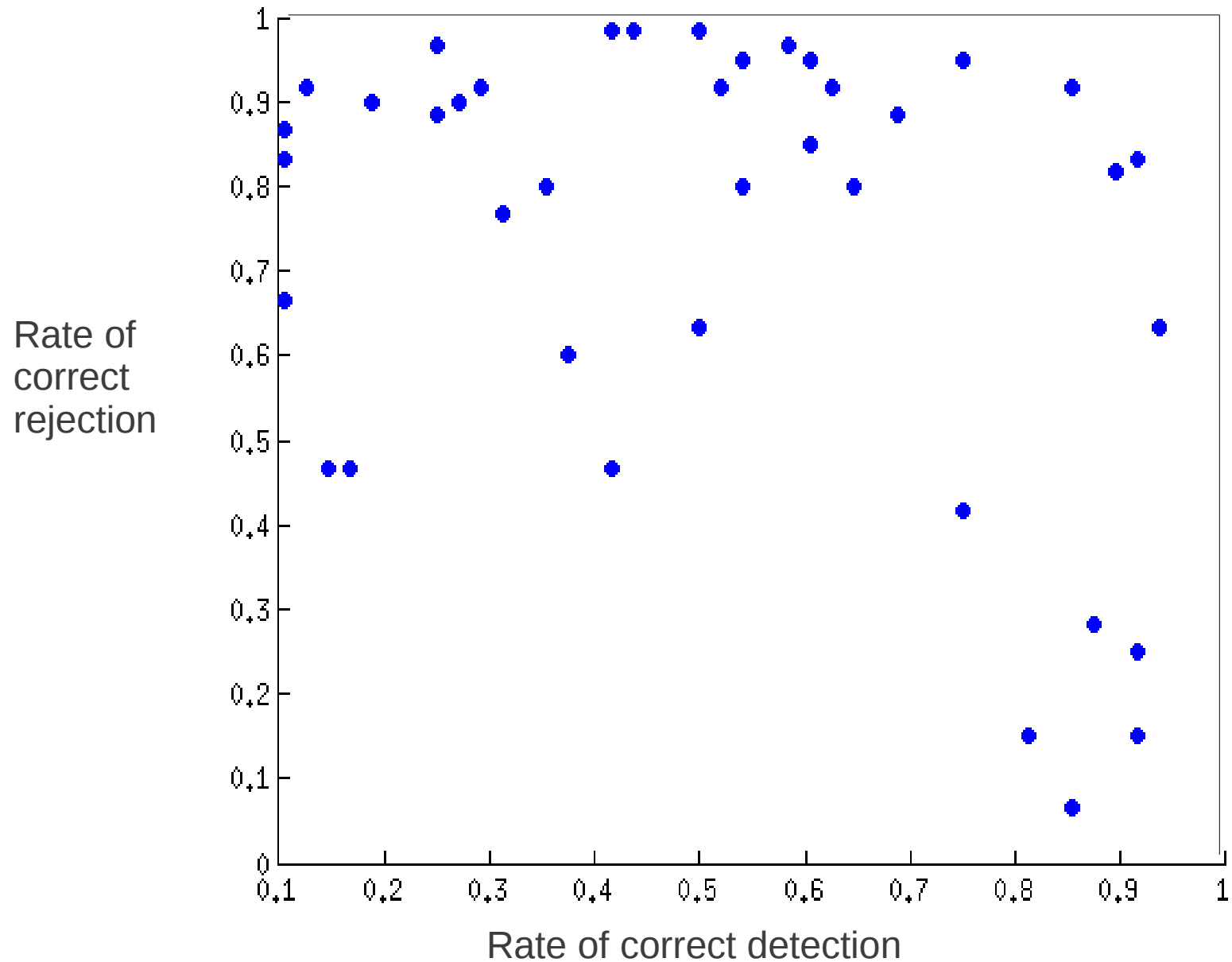


A B C D

Motivation

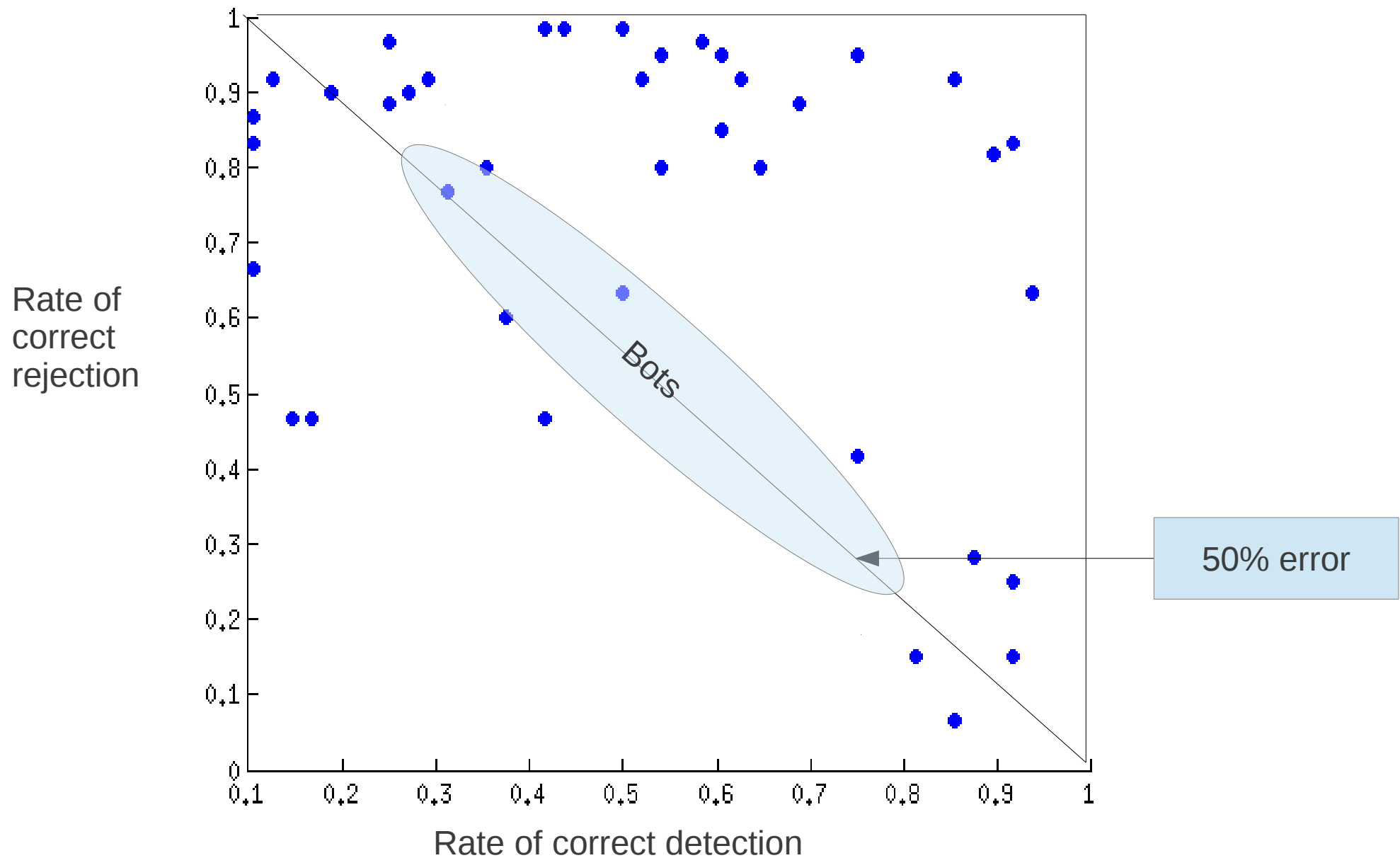
Distribution of Human Expertise – Task: Finding bluebirds

Experiment #1



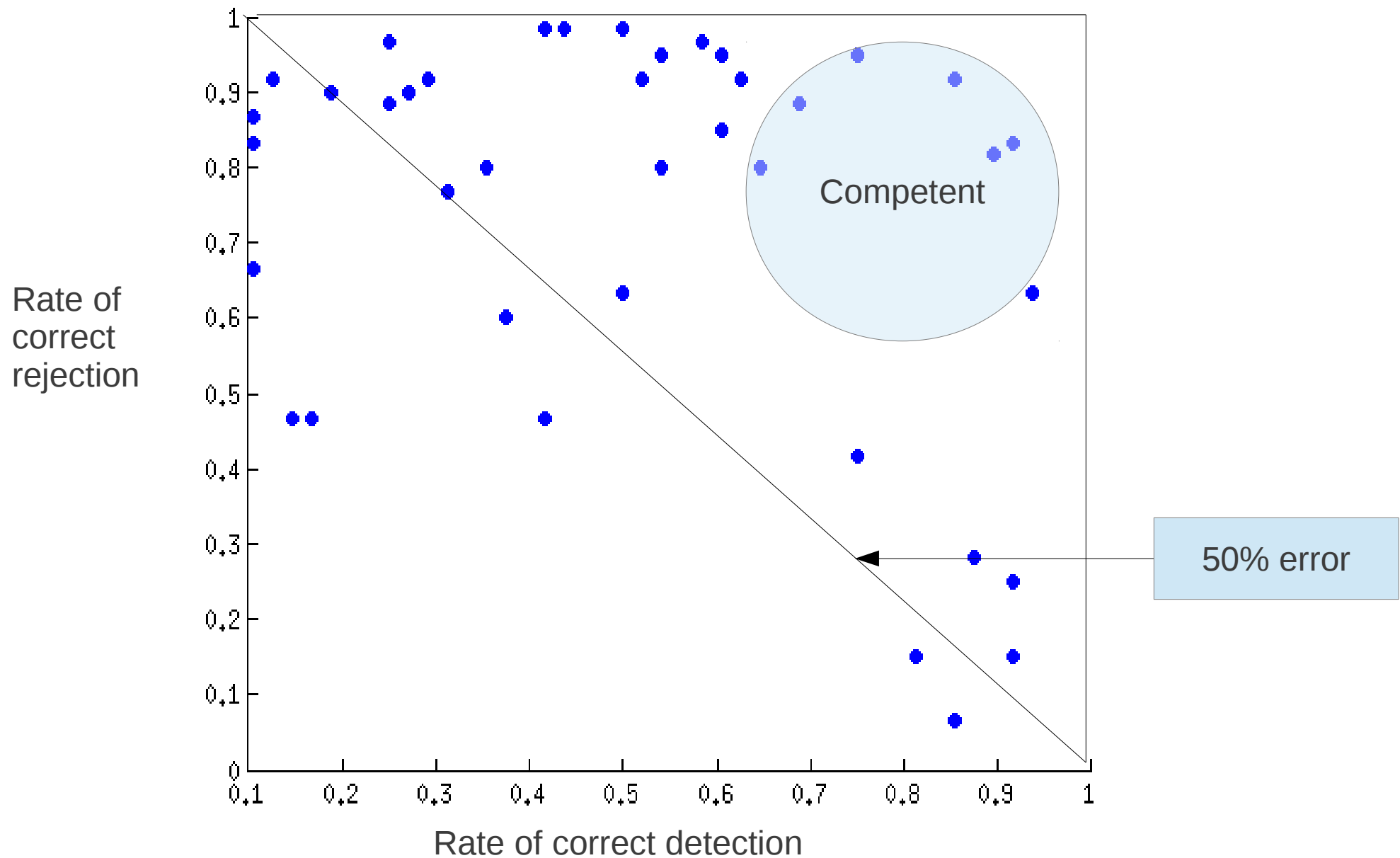
Distribution of Human Expertise – Task: Finding bluebirds

Experiment #1



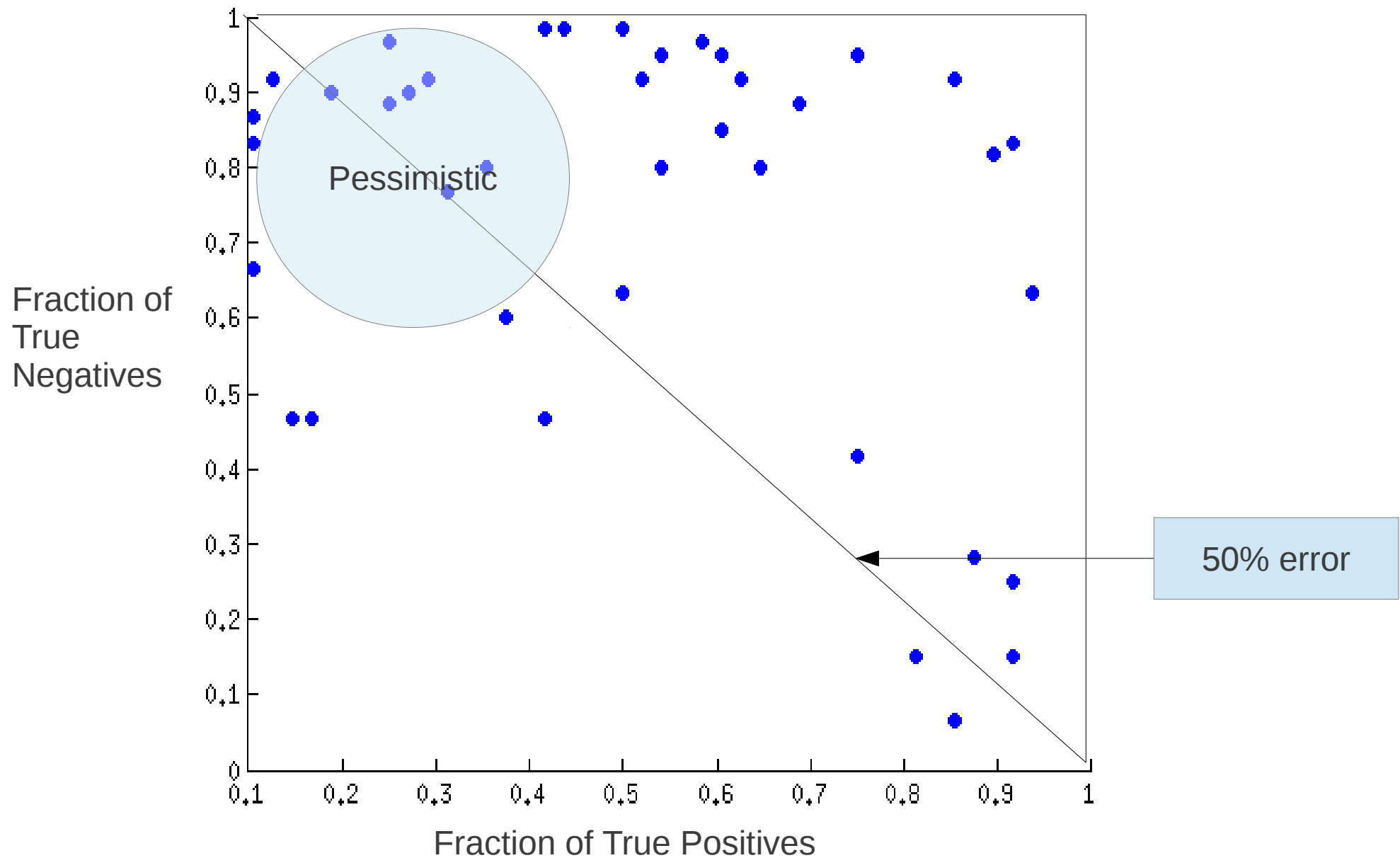
Distribution of Human Expertise – Task: Finding bluebirds

Experiment #1



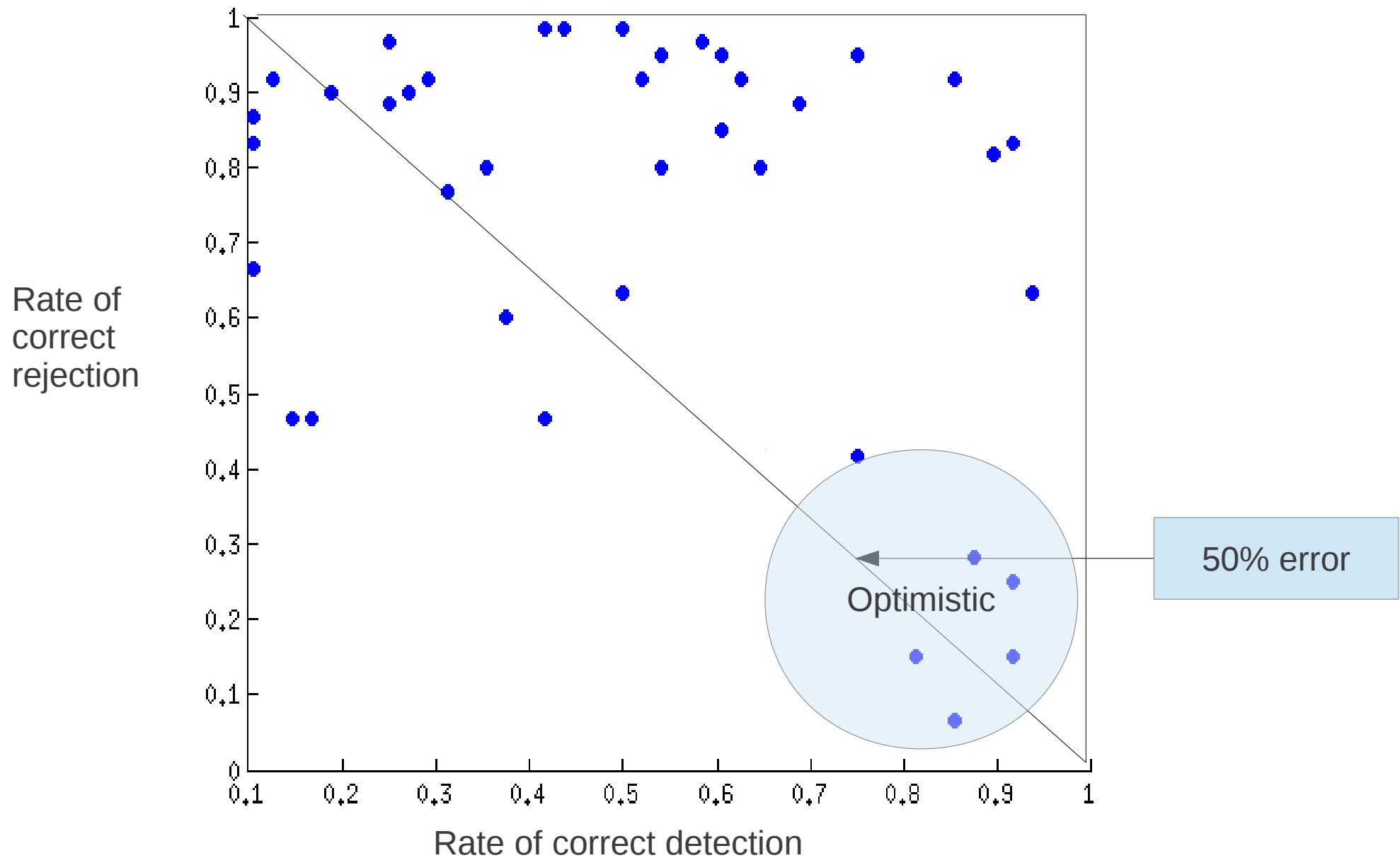
Distribution of Human Expertise – Task: Finding bluebirds

Experiment #1



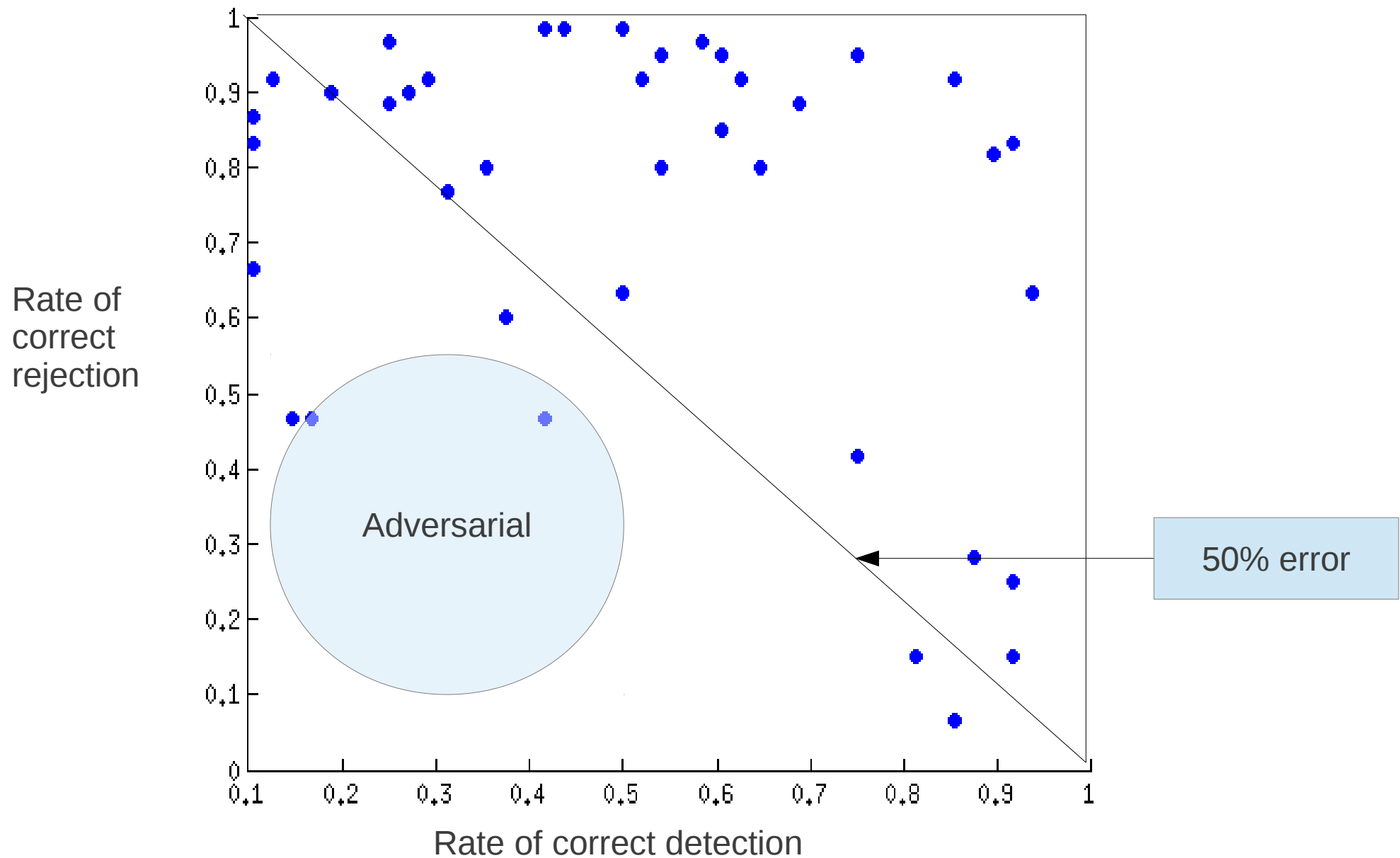
Distribution of Human Expertise – Task: Finding bluebirds

Experiment #1



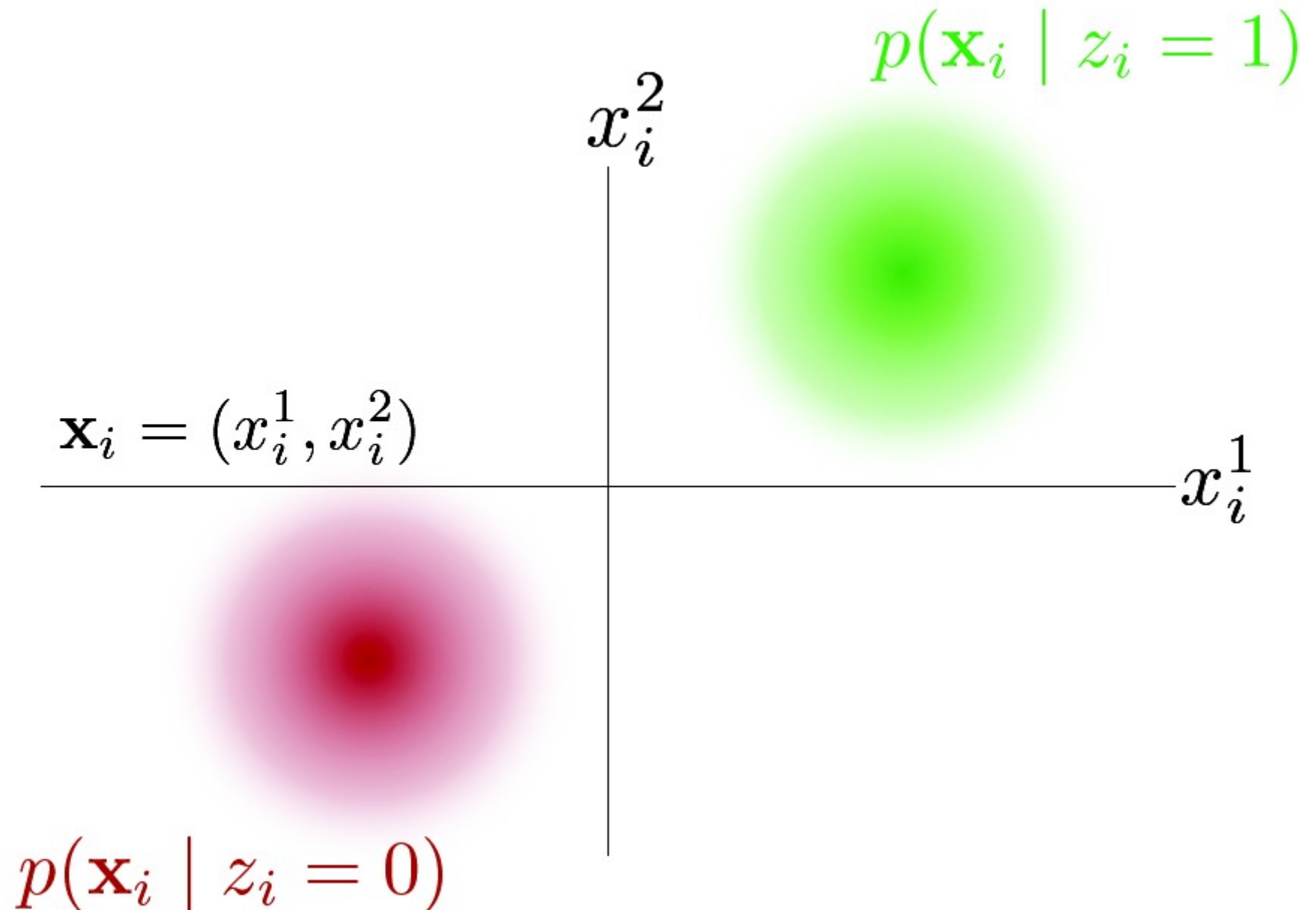
Distribution of Human Expertise – Task: Finding bluebirds

Experiment #1

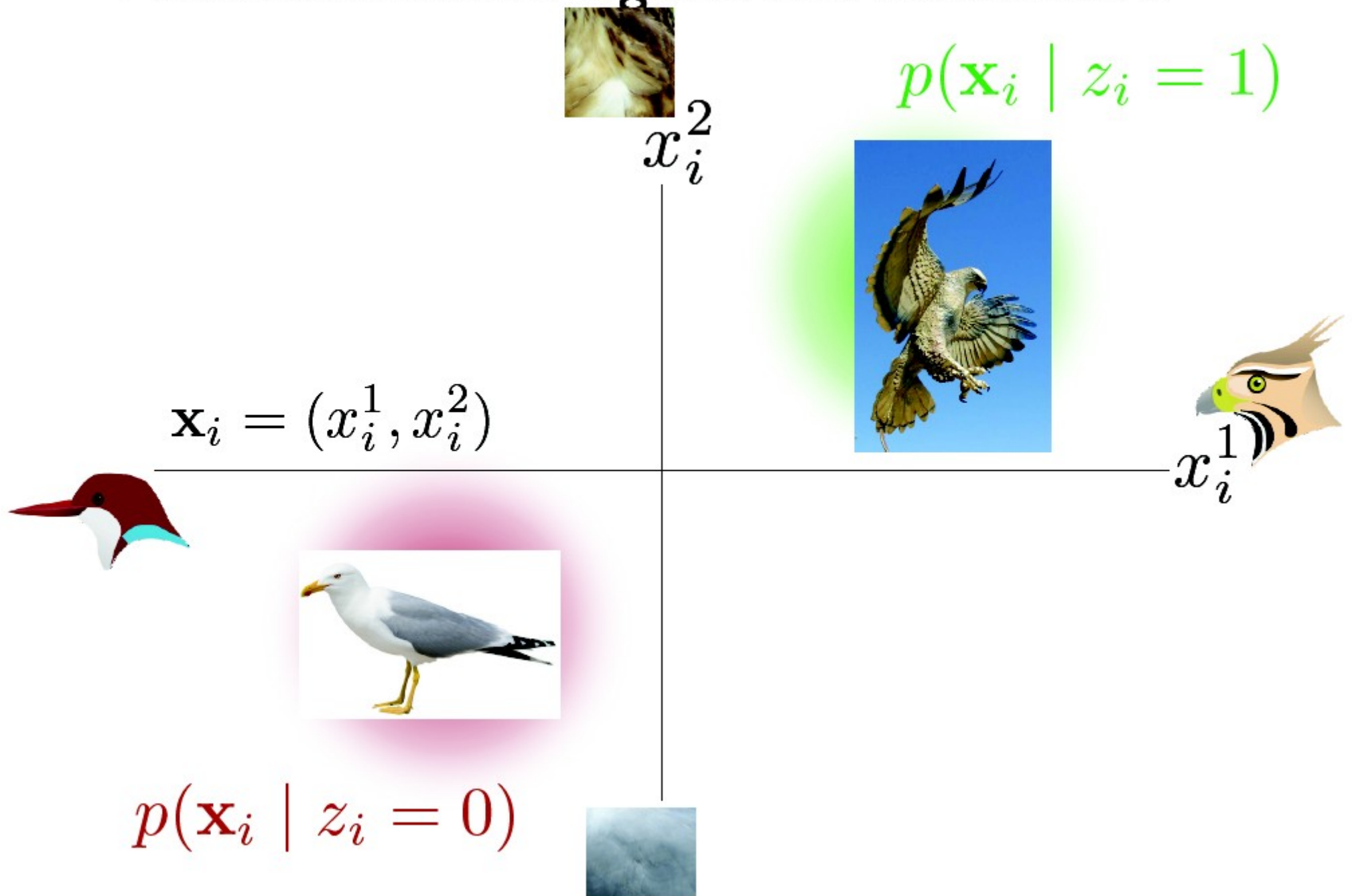


The Idea

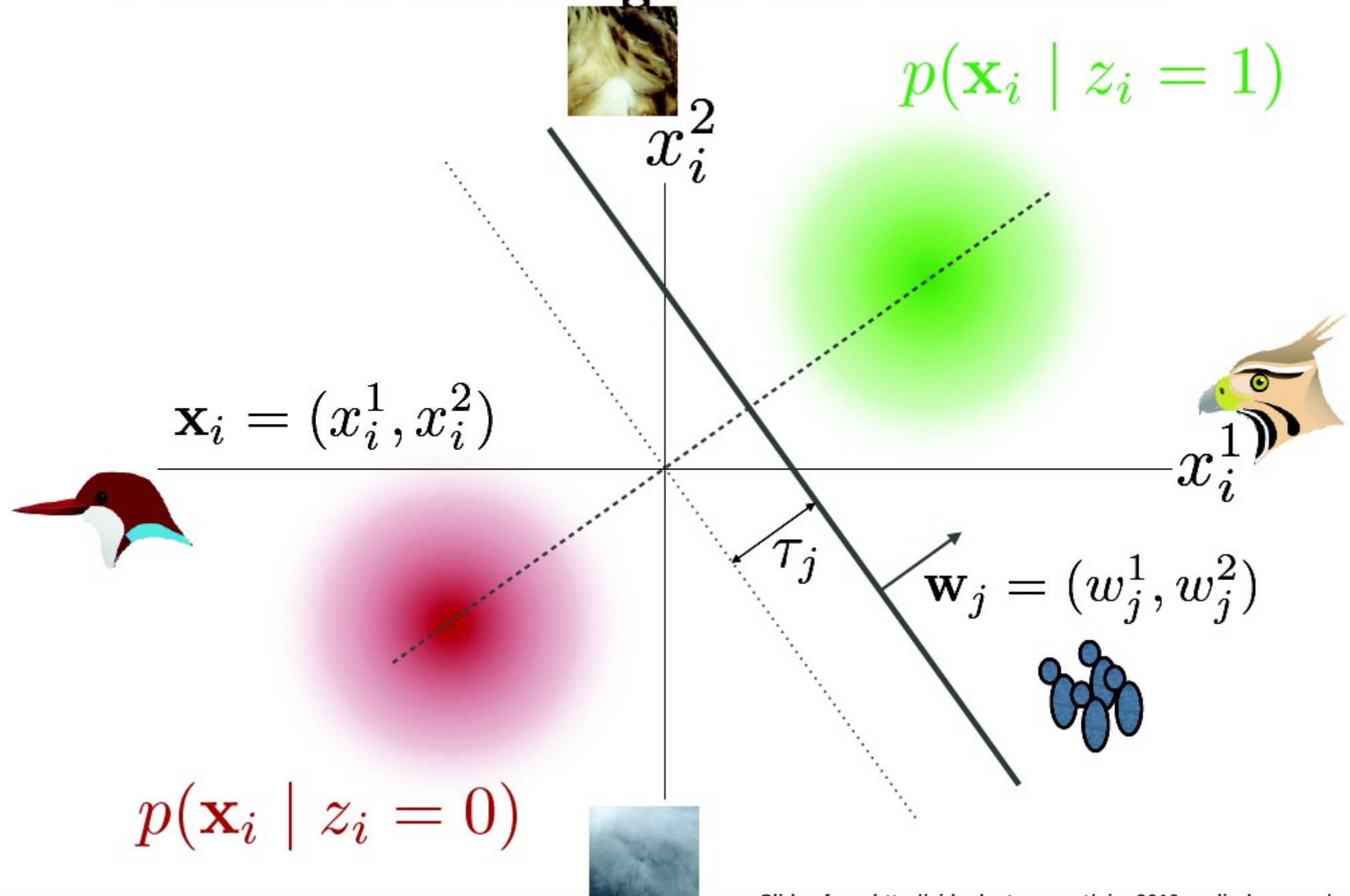
Multidimensional signals and annotators



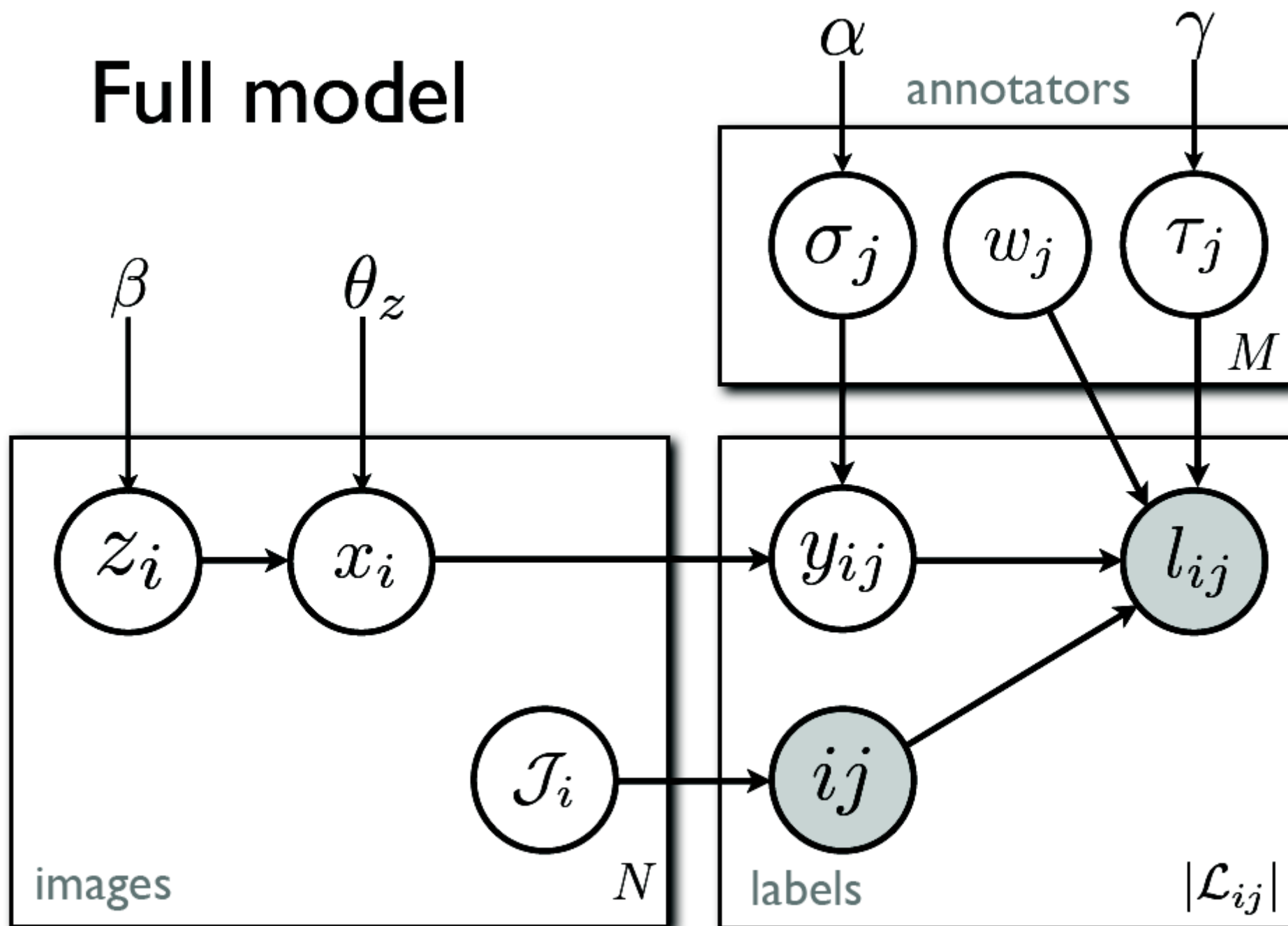
Multidimensional signals and annotators



Multidimensional signals and annotators

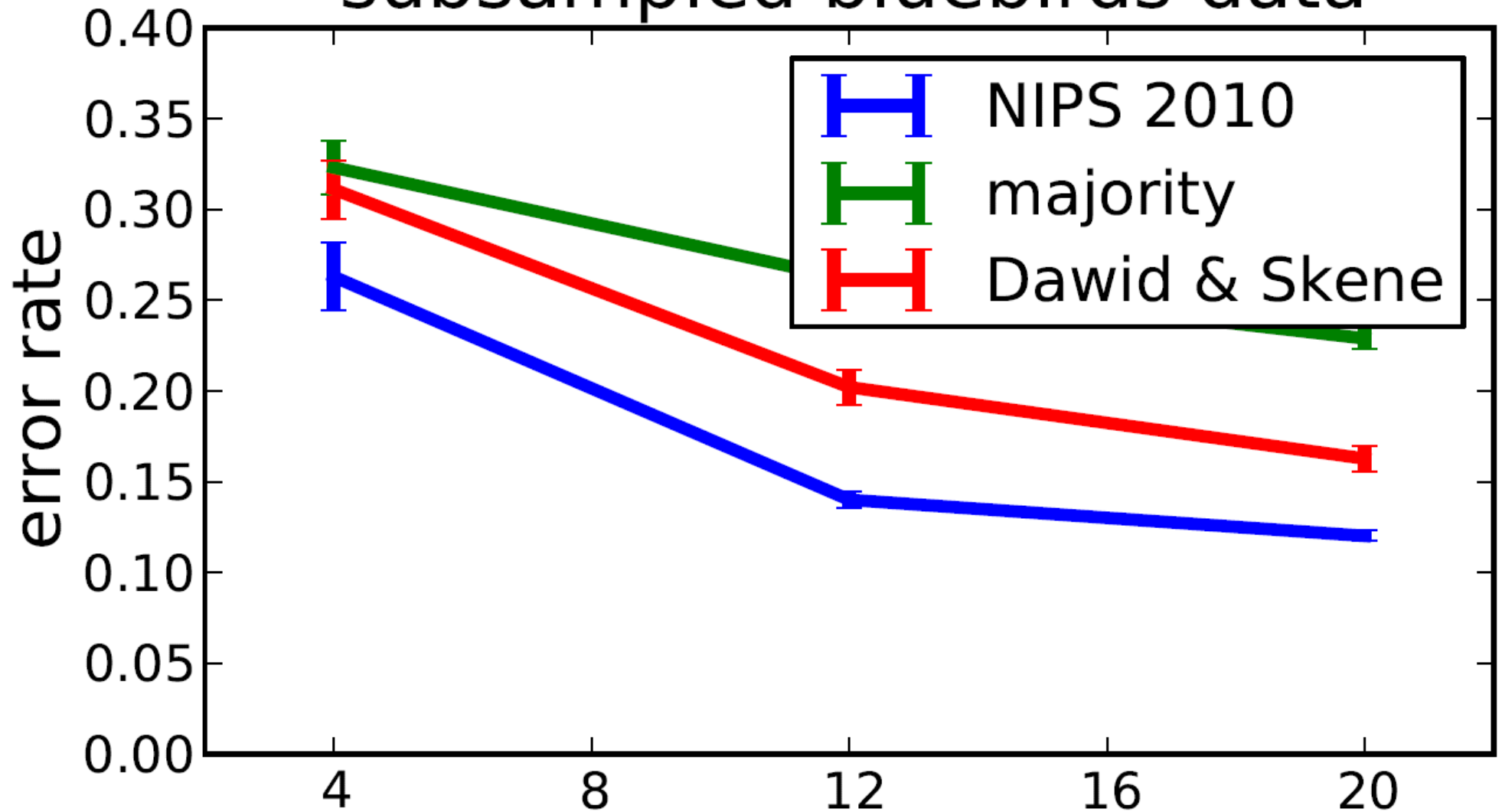


Full model

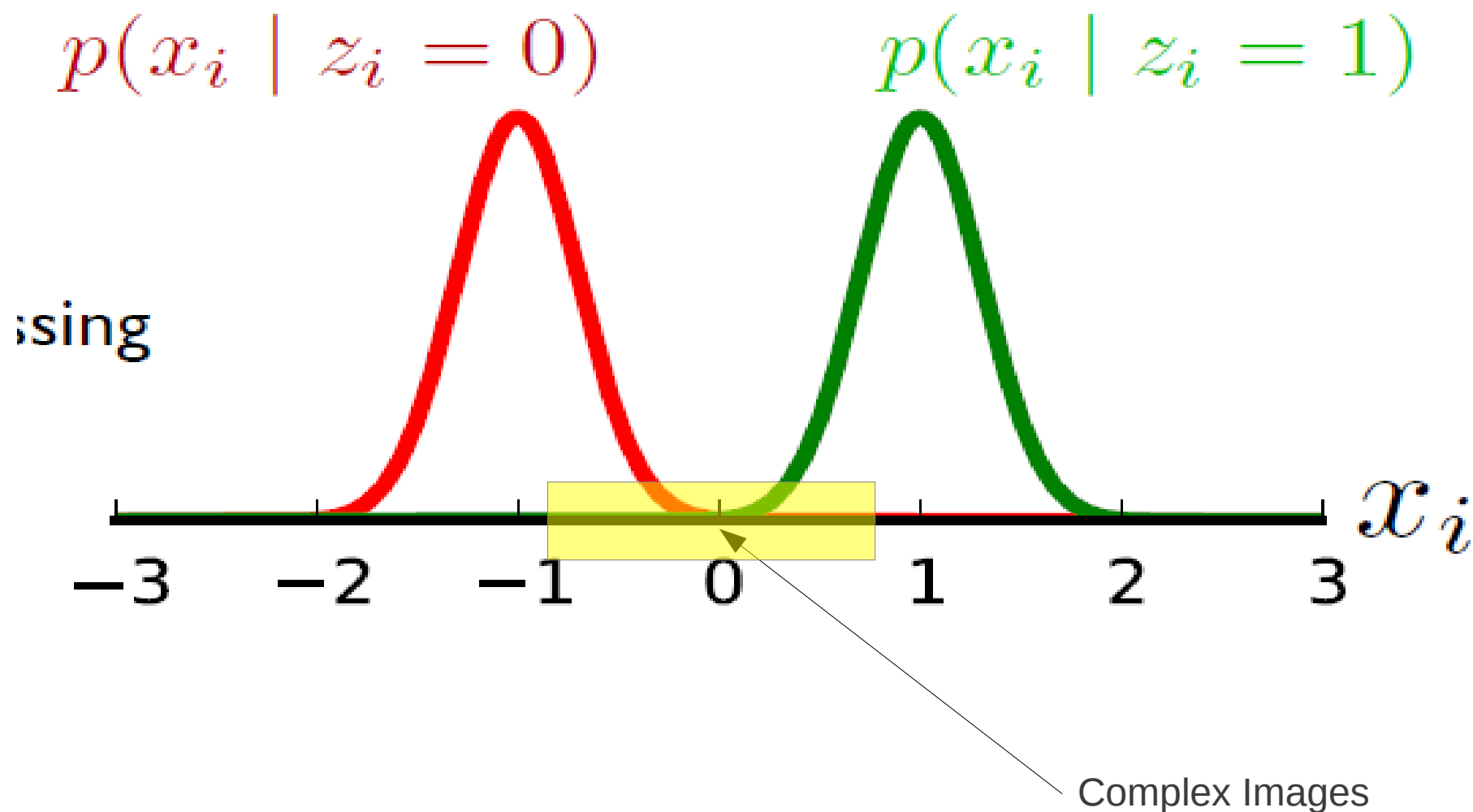


Error Rate for Bluebirds dataset

subsampled bluebirds data

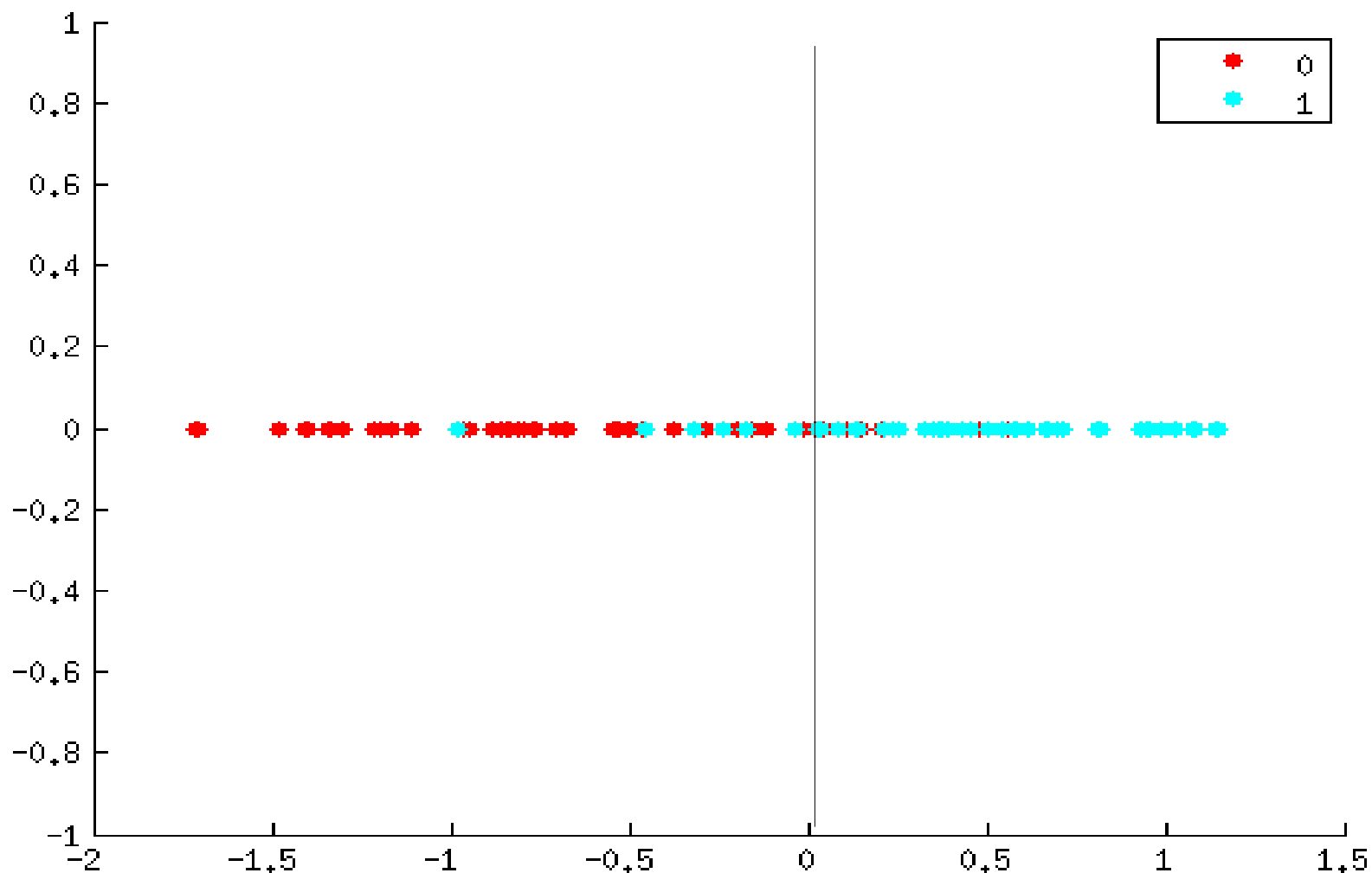


Estimating Image Difficulty



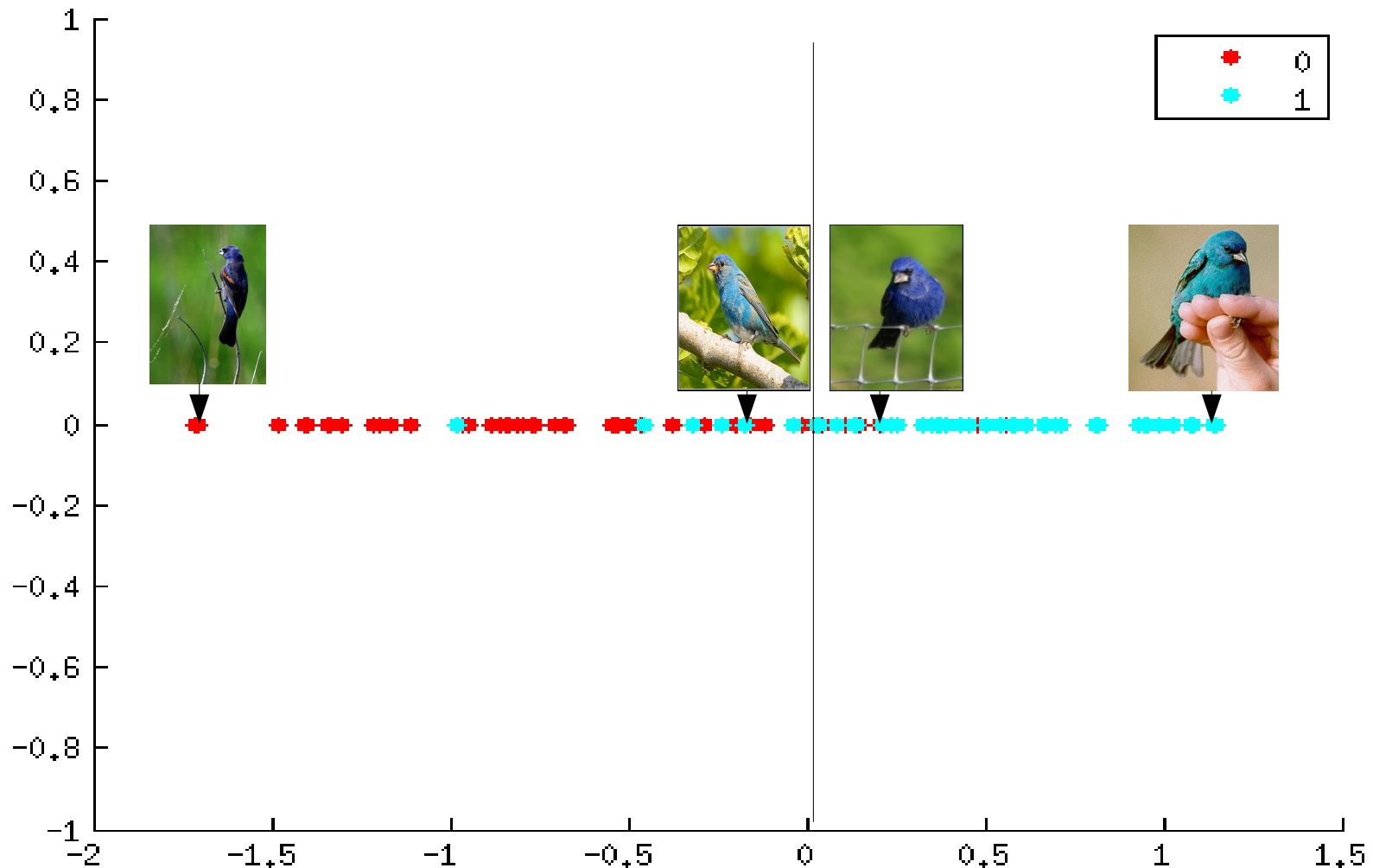
1D clusters from learned X_i values

Dataset: Bluebirds



1D clusters from learned X_i values

Dataset: Bluebirds



How do these learned image complexities compare with vision-based techniques?

Vision-based measure:

Predicted time* to label an image as a measure of image complexity

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

Approach:

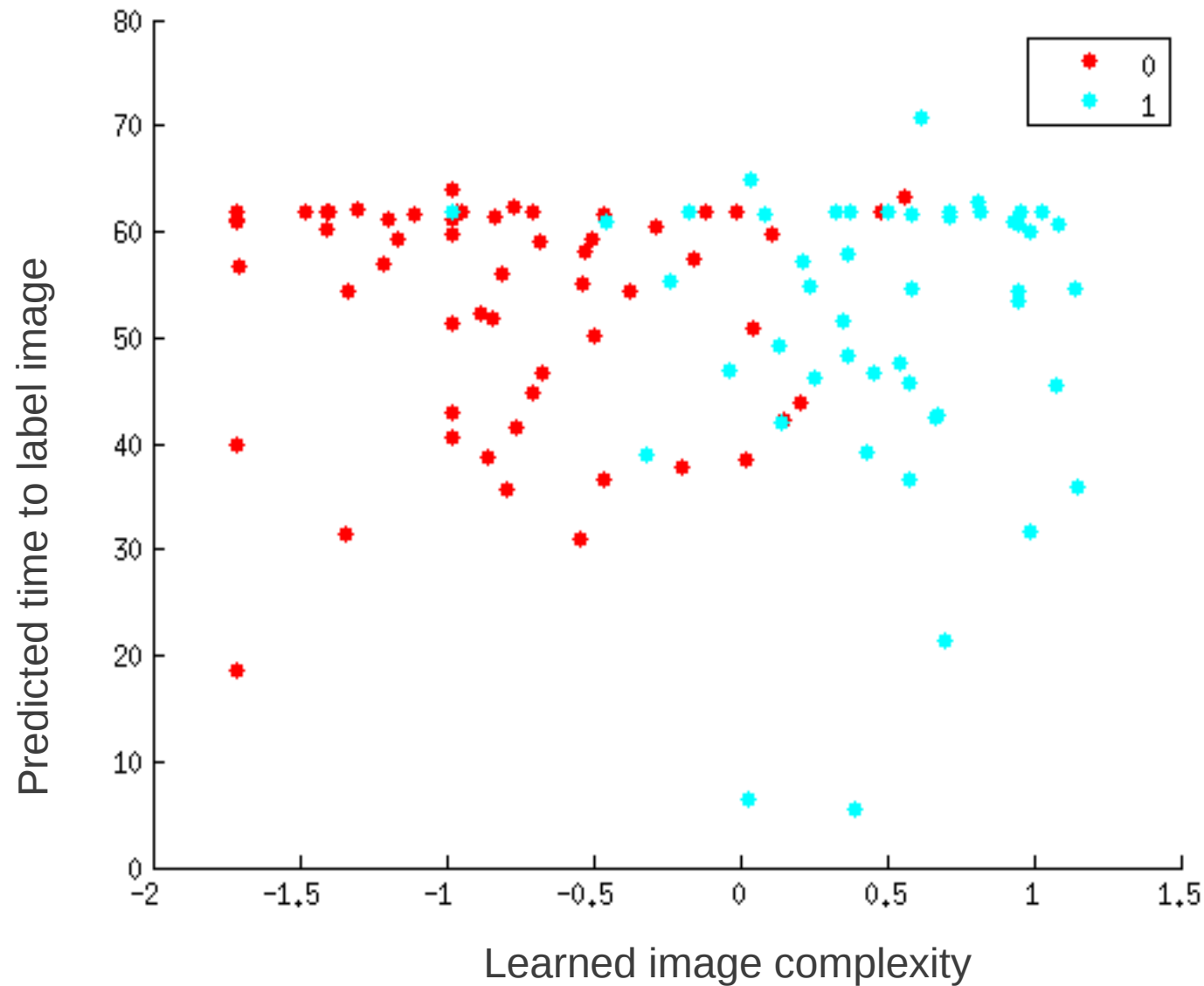
Extract 2804-d feature vectors for MSRC dataset

- Pyramid of HoG
- Color histogram
- Grayscale histogram
- Spatial pyramid of edge density (Canny edge)

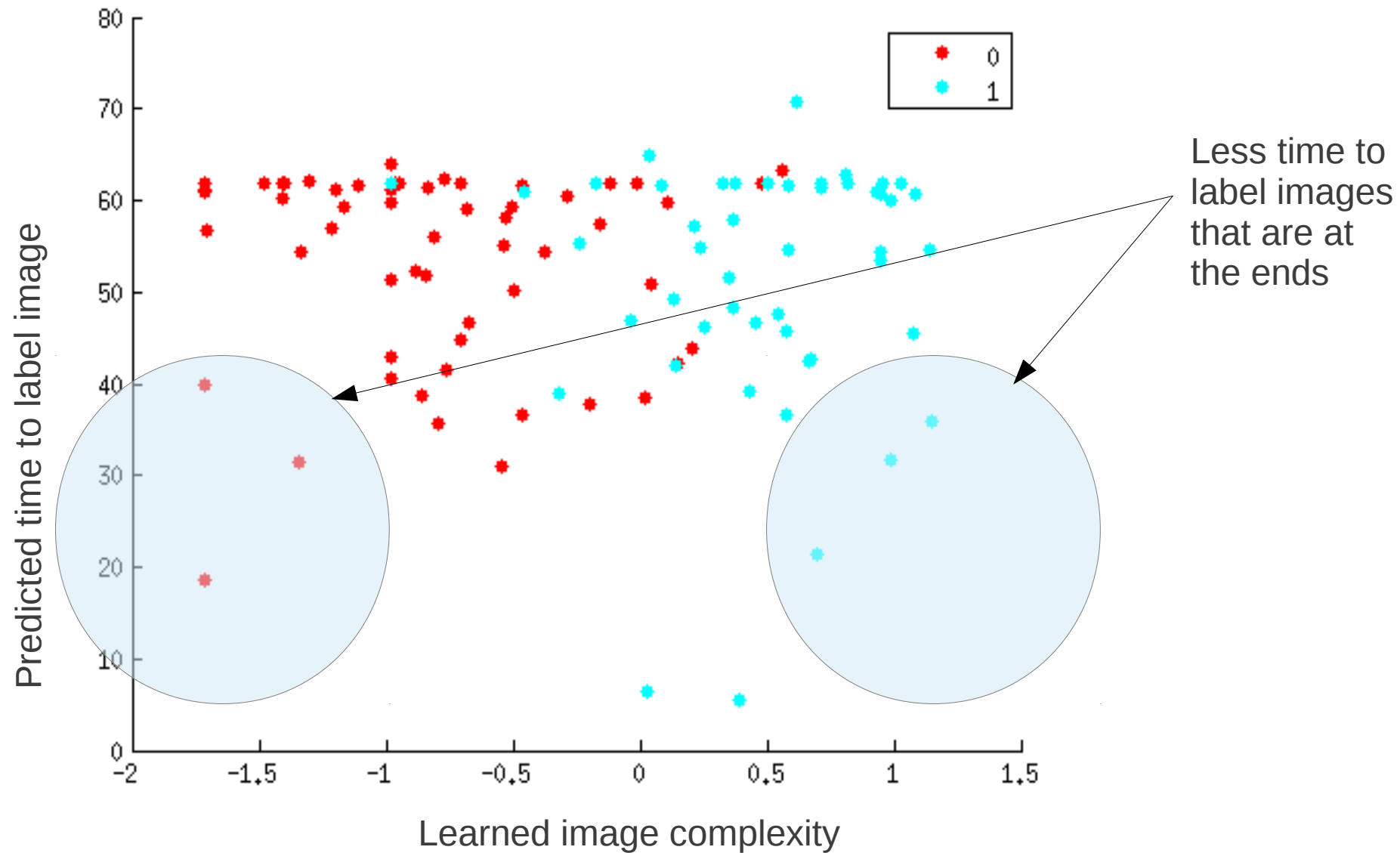
Train a regressor on top 200 features selected using ReliefF

Predict time to label images for bluebirds dataset

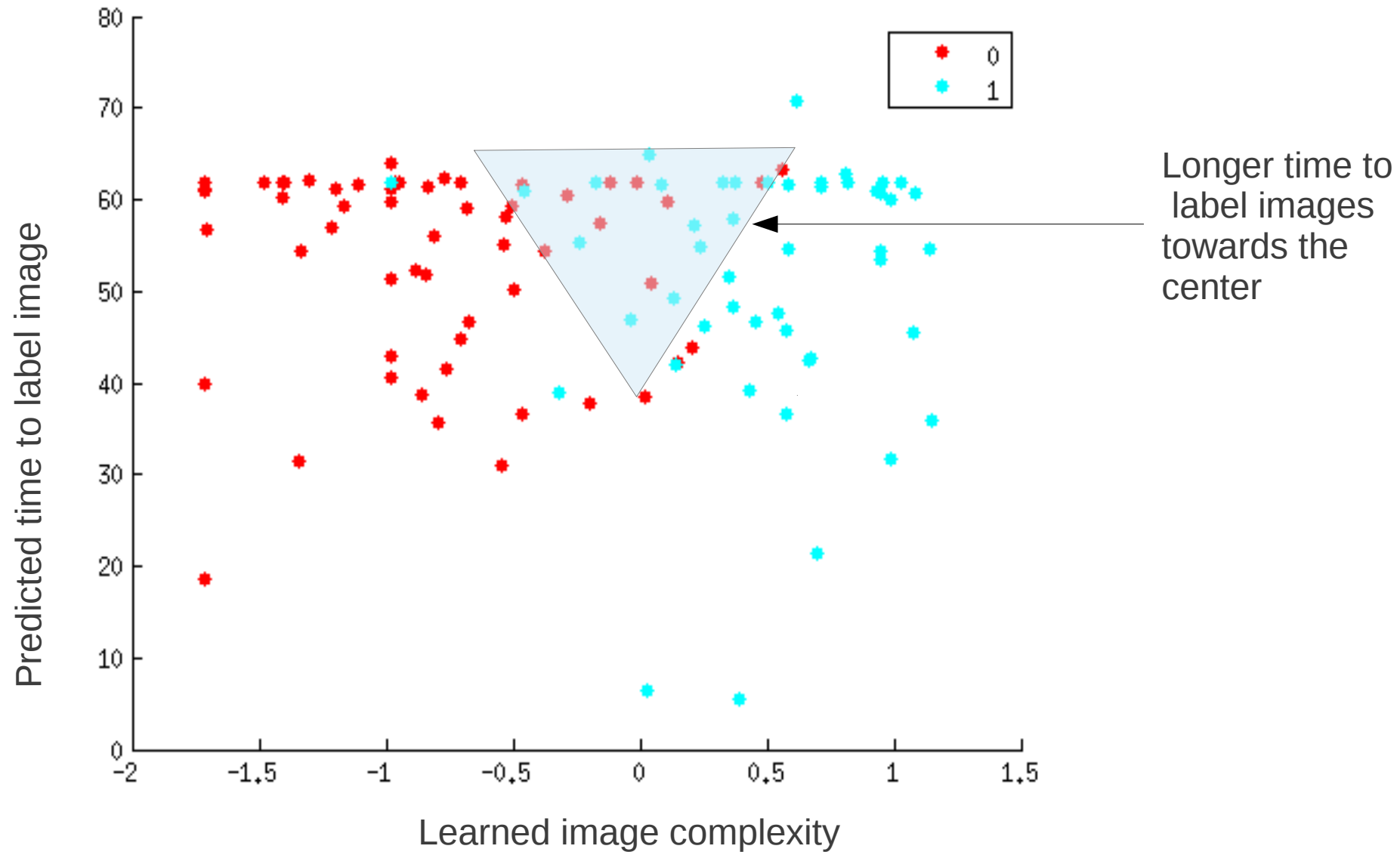
Vision-based complexity (vs) Learned Image Complexity



Vision-based complexity (vs) Learned Image Complexity



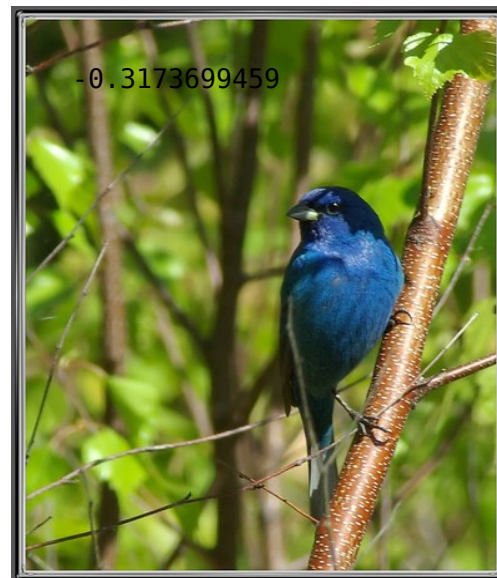
Vision-based complexity (vs) Learned Image Complexity



Qualitative Comparison

Complex Images – Examples

False negatives

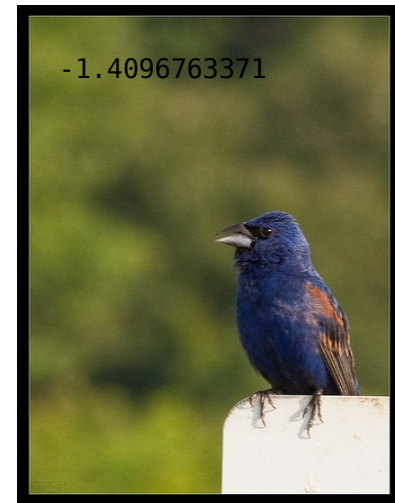


Complex Images – Examples

False positive



Easy Images – Examples True negatives



Easy Images – Examples True positives



Task: Finding ducks

Mallard



American Black Duck



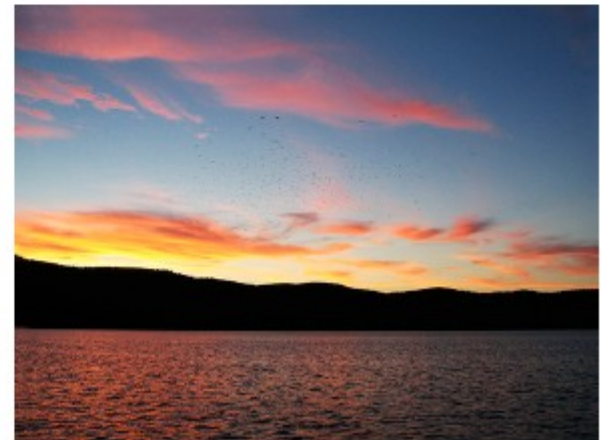
Canada Goose



Red Necked Grebe



Non-bird



DUCKS

Mallard



American Black Duck



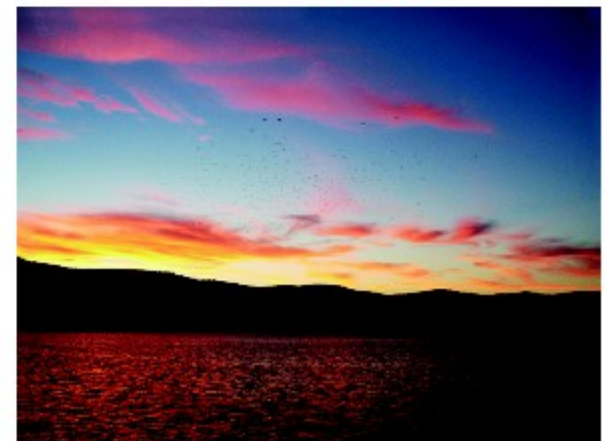
Canada Goose



Red Necked Grebe

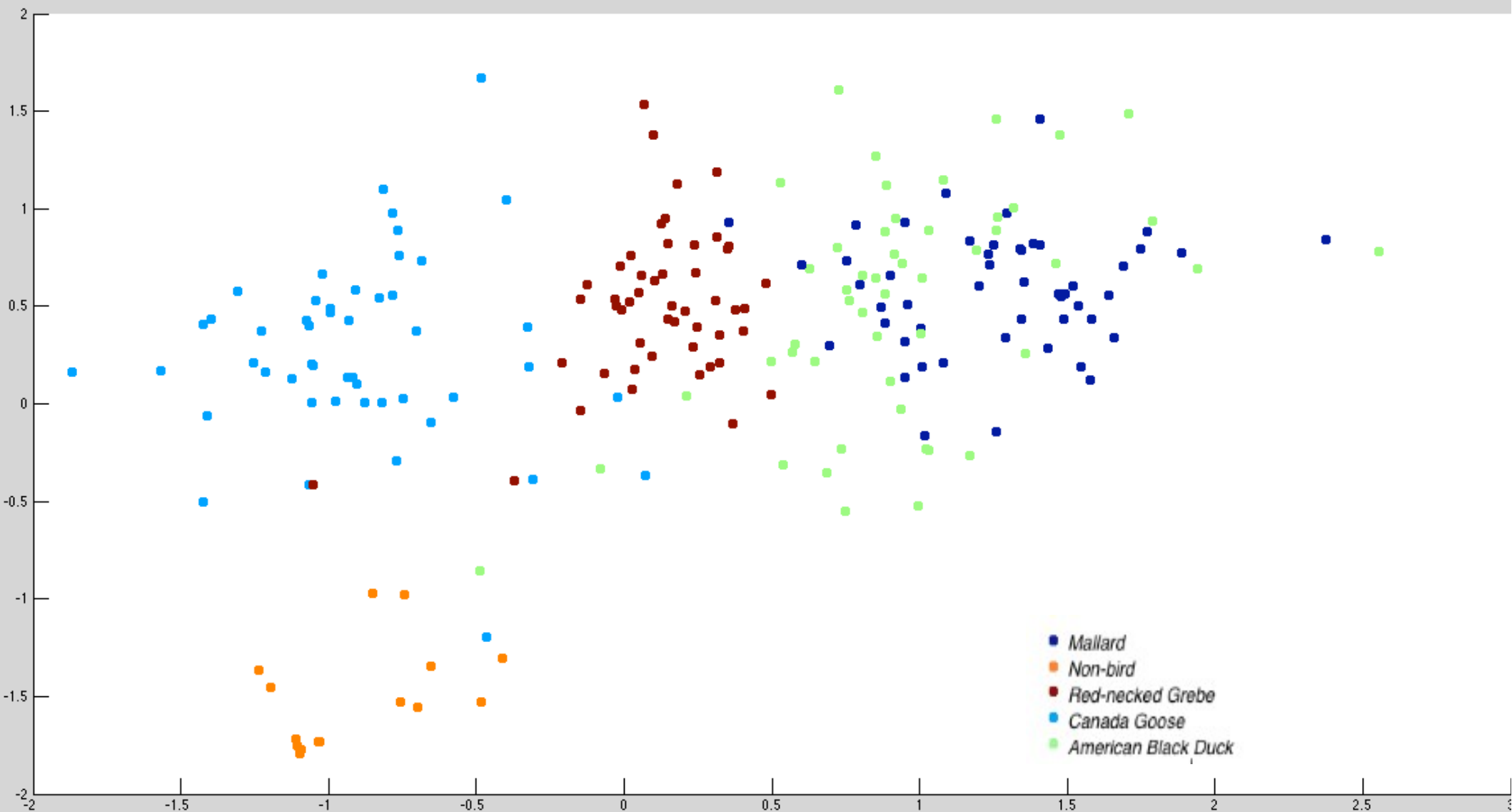


Non-bird



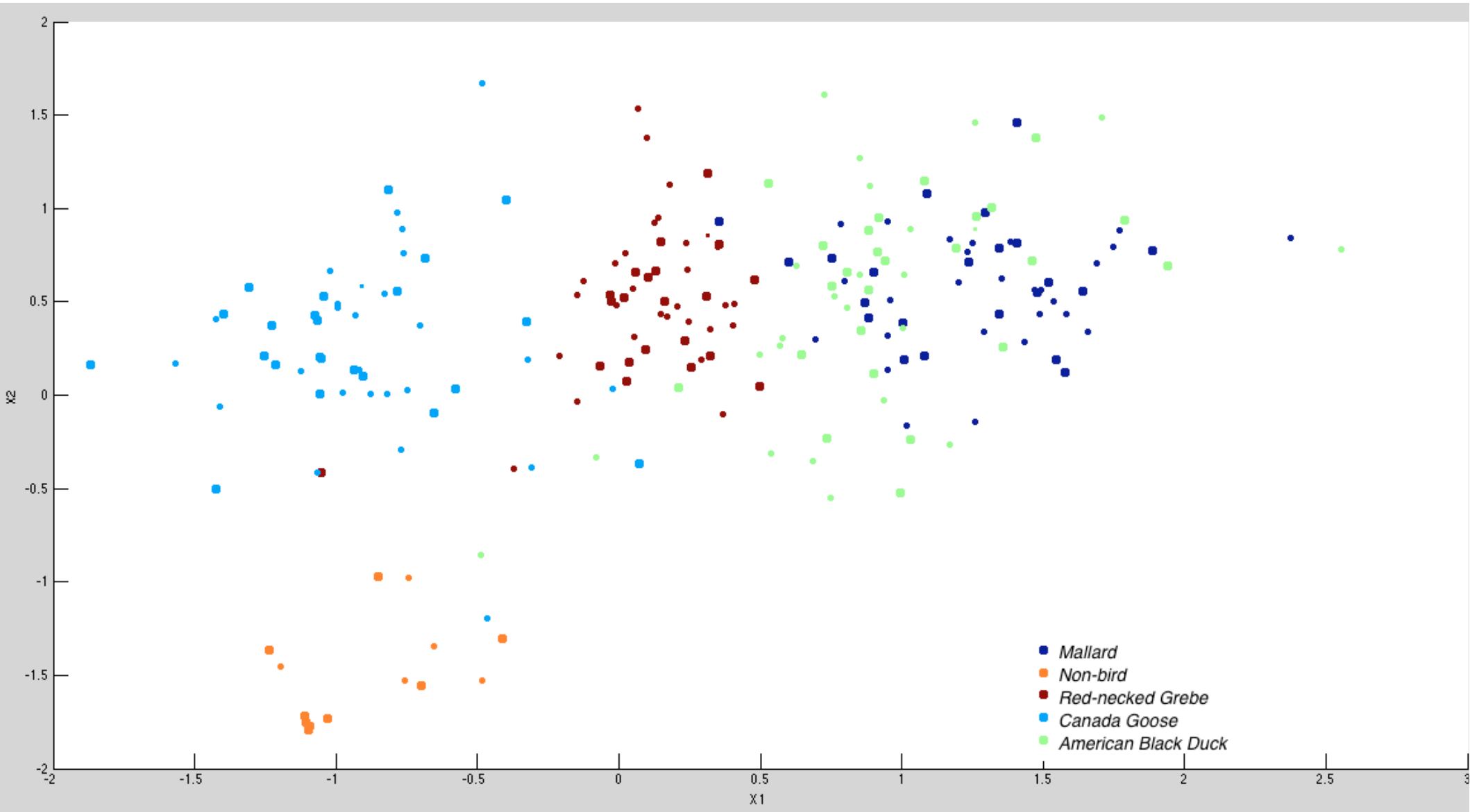
2D clusters from learned X_i values

Dataset: Ducks



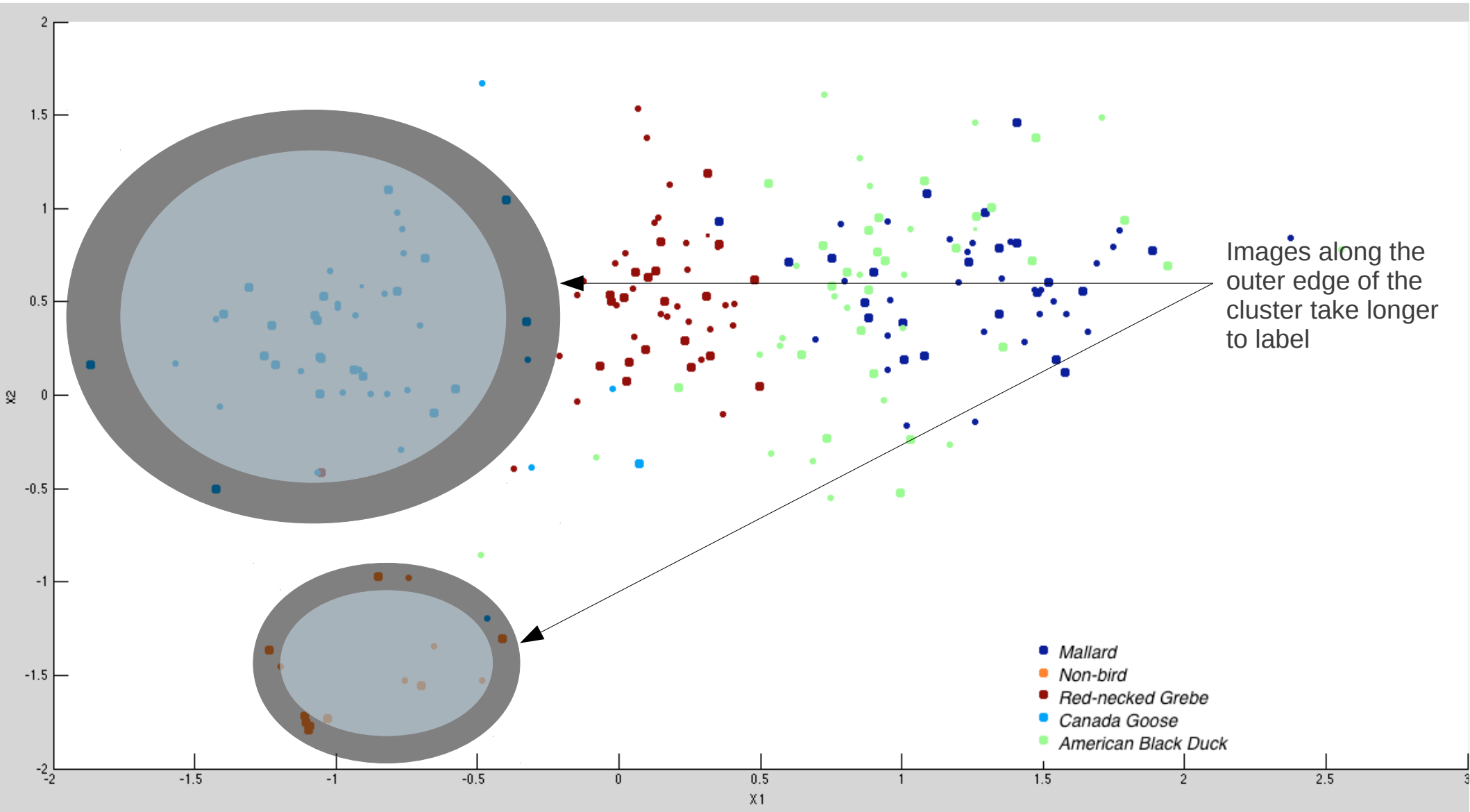
Recreated in MATLAB

Vision-based complexity (vs) Learned Image Complexity



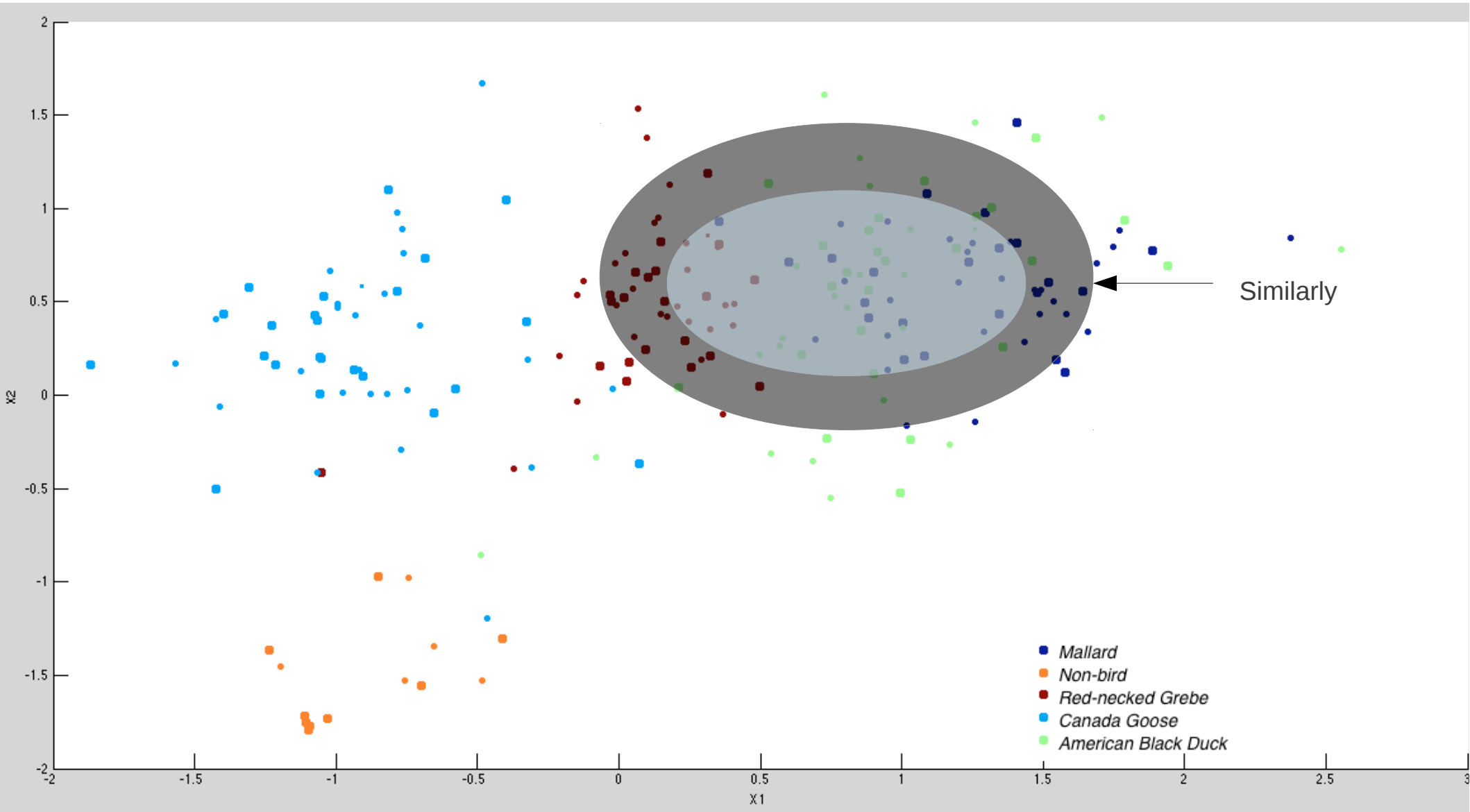
Recreated in MATLAB: **Size of point is proportional to the predicted time needed to label it**

Vision-based complexity (vs) Learned Image Complexity



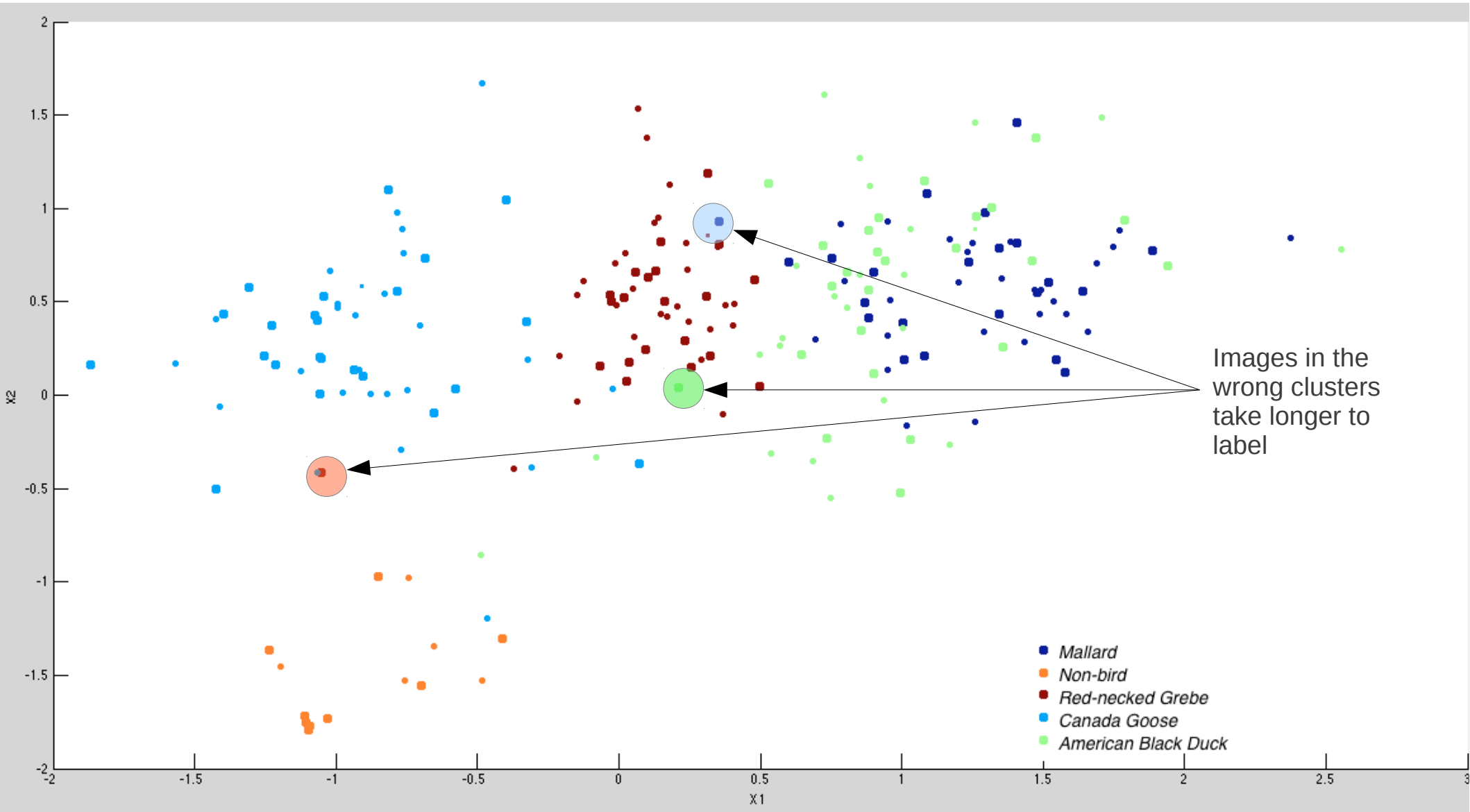
Recreated in MATLAB: Size of point is proportional to the predicted time needed to label it

Vision-based complexity (vs) Learned Image Complexity



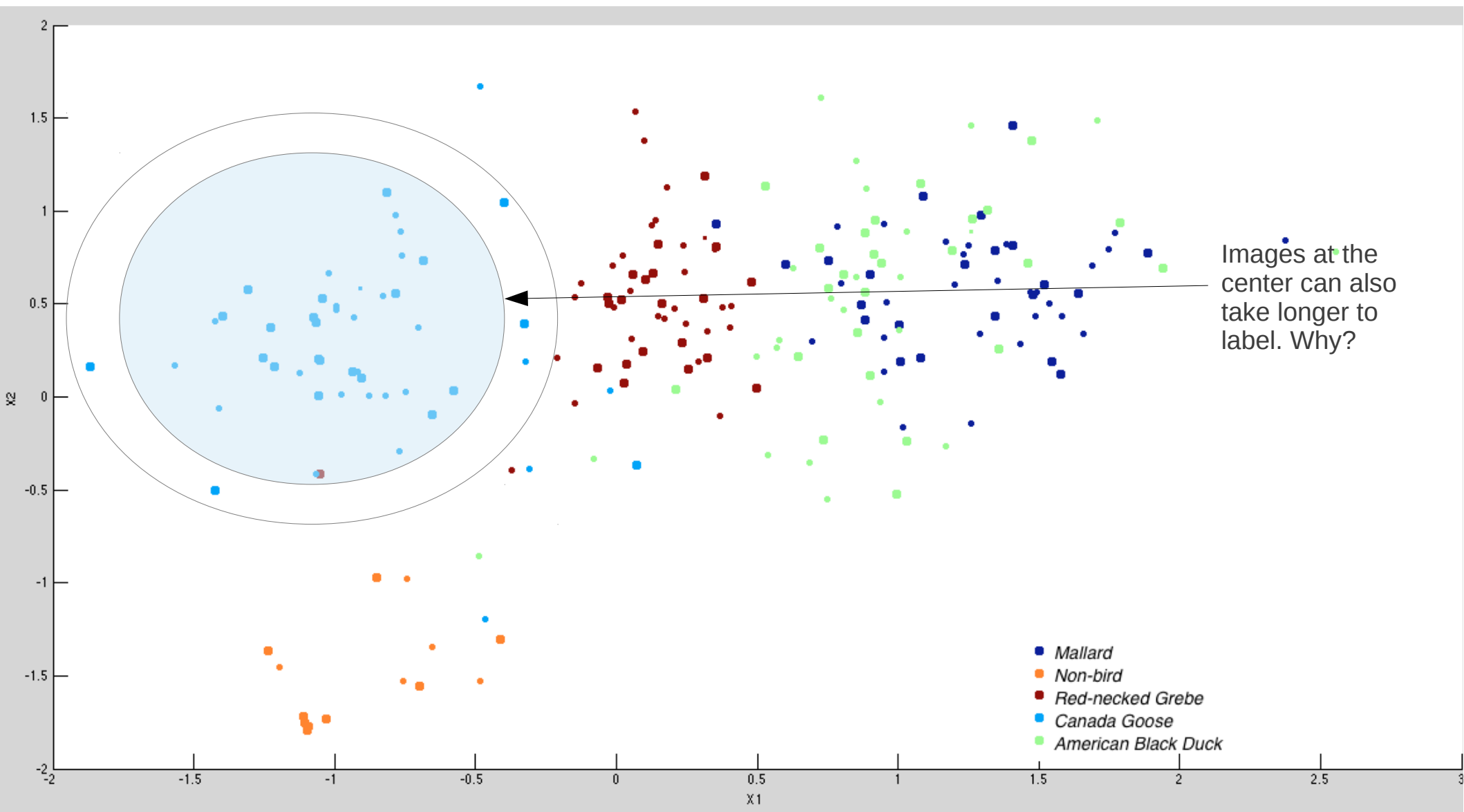
Recreated in MATLAB: Size of point is proportional to the predicted time needed to label it

Vision-based complexity (vs) Learned Image Complexity



Recreated in MATLAB: **Size of point is proportional to the predicted time needed to label it**

Vision-based complexity (vs) Learned Image Complexity



Recreated in MATLAB: **Size of point is proportional to the predicted time needed to label it**

Discussion

Is vision-based image complexity a good indicator of difficulty in labeling an image?

What are the other factors?

Discussion

Is vision-based image complexity a good indicator of difficulty in labeling an image?

What are the other factors?

Bird pose

Occlusions

Lighting

Discussion

1. The authors experiment only with a **2-dimensional** model of human expertise

How would this model perform by increasing the number of intrinsic dimensions?

Extending this approach
to a video dataset
YouTube corpus

Example YouTube video with descriptions

<http://youtu.be/FYyqIJ36dSU>

A french bulldog is **playing** with a big ball

A small dog **chases** a big ball.

A French bulldog is **running** fast and **playing** with a blue yoga ball all by himself in a field.

The little dog **pushed** a big blue ball.

A dog is **playing** with a very large ball.

A dog **chases** a giant rubber ball around

A dog is **playing** with ball

Approach

YouTube corpus is not cut out for this task.

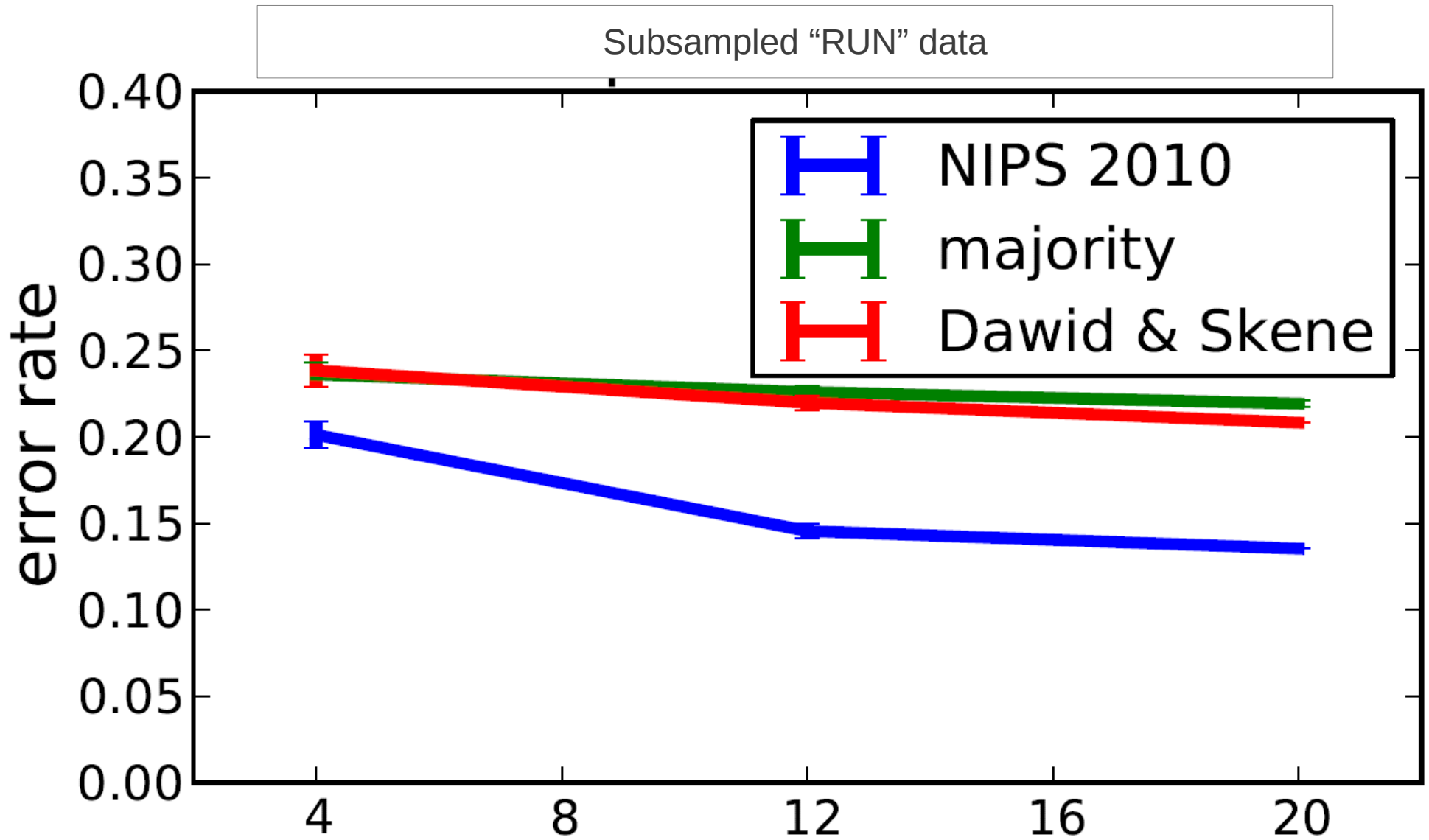
Consider predicting the presence of the activity “**run**”

1. Selected **50 videos where “run”** was the predicted activity using majority voting
2. Selected **30 videos where “play”** was the predicted activity using majority voting
3. Selected **20 videos where “walk”** was the predicted activity using majority voting

Ground Truth Labels were assigned accordingly

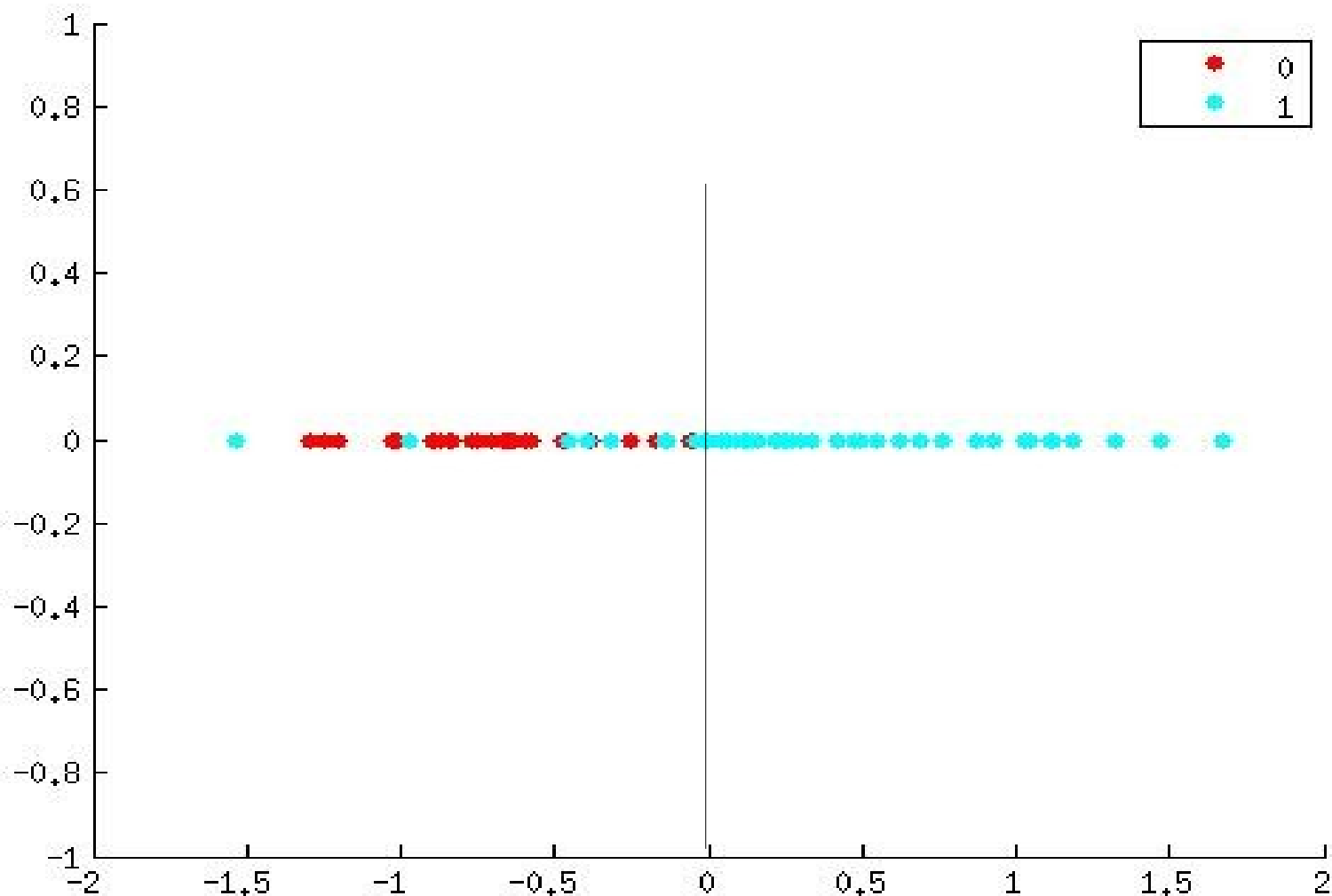
Each video had variable number of annotators. Picked the **20 most frequent annotators**.

Results



1D clusters from learned X_i values

Dataset: YouTube videos



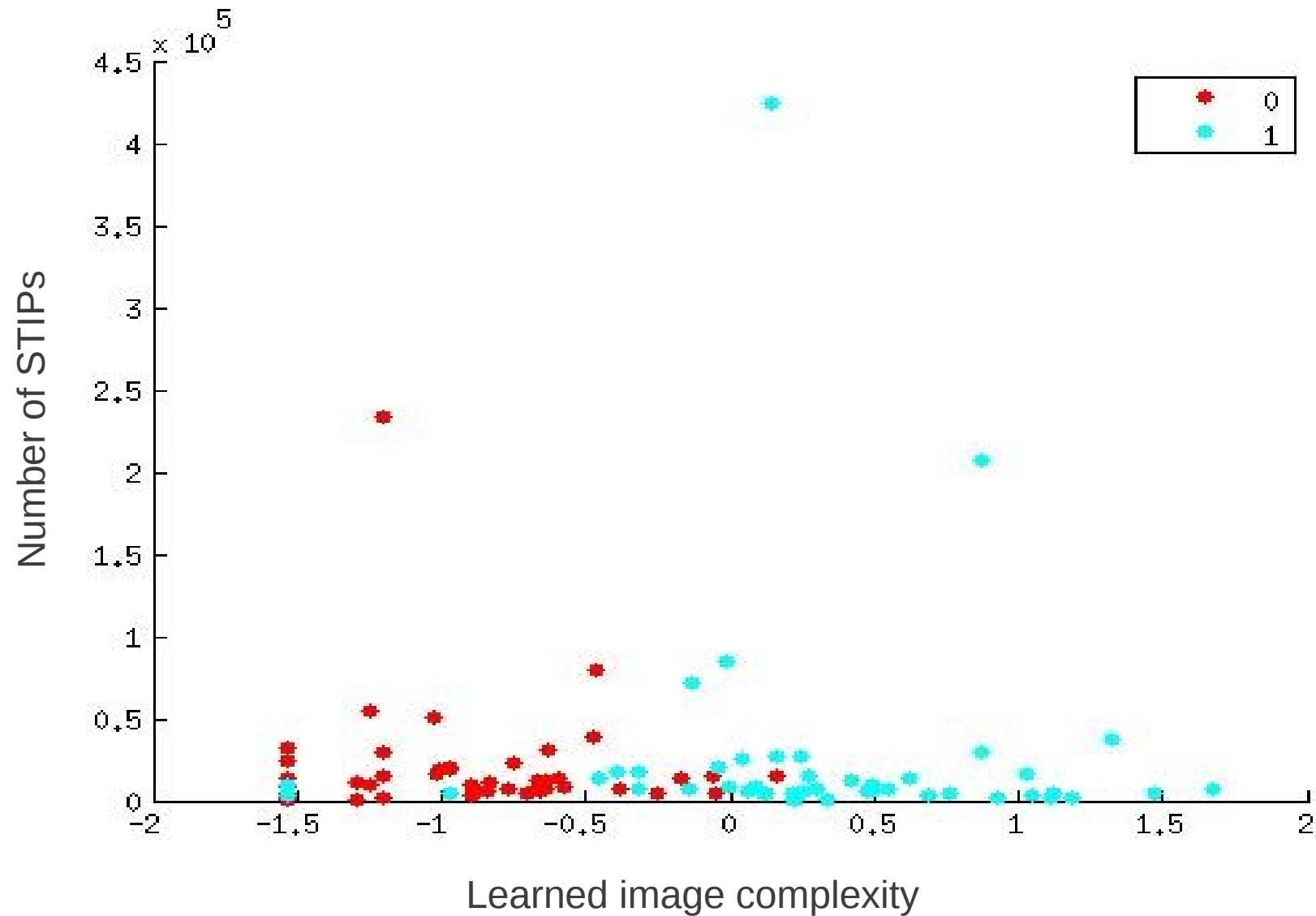
How do these learned video complexities compare with vision-based techniques?

Vision-based measure:

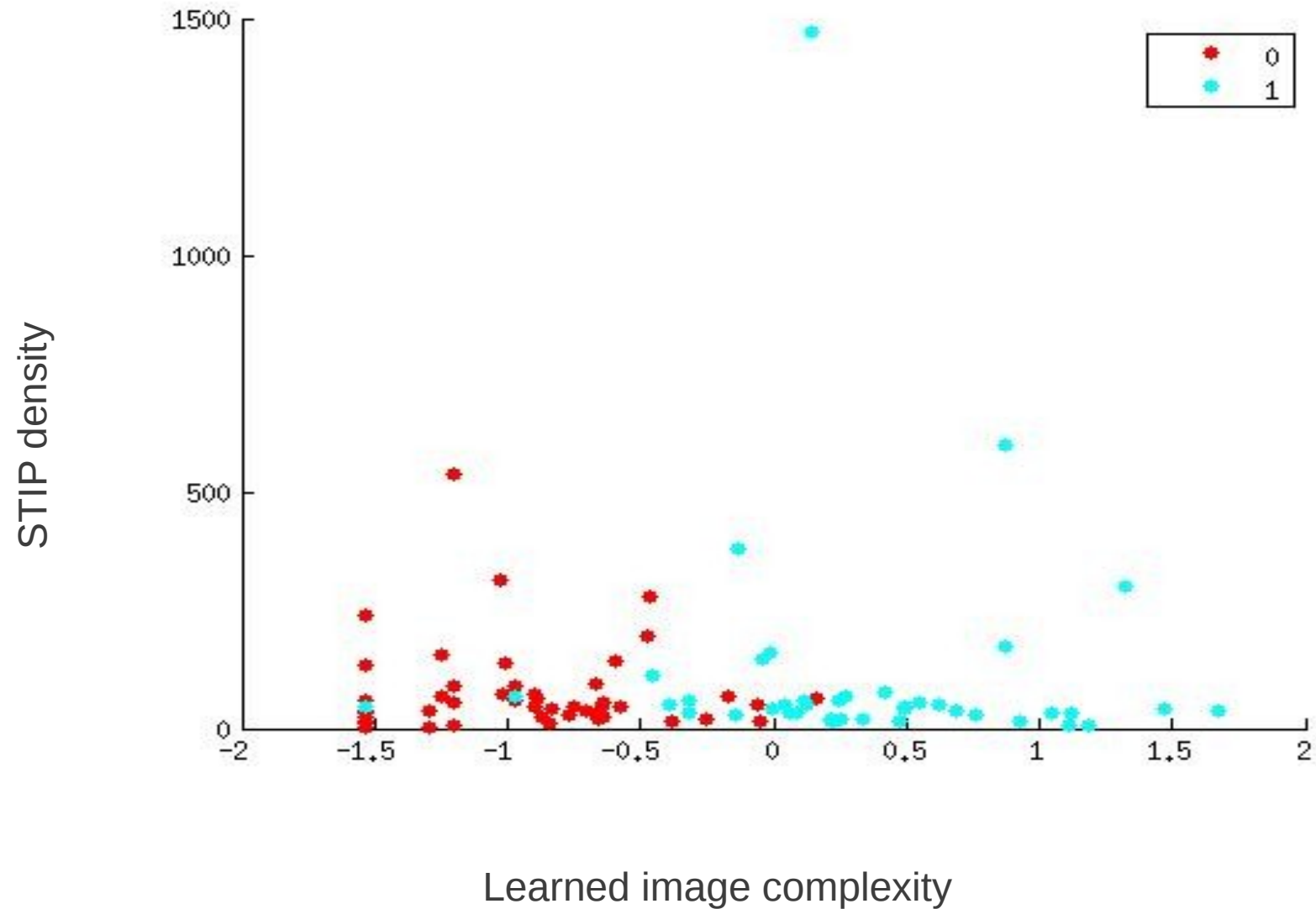
Number of STIPS in the video
STIP density

***Learning Realistic Human Actions from Movies. I. Laptev, M. Marszałek, C. Schmid and B. Rozenfeld. CVPR 2008.**

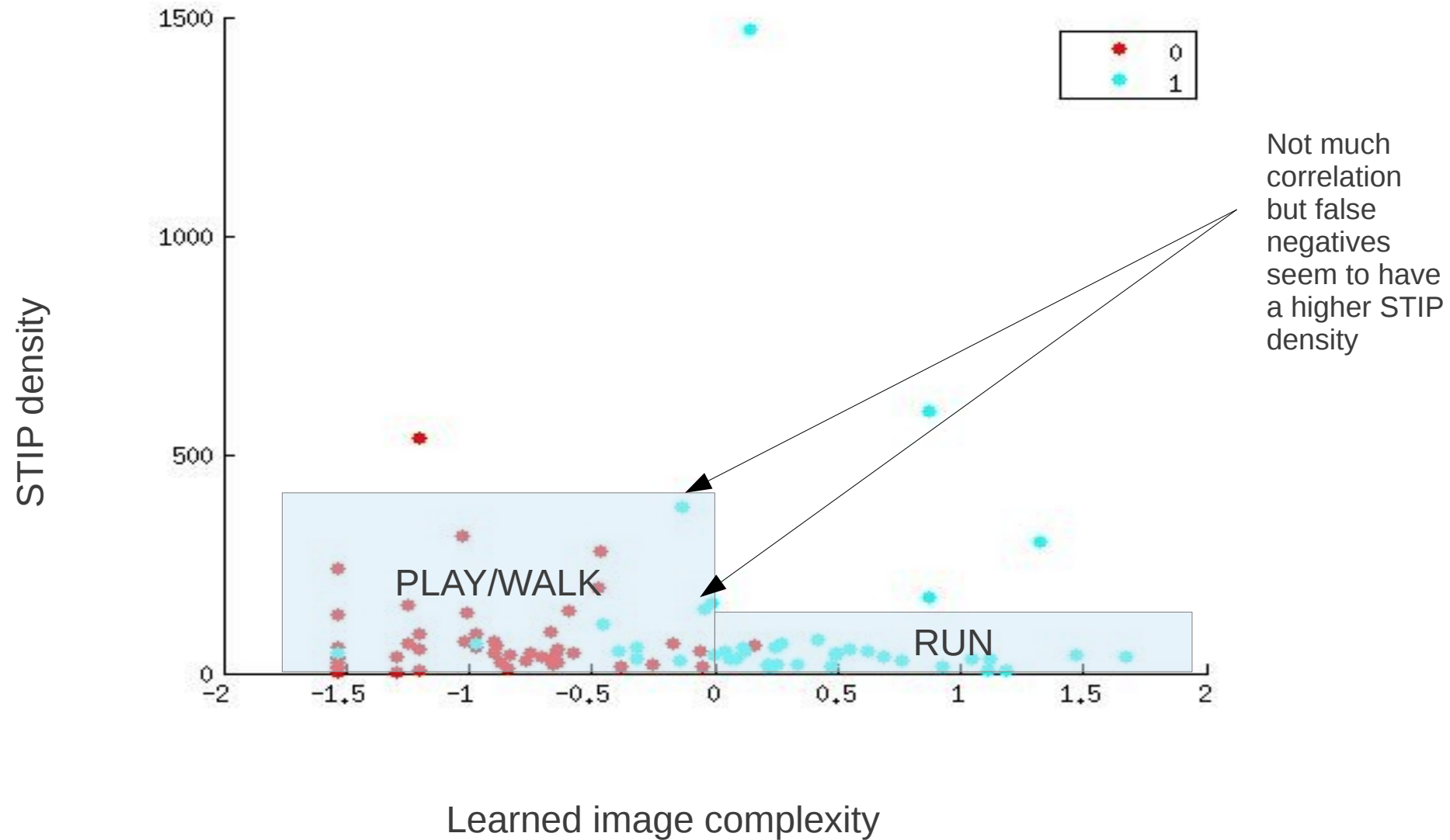
Vision-based complexity (vs) Learned Image Complexity



Vision-based complexity (vs) Learned Image Complexity



Vision-based complexity (vs) Learned Image Complexity



Discussion

How can we quantify the complexity of a video?

STIP density?

Video length?

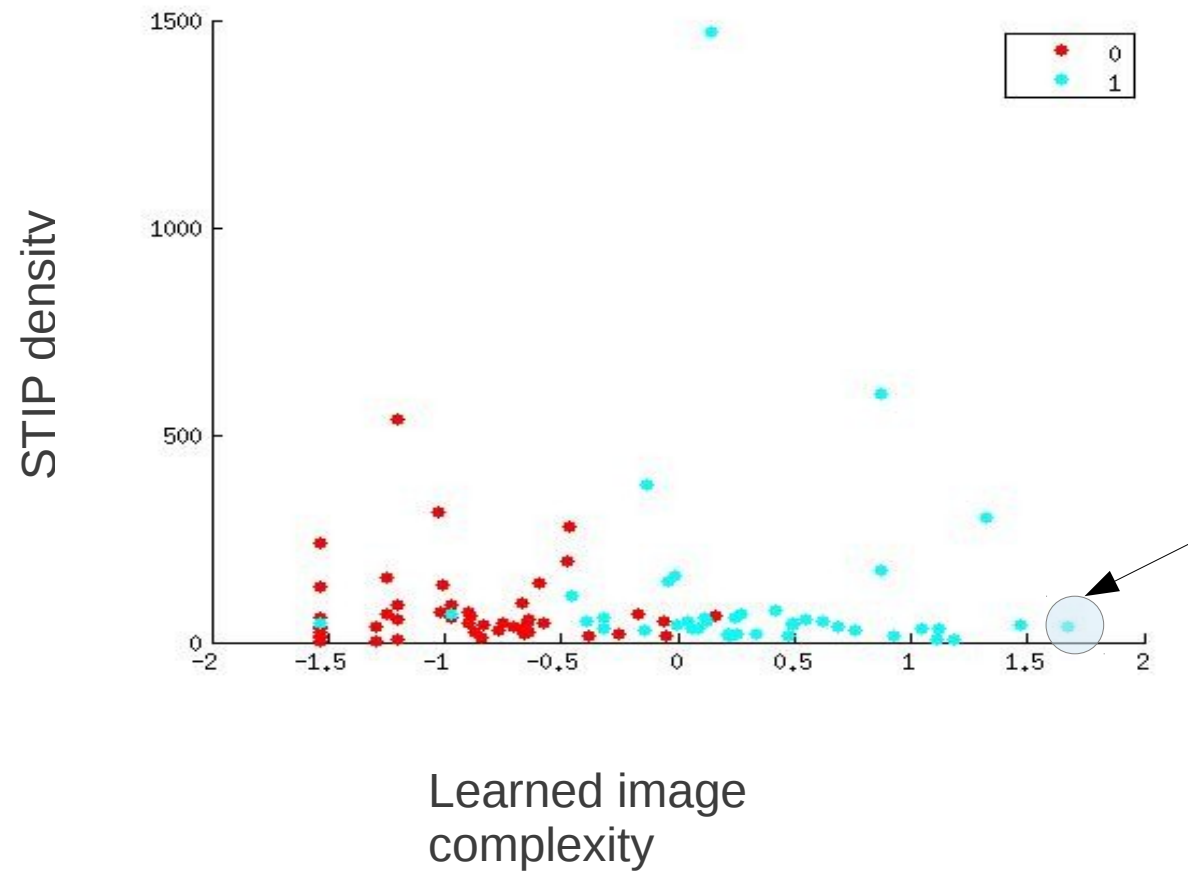
Variety in STIPS?

Confusion amongst multiple annotators?

How can we quantify the effort involved in labeling a video?

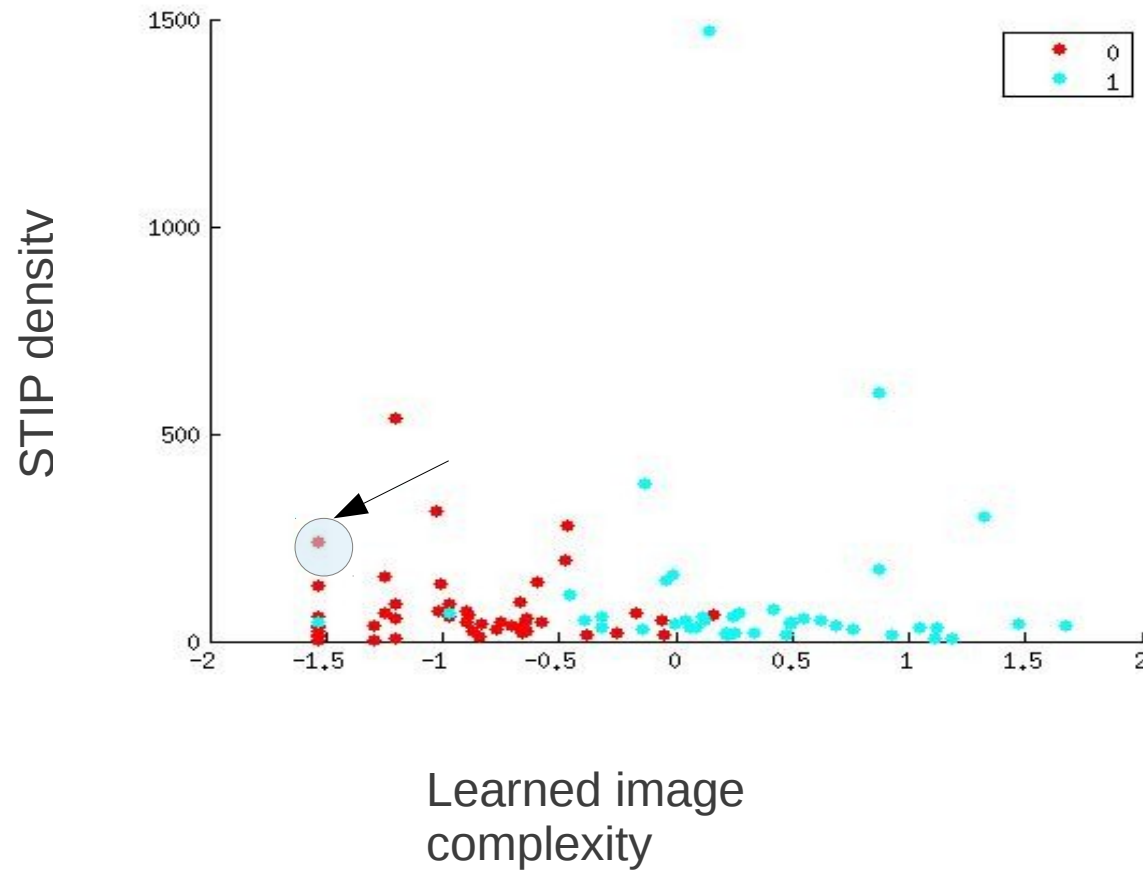
How do these relate to video ambiguity?

Qualitative Comparison – True positive



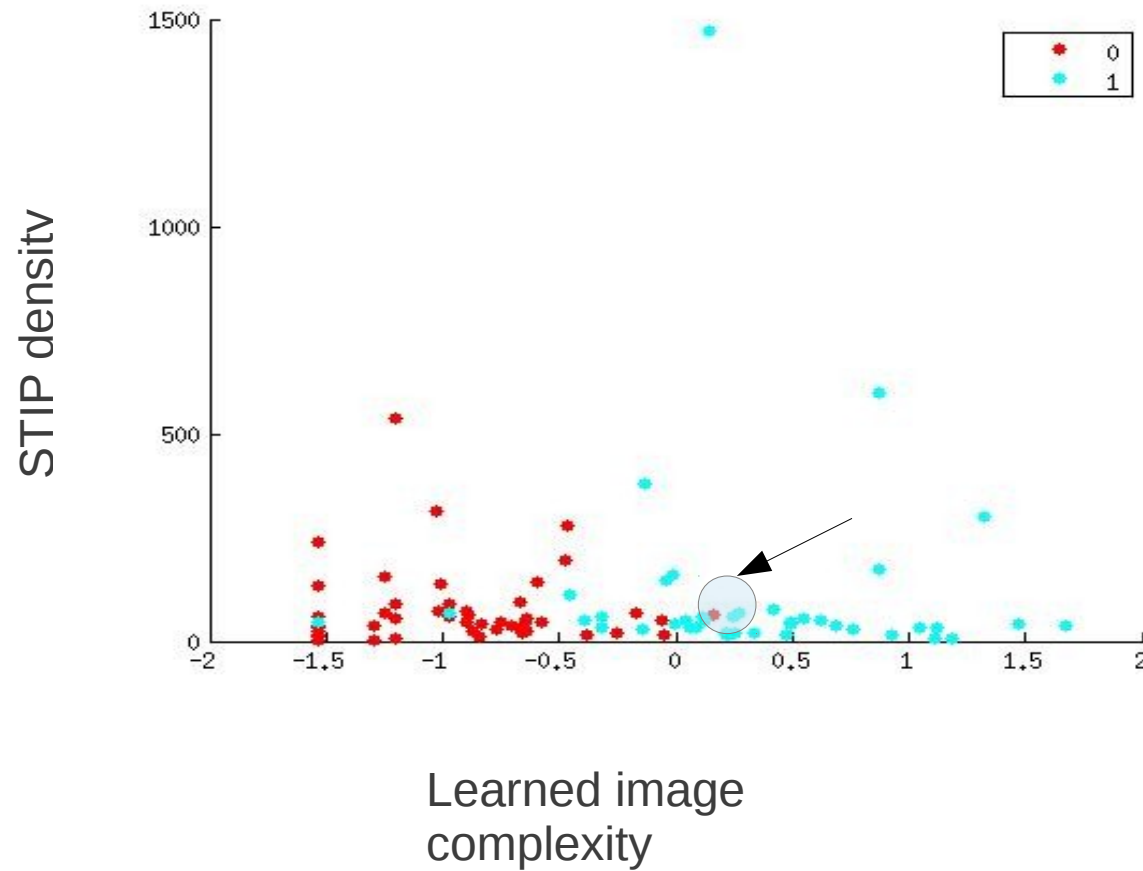
http://youtu.be/NKm8c_7mgx4

Qualitative Comparison – True negative



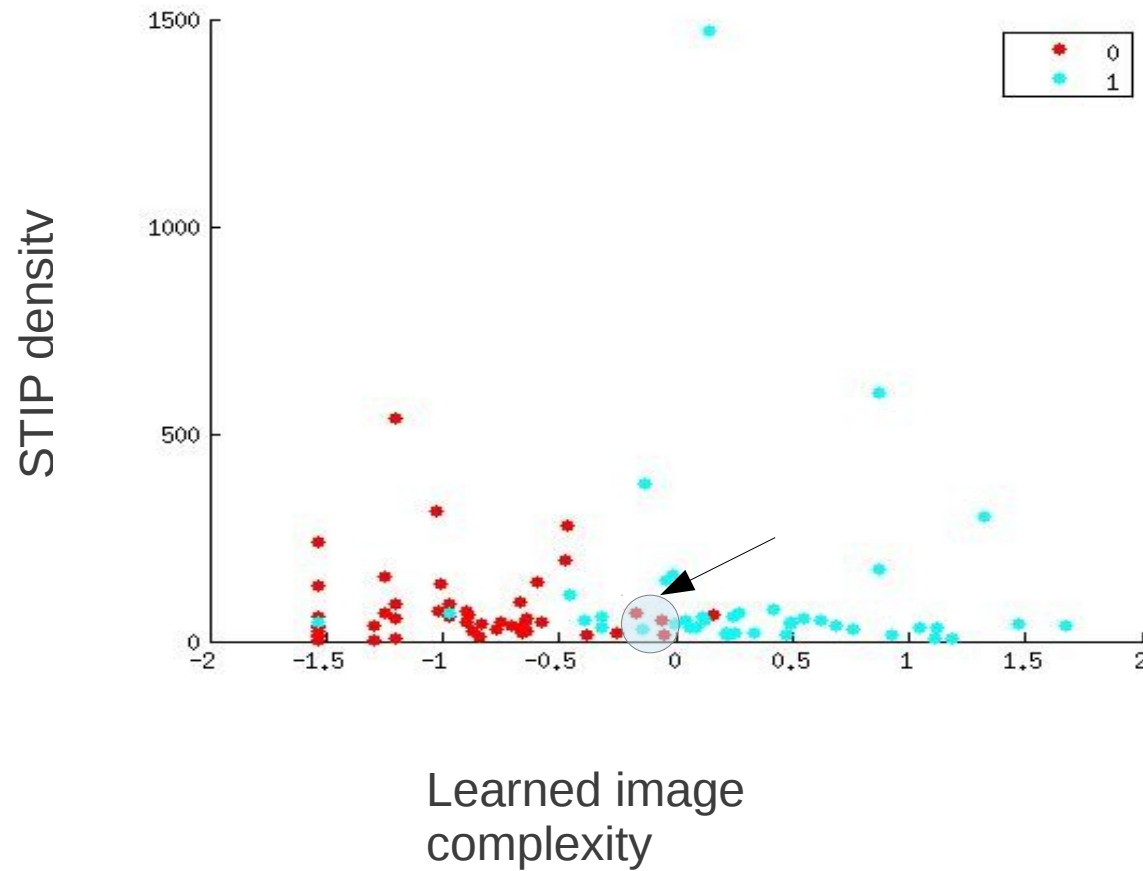
<http://youtu.be/abiezv1p7SY>

Qualitative Comparison – False positive



http://youtu.be/1l9Hx1kX_tQ

Qualitative Comparison – False negative



<http://youtu.be/8miosT-Fs1k>

Strengths

1. Each annotator is modeled as a multi-dimensional entity – competence, expertise, bias
2. Can be extended to any domain to estimate the ground truth with least error
3. Models image complexities without even seeing the image
4. The model discovers groups of annotators with varying skill sets.

Discussion

1. Image difficulties are learned from human annotations only, which is great!

But would the model perform better if image difficulty was incorporated as a known parameter (using some vision-based technique) into the graphical model?

