

# Scene Classification

**Danna Gurari**

The University of Texas at Austin

Fall 2019



<https://www.ischool.utexas.edu/~dannag/Courses/CrowdsourcingForCV/CourseContent.html>

This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/)

# Review

- Last week:
  - Object recognition applications
  - Object recognition datasets: key steps in creating them
  - Object recognition datasets: scaling up their size with *crowdsourcing*
  - Scaling up community working on object recognition with *workshop challenges*
- Assignments (Canvas)
  - Reading assignment 2 due yesterday
  - Lab assignment 1 due next week
- Questions?



# Today's Topics

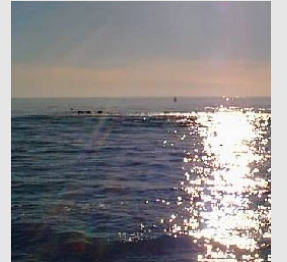
- Scene classification applications
- Scene classification datasets: key steps in creating them
- Scene classification datasets: scaling up with *crowdsourcing* and *challenges*
- Class discussion (chosen by YOU 😊)
- Lab: Javascript

# Today's Topics

- Scene classification applications
- Scene classification datasets: key steps in creating them
- Scene classification datasets: scaling up with *crowdsourcing* and *challenges*
- Class discussion (chosen by YOU 😊)
- Lab: Javascript

# Today's Focus: Scene Classification

Input:



Label:

Kitchen

Store



Coast

*"... a place in which a human can act within, or a place to which a human being could navigate."*

- Xiao et al; 2010

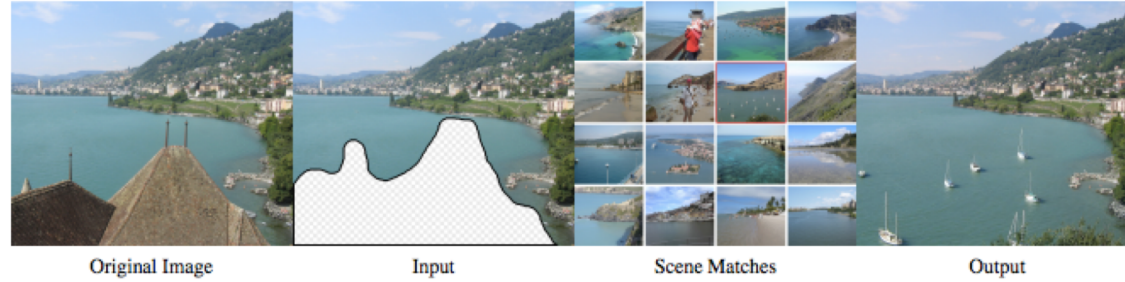
# Why Scene Classification?



- Object Recognition
  - e.g., What would you expect (or not expect) to find in the scene [now, earlier, later]?
- Activity Recognition/Prediction
  - e.g., What would you expect people to do (or not do) in the scene [now, earlier, later]?

# Why Scene Classification?

Idea:



Example:



Example:

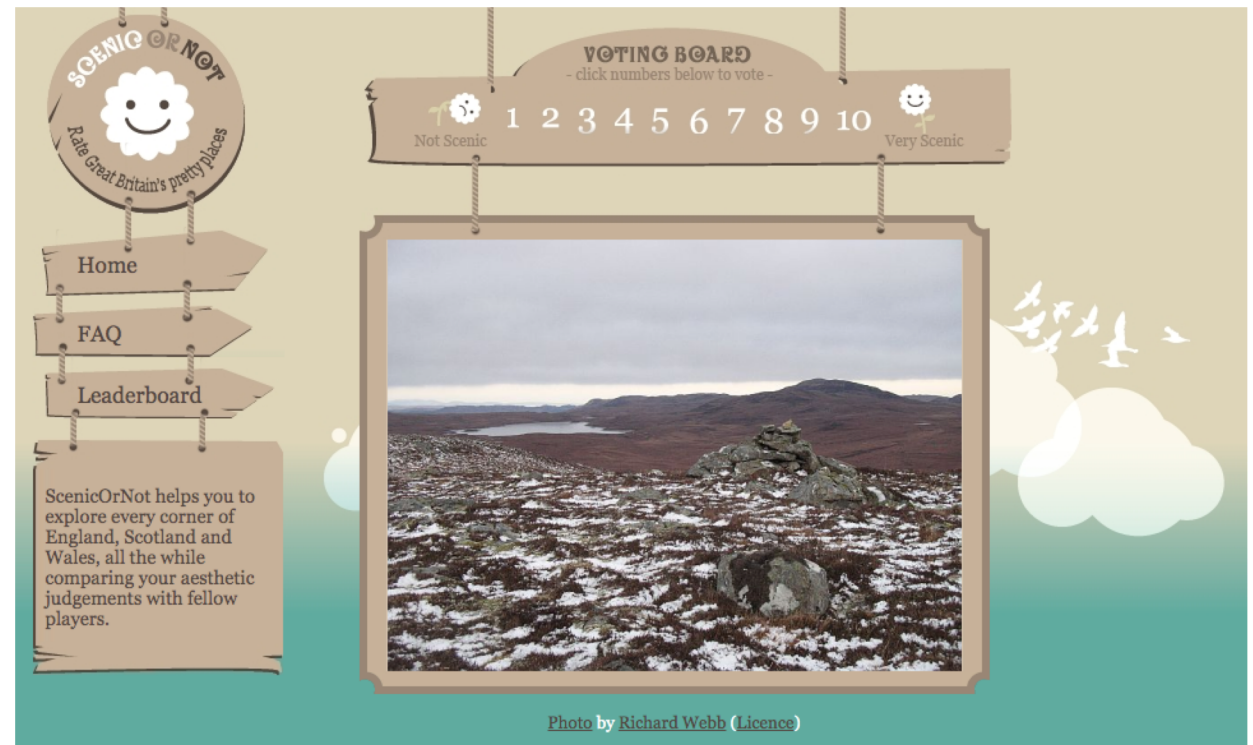


# Why Scene Classification?

Urban planning, since people's *well-being* is correlated with *scenic* places

Demo:

<http://scenicornot.datasciencelab.co.uk/>



Dataset: <http://scenicornot.datasciencelab.co.uk/>

Chanuki Illushka Seresinhe et al. Happiness is greater in more scenic locations. *Scientific reports*, 2019.

<https://www.economist.com/science-and-technology/2017/07/20/computer-analysis-of-what-is-scenic-may-help-town-planners>

# Scene vs Object Recognition

How is scene classification distinct from object recognition?

# Today's Topics

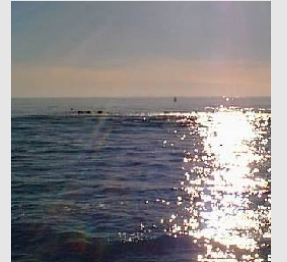
- Scene classification applications
- Scene classification datasets: key steps in creating them
- Scene classification datasets: scaling up with *crowdsourcing* and *challenges*
- Class discussion (chosen by YOU 😊)
- Lab: Javascript



# Recall: Need Datasets to **Train** & Evaluate Algorithms

## 1. Create Training Data

Input:



Label:

Kitchen

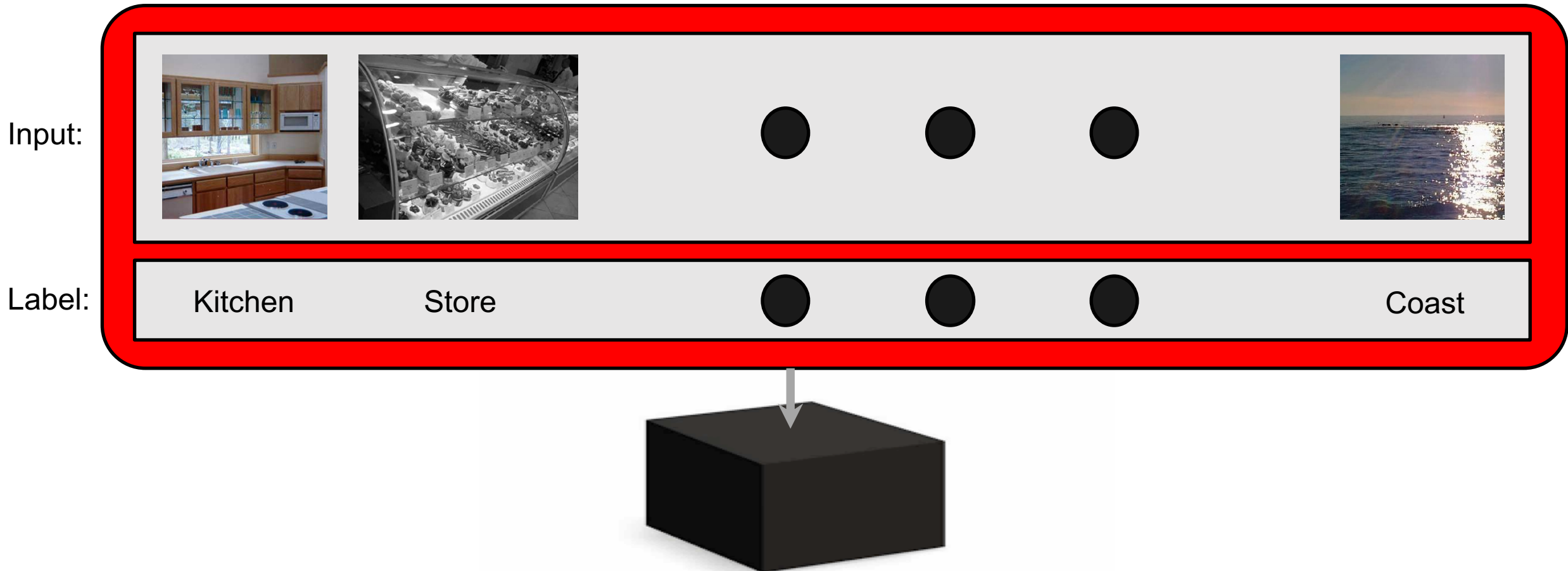
Store



Coast

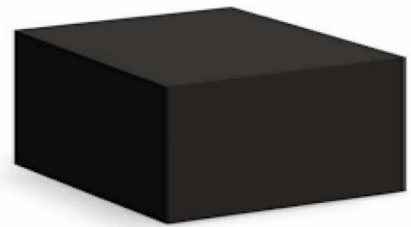
# Recall: Need Datasets to **Train** & Evaluate Algorithms

## 2. Train Prediction System

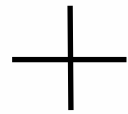


# Recall: Need Datasets to Train & Evaluate Algorithms

## 3. Apply Prediction System to Novel Images



Prediction Model



Input:



Predicted Label:

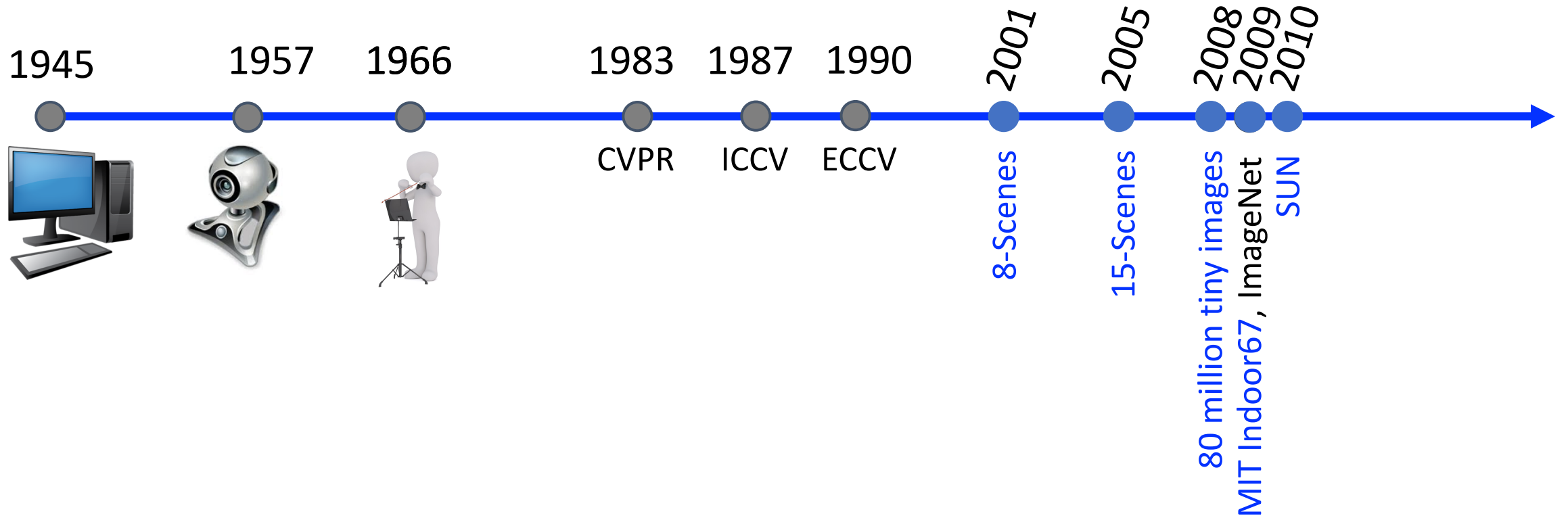
Highway

Mountain

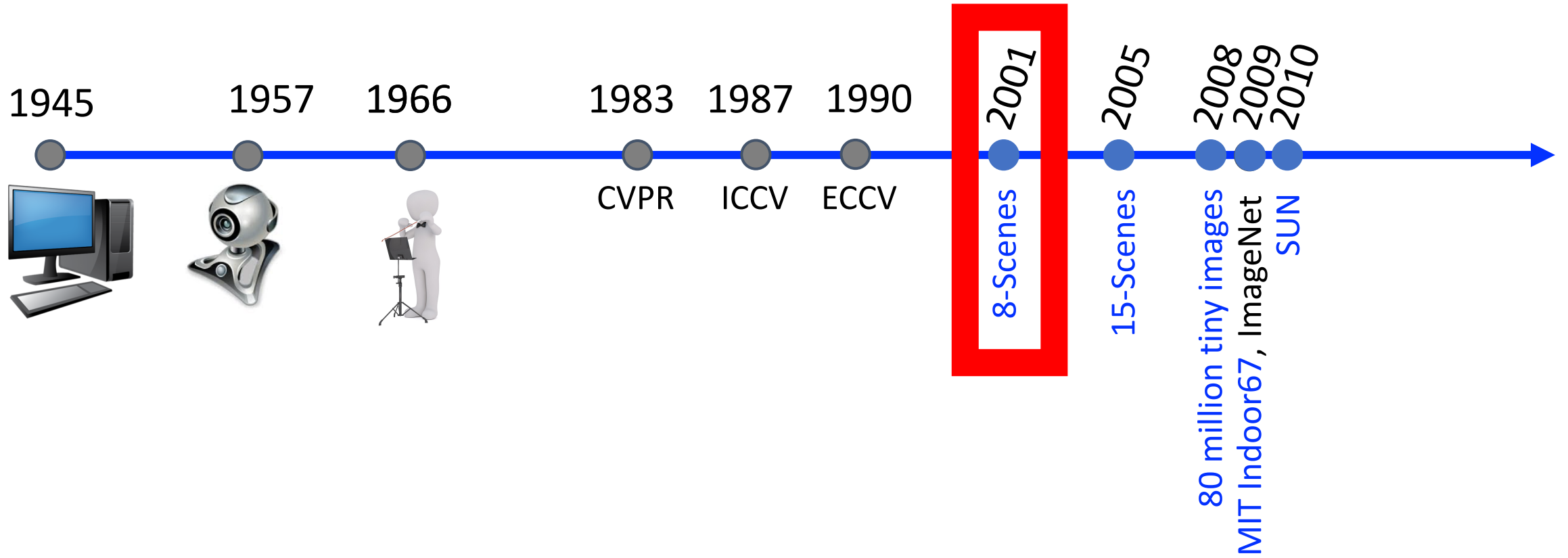


Office

# Scene Classification Datasets



# Scene Classification Datasets



# Scene Classification Datasets: 8-Scenes

**Taxonomy Source:** unclear

**Image Source:** COREL stock photo library, personal photographs, downloaded from Internet

Coast



Fields



Forests



Mountains



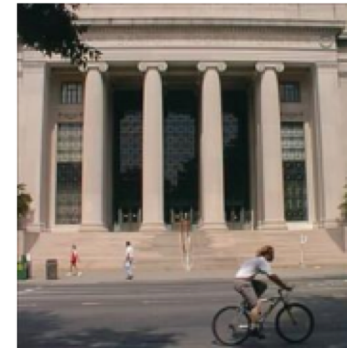
Highways



Streets



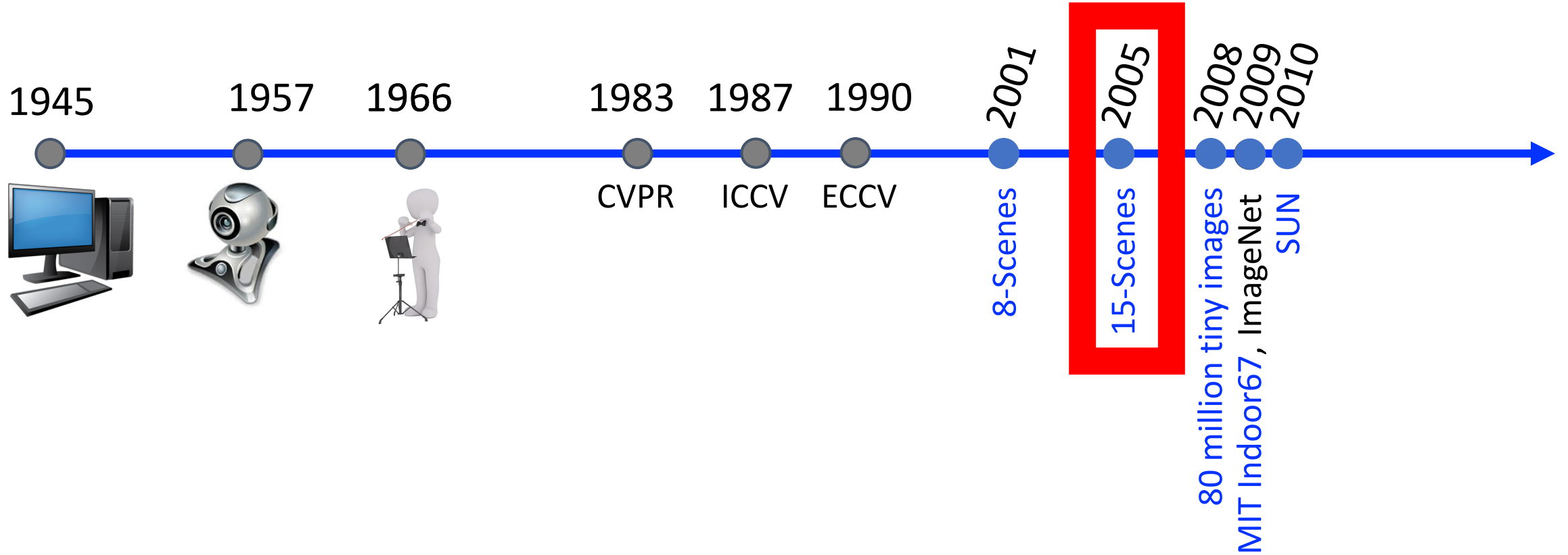
Inside City



Skyscrapers



# Scene Classification Datasets





# Scene Classification Datasets: 15-Scenes

**Taxonomy Source:** unclear

**Image Source:** COREL stock photo library, personal photographs, downloaded from Internet (contains 8-scenes dataset)



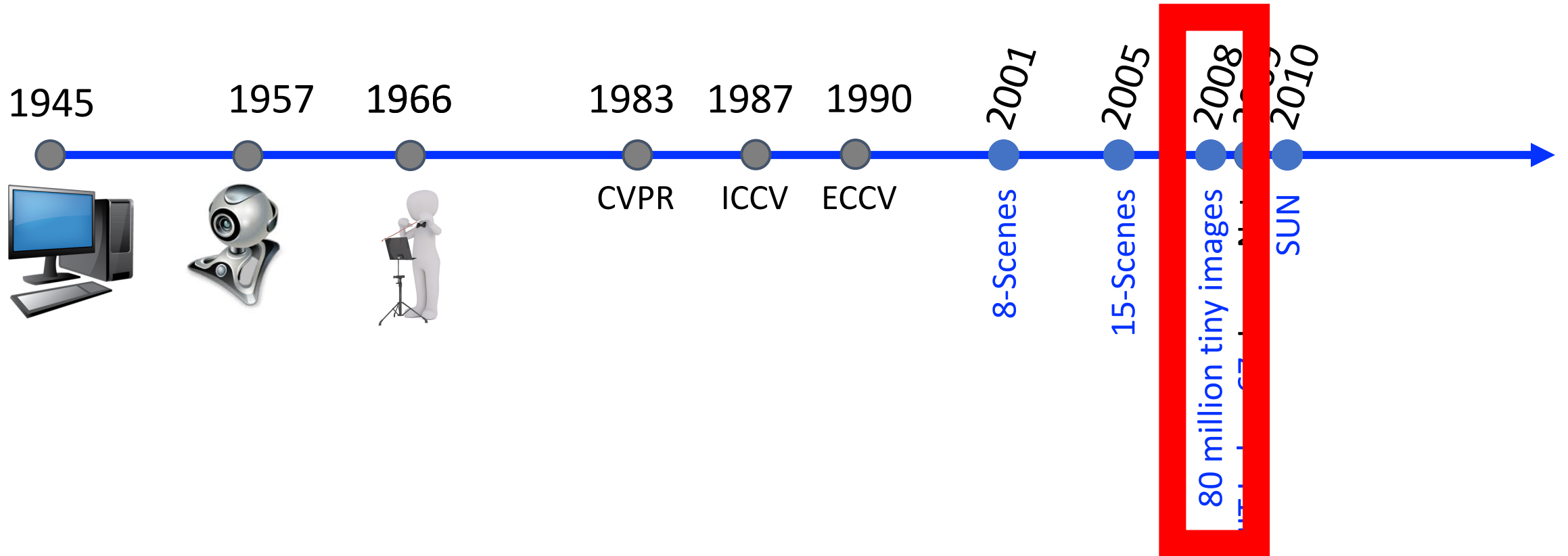
Dataset: <https://www.kaggle.com/zaiyankhan/15scene-dataset>

Fei Fei Li and Pietro Perona. A Bayesian Hierarchical Model for Learning Natural Scene Categories. CVPR 2005.

Svetlana Labeznik et al. Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories. CVPR 2005.



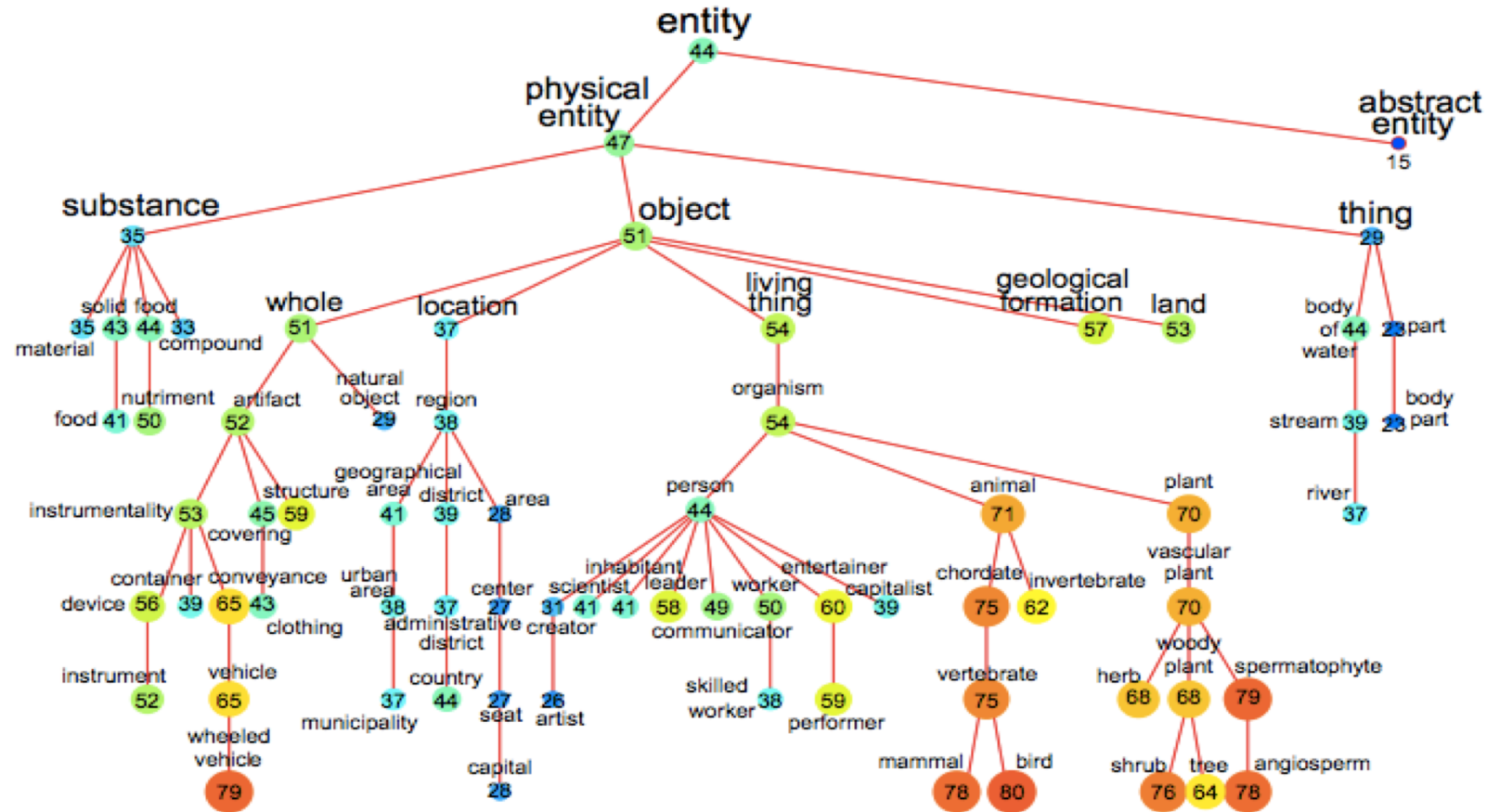
# Scene Classification Datasets



# Scene Classification Datasets: 80 Million Tiny Images

## 1. Category Selection

75,000 non-abstract nouns  
from WordNet



# Scene Classification Datasets: 80 Million Tiny Images

## 1. Category Selection

75,000 non-abstract nouns  
from WordNet



## 2. Image Collection

Images downloaded for 8  
months from 7 online  
image search engines to  
32x32 resolution



(Adapted from slides by Antonio Torralba)

# Scene Classification Datasets: 80 Million Tiny Images

## 1. Category Selection

75,000 non-abstract nouns  
from WordNet

## 2. Image Collection

Images downloaded for 8  
months from 7 online

in

32x32 resolution

**Why “tiny” images?**

# Scene Classification Datasets: 80 Million Tiny Images

256x256



## Why “tiny” images?

Idea: What resolution does a human need to recognize a scene?

### Study:

- 6 participants
- 585 color images
- Classify as 1 of 15 scene categories
- Images presented at 5 possible resolutions ( $8^2$ ,  $16^2$ ,  $32^2$ ,  $64^2$ ,  $256^2$ )

# Scene Classification Datasets: 80 Million Tiny Images

## 1. Category Selection

75,000 non-abstract nouns  
from WordNet

## 2. Image Collection

Images downloaded for 8  
months from 7 online  
image search engines to  
32x32 resolution



# Scene Classification Datasets: 80 Million Tiny Images



# Scene Classification Datasets: 80 Million Tiny Images

## 1. Category Selection

75,000 non-abstract nouns  
from WordNet

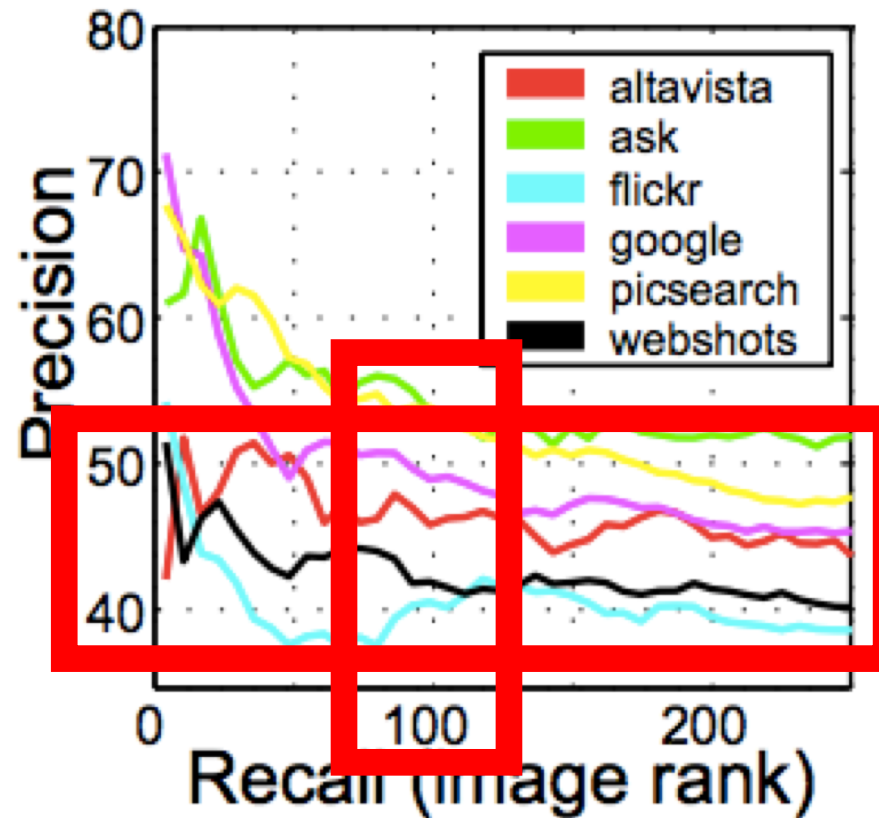
## 2. Image Collection

Images downloaded for 8  
months from 7 online  
image search engines to  
32x32 resolution

**Why no human review?**



# Scene Classification Datasets: 80 Million Tiny Images

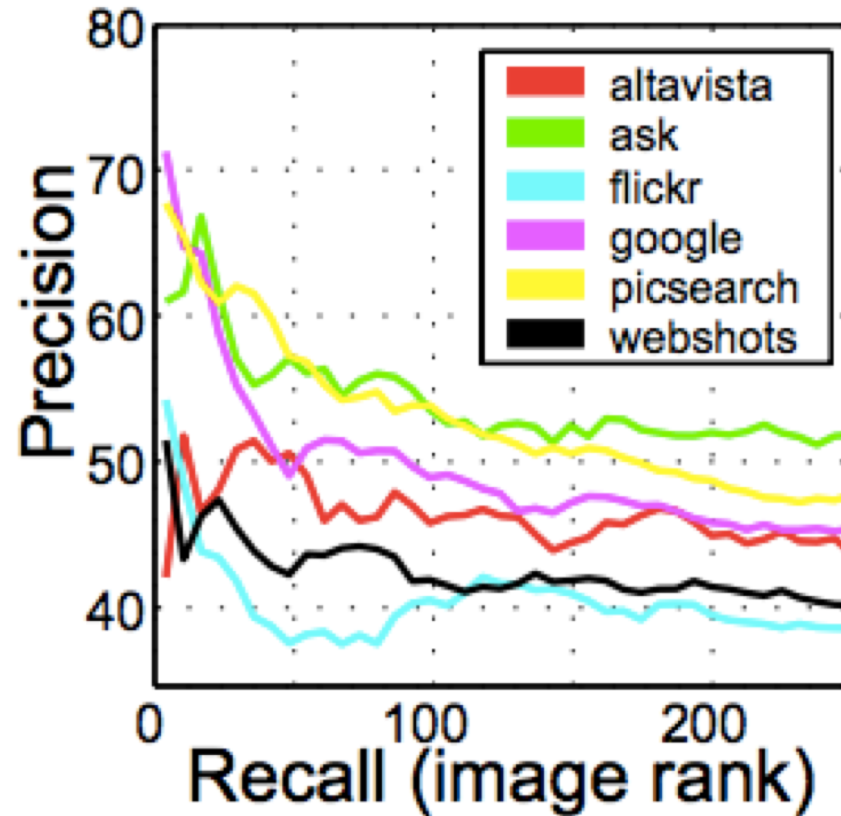


**Result of no human review?**

For each word, examined % of correct queries up to 250 words

# Scene Classification Datasets: 80 Million Tiny Images

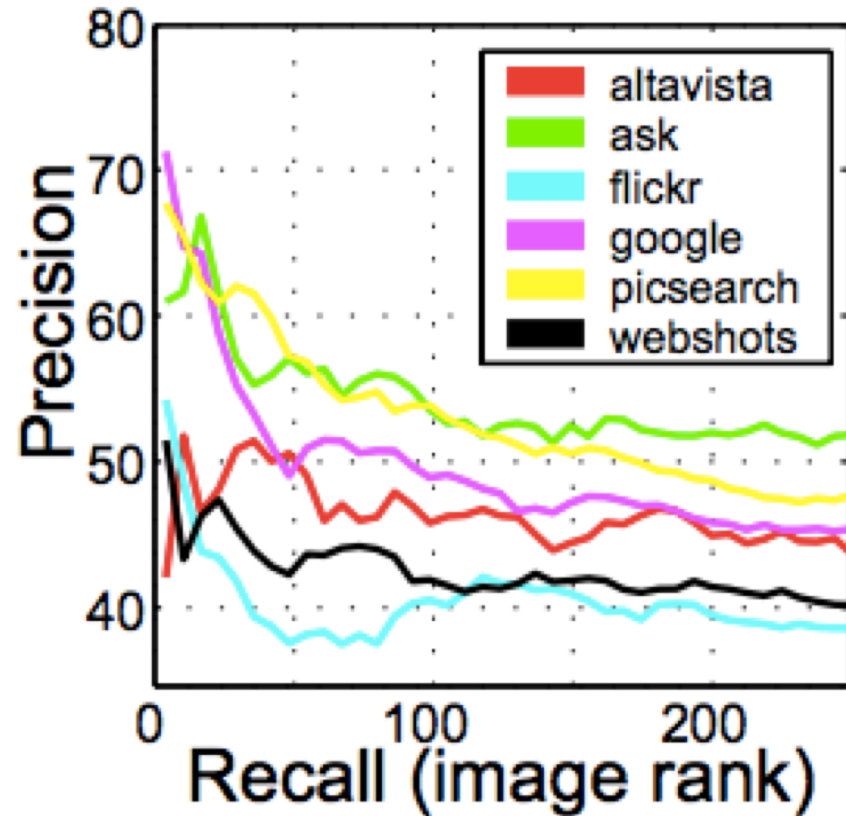
**Most accurate website?**



**Result of no human review?**

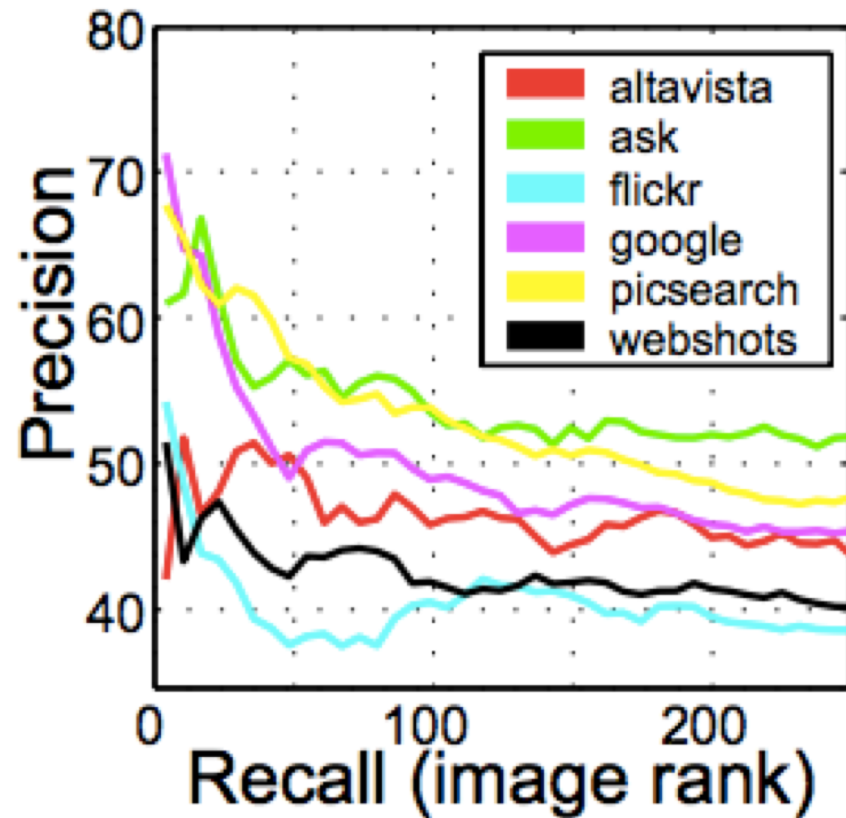
# Scene Classification Datasets: 80 Million Tiny Images

**Least accurate website?**



**Result of no human review?**

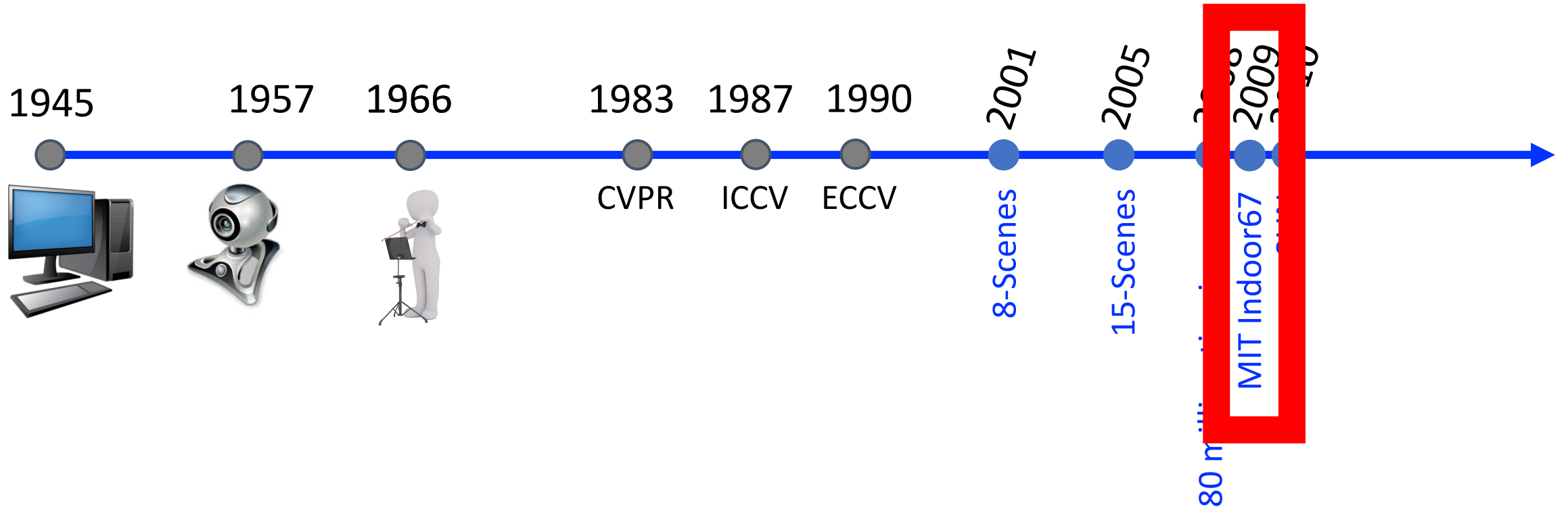
# Scene Classification Datasets: 80 Million Tiny Images



**Result of no human review?**

**Dataset is noisy!**

# Scene Classification Datasets



# MIT Indoor67

## 1. Category Selection

67 categories for 5 domains





# Scene Classification Datasets: MIT Indoor67

## 1. Category Selection

67 categories for 5 domains

## 2. Image Collection

Images downloaded from  
2 image search tools,  
1 online photo sharing sites,  
and 1 vision dataset

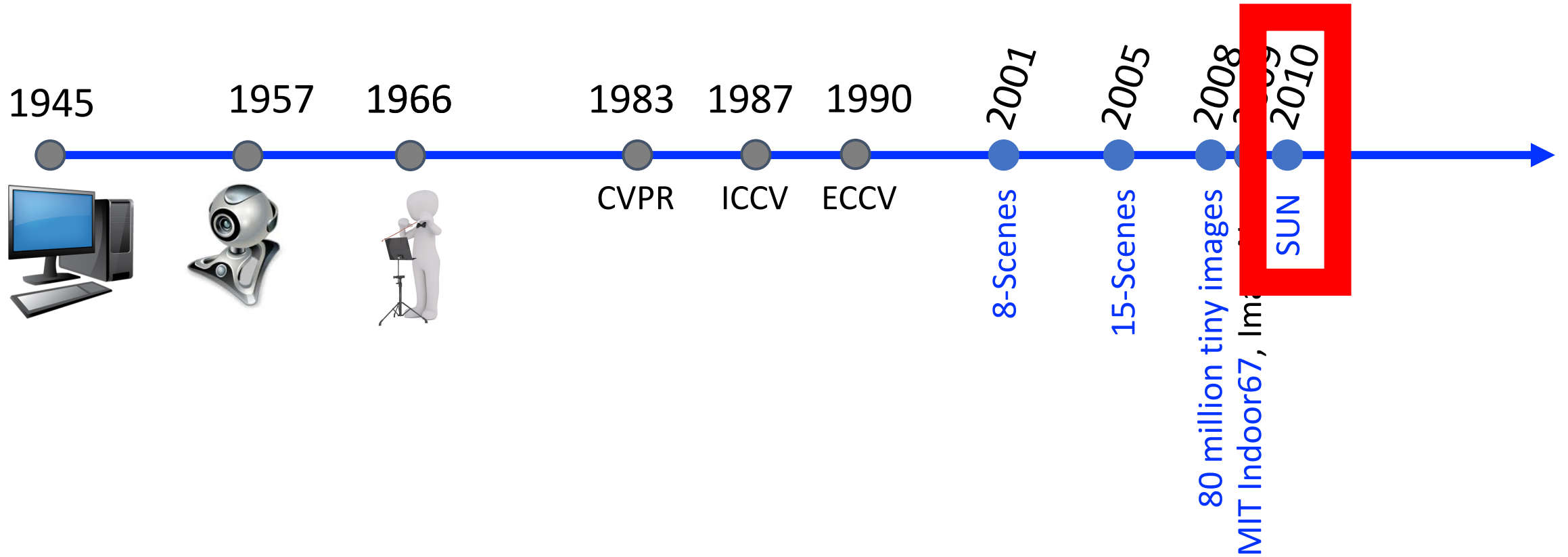
Google  
Image Search

flickr  
GAMMA

LabelMe

altavista

# Scene Classification Datasets





# Scene Classification Datasets: SUN

## 1. Category Selection

- From 70,000 categories in “Tiny Images” (WordNet), chose 908 categories describing scenes, places, and environments, excluding:

- 1) names of specific places (e.g., New York)
- 2) non-navigable scenes
- 3) “mature” data

- Extra categories; e.g., mission, jewelry store



# Scene Classification Datasets: SUN

## 1. Category Selection

- From 70,000 categories in “Tiny Images” (WordNet), chose 908 categories describing scenes, places, and environments, excluding:
  - 1) names of specific places (e.g., New York)
  - 2) non-navigable scenes
  - 3) “mature” data
- Extra categories; e.g., mission, jewelry store

## Category Validation Experiment:

- 7 subjects wrote every 30 minutes the name of the scene category for their location
- All resulting 52 categories were in SUN

# Scene Classification Datasets: SUN

## 1. Category Selection

- From 70,000 categories in “Tiny Images” (WordNet), chose 908 categories describing scenes, places, and environments, excluding:
  - 1) names of specific places (e.g., New York)
  - 2) non-navigable scenes
  - 3) “mature” data
- Extra categories; e.g., mission, jewelry store

## 2. Image Collection

- Downloaded from search engines
- Automatically discarded images that are:
  - 1) not color
  - 2) less than 200x200
  - 3) very blurry or noisy
  - 4) aerial views
  - 5) duplicates



(Adapted from slides by Antonio Torralba)

# Scene Classification Datasets: SUN

## 1. Category Selection

- From 70,000 categories in “Tiny Images” (WordNet), chose 908 categories describing scenes, places, and environments, excluding:
  - 1) names of specific places (e.g., New York)
  - 2) non-navigable scenes
  - 3) “mature” data
- Extra categories; e.g., mission, jewelry store

## 2. Image Collection

- Downloaded from search engines
- Automatically discarded images that are:
  - 1) not color
  - 2) less than 200x200
  - 3) very blurry or noisy
  - 4) aerial views
  - 5) duplicates

## 3. Human Verification

- 9 in-house people reviewed & discarded irrelevant images
- Result is 130,519 images spanning 397 categories with >99 images per category







# Scene Classification Datasets: SUN

## 3. Human Verification



- 9 in-house people reviewed & discarded irrelevant images
- Result is 130,519 images spanning 397 categories with >99 images per category

What are poorly represented categories?

# Scene Classification Datasets: SUN

## Dataset Validation Experiment: Crowdsourcing

### 3. Human Verification

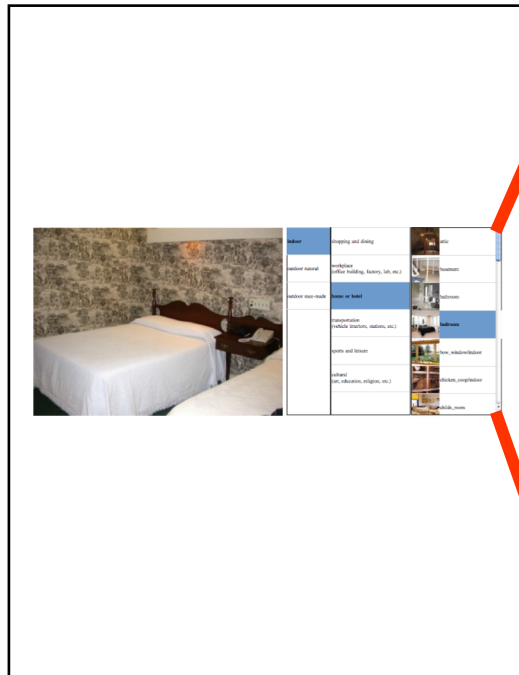
- 9 in-house people reviewed & discarded irrelevant images
- Result is 130,519 images spanning 397 categories with >99 images per category



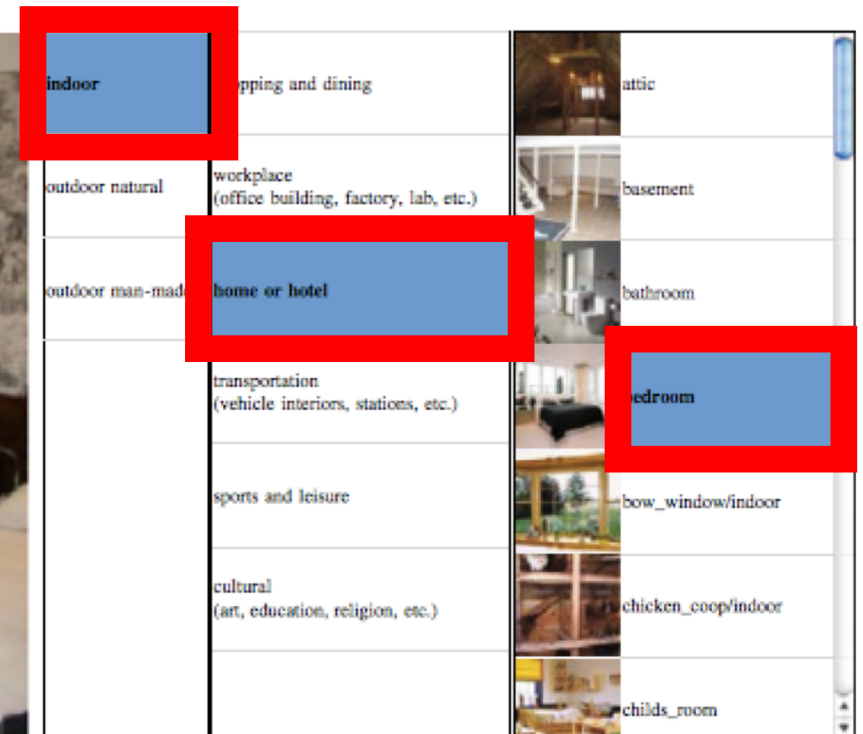
# Scene Classification Datasets: SUN

## Dataset Validation Experiment: Crowdsourcing

### 1. Task Design

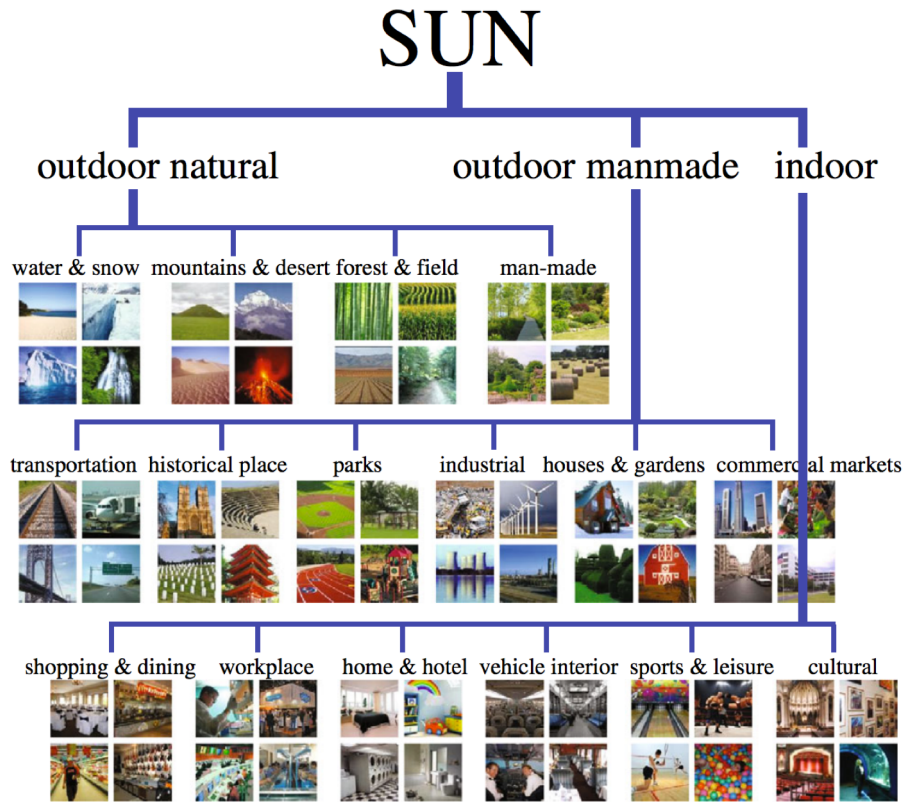


User interface: Forced-choice 3 level hierarchy

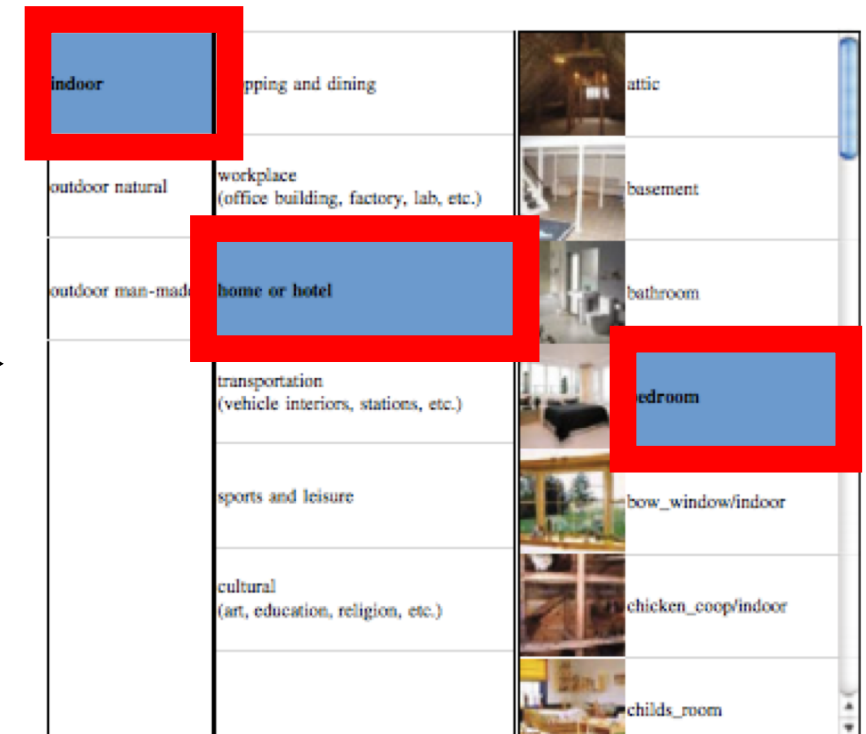


# Scene Classification Datasets: SUN

## Dataset Validation Experiment: Crowdsourcing



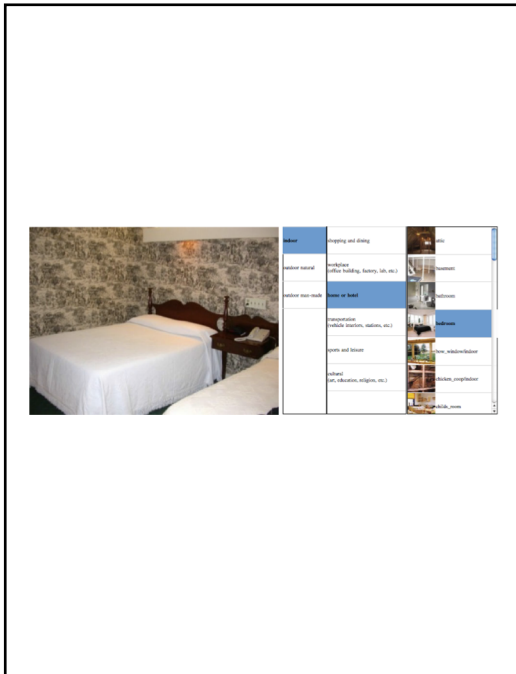
User interface: Forced-choice 3 level hierarchy



# Scene Classification Datasets: SUN

## Dataset Validation Experiment: Crowdsourcing

### 1. Task Design



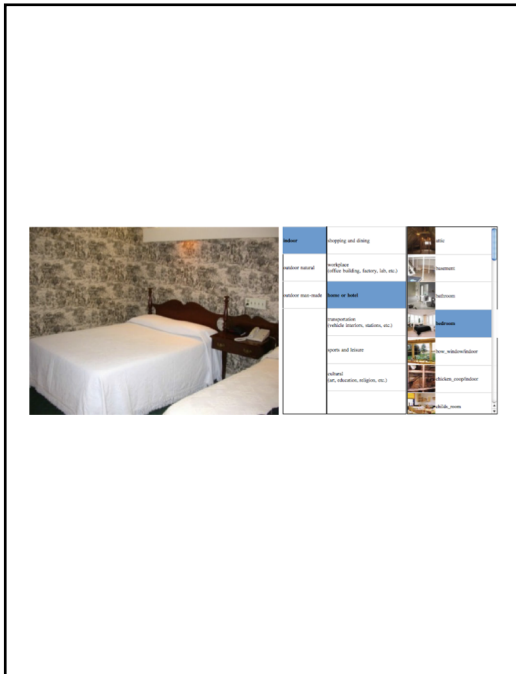
### 2. Crowdsourcing Properties

- Crowdsourcing Platform: AMT
- Worker Requirements:
  - 1) U.S. workers
  - 2) Completed at least 100 HITs
  - 3) > 95% accuracy on first level hierarchy voting
- Worker Pay: \$0.03 per HIT (61 s avg)

# Scene Classification Datasets: SUN

## Dataset Validation Experiment: Crowdsourcing

### 1. Task Design



### 2. Crowdsourcing Properties

- Crowdsourcing Platform: AMT
- Worker Requirements:
  - 1) U.S. workers
  - 2) Completed at least 100 HITs
  - 3) > 95% accuracy on first level hierarchy voting
- Worker Pay: \$0.03 per HIT (61 s avg)

### 3. Crowdsourcing Validation

- Test Data
- 7,940 jobs: 20 examples per category

- Crowd Performance
- 58.6% leaf accuracy
- Author Performance
- 68.5% leaf accuracy

# Scene Classification Datasets: SUN

## 1. Category Selection

- From 70,000 categories in “Tiny Images” (WordNet), chose 908 categories describing scenes, places, and environments, excluding:
  - 1) names of specific places (e.g., New York)
  - 2) non-navigable scenes
  - 3) “mature” data
- Extra categories; e.g., mission, jewelry store

Ran Validation Experiment

## 2. Image Collection

- Downloaded from search engines
- Automatically discarded images that are:
  - 1) not color
  - 2) less than 200x200
  - 3) very blurry or noisy
  - 4) aerial views
  - 5) duplicates

## 3. Human Verification

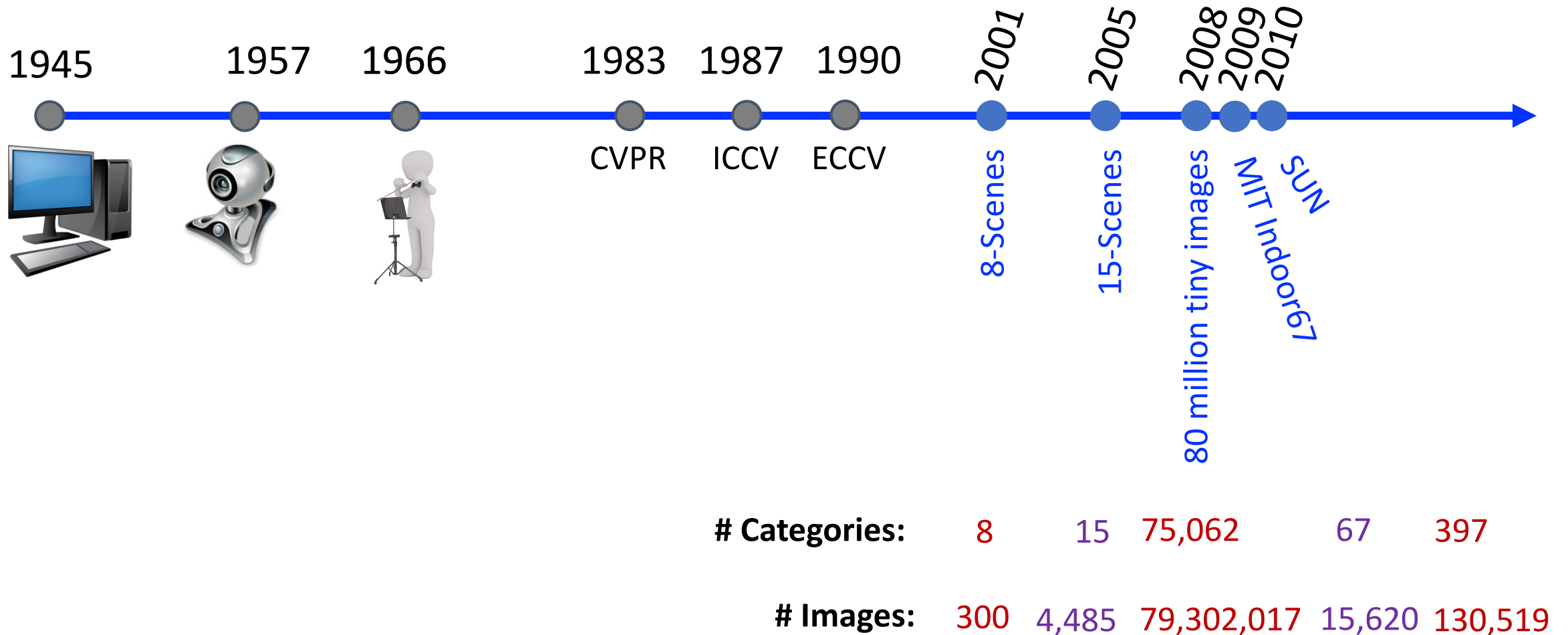
- 9 in-house people reviewed & discarded irrelevant images
- Result is 130,519 images spanning 397 categories with >99 images per category

Ran Validation Experiment

# Scene Classification Datasets: SUN Image Browser

Demo: <https://groups.csail.mit.edu/vision/SUN/>

# Scene Classification Datasets



# Scene Classification Datasets: Summary

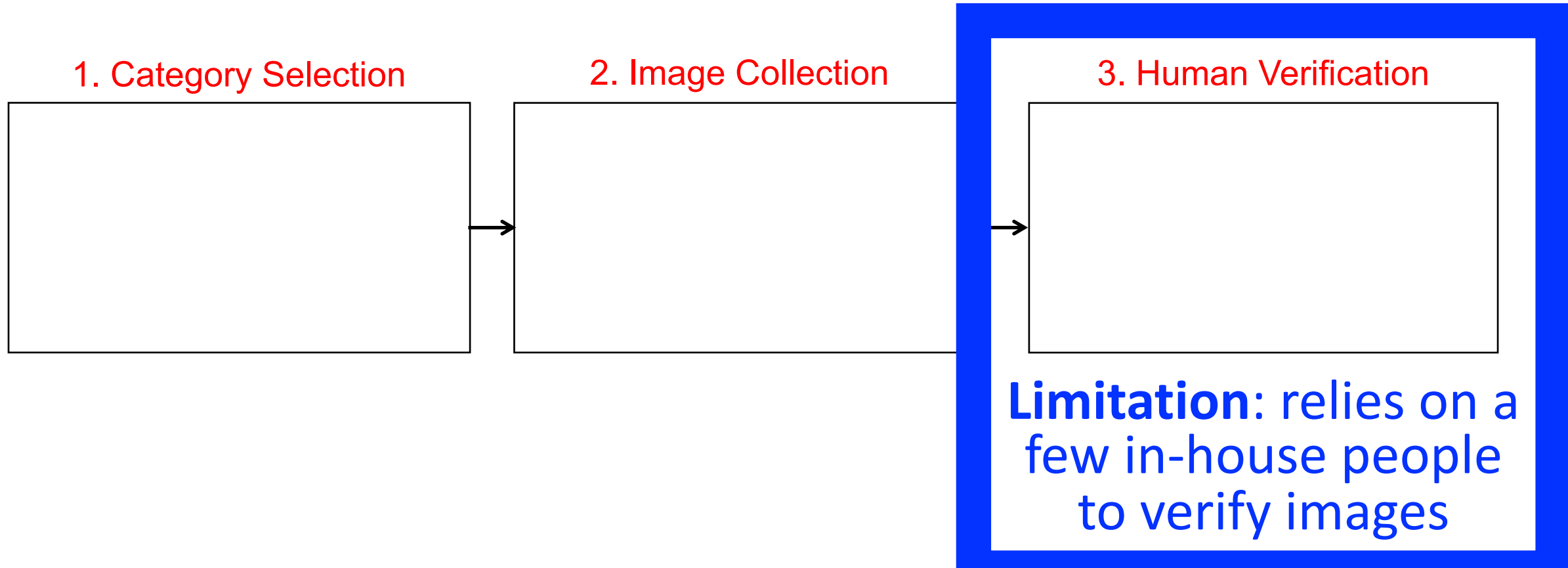
- Key steps in creating dataset:





# Scene Classification Datasets: Summary

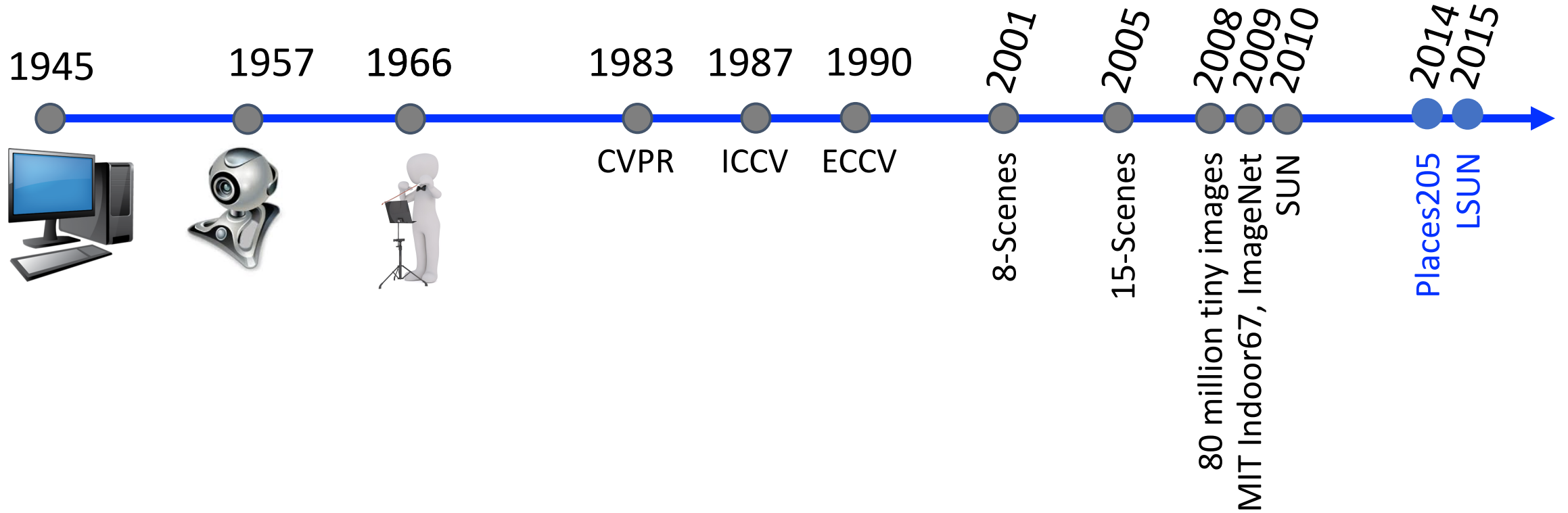
- Key steps in creating dataset:



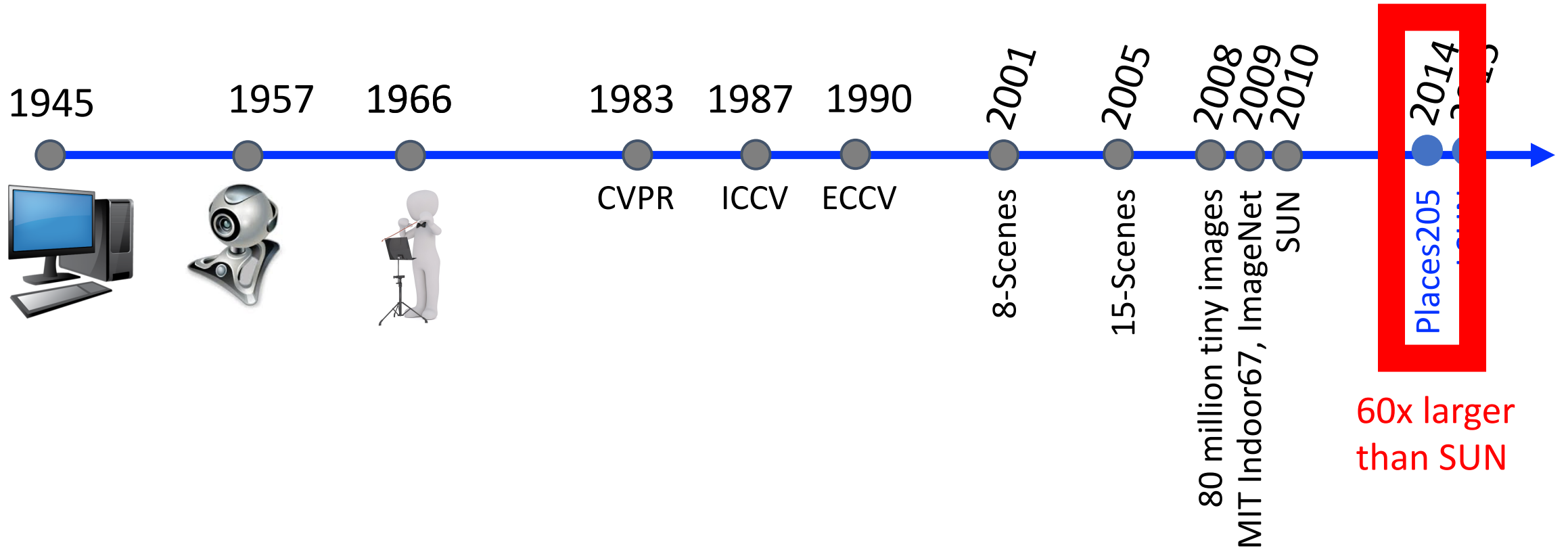
# Today's Topics

- Scene classification applications
- Scene classification datasets: key steps in creating them
- Scene classification datasets: scaling up with *crowdsourcing* and *challenges*
- Class discussion (chosen by YOU 😊)
- Lab: Javascript

# Scene Classification Datasets



# Scene Classification Datasets

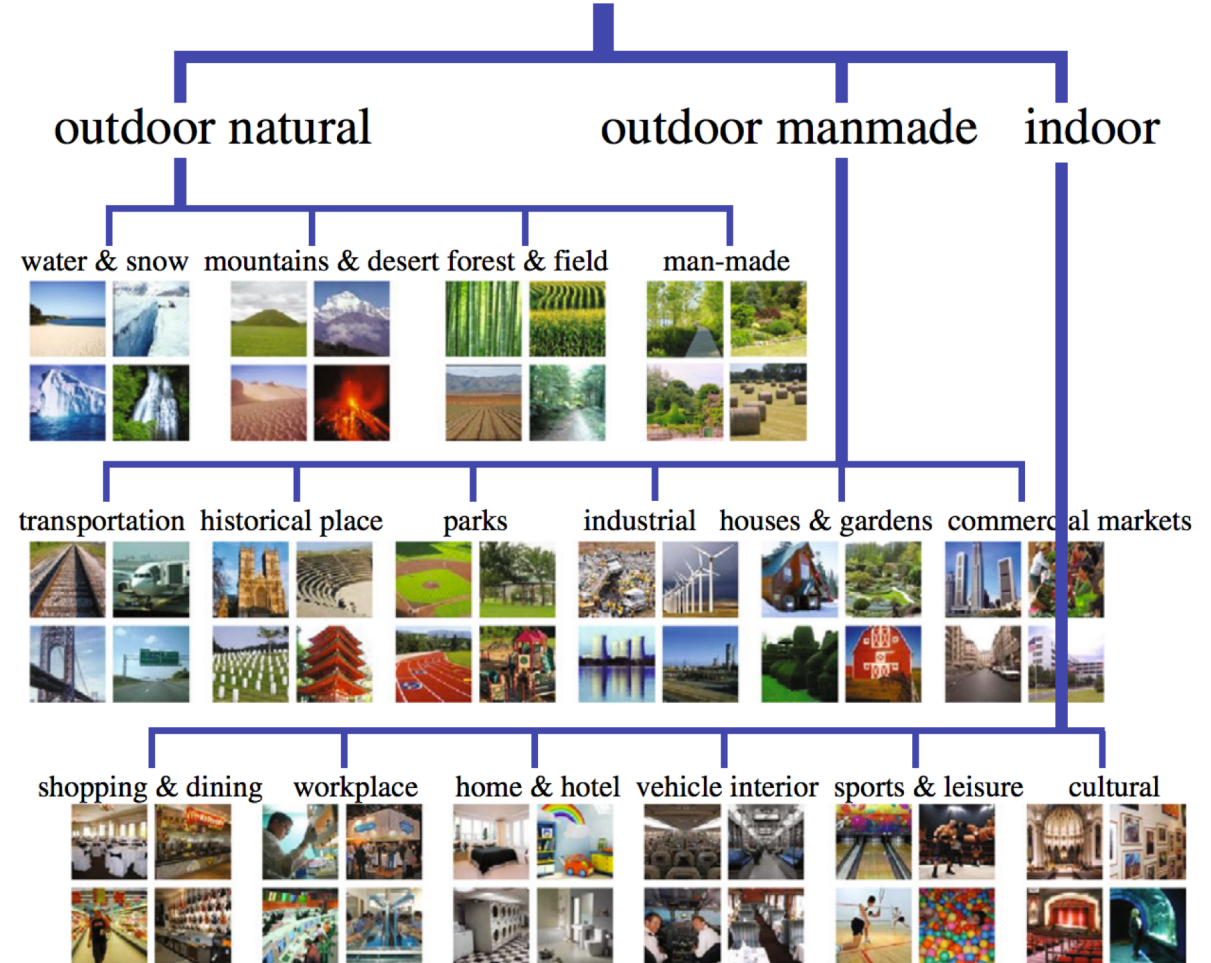


# Scene Classification Datasets: Places205

## SUN

### 1. Category Selection

Same taxonomy as SUN



# Scene Classification Datasets: Places205

## 1. Category Selection

Same taxonomy as SUN

## 2. Image Collection

- Downloaded images from three search engines; query terms were 696 common adjectives (messy, spare, sunny, desolate, etc) with each scene category
- Automatically discarded images that are:
  - 1) not color
  - 2) less than 200x200



# Scene Classification Datasets: Places205

## 1. Category Selection

Same taxonomy as SUN

## 2. Image Collection

- Downloaded images from three search engines; query terms were 696 common adjectives (messy, spare, sunny, desolate, etc) with each scene category
- Automatically discarded images that are:
  - 1) not color
  - 2) less than 200x200

## 3. Human Verification

- AMT crowd workers identified (ir)relevant images for batches of 750 images
- Result is 7,076,580 images spanning 476 categories

# Places205

## User interface: Instructions

### 1. Task Design

#### Instructions:



#### Interface:



### Examples

**Start** **Is this a cliff scene?**

**Definition:** high, steep or overhanging face of rock.

**Task**

For each of the **810** images, answer yes or no to the above question. Only answer **Yes** to **real photos**. Always answer **No** to **cartoon, drawing, CG rendering**, or real photos with a **large text overlay** on the photo. Here are some examples:

No Single Object No Text Overlay No Drawing No Screenshot No Graphics No Bad Photo

Not Only Logo No Magazine/Newspaper No No Yes Yes

### Instructions



# Places205

User interface: Task

# Tasks left

## 1. Task Design

Instructions:



Interface:



Instruction **Is this a cliff scene?** Submit (790 images left)

Definition: a high, steep or overhanging face of rock

Current Task: press a key on keyboard

Completed Tasks

No No No



Yes



Next Tasks

No



# Scene Classification Datasets: Places205

## 1. Task Design

**Instructions:**



**Interface:**



## 2. Crowdsourcing Platform



# Scene Classification Datasets: Places205

## 1. Task Design

**Instructions:**



**Interface:**



## 2. Crowdsourcing Platform



## 3. Quality Control

- Run images through crowd twice with default "yes" and then default "no answer"
- "Honeypot"
  - labelled at least 90% on control set correctly, where it includes 30 known positive and negative labelled images per "HIT"

# Scene Classification Datasets: Places205 Summary

## 1. Category Selection

Same taxonomy as SUN

## 2. Image Collection

- Downloaded images from three search engines; query terms were 696 common adjectives (messy, spare, sunny, desolate, etc) with each scene category
- Automatically discarded images that are:
  - 1) not color
  - 2) less than 200x200

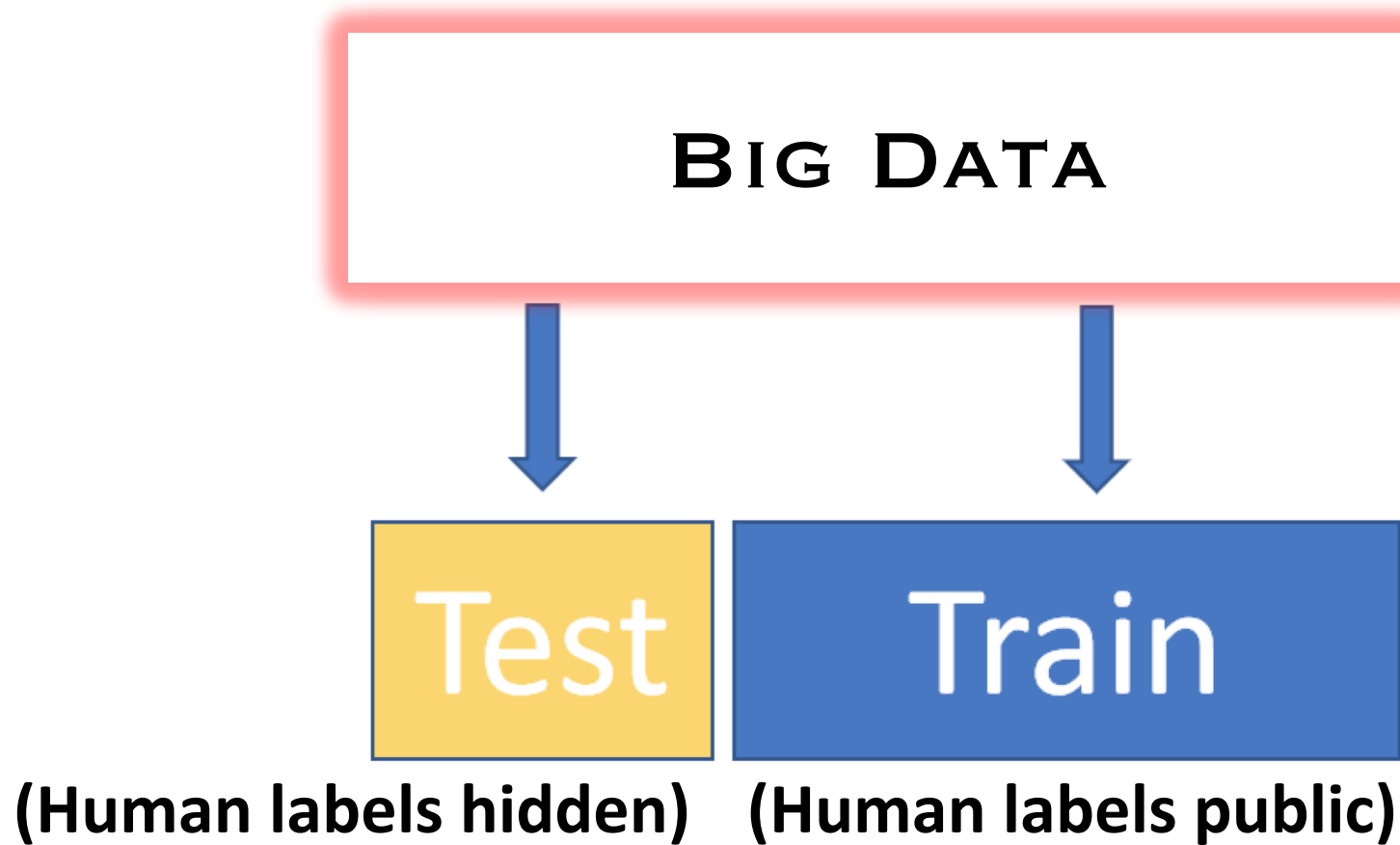
## 3. Human Verification

- AMT crowd workers identified (ir)relevant images for batches of 750 images
- Result is 7,076,580 images spanning 476 categories

# Scene Classification: Places Challenge

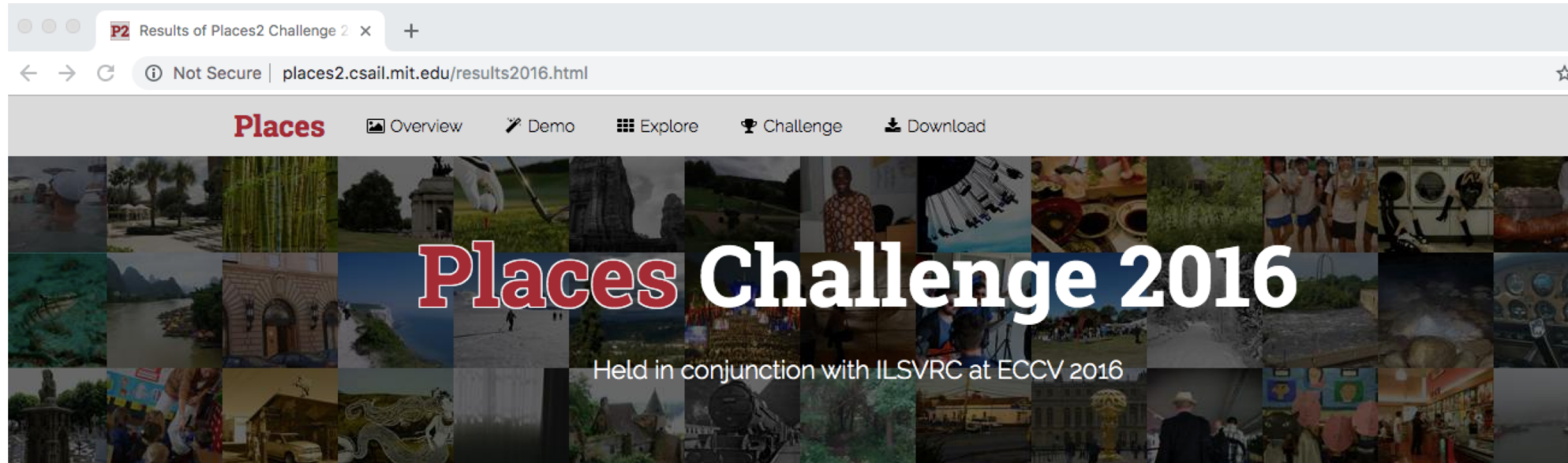


# Scene Classification: Places Challenge (Recall)



**Winner: highest scoring method on the hidden test set**

# Scene Classification: Places Challenge



## Results

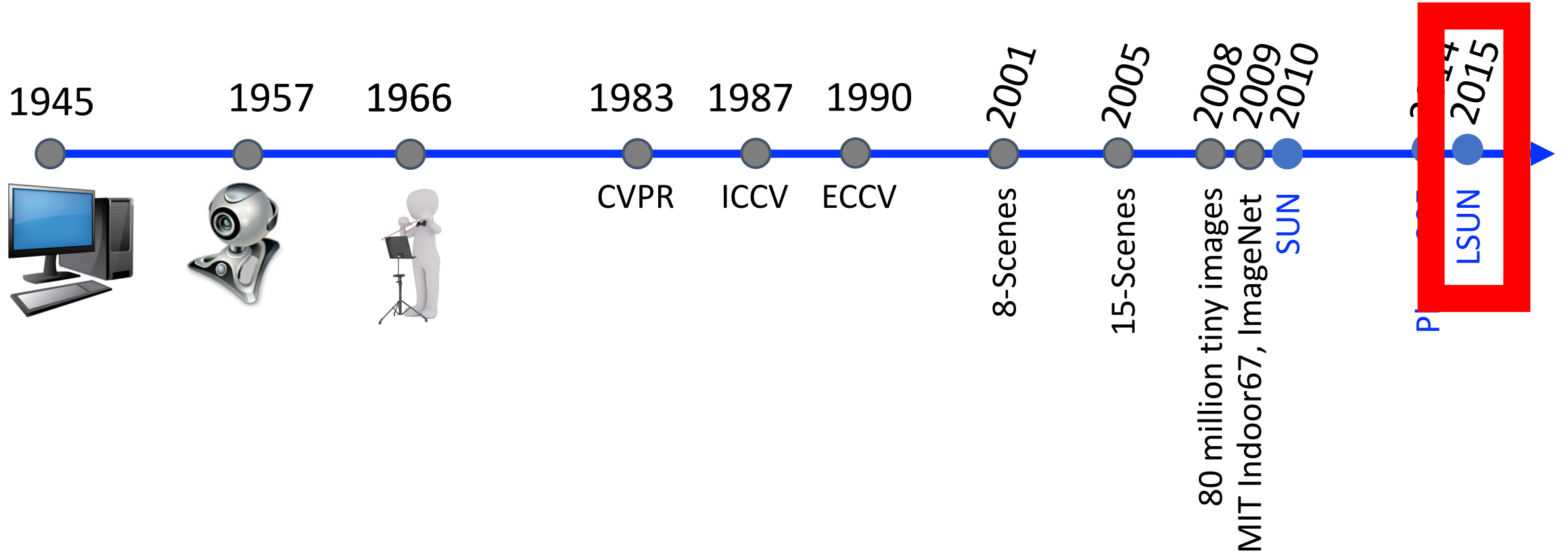
### Contents:

- Summary: There are totally **92** valid submissions from **27** teams. Hikvision won the 1st place with **0.0901** top-5 error, MW won the 2nd place with **0.1019** top-5 error, and Trimps-Soushen won the 3rd place with **0.1030** top-5 error. Congratulations to all the teams. See below for the leaderboard and the team information.
- Rule: Each teams can only use the provided data in Places2 Challenge 2016 to train their networks. Standard pre-trained CNN models trained on Imagenet-1.2million and previous Places are allowed to use. Each teams can submit at most 5 prediction results. Ranks are based on the top-5 classification error of each submission.
- [Scene classification with provided training data](#)
- [Team information](#)

Demo: <http://places2.csail.mit.edu/results2016.html>



# Scene Classification Datasets





# Scene Classification Datasets: LSUN

## 1. Category Selection

10 scene categories from  
SUN

# Scene Classification Datasets: LSUN

## 1. Category Selection

10 scene categories from SUN

## 2. Image Collection

- Downloaded images from Google Images; query terms were 696 common adjectives (messy, spare, sunny, desolate, etc) with each scene category for all 3-day time spans since 2009
- Automatically discarded images that are  $< 256 \times 256$



# Scene Classification Datasets: LSUN

## 1. Category Selection

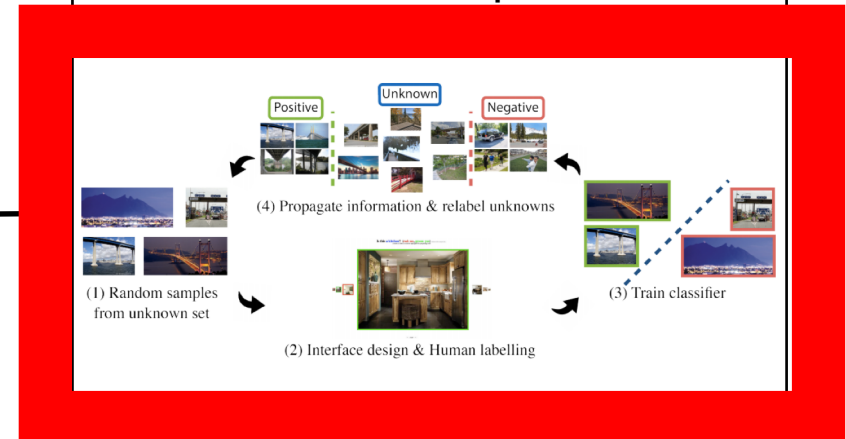
10 scene categories from SUN

## 2. Image Collection

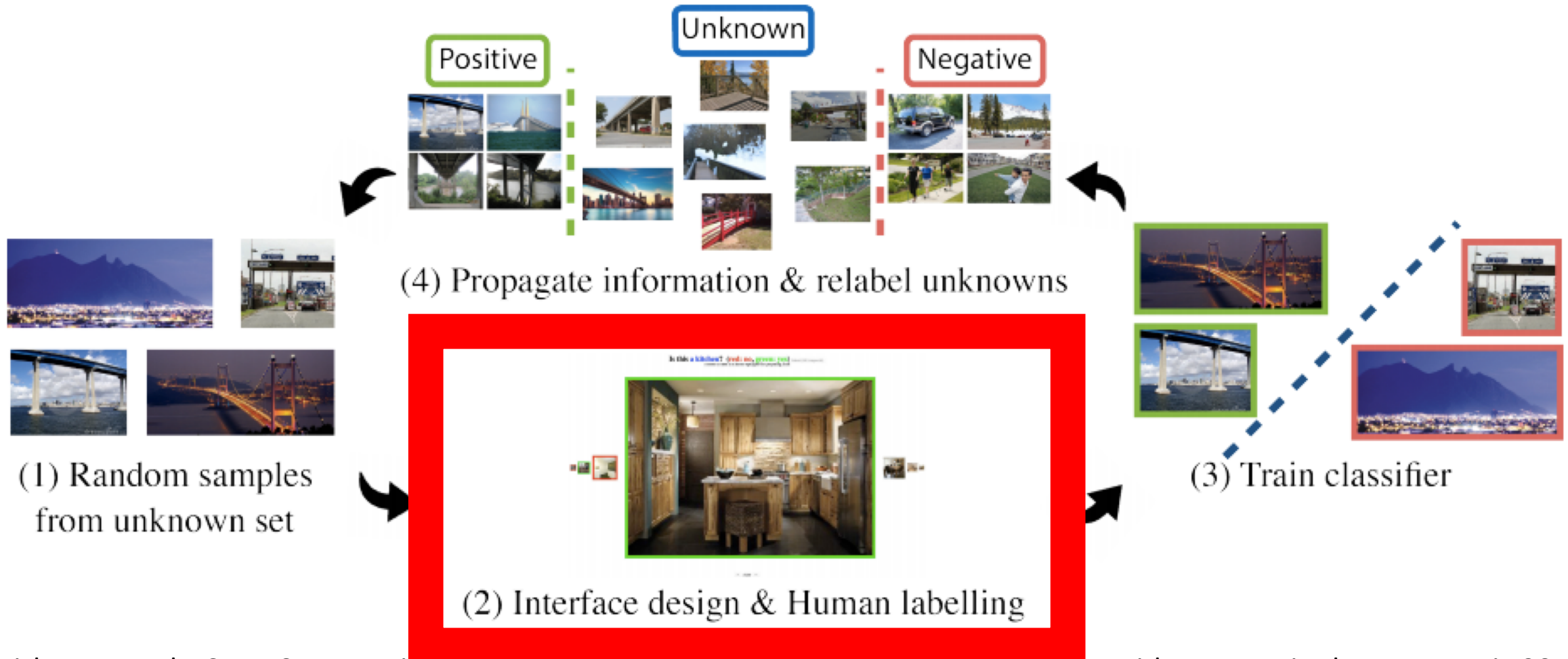
- Downloaded images from Google Images; query terms were 696 common adjectives (messy, spare, sunny, desolate, etc) with each scene category for all 3-day time spans since 2009
- Automatically discarded images that are  $< 256 \times 256$

## 3. Label Verification

- Human in the loop



# Scene Classification Datasets: LSUN Label Verification with Humans in the Loop



# LSUN

## 1. Task Design

Instructions:

- For consistency, include examples commonly leading to (experimentally observed) crowd disagreement; e.g., occlusion
- General categories: e.g., cartoons

Interface:

## Recall User Interface for Creating “Places”

The screenshot shows a user interface for a task. At the top, there is an "Instruction" box containing the question "Is this a cliff scene?" and a "Submit (790 images left)" button. Below the question is a "Definition: a high, steep or overhanging face of rock." and a "Current Task: press 'space' to toggle answer & arrows to move to previous/next" instruction. The main area displays a large image of a river flowing through a rocky, forested landscape. The word "Yes" is written in green above the image. To the left of the main image is a "Completed Tasks" section showing three smaller images, each with the word "No" above it. To the right is a "Next Tasks" section showing three smaller images, each with the word "No" above it. The main image and the "Completed Tasks" section are highlighted with a red border, and the "Next Tasks" section is also highlighted with a red border.

# Scene Classification Datasets: LSUN Label Verification with Humans in the Loop

## 1. Task Design

Instructions:

- For consistency, include examples commonly leading to (experimentally observed) crowd disagreement; e.g., occlusion
- General categories: e.g., cartoons

Interface:

## 2. Crowdsourcing Platform

amazon mechanical turk™  
Artificial Artificial Intelligence

# Scene Classification Datasets: LSUN Label Verification with Humans in the Loop

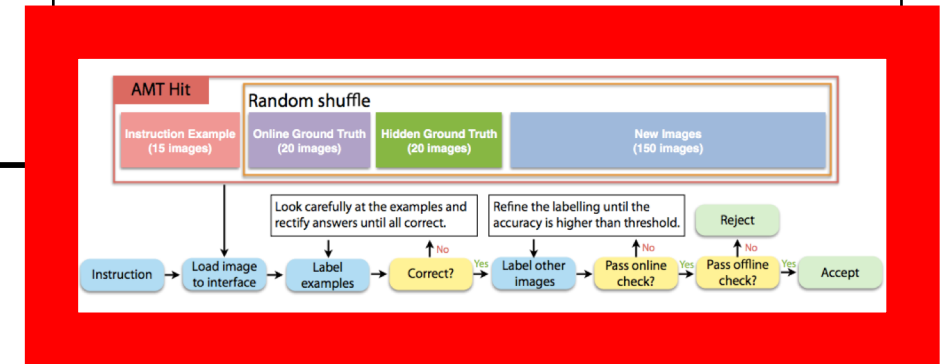
## 1. Task Design

Instructions:  
- For consistency, include examples commonly leading to (experimentally observed) crowd disagreement; e.g., occlusion  
- General categories: e.g., cartoons  
Interface:

## 2. Crowdsourcing Platform

amazon mechanical turk™  
Artificial Artificial Intelligence

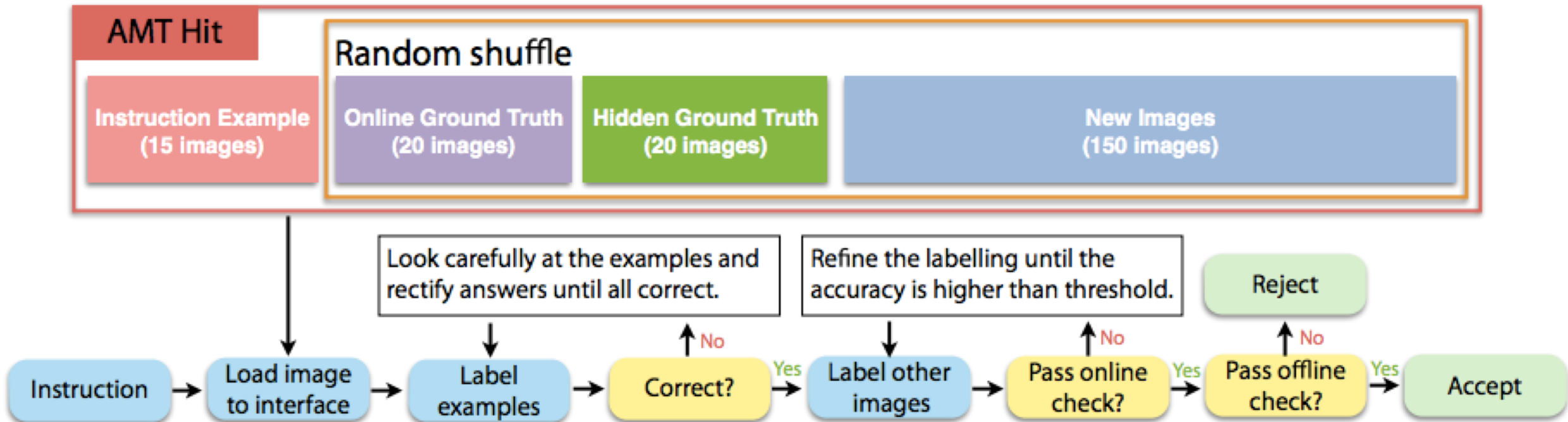
## 3. Quality Control





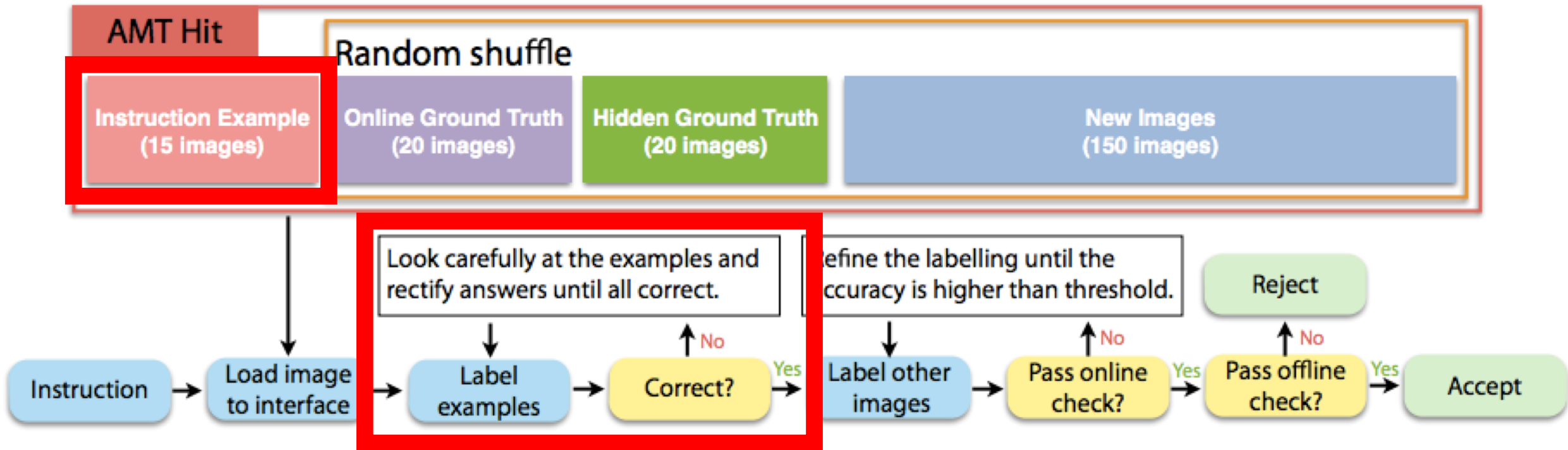
# Scene Classification Datasets: LSUN Label Verification with Humans in the Loop

Crowdsourcing Quality Control:



# Scene Classification Datasets: LSUN Label Verification with Humans in the Loop

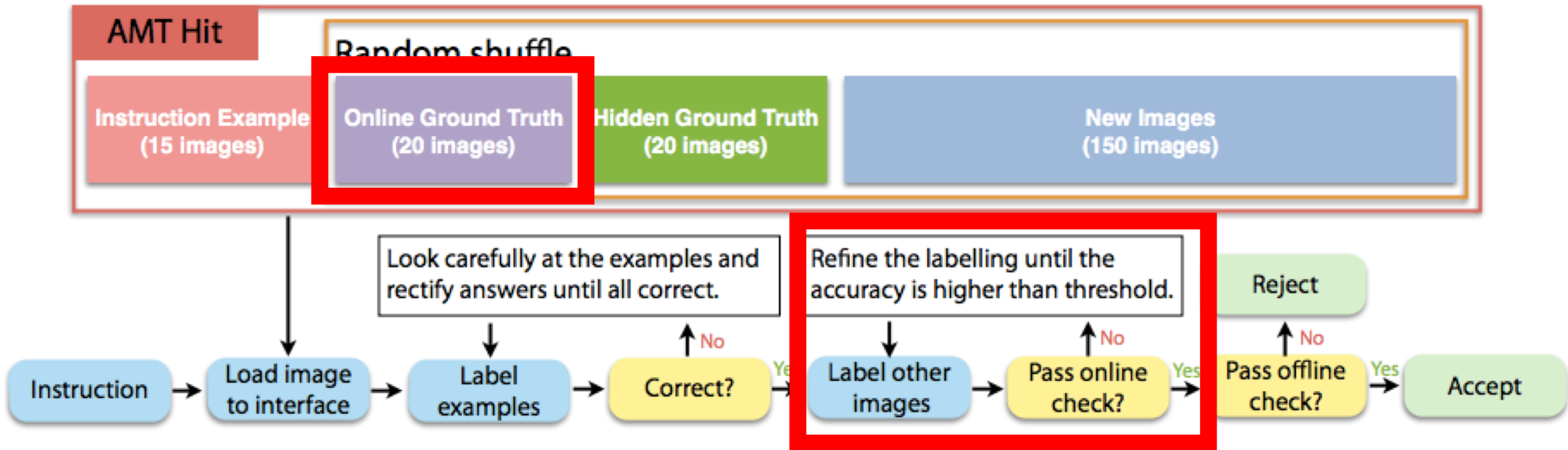
Crowdsourcing Quality Control:



Seed initial images with typical categories and common mistakes; pop-up box (i.e., tutorial) requiring mistake fixed

# Scene Classification Datasets: LSUN Label Verification with Humans in the Loop

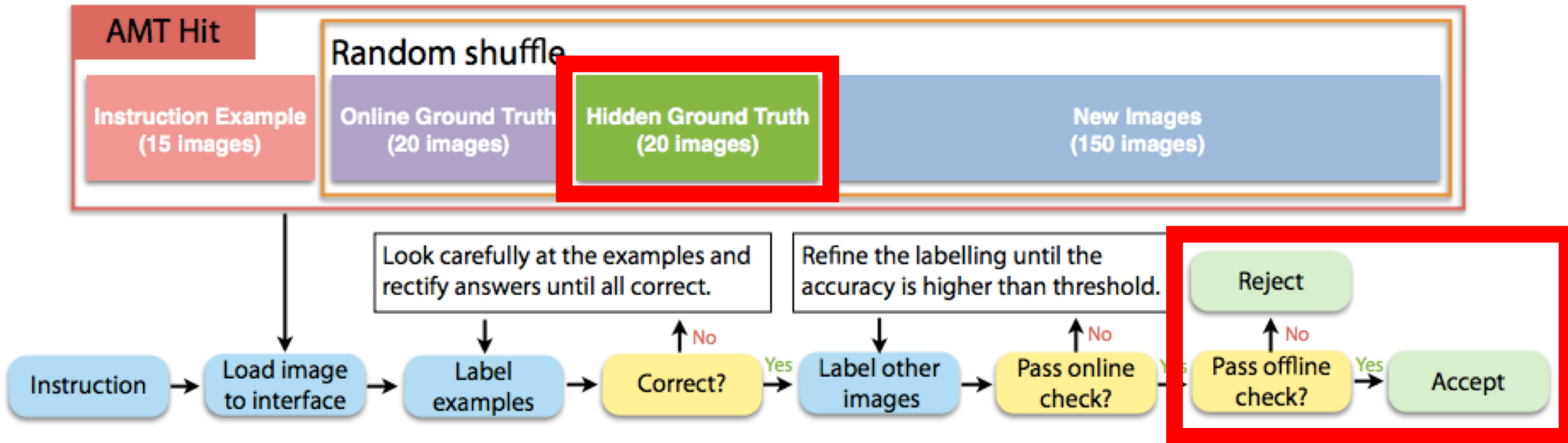
Crowdsourcing Quality Control:



Crowd worker can submit results when  $>90\%$  of “honeypot” examples are correct

# Scene Classification Datasets: LSUN Label Verification with Humans in the Loop

Crowdsourcing Quality Control:



Accept crowd worker's results when >85% of "honeypot" examples are correct

# Scene Classification Datasets: LSUN Summary

## 1. Category Selection

10 scene categories from SUN

## 2. Image Collection

- Downloaded images from Google Images; query terms were 696 common adjectives (messy, spare, sunny, desolate, etc) with each scene category for all 3-day time spans since 2009

- Automatically discarded images that are  $< 256 \times 256$

## 3. Label Verification

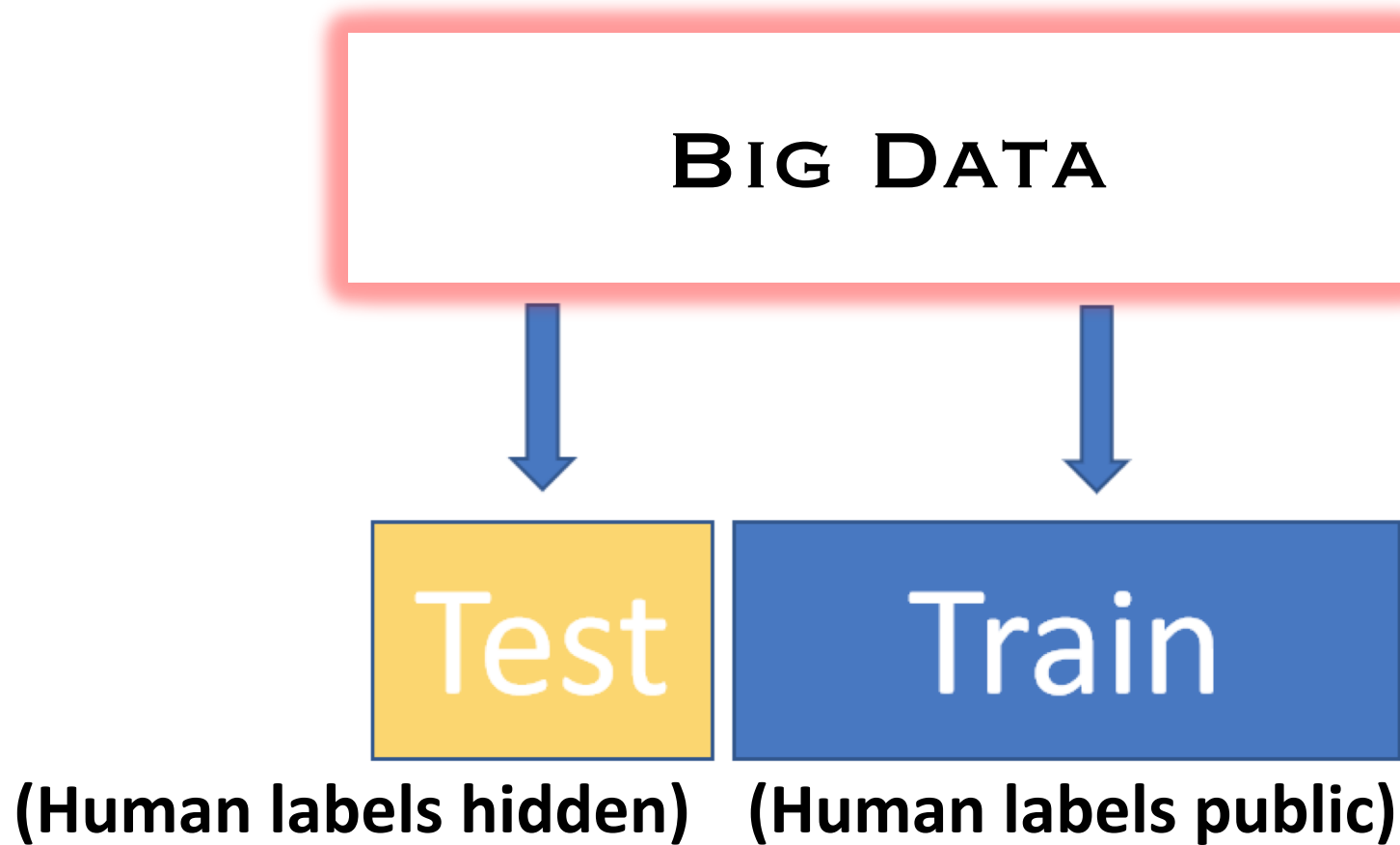
- Human in the loop



# Scene Classification Datasets: LSUN Challenge



# Scene Classification Datasets: LSUN Challenge



**Winner: highest scoring method on the hidden test set**



# Scene Classification Datasets: LSUN Challenge

jointscene.csail.mit.edu



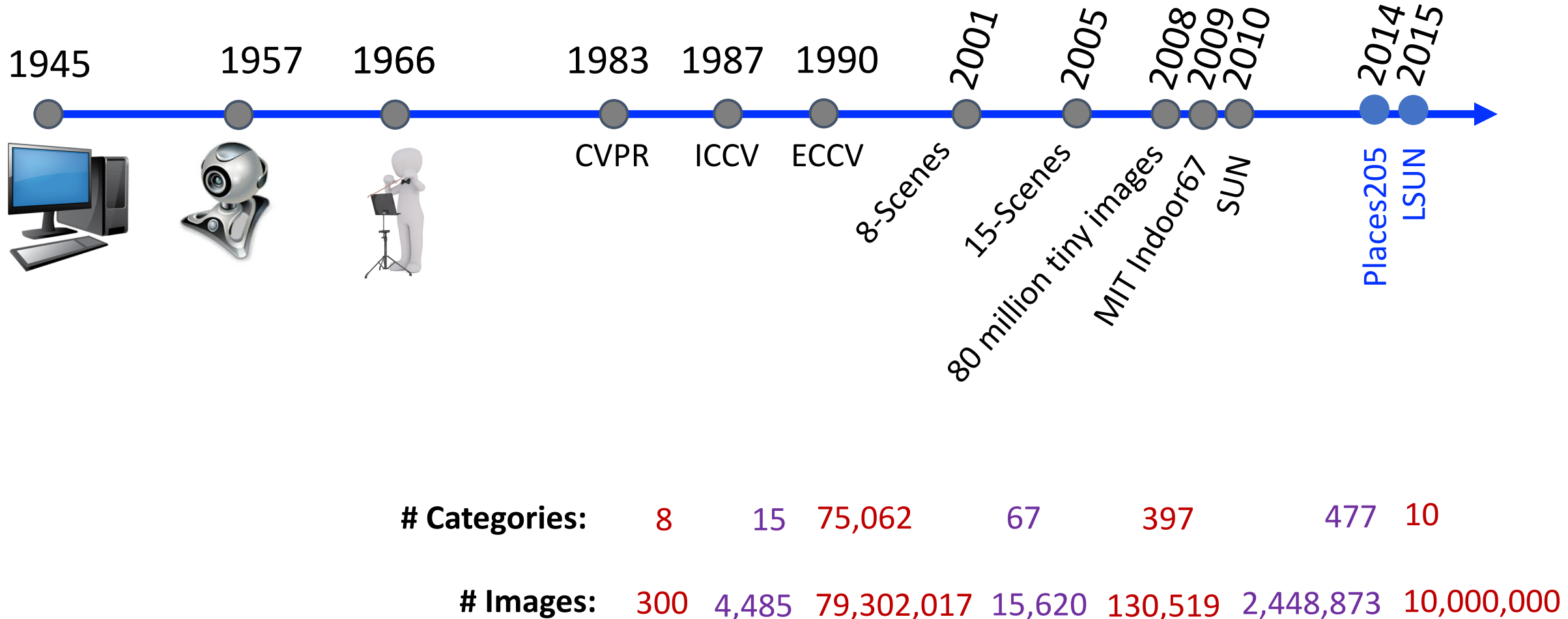
## Morning Session: Scene Understanding Workshop (SUNw'17)

Organizers: Bolei Zhou, Aditya Khosla, Jianxiong Xiao, James Hays

## Afternoon Session: Large SUN Challenge (LSUN'17)

Organizers: Fisher Yu, Peter Kotschieder, Shuran Song, Ming Jiang, Yinda Zhang, Catherine Qi Zhao, Thomas Funkhouser, Jianxiong Xiao

# Scene Classification Datasets



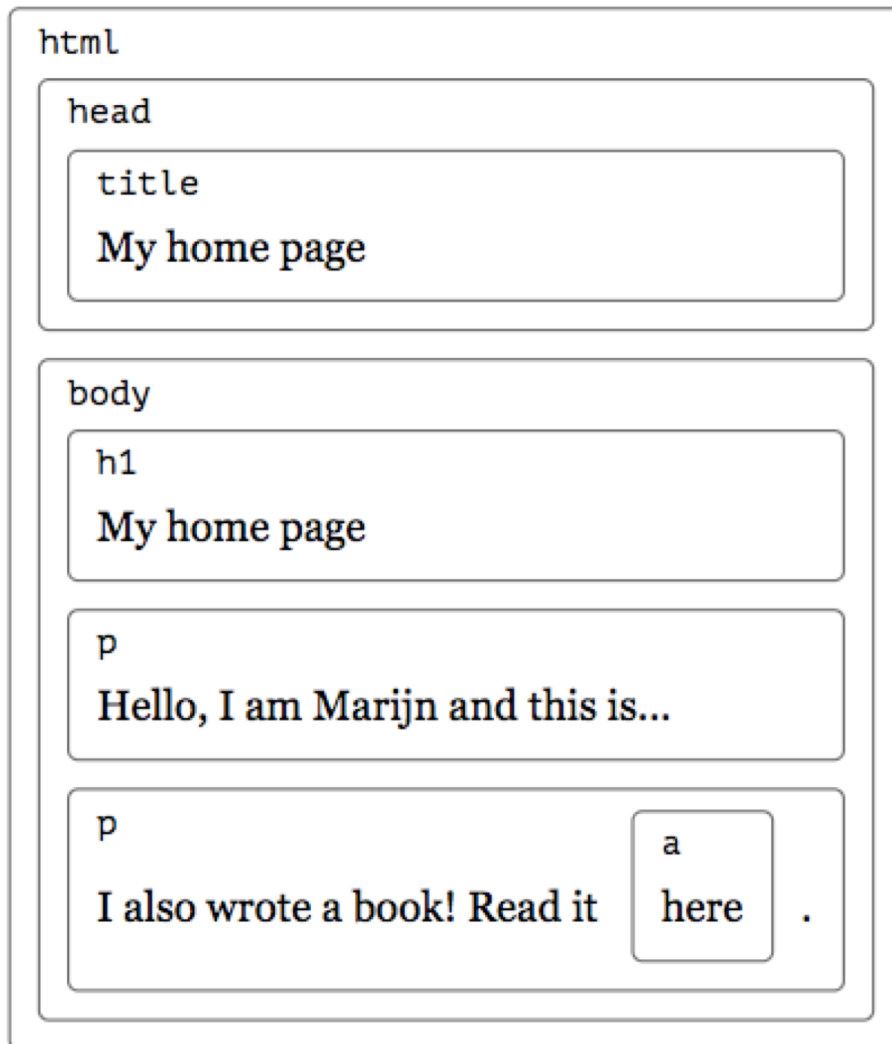
# Today's Topics

- Scene classification applications
- Scene classification datasets: key steps in creating them
- Scene classification datasets: scaling up with *crowdsourcing* and *challenges*
- Class discussion (chosen by YOU 😊)
- Lab: Javascript

# Today's Topics

- Scene classification applications
- Scene classification datasets: key steps in creating them
- Scene classification datasets: scaling up with *crowdsourcing* and *challenges*
- Class discussion (chosen by YOU 😊)
- Lab: Javascript

# Document Object Model (DOM)

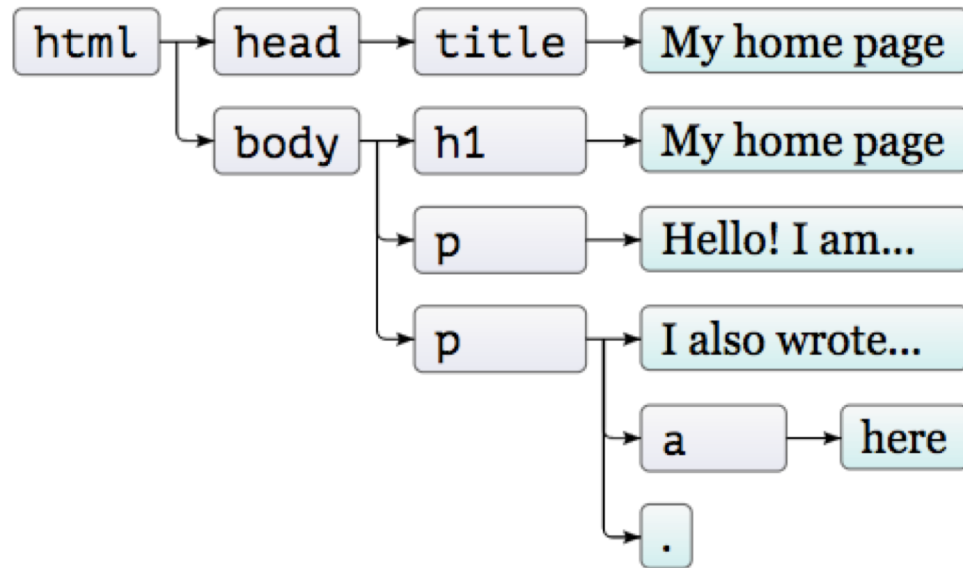


- Web browsers parse html into a DOM
- JavaScript programs interact with the html using the DOM

```
<!doctype html>
<html>
  <head>
    <title>My home page</title>
  </head>
  <body>
    <h1>My home page</h1>
    <p>Hello, I am Marijn and this is my home page.</p>
    <p>I also wrote a book! Read it
      <a href="http://eloquentjavascript.net">here</a>.</p>
  </body>
</html>
```

# Document Object Model (DOM)

- Web browsers parse html into a DOM
- JavaScript programs interact with the html using the DOM



```
<!doctype html>
<html>
  <head>
    <title>My home page</title>
  </head>
  <body>
    <h1>My home page</h1>
    <p>Hello, I am Marijn and this is my home page.</p>
    <p>I also wrote a book! Read it
      <a href="http://eloquentjavascript.net">here</a>.</p>
  </body>
</html>
```