

# Semantics of Visual Discrimination

John R. Smith  
IBM T. J. Watson Research Center  
[jsmith@us.ibm.com](mailto:jsmith@us.ibm.com)



February 2015

# Many Thanks To ...

## IBM T. J. Watson - Multimedia Research Team



Liangliang Cao



Michele Merler



Noel Codella



Matthew Hill



Quoc-Bao Nguyen



Wei Liu



Rosario Uceda-Sosa

## IBM T. J. Watson - Exploratory Computer Vision Team



Sharath Pankanti



Rogerio Feris



Quanfu Fan



Lisa Brown



Nalini Ratha



Chiao-Fe Shu



Chung-Ching Lin

## IBM Research Collaborators



Gang Hua



Shih-Fu Chang



John Kender



Daniel Ellis



Felix Yu

# Big News Event in 2005

How many cameras?



# Similar Big News Event Eight Years Later

How many cameras?



Cannot count them all ... image and video is here!

# Massive Multimedia is the Biggest Wave of All!



## Safety / Security



10s millions cameras

## Medical

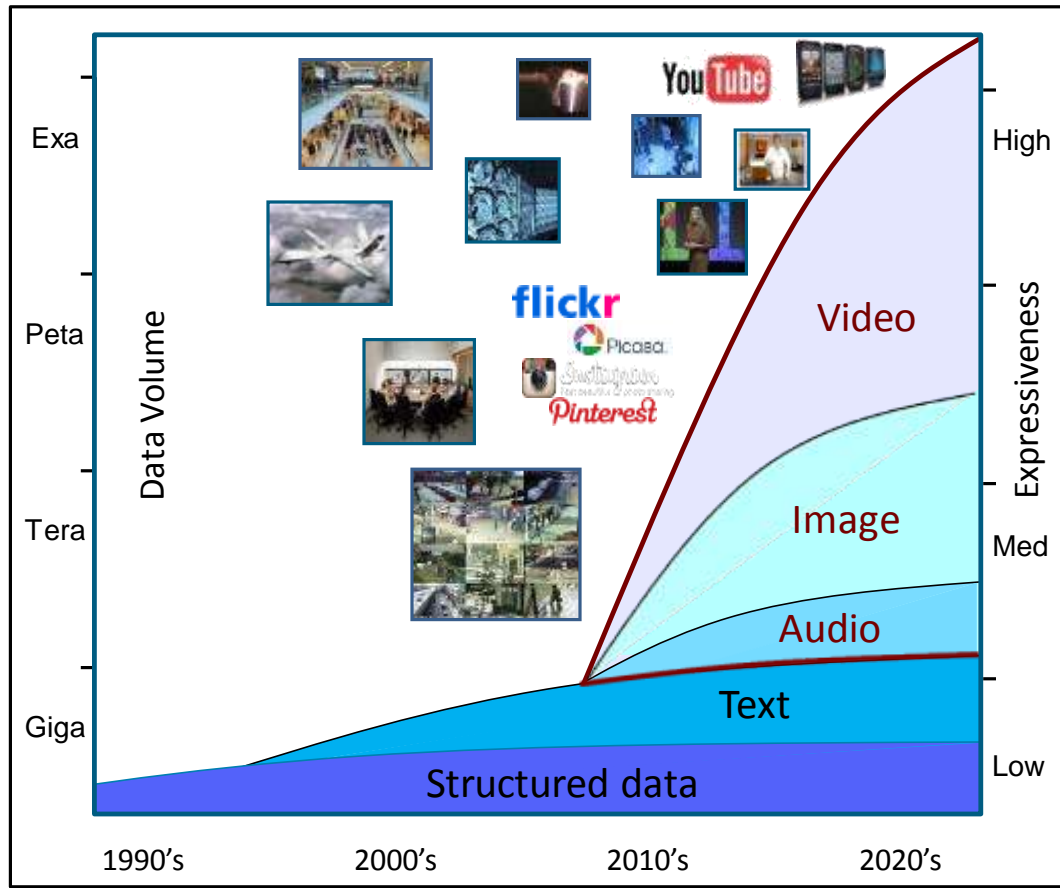


1B medical images

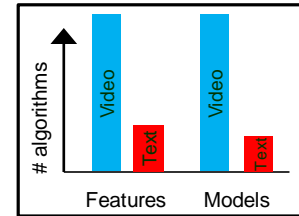
## Customer



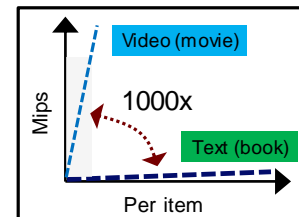
1B camera phones



## >> Sophistication



## >> Computation



## Media



100 video hrs/minute

## Wide Area Imagery



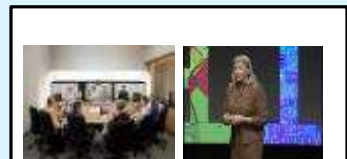
100's TB per day

## Digital Marketing



12% of video views

## Enterprise Video



Used by 1/3 of enterprises

## Key Messages:

**Image and video data is growing in volume and importance**

**Manual analysis is not scalable ... cannot keep up  
Vision system development has historically required deep expertise**

**Expert Analysis → Data-driven Visual Learning**

**Increasing availability of labeled training data  
Emergence of increasingly sophisticated visual learning algorithms**

**Focus needed on Visual Semantic Modeling**

**How to best model visual world (concepts and relationships)  
How to combine visual semantic modeling with learning systems**

# Use Cases

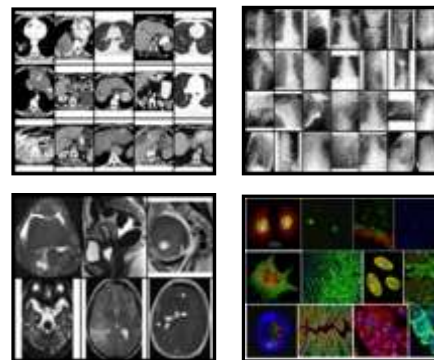
# Multiple Industry Problems Require Large-Scale Visual Content Analysis

## Safety and Security (IBM Intelligent Video Analytics)

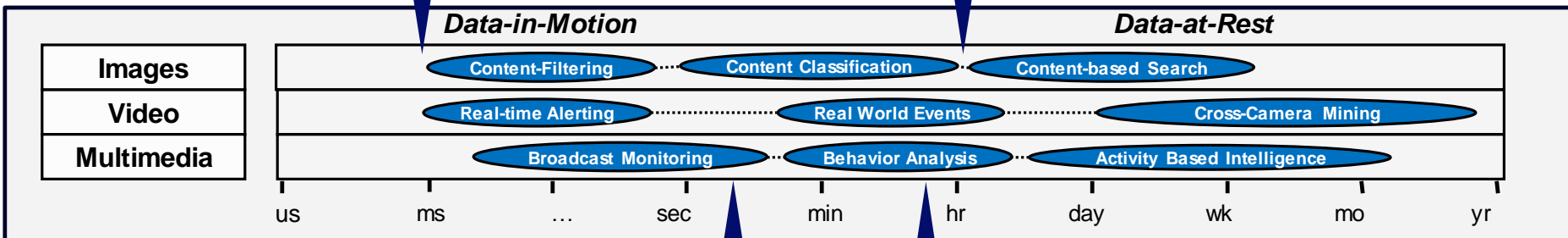


**Volume:** 10K's managed cameras per city  
**Velocity:** real-time alerts, 20M video events/day  
**Variety:** street scenes, rail stations, crowds, people, environmental conditions  
**Veracity:** analysis of complex activities (trip wires, abandoned objects)

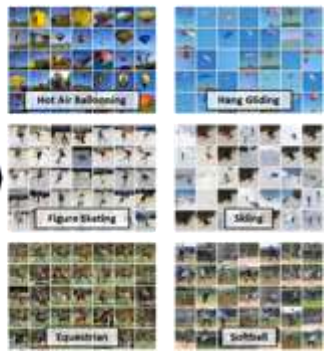
## Healthcare Cognitive Systems (IBM Watson Research)



**Volume:** 1B medical images per year (growing 20-40% /yr)  
**Velocity:** 50K radiology images per day per radiology dept.  
**Variety:** images, video, text, patent records, cases, scientific literature, ontologies/semantics  
**Veracity:** subjective interpretation across millions of categories (modalities, body views, organ systems, pathologies, anomalies)



## Enterprise Content Management (IBM IMARS)



**Volume:** 70PB broadcast/yr, 40K hrs per news archive  
**Velocity:** 100 video-hrs/min to YouTube  
**Variety:** mobile, user generated, professional  
**Veracity:** robust content extraction for objects, places, scenes, activities, people

## Retail and Mobile Commerce (IBM System V)



**Volume:** 500B consumer photos/yr  
**Velocity:** 100M customers per week for large retailers  
**Variety:** transient and dynamic content  
**Veracity:** predicting consumer attributes from diverse sources including visual data (images and video)

# Safety and Security:

## Urban Surveillance



**Trip Wires**

**Object Detection**

**Alerts**

**Offline Training**

**Forensic Search**

**Object Tracking**

- Expert developed algorithms
- Limited object detection and tracking
- Limited robustness to environment
- Manual tuning for new deployments

## Semantic Extraction



**Machine Learning**

**Visual Recognition**

**Attributes**

**Pattern Discovery**

**Activity Detection**

- Large-scale data-driven visual learning
- Semantic extraction including attributes
- Integrated analysis across cameras
- Self-configuration, tuning, adaptation

# IBM Intelligent Video Analytics (IVA) for Smarter Cities



Vehicle detection/classification



Person and face attributes



Activities and behaviors

## Visual Learning



## Traditional Computer Vision



Real world metrics (speed, size)



Trip wires and safety regions



Object detection and tracking

## Media-in-the-Wild:

### Traditional Media Archives

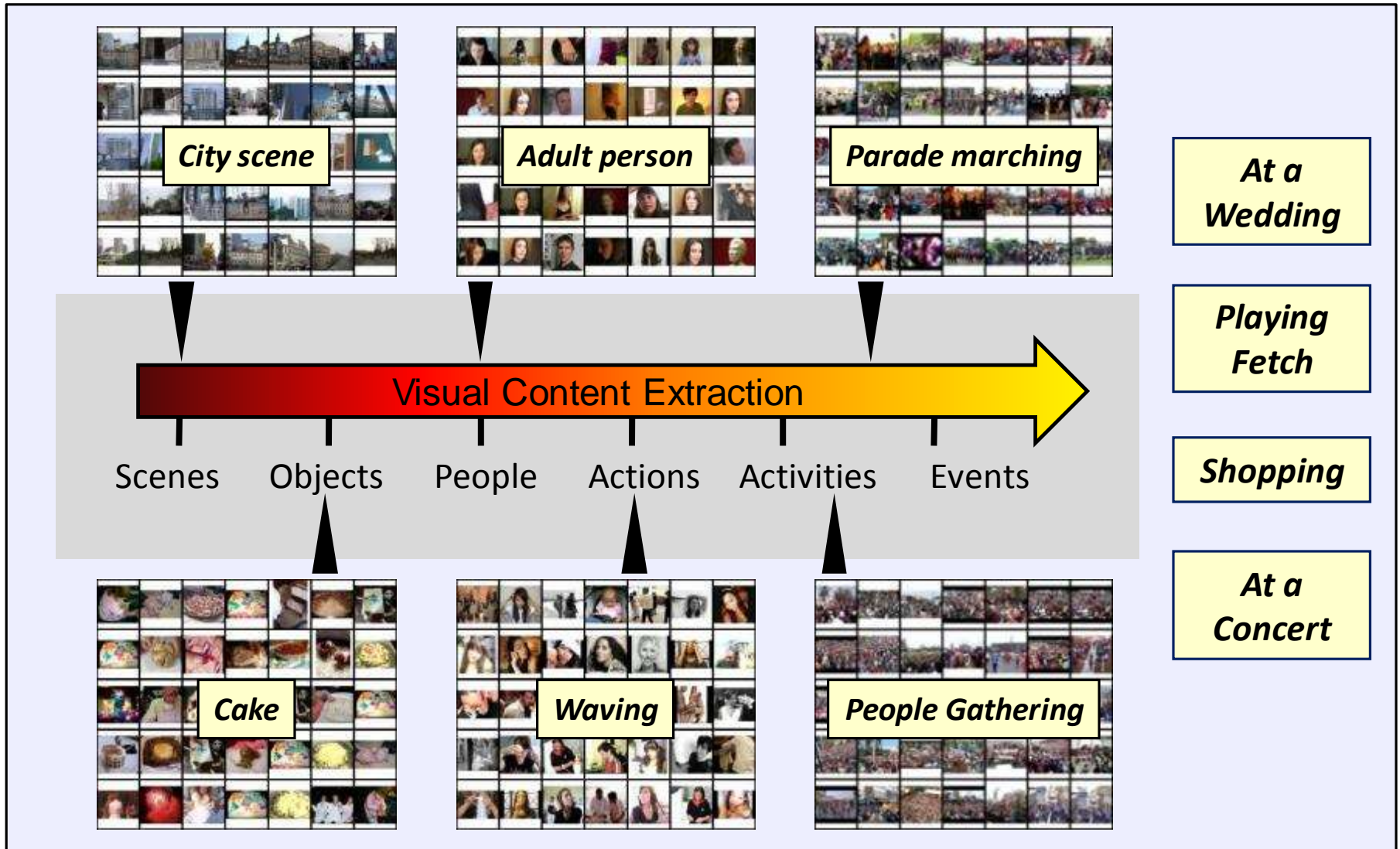
- Limited indexing and search of *news archives, TV programs, movies*
- Relies heavily on related information (e.g., *metadata, speech transcript, user tags*)



### User Generated Content

- Want to understand diverse *user generated content* at a semantic level (e.g., *sports, activities, life events*)
- Extraction of visual insights
- Discovery of patterns across *users, segments, time, geographies*

**Visual search uses thousands of Visual Classifiers across dozens of Facets (*Objects, Scenes, Locations, Activities, People, Events, etc.*)**



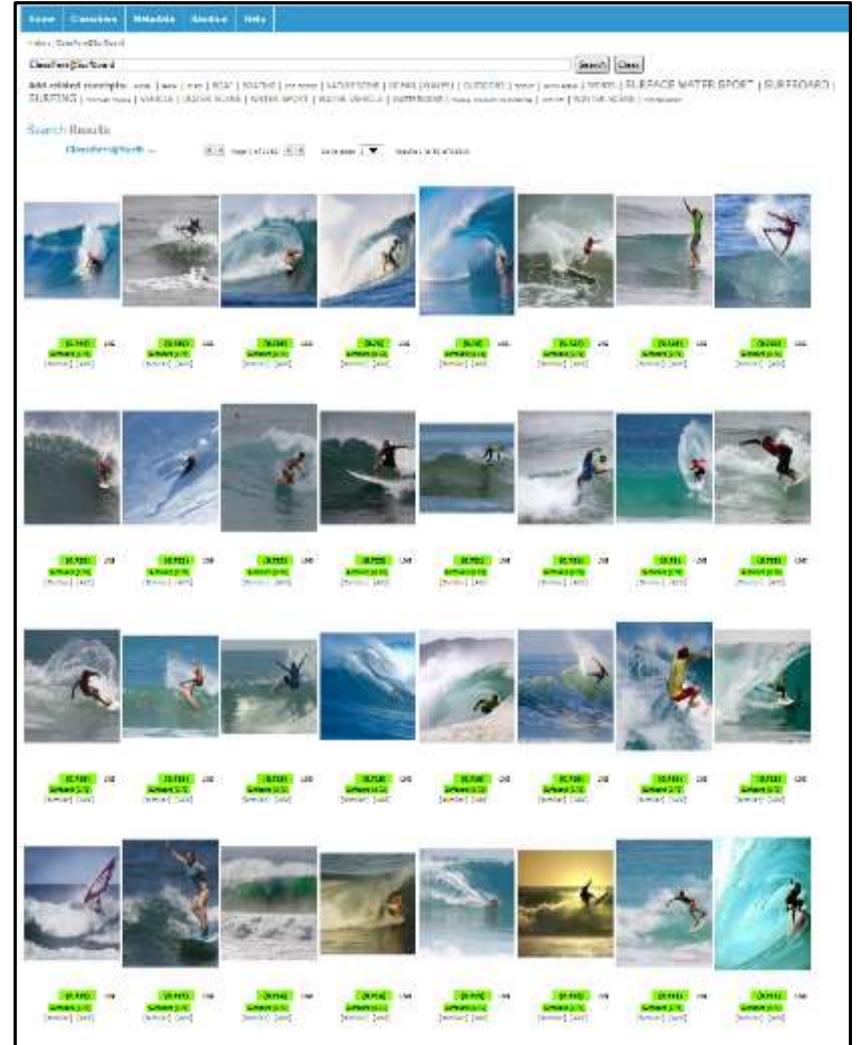
# Example: Visual Recognition for Image Semantic Indexing

Demo

## Visual Semantic Label Categories



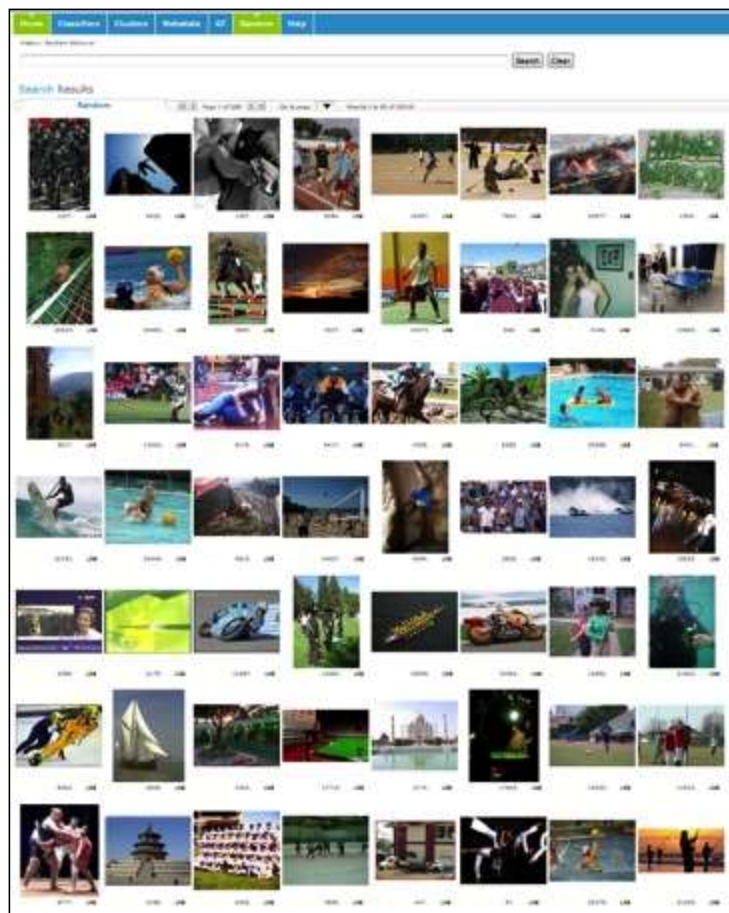
## E.g., *Surfing*







# Machine Learning Shifts Effort to Organizing Image Data for Training



**Accurate recognition of *sports* by training 150 categories**

# Semantic Searches is supported by Automatically Extracted Semantic Labels (Boolean = AND/OR)



## Car AND Street Scene

The screenshot shows a search interface with a navigation bar (Home, Classifiers, Outlets, Analytics, Near-Overlaid, Metadata, Results, Help) and a search bar containing the query 'Car AND Street Scene'. Below the search bar, there are several search results displayed as a grid of images. Each image is accompanied by a small green and red label indicating the detected semantic categories. The results include various types of vehicles (cars, trucks, vans) in different street settings.

## Beach OR Sunset

The screenshot shows a search interface with a navigation bar (Home, Classifiers, Outlets, Analytics, Near-Overlaid, Metadata, Results, Help) and a search bar containing the query 'Beach OR Sunset'. Below the search bar, there are several search results displayed as a grid of images. Each image is accompanied by a small green and red label indicating the detected semantic categories. The results include various beach scenes (sand, water, waves) and sunset scenes (orange sky, lighthouse, beach at dusk).

# Semantic Indexing applies to Video Content Extraction and Search



## Boating

The screenshot shows a search interface for the term "Boating". The search results are displayed as a grid of 84 small video thumbnails, arranged in 7 rows and 12 columns. Each thumbnail includes a small green "SEARCH" button and a "1:00" duration indicator. The thumbnails depict various boating activities, including people on boats, kayakers, and sailboats on a lake.

Other activities: *Running, Skiing, ...*

## Parade AND Urban Scene

The screenshot shows a search interface for the combined terms "Parade AND Urban Scene". The search results are displayed as a grid of 108 small video thumbnails, arranged in 9 rows and 12 columns. Each thumbnail includes a small green "SEARCH" button and a "1:00" duration indicator. The thumbnails depict various scenes from parades and urban environments, including people in costumes, street scenes, and buildings.

Other Combining Functions: *AND, OR, X, MIN*

# On-line photos and videos have enormous potential as a rich source of information about consumers

---



**1/4 Trillion**

*photos* hosted on Facebook from 1.3 Billion users



**100 Hours**

of video are uploaded to YouTube *every minute*



**51%**

of the Class of 2015 (high school) use Instagram *daily*



**80%**

of Pinterest users are *female*



. . . and growing

# Growing Amount of Consumer Images and Video is Source for Insights



## Healthcare:

### Radiology Image Analysis

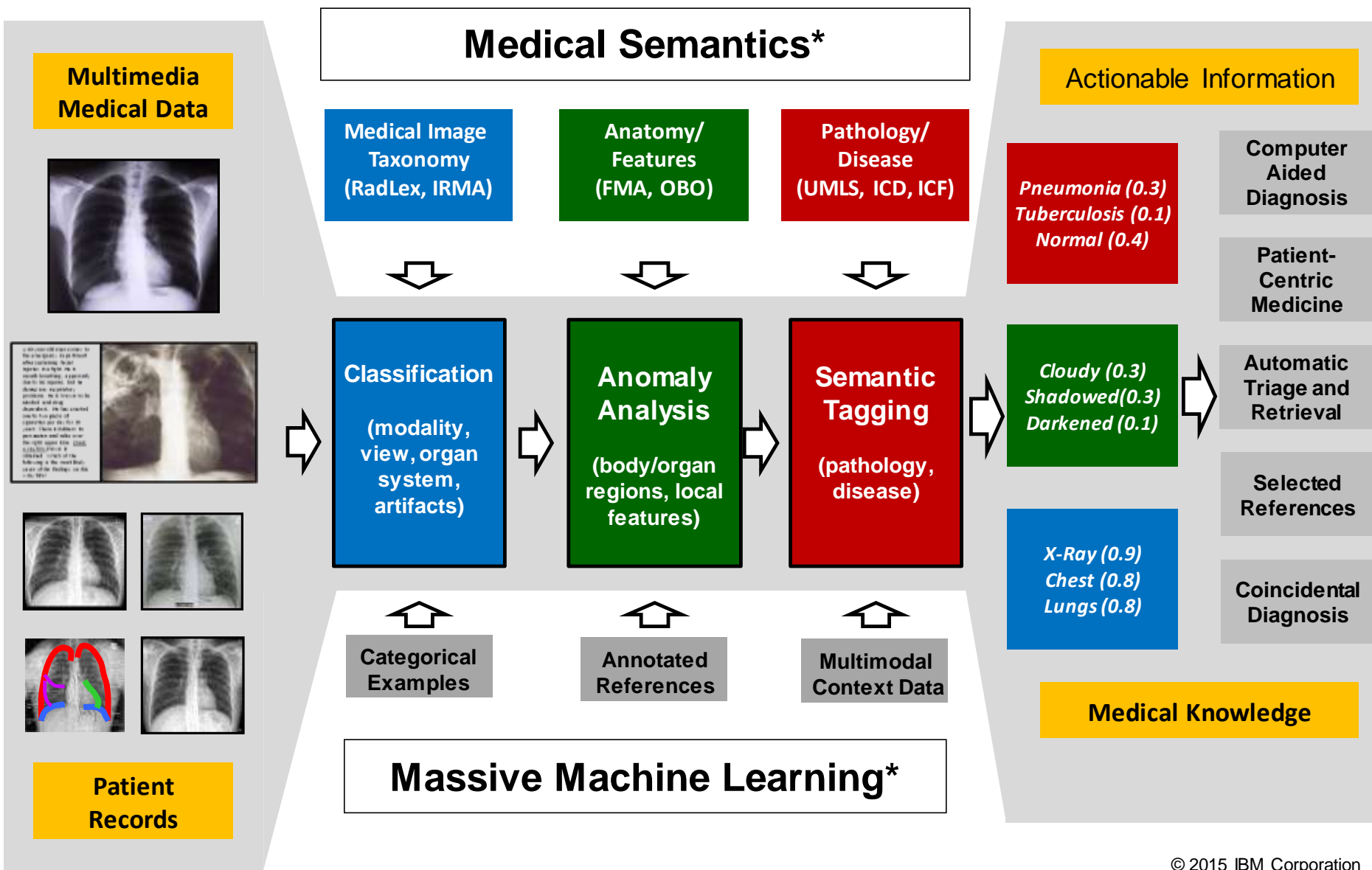
- Specialized algorithms per medical *modality* and *disease*
- Limited scaling and coverage
- Bottleneck in algorithm development by computer vision experts



### Medical Insights

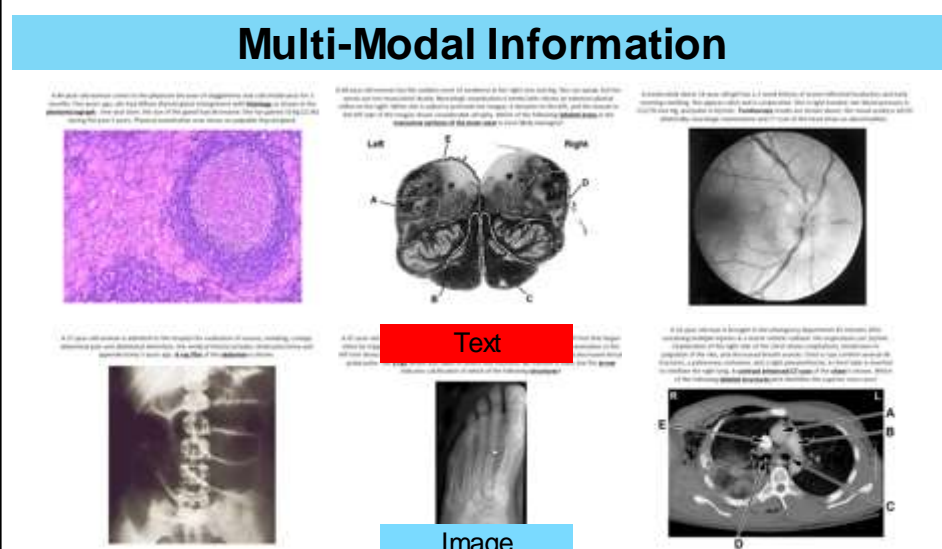
- Diverse visual extraction across *modality, anatomy, pathology*
- Address wide spectrum across patient image data, images and figures in medical literature, visual knowledge repositories
- Beyond radiology to analyze images broadly in medicine

# Multi-modal Analysis helps Build Medical Knowledge and Aid Diagnosis



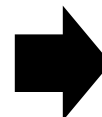
# Huge Diversity of Imagery in Medicine Spans Modality, Anatomy, Disease

## Multi-Modal Information


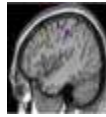



Text

Image

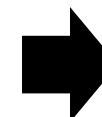


## Visual Classification

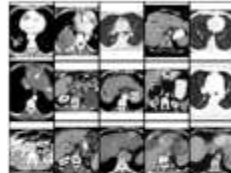
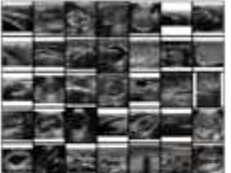
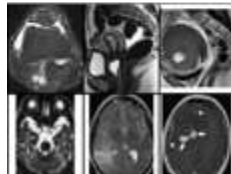
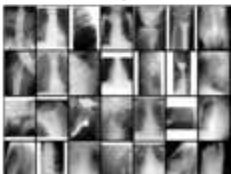
|   |     | Modality     |        |   |
|---|-----|--------------|--------|---|
|  | ... | X-ray        | CT     | MRI   |
|   |     | Region       |        |   |
|   |     | Pelvis       | Head   | Chest ...  |
|   |     | Disease      |        |   |
|  | ... | Tuberculosis | Cancer | Collapse  |

## NIH PubMed Medical Modality Classification:

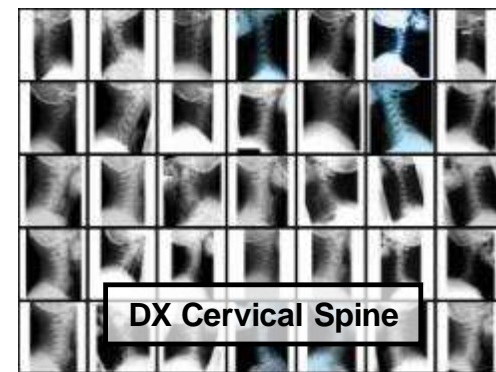
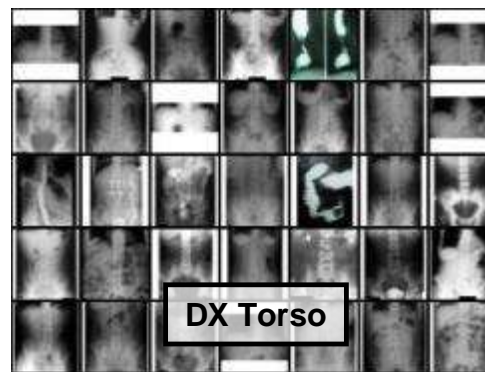
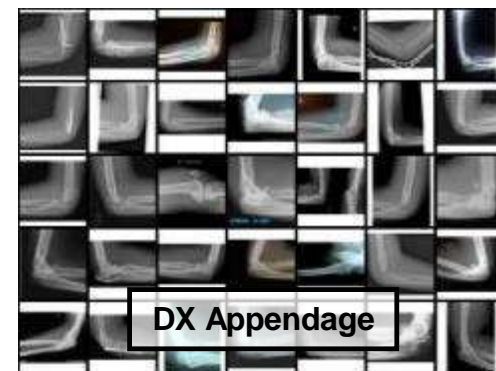
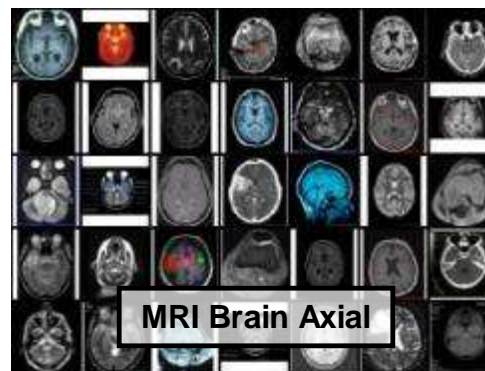
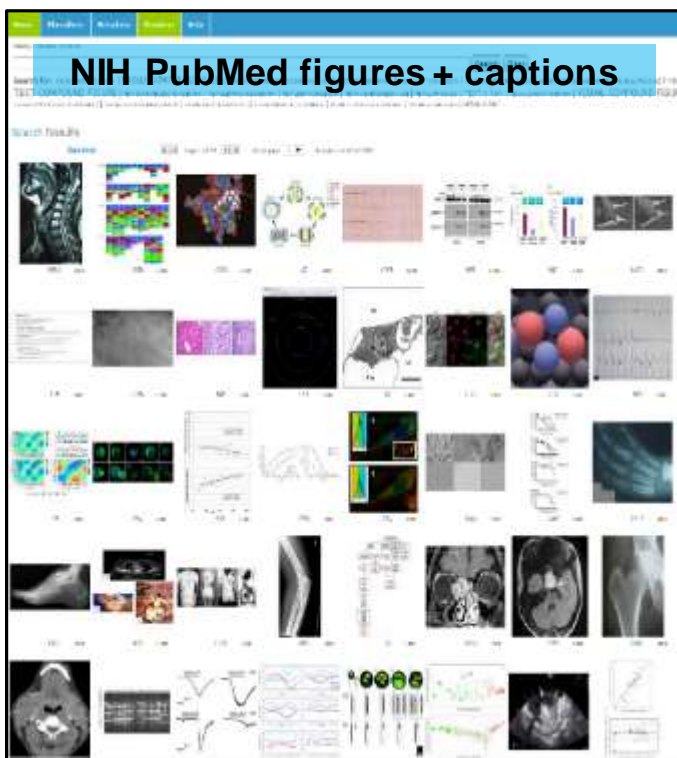
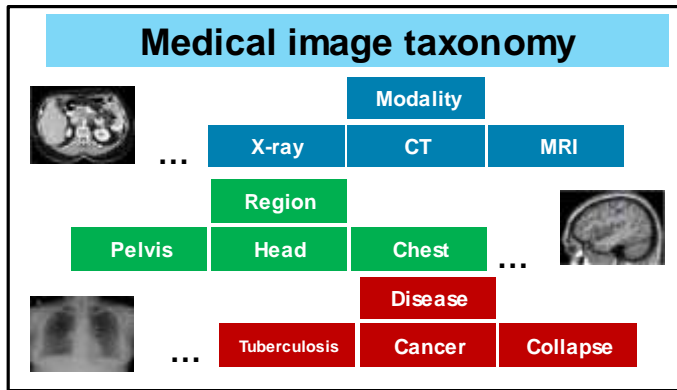
- **Extracting knowledge from NLM/NIH PubMed** 14.2 million articles (3.8 million *free* full-text articles)
- Rich **contextual information**: text + figures + captions
- Challenge: given an unknown PubMed image determine medical category automatically
- Automatically classifying millions of published medical images by **modality, region, disease** along with associated text builds multi-modal knowledge



## Visual Recognition

| CT  | Echo  |
|---|---|
|   |   |
| MRI   | Xray  |
|  |  |

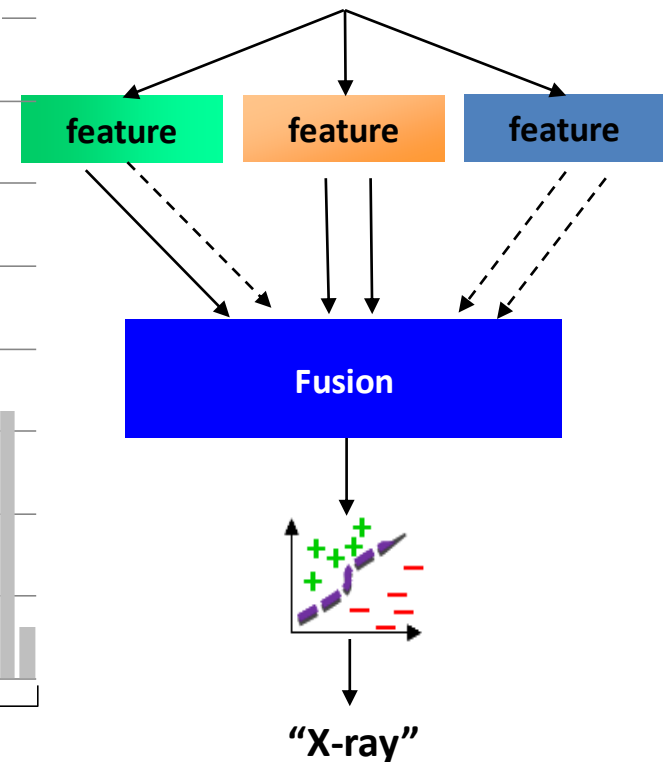
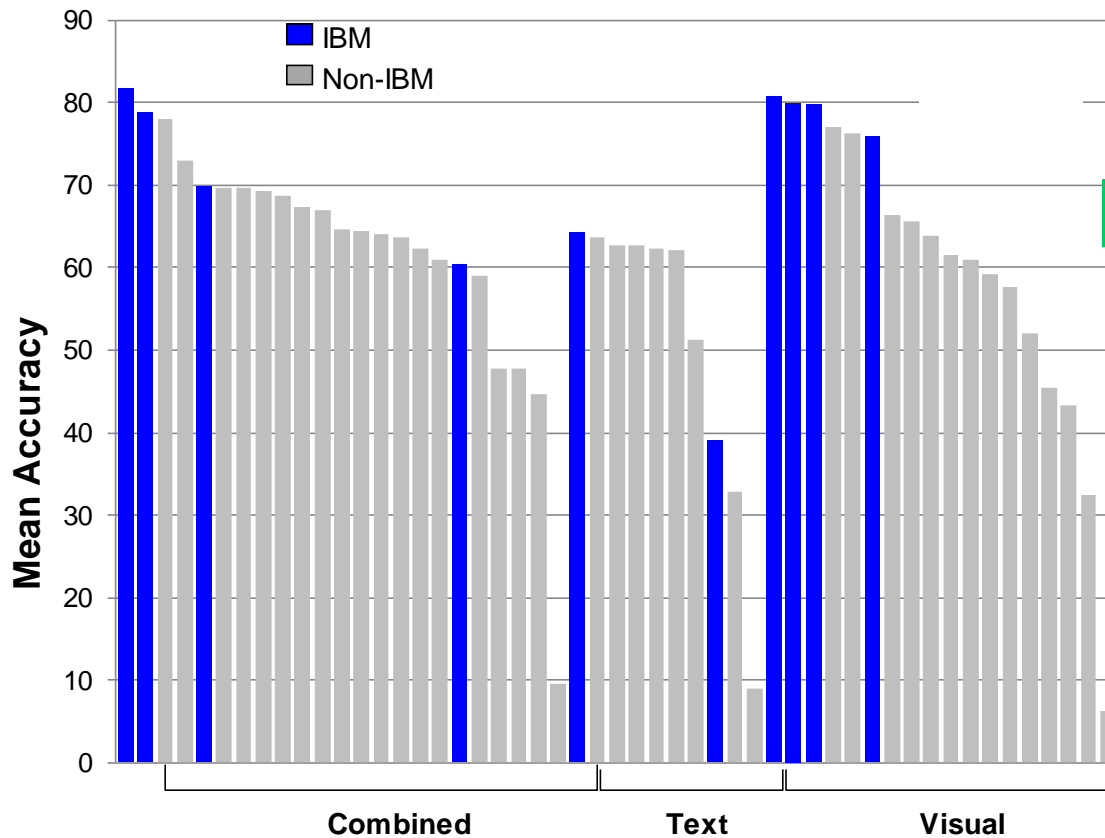
# Medical Knowledge Management – Automatic Medical Image Classification



# ImageCLEF Medical Image Classification\*

**Goal:** automatically classify images into correct medical category

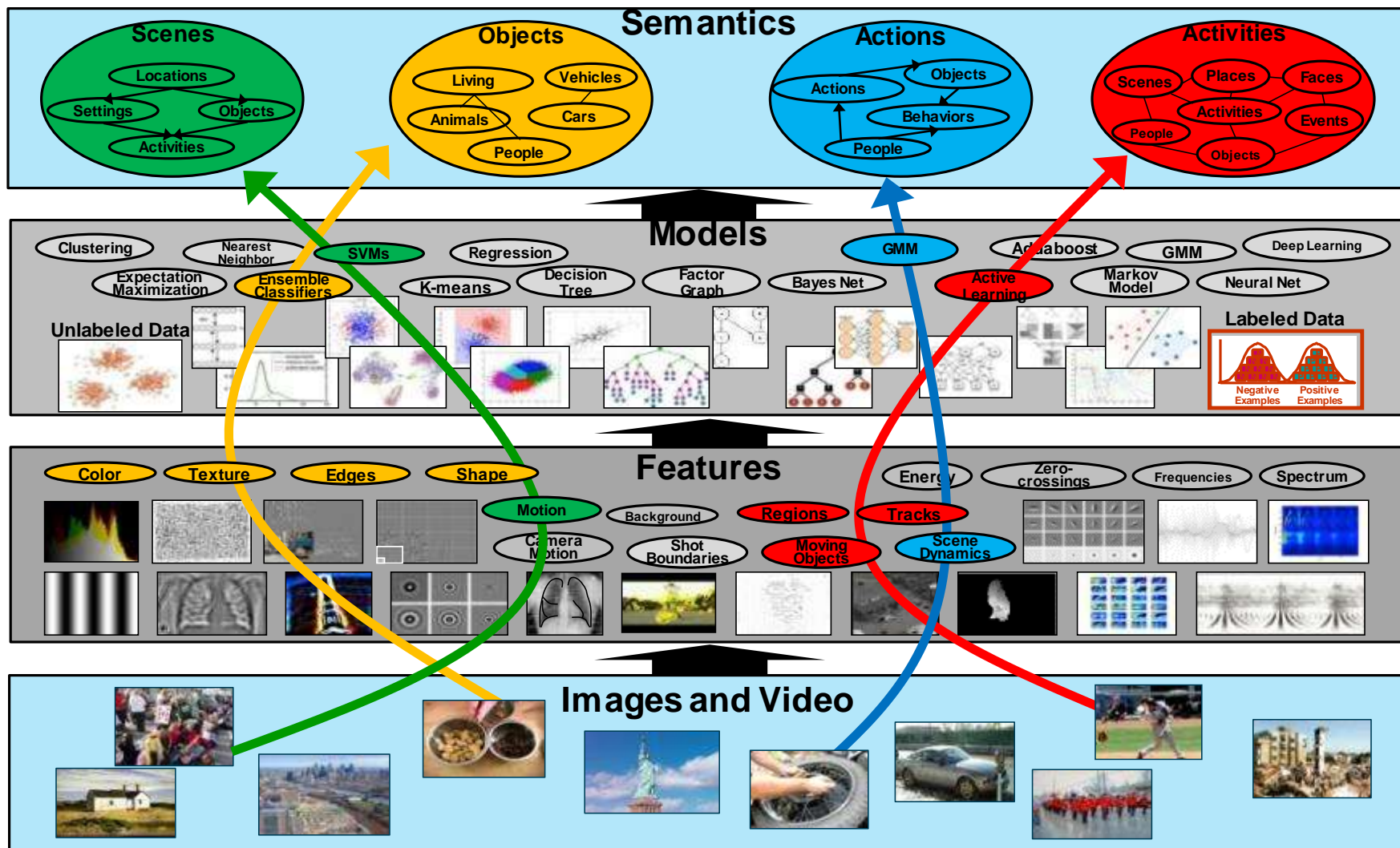
- Data-set of NIH PubMed articles with 305,000 images
- IBM achieved #1 performance in 2012 and 2013
- Top performance in each task (visual, text, combined)



\* <http://imageclef.org/2013/medical>

# **Visual Recognition Technical Foundation**

# Multi-layer Learning Architecture for Image and Video Analysis

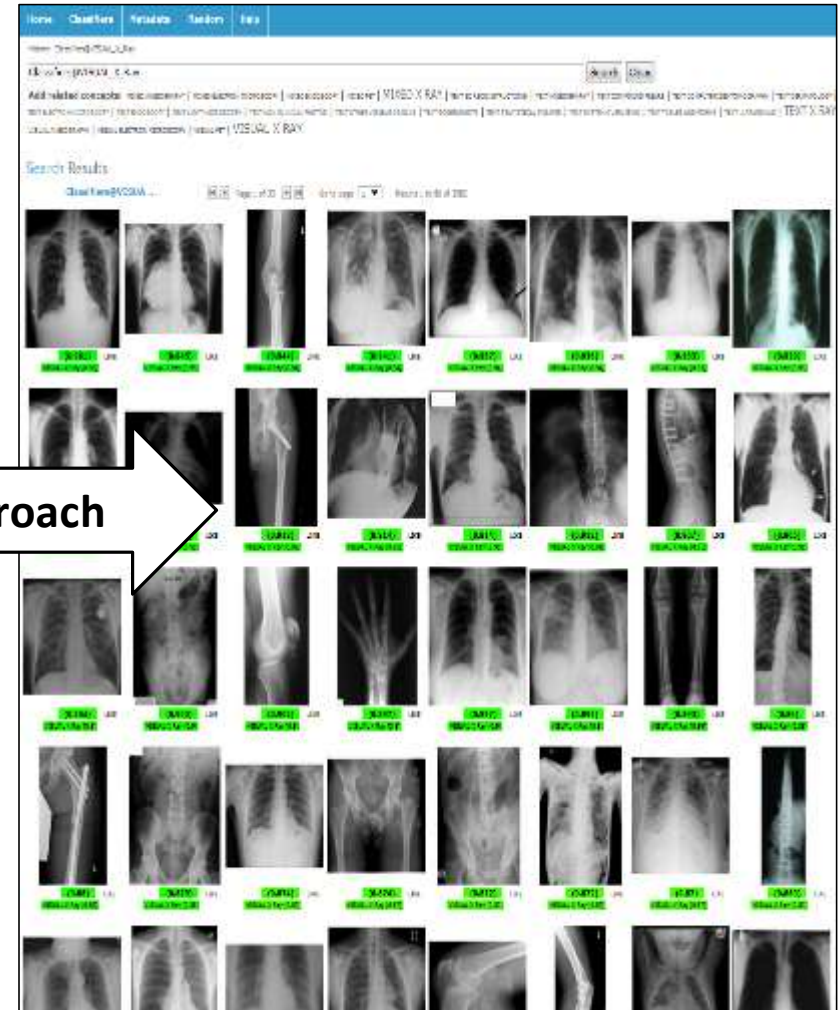


• **Visual Recognition** = predicting “semantic” labels for unknown images and video

## Natural Photos and Video



## Medical Modalities & Viewpoints



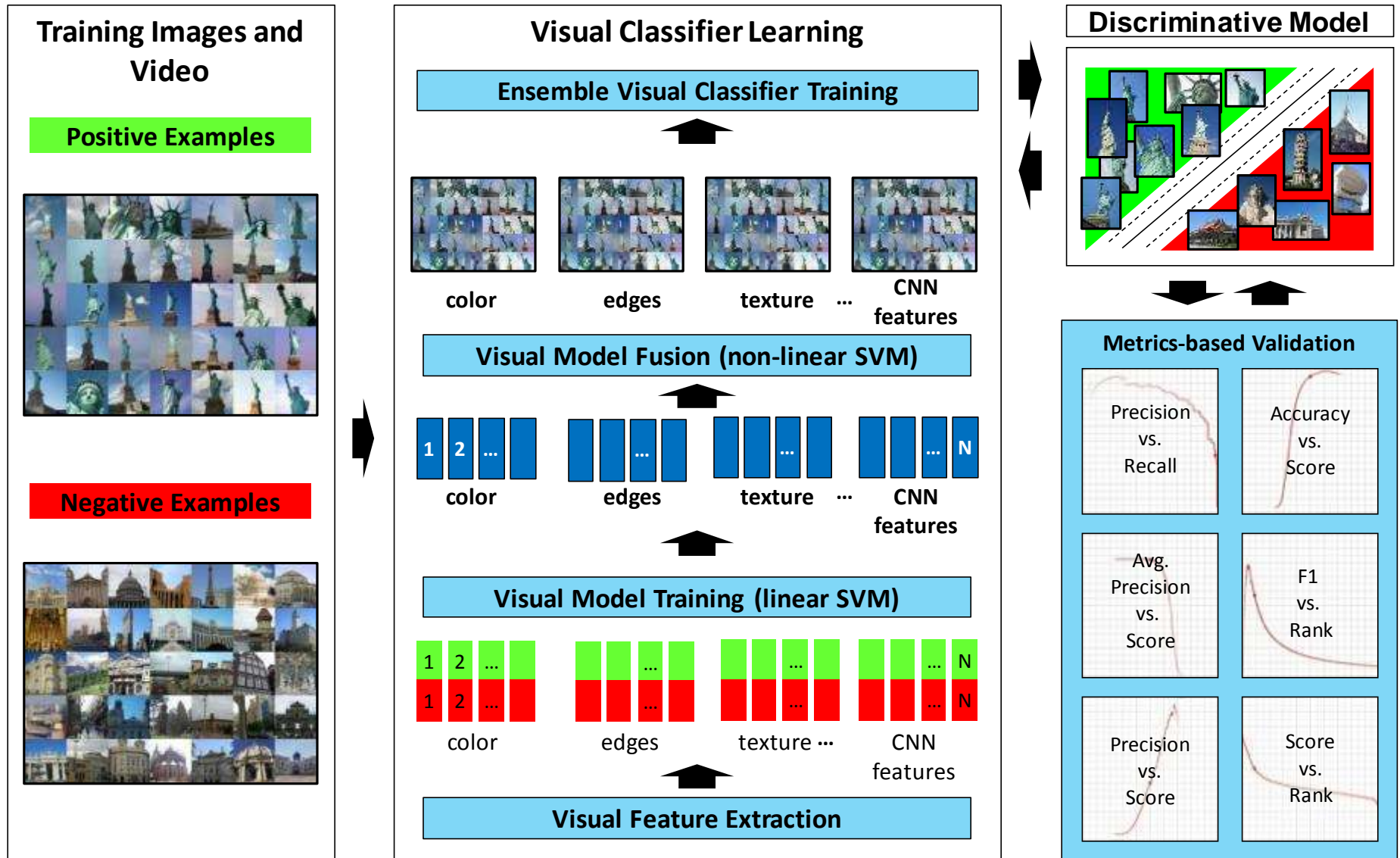
Same Approach

# Rich Set of Visual Features is Needed to Learn Semantic Discrimination





# Visual Classifier Learning Allows Metrics-based Optimization



# Visual Recognition Example

## Upload Unknown Images

IMARS

HOME RESULTS CONTACT

Max labels per image: 4 Max labels per facet: 4 Threshold: 0.5 RECOGNIZED

Grid of 25 image thumbnails with file sizes:

- 18.6 KB, 8 KB, 4.2 KB, 24.8 KB, 37.6 KB, 21.8 KB
- 18.9 KB, 25.9 KB, 26.5 KB, 20 KB, 28.2 KB, 28.7 KB
- 22.8 KB, 25.2 KB, 21.3 KB, 16.8 KB, 38.2 KB, 21.8 KB
- 16.4 KB, 20.4 KB, 19 KB, 24.4 KB, 13.1 KB, 21.6 KB
- 18.8 KB, 16.9 KB, 7.1 KB, 28.9 KB, 38.4 KB, 21.2 KB
- 32.8 KB, 28.7 KB, 20.8 KB, 26.3 KB, 21.3 KB



## Visual Recognition Results

IMARS

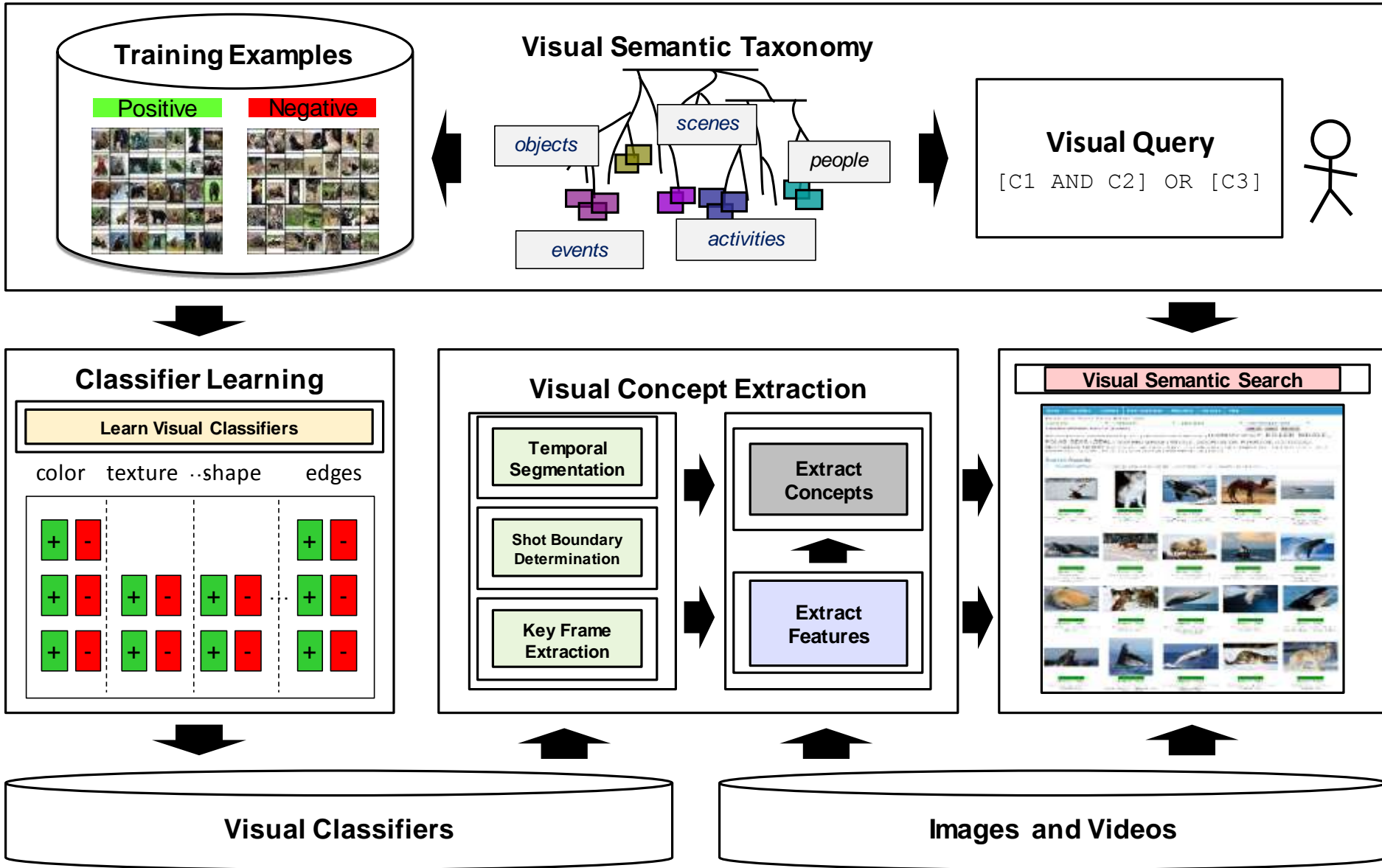
HOME RESULTS CONTACT

Grid of 25 image thumbnails with associated labels and scores:

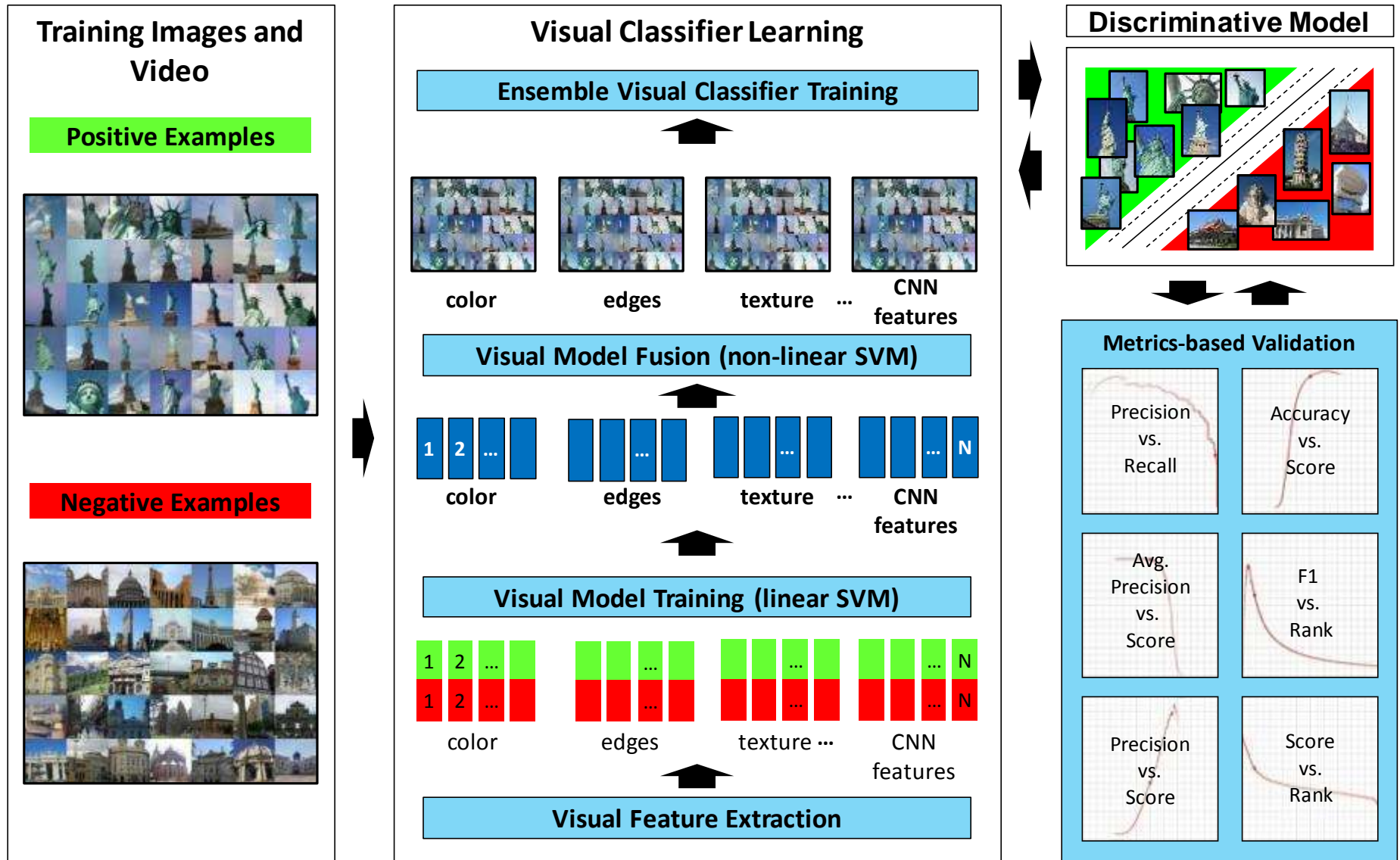
- Person: 0.691, Head Color: 0.690, Person View: 0.685, Person: 0.685
- Tennis Sport: 0.687, Team Field Sport: 0.687, Ball Sport: 0.682, Baseball: 0.648, Sports Activity: 0.643
- Stadium: 0.718, Crowd: 0.711, Head Color: 0.705, Demonstration: 0.659, Banner: 0.653
- Scene: 0.718, Indoor: 0.708, Room: 0.698, Object: 0.698, State of People: 0.615
- Monument: 0.693, Bus: 0.735, Landmark: 0.838, Air Sport: 0.689, Statue Statue: 0.641
- Indoor: 0.760, Scene: 0.636, Connection: 0.636, Head Color: 0.581, Sports: 0.511, Team Sport: 0.510
- Person: 0.691, Head Color: 0.690, Person View: 0.685, Person: 0.685
- Tennis Sport: 0.687, Team Field Sport: 0.687, Ball Sport: 0.682, Baseball: 0.648, Sports Activity: 0.643
- Stadium: 0.718, Crowd: 0.711, Head Color: 0.705, Demonstration: 0.659, Banner: 0.653
- Scene: 0.718, Indoor: 0.708, Room: 0.698, Object: 0.698, State of People: 0.615
- Monument: 0.693, Bus: 0.735, Landmark: 0.838, Air Sport: 0.689, Statue Statue: 0.641
- Indoor: 0.760, Scene: 0.636, Connection: 0.636, Head Color: 0.581, Sports: 0.511, Team Sport: 0.510
- Person: 0.691, Head Color: 0.690, Person View: 0.685, Person: 0.685
- Tennis Sport: 0.687, Team Field Sport: 0.687, Ball Sport: 0.682, Baseball: 0.648, Sports Activity: 0.643
- Stadium: 0.718, Crowd: 0.711, Head Color: 0.705, Demonstration: 0.659, Banner: 0.653
- Scene: 0.718, Indoor: 0.708, Room: 0.698, Object: 0.698, State of People: 0.615
- Monument: 0.693, Bus: 0.735, Landmark: 0.838, Air Sport: 0.689, Statue Statue: 0.641
- Indoor: 0.760, Scene: 0.636, Connection: 0.636, Head Color: 0.581, Sports: 0.511, Team Sport: 0.510
- Person: 0.691, Head Color: 0.690, Person View: 0.685, Person: 0.685
- Tennis Sport: 0.687, Team Field Sport: 0.687, Ball Sport: 0.682, Baseball: 0.648, Sports Activity: 0.643
- Stadium: 0.718, Crowd: 0.711, Head Color: 0.705, Demonstration: 0.659, Banner: 0.653
- Scene: 0.718, Indoor: 0.708, Room: 0.698, Object: 0.698, State of People: 0.615
- Monument: 0.693, Bus: 0.735, Landmark: 0.838, Air Sport: 0.689, Statue Statue: 0.641
- Indoor: 0.760, Scene: 0.636, Connection: 0.636, Head Color: 0.581, Sports: 0.511, Team Sport: 0.510
- Person: 0.691, Head Color: 0.690, Person View: 0.685, Person: 0.685
- Tennis Sport: 0.687, Team Field Sport: 0.687, Ball Sport: 0.682, Baseball: 0.648, Sports Activity: 0.643
- Stadium: 0.718, Crowd: 0.711, Head Color: 0.705, Demonstration: 0.659, Banner: 0.653
- Scene: 0.718, Indoor: 0.708, Room: 0.698, Object: 0.698, State of People: 0.615
- Monument: 0.693, Bus: 0.735, Landmark: 0.838, Air Sport: 0.689, Statue Statue: 0.641
- Indoor: 0.760, Scene: 0.636, Connection: 0.636, Head Color: 0.581, Sports: 0.511, Team Sport: 0.510

Demo

# Visual Indexing and Search using Semantic Faceted Taxonomy



# Visual Classifier Learning Allows Metrics-based Optimization



# Different Approaches can be used for Managing and Selecting Negatives for Visual Discriminative Learning

## Positive Examples

## Negative Examples



*Whale*



*Person, Dog, Whale*

Common background  
of all images  
(unlabeled positive  
images are included as  
negatives)



*Whale*



*Zebra, Nature, Sunset*

Traditional taxonomy  
mutually exclusivity  
(unlabeled positive  
images are included as  
negatives)



*Whale*



*Zebra, Cat, Koala*

Faceted classification  
ensures mutual  
exclusivity within facets  
(allows correct positives  
and negatives)

# What label?



- Man
- Dog
- Frisbee
- Beach
- Playing

## FACETS:

[Person]

[Animal]

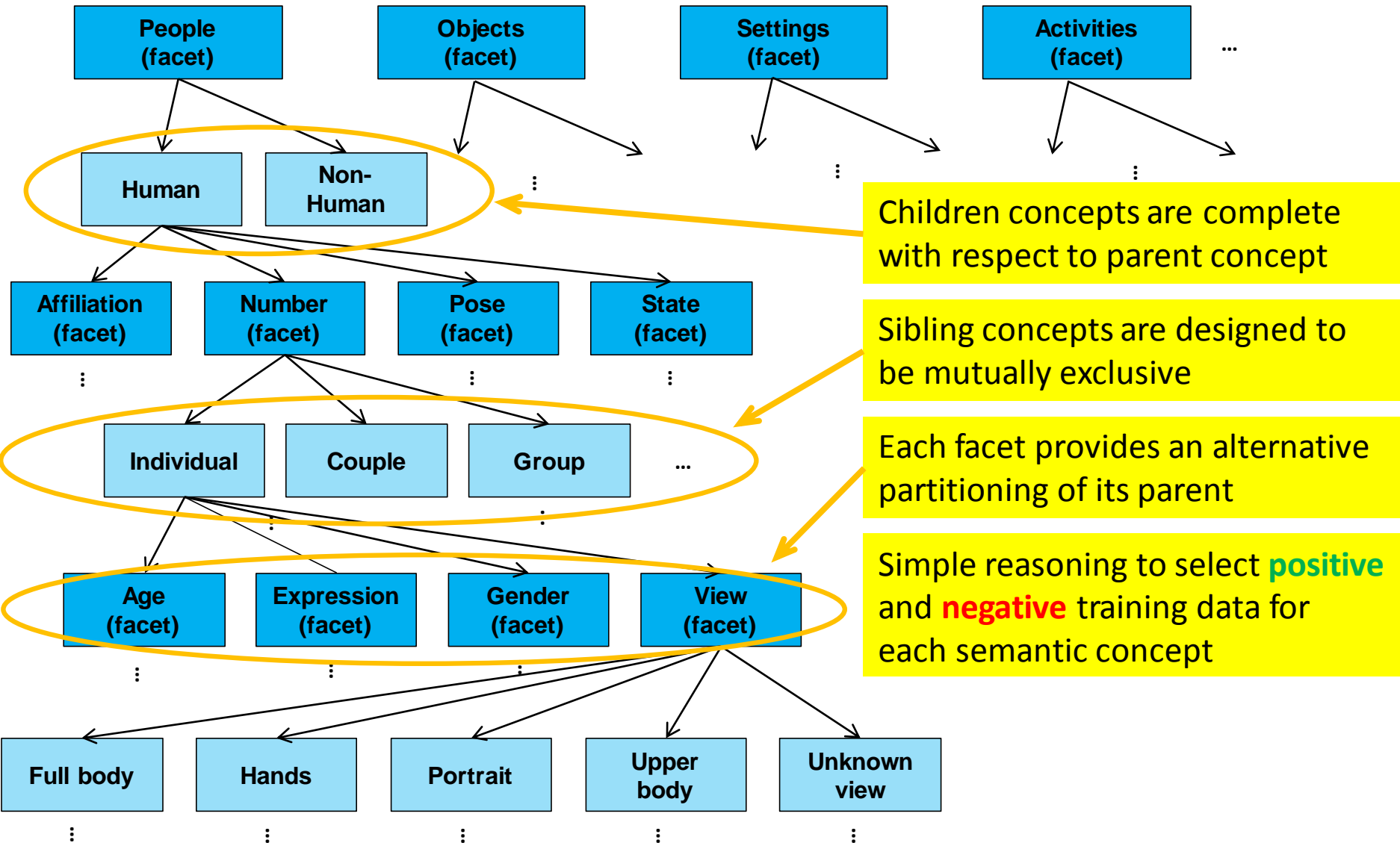
[Object]

[Setting]

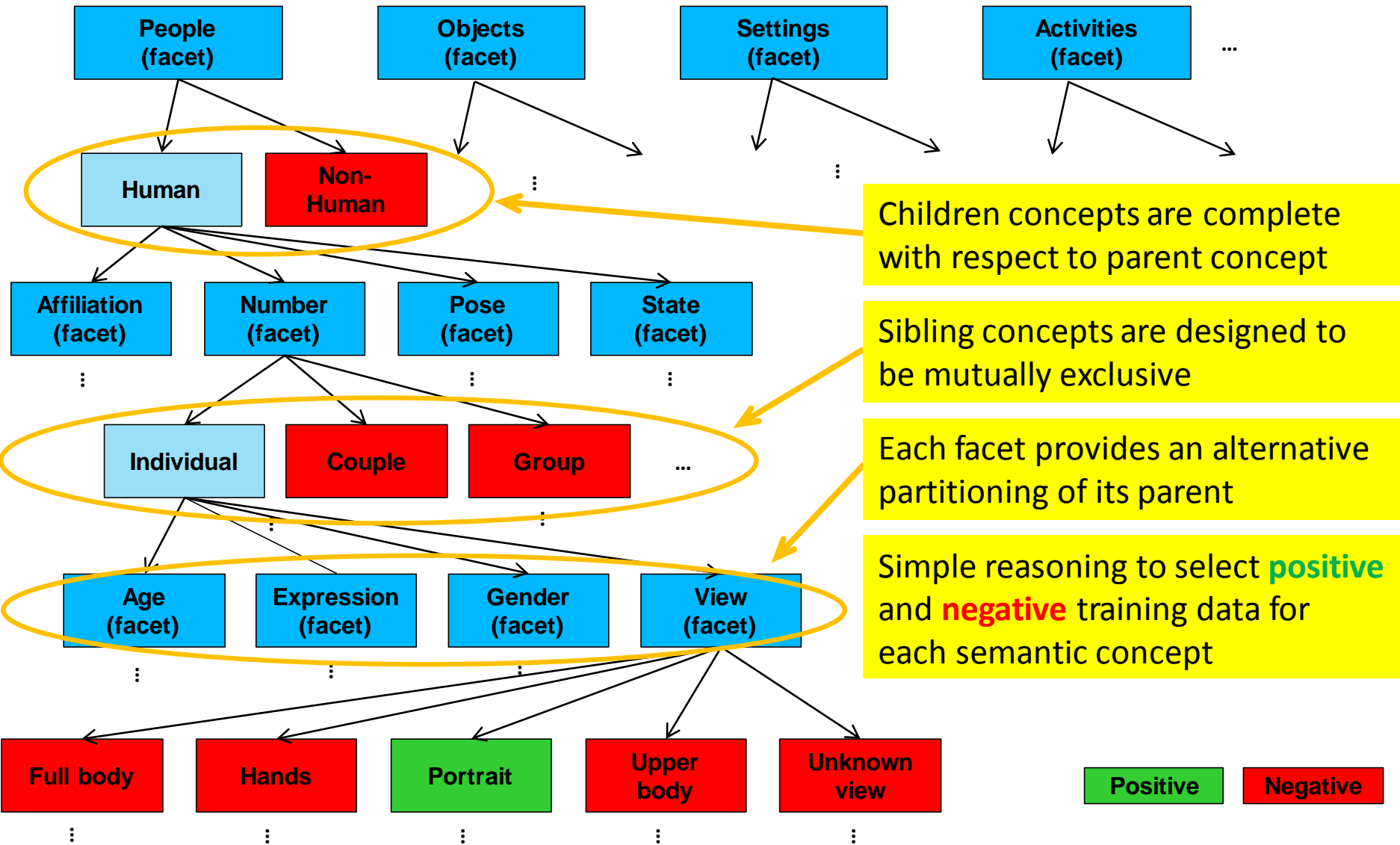
[Activity]

**Need to learn and assign labels from multiple facets!**

# Facets can be Nested in Hierarchical Faceted Classification Scheme



# Facets can be Nested in Hierarchical Faceted Classification Scheme



# Key Challenge in Visual Discriminative Learning is Managing and Selecting of Negative Training Examples

## Training Setup ?

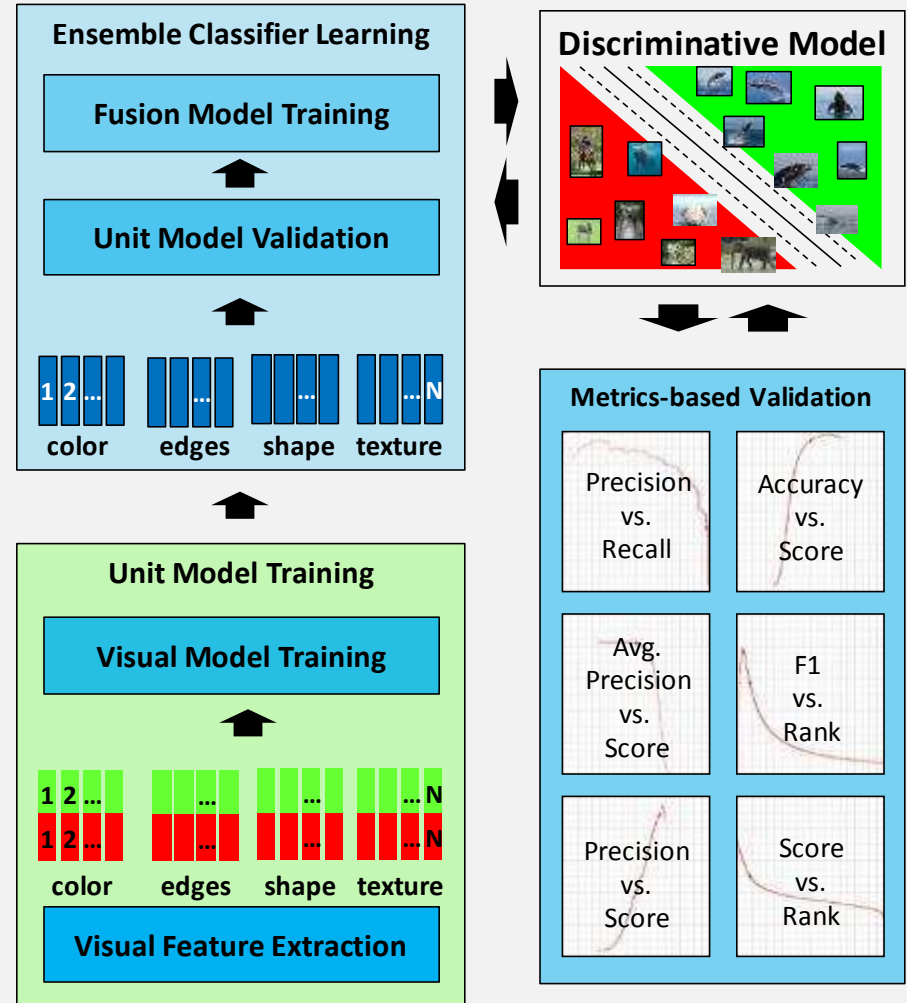
### Positive Examples



### Negative Examples ???



## Semantic Learning



# Approaches for Labeling Training Images for Discriminative Learning

$N$  Images

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|
| A | N | Y | N | N | N | N | Y | N | Y | N  | N  | N  |
| B | N | N | Y | N | N | N | N | N | N | N  | N  | N  |
| C | N | N | N | N | N | N | Y | N | N | Y  | N  | N  |
| D | Y | N | N | N | N | N | N | N | N | N  | N  | Y  |
| E | N | N | N | N | N | Y | N | Y | N | N  | N  | N  |
| F | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| G | N | N | N | Y | Y | N | N | N | N | N  | N  | N  |
| H | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| I | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| J | N | Y | N | N | N | N | N | N | N | N  | Y  | N  |
| K | N | N | N | N | N | N | Y | N | N | N  | N  | N  |
| L | N | N | N | N | N | N | N | N | Y | N  | N  | N  |
| M | N | N | N | Y | N | N | N | N | N | N  | N  | N  |

$K$  Categories

## Approaches for Labeling Training Images for Discriminative Learning

$N$  Images

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|
| A | N | Y | N | N | N | N | Y | N | Y | N  | N  | N  |
| B | N | N | Y | N | N | N | N | N | N | N  | N  | N  |
| C | N | N | N | N | N | N | Y | N | N | Y  | N  | N  |
| D | Y | N | N | N | N | N | N | N | N | N  | N  | Y  |
| E | N | N | N | N | N | Y | N | Y | N | N  | N  | N  |
| F | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| G | N | N | N | Y | Y | N | N | N | N | N  | N  | N  |
| H | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| I | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| J | N | Y | N | N | N | N | N | N | N | N  | Y  | N  |
| K | N | N | N | N | N | N | Y | N | N | N  | N  | N  |
| L | N | N | N | N | N | N | N | N | Y | N  | N  | N  |
| M | N | N | N | Y | N | N | N | N | N | N  | N  | N  |

K Categories

- **Exhaustive:** label all  $N$  images by all  $K$  categories  $\rightarrow N * K$  is large (does not scale!)
- Assume a very “efficient” person can label 100K images per day per concept
- 1M image training set
- 10K concepts
- $\rightarrow$  0.1 concept per day per person
- $\rightarrow$  **Need 100K person-days**

## Approaches for Labeling Training Images for Discriminative Learning

**$N$  Images**

**$K$  Categories**

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|
| A | N | Y | N |   |   |   | Y | N | Y |    |    | N  |
| B |   | N | Y | N |   |   |   |   |   | N  |    |    |
| C | N |   |   |   |   |   | Y | N |   | Y  | N  | N  |
| D | Y | N |   |   | N |   |   |   |   |    |    | Y  |
| E | N |   |   |   |   | Y | N | Y |   |    |    |    |
| F |   | N | N |   |   |   |   |   |   |    |    |    |
| G | N |   | N | Y | Y | N | N |   |   |    |    |    |
| H |   |   |   |   | N | N |   |   |   |    |    |    |
| I | N | N | N |   |   |   |   |   |   |    |    | N  |
| J | N | Y | N |   |   |   |   |   | N | N  | Y  | N  |
| K |   |   |   |   |   | N | Y | N |   |    | N  | N  |
| L | N |   |   |   |   |   |   | N | Y | N  | N  | N  |
| M |   | N | N | Y |   |   |   |   | N |    |    |    |

- **Exhaustive:** label all  $N$  images by all  $K$  categories  $\rightarrow N * K$  is large (does not scale!)

- **Binary:** label  $p$  positives ( $p \ll M$ ) and  $n$  negatives ( $n \ll N$ ) for each of  $K$  categories  $\rightarrow (p+n) * K \ll N * K$  (much better, still a lot)

Assume need  $2K$  labeled images per concepts

“Efficient” person can label 50 concepts per day

$\rightarrow$  **Need 200 person-days**

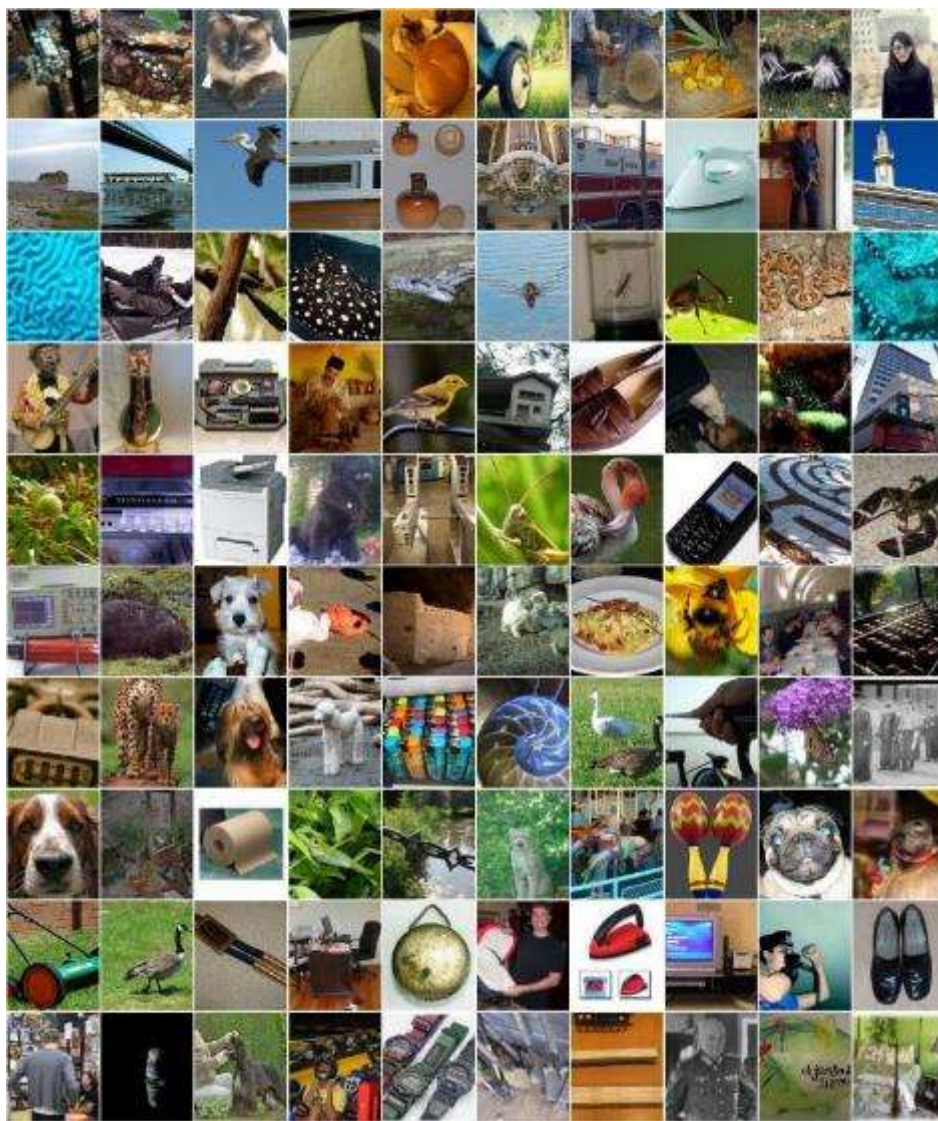
## Approaches for Labeling Training Images for Discriminative Learning

$N$  Images

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|
| A |   | Y |   |   |   |   | Y |   | Y |    |    |    |
| B |   |   | Y |   |   |   |   |   |   |    |    |    |
| C |   |   |   |   |   |   | Y |   |   | Y  |    |    |
| D | Y |   |   |   |   |   |   |   |   |    |    | Y  |
| E |   |   |   |   |   | Y |   | Y |   |    |    |    |
| F |   |   |   |   |   |   |   |   |   |    |    |    |
| G |   |   |   | Y | Y |   |   |   |   |    |    |    |
| H |   |   |   |   |   |   |   |   |   |    |    |    |
| I |   |   |   |   |   |   |   |   |   |    |    |    |
| J |   | Y |   |   |   |   |   |   |   |    | Y  |    |
| K |   |   |   |   |   |   | Y |   |   |    |    |    |
| L |   |   |   |   |   |   |   |   | Y |    |    |    |
| M |   |   |   | Y |   |   |   |   |   |    |    |    |

$K$  Categories

- **Exhaustive:** label all  $N$  images by all  $K$  categories  $\rightarrow N * K$  is large (does not scale!)
  - **Binary:** label  $p$  positives ( $p \ll M$ ) and  $n$  negatives ( $n \ll N$ ) for each of  $K$  categories  $\rightarrow (p+n) * K \ll N * K$  (much better, still a lot)
  - **Multi-Class:** label  $p$  positives for each of  $K$  categories and infer negatives  $\rightarrow p * K$
- $\rightarrow$  Need 100 person-days**

**Problem:** ImageNet ILSVRC Classification and Localization (CLS-LOC) Task

- 1,000 **object** categories that correspond to WordNet synsets
- No hierarchical “overlap” between synsets
- For any pair of synsets  $i$  and  $j$ ,  $i$  is not an ancestor of  $j$  in the WordNet hierarchy
- Every image has a **single image-level label** specifying the presence of **one object** category
- **No explicitly labeled negative** training images

## Example Semantic Labeling Problems in ImageNet 2015



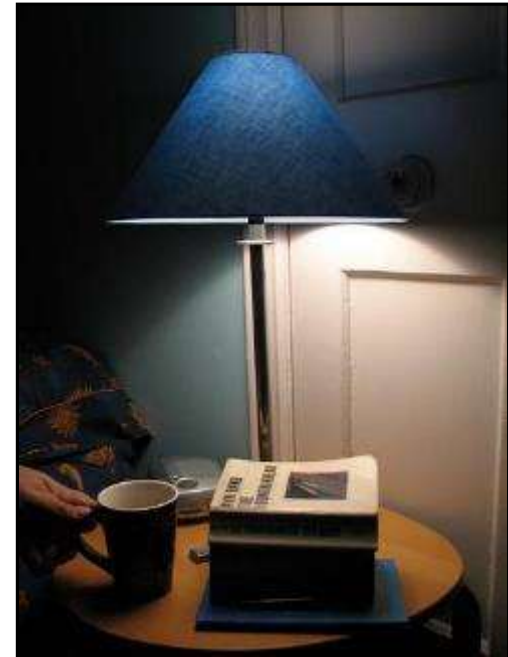
hay

→ Not horse

coffee mug

→ Not lamp shade

→ Not table lamp



# ImageNet CLS-LOC Training Images for Discriminative Learning

**$N$  Images**

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|
| A | N | Y | N | N | N | N | Y | N | Y | N  | N  | N  |
| B | N | N | Y | N | N | N | N | N | N | N  | N  | N  |
| C | N | N | N | N | N | N | Y | N | N | Y  | N  | N  |
| D | Y | N | N | N | N | N | N | N | N | N  | N  | Y  |
| E | N | N | N | N | N | Y | N | Y | N | N  | N  | N  |
| F | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| G | N | N | N | Y | Y | N | N | N | N | N  | N  | N  |
| H | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| I | N | N | N | N | N | N | N | N | N | N  | N  | N  |
| J | N | Y | N | N | N | N | N | N | N | N  | Y  | N  |
| K | N | N | N | N | N | N | Y | N | N | N  | N  | N  |
| L | N | N | N | N | N | N | N | N | Y | N  | N  | N  |
| M | N | N | N | Y | N | N | N | N | N | N  | N  | N  |

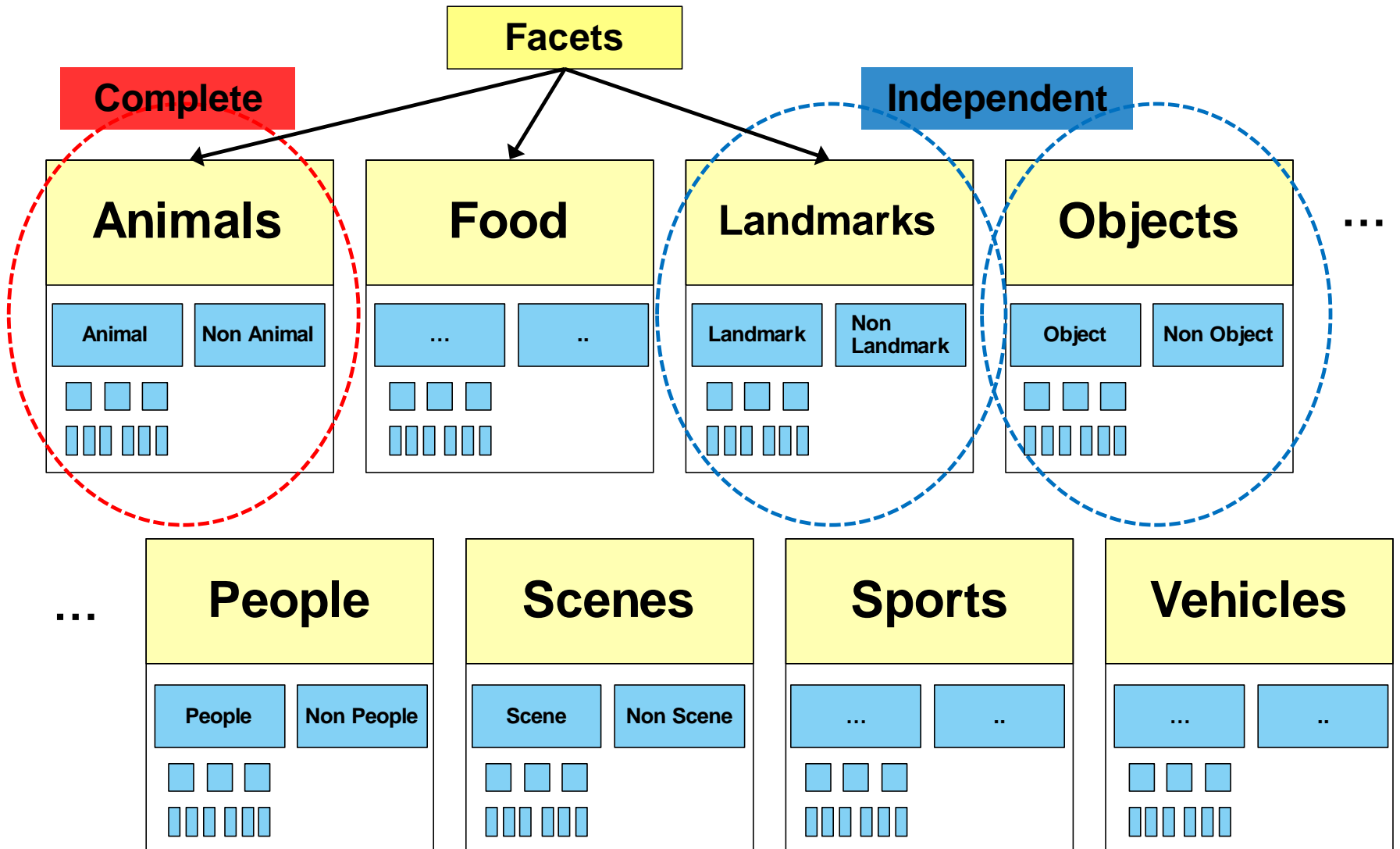
**$K$  Categories**

- **Exhaustive:** assign  $K$  labels to  $N$  images  $\rightarrow N * K$  is large (100K person-days for 10K labels and 1M images)
- **Binary:** label  $p$  positives ( $p \ll M$ ) and  $n$  negatives ( $n \ll N$ ) for each of  $K$  categories  $\rightarrow (p+n)*K \ll N*K$  (how to get negatives?)

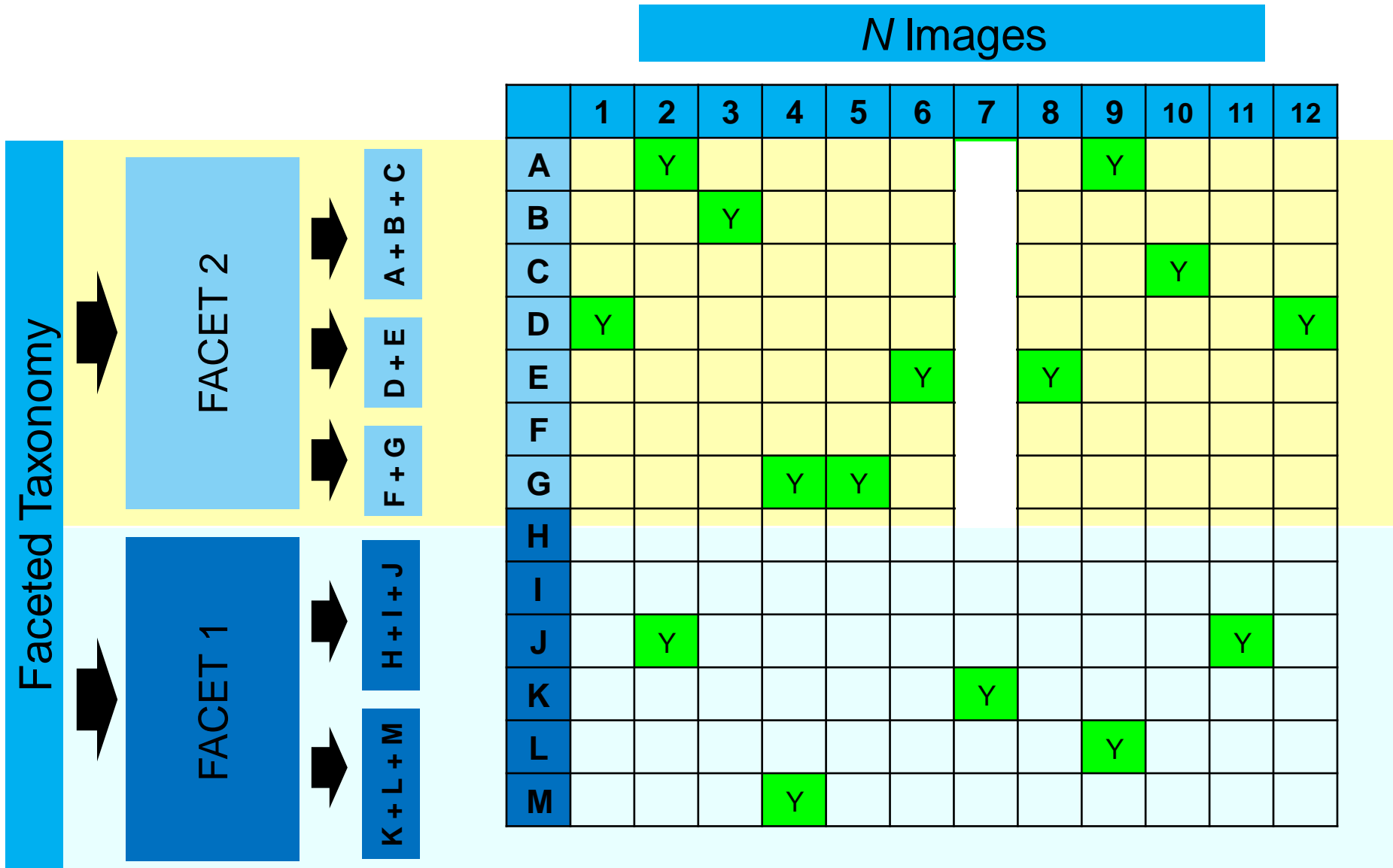
- **Multi-Class:** label  $p$  positives for each of  $K$  categories and infer negatives  $\rightarrow p*K$

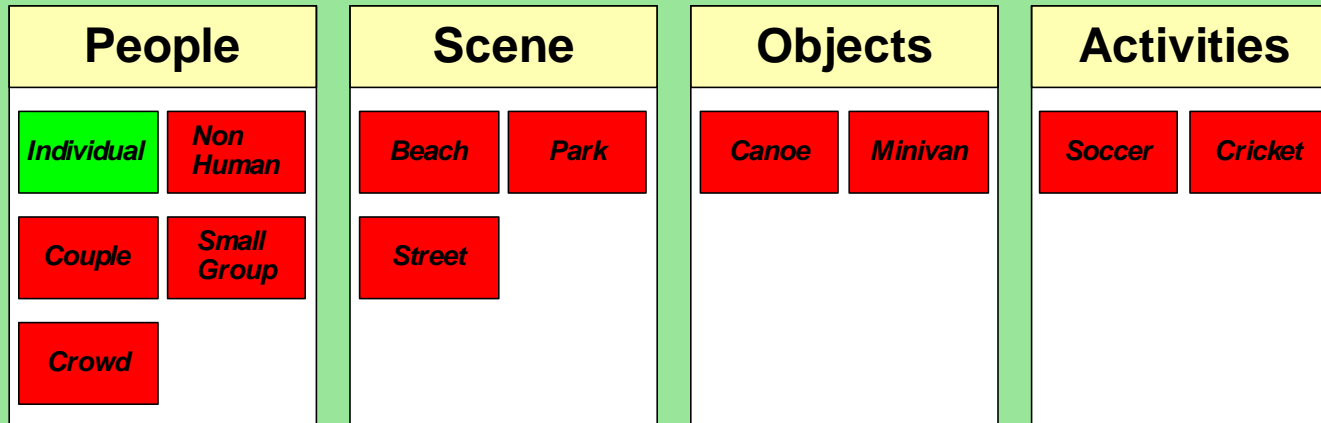
**$\rightarrow$  Need 100 person-days**

# Better Managing Visual Semantics using a Faceted Taxonomy



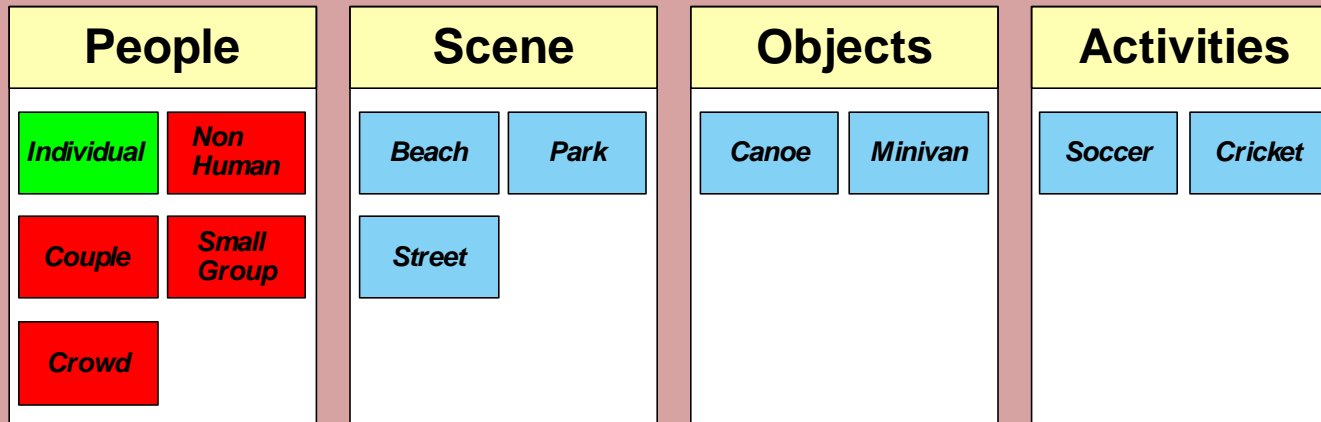
# Faceted Taxonomy Authoring and Maintenance ( $p \times K_i$ labels per FACET $i$ ):





## Traditional Taxonomy

- Explicitly labeled positives
- Negatives inferred based on assumption of mutually exclusivity



## Faceted Hierarchy

- Explicitly labeled positives
- Negatives inferred based on mutually exclusivity within facets only

# Results: Semantic Concepts are not Mutually Exclusive across Facets



[SCENE] *Beach*



[OBJECT] *Canoe*



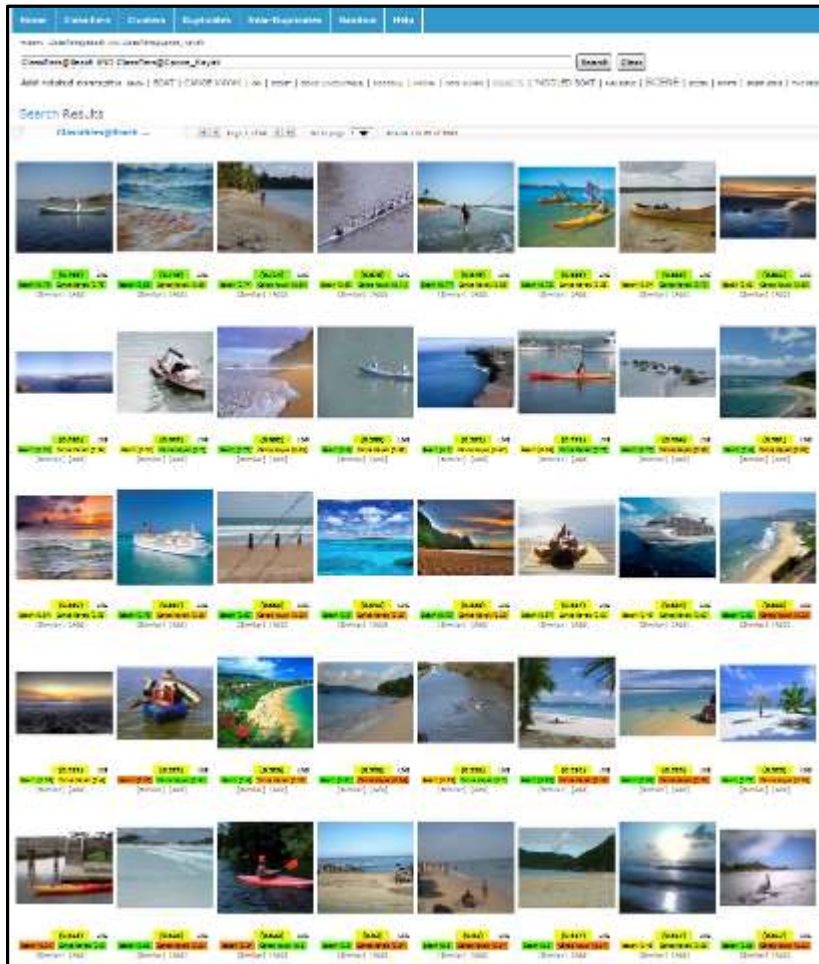
[SCENE] *Park*



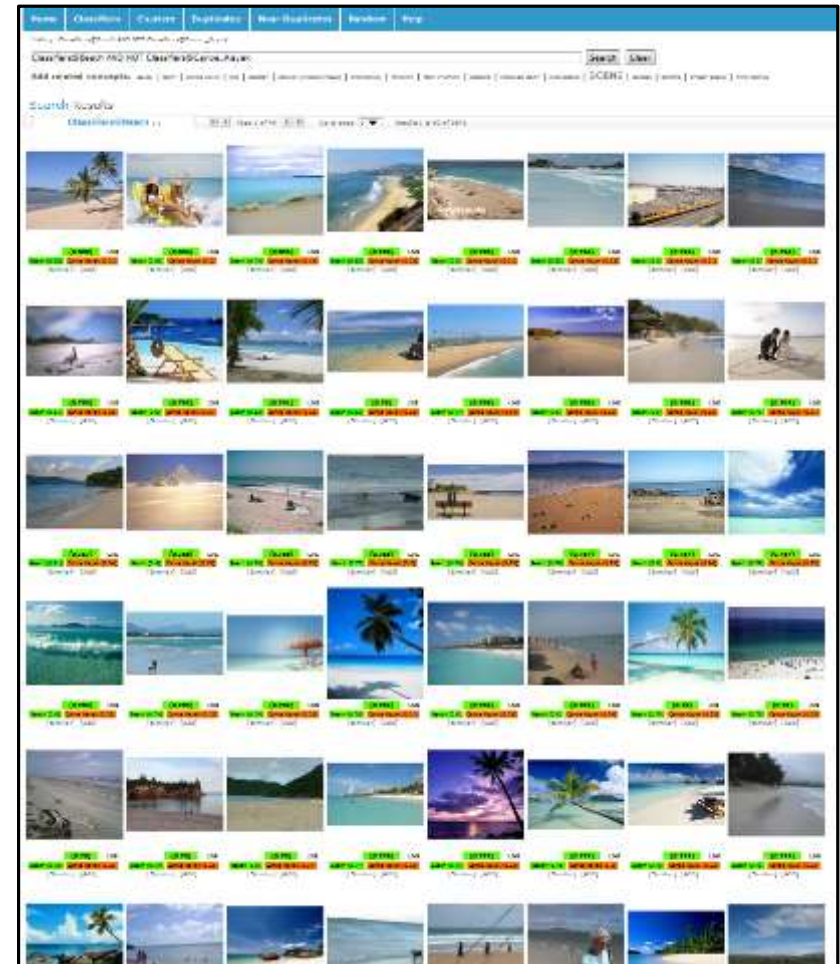
[PEOPLE] *Non Human*

# Example Search Results using Visual Semantic Faceted Hierarchy

## Beach AND Canoe



## Beach AND NOT Canoe

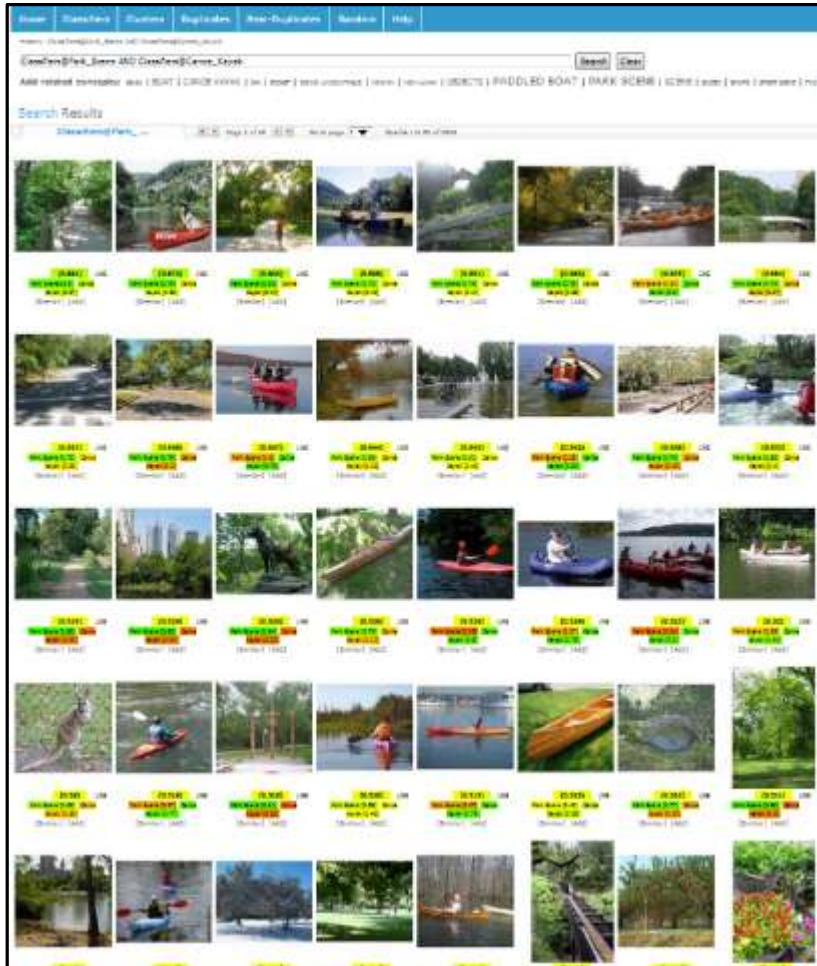


**Beach** and **Canoe** can be combined across facets → not mutually exclusive

**Beach** and **NOT Canoe** can also be combined across facets

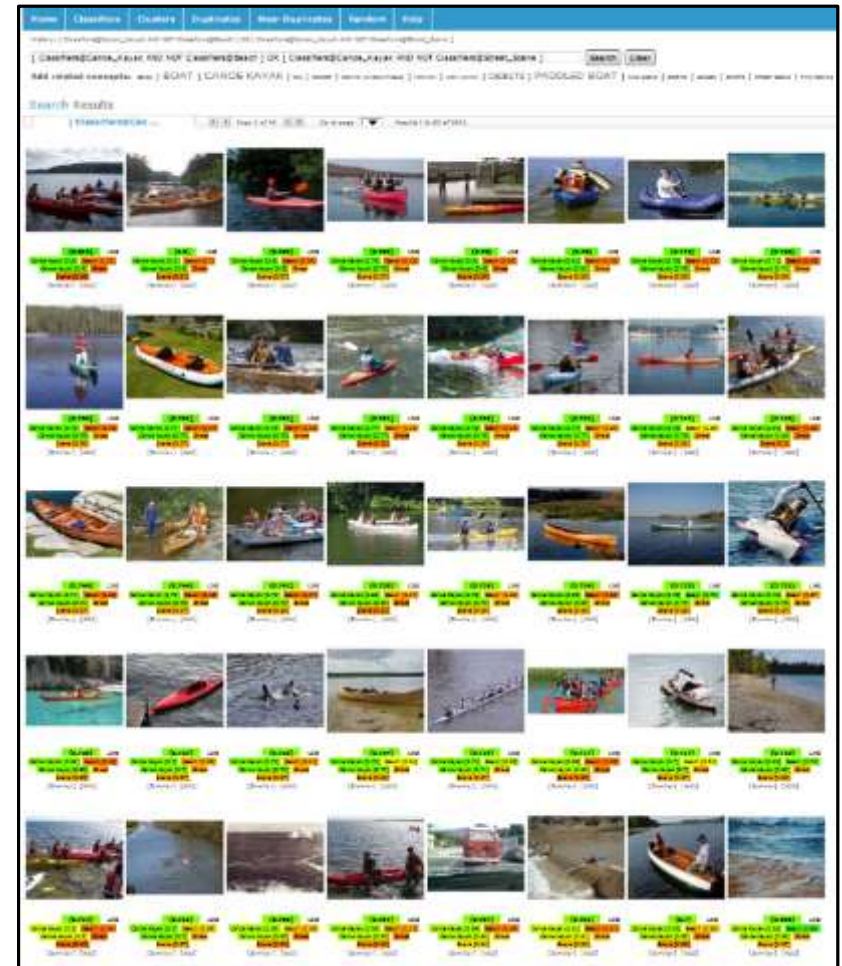
# Example Search Results using Visual Semantic Faceted Hierarchy

**Park AND Canoe**



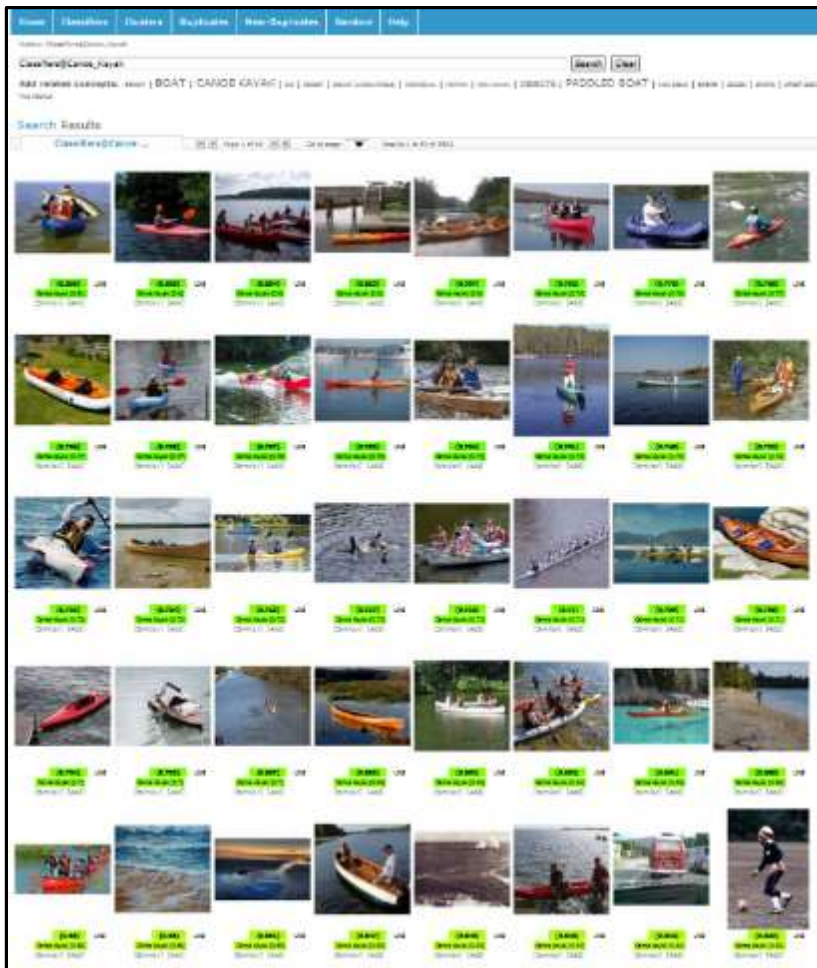
**Park AND Canoe** combines discriminative semantic concepts across facets

**[ Canoe AND NOT Beach ]  
OR [ Canoe AND NOT Street ]**

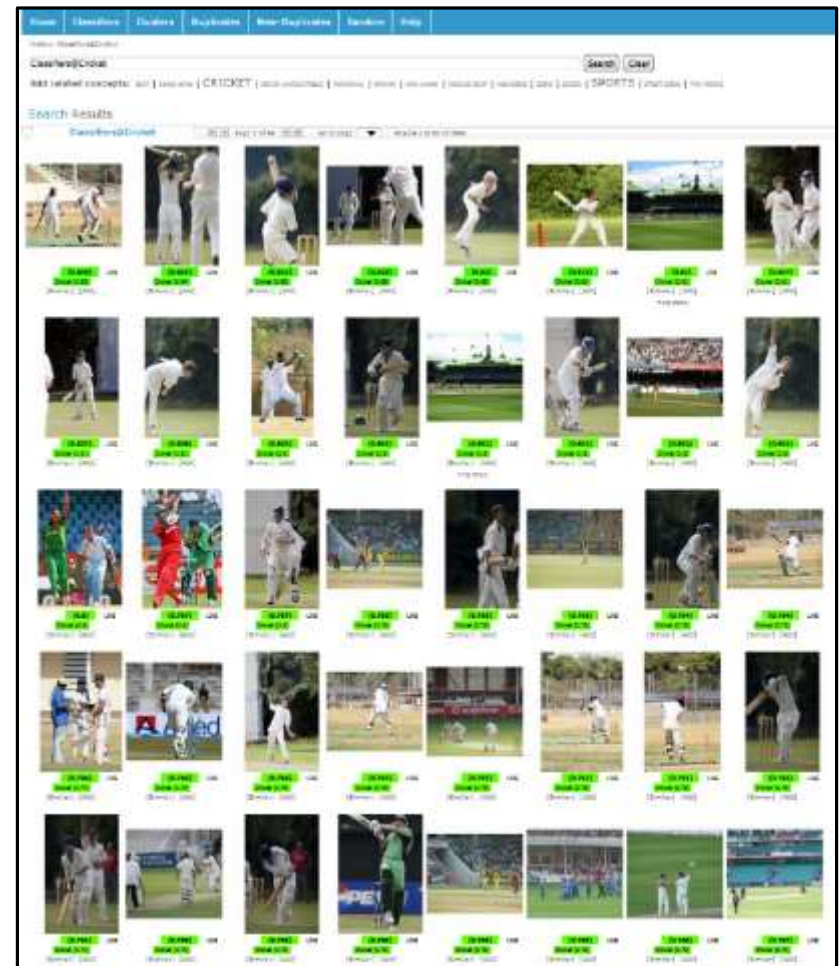


**[ Canoe AND NOT Beach ] OR [ Canoe AND NOT Street ]** → approx. **Park AND Canoe**

## [OBJECT] Canoe



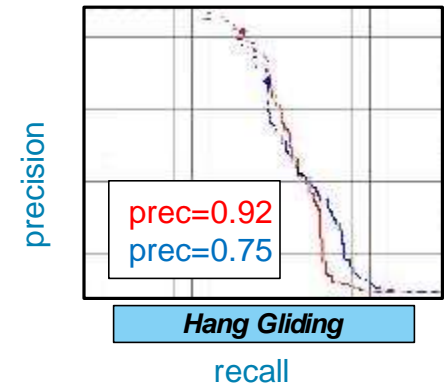
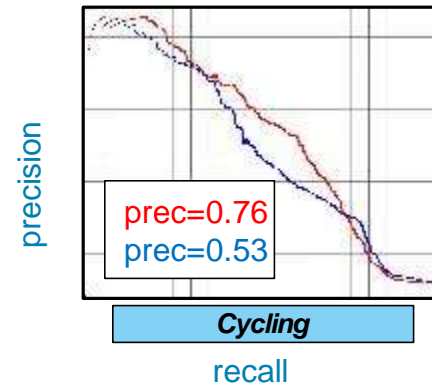
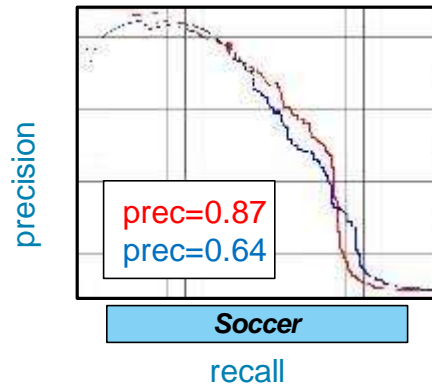
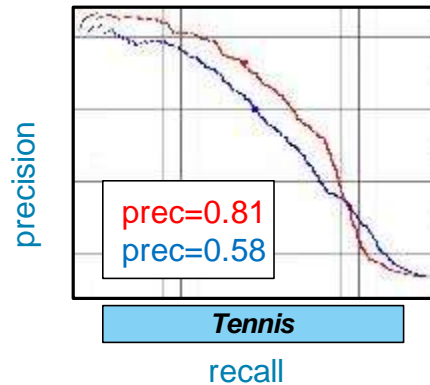
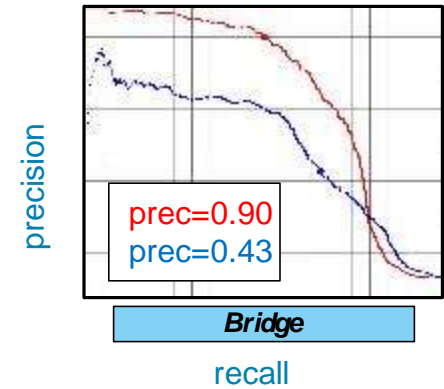
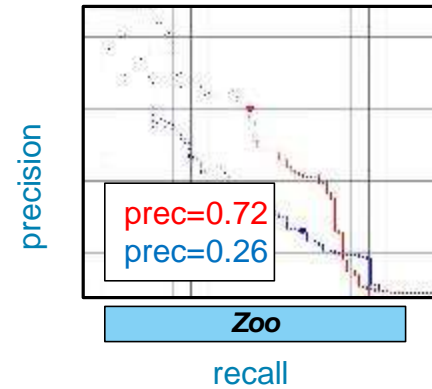
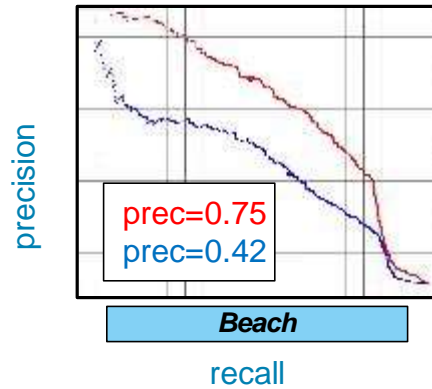
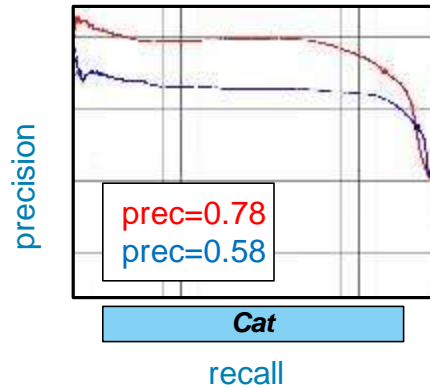
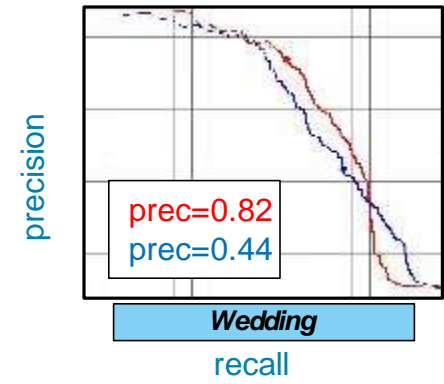
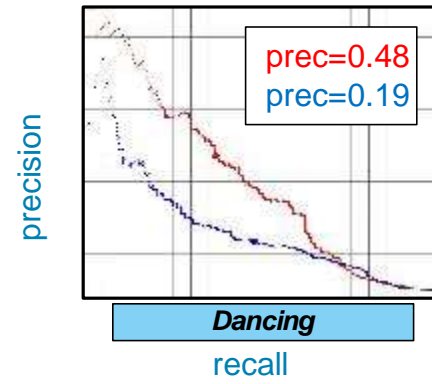
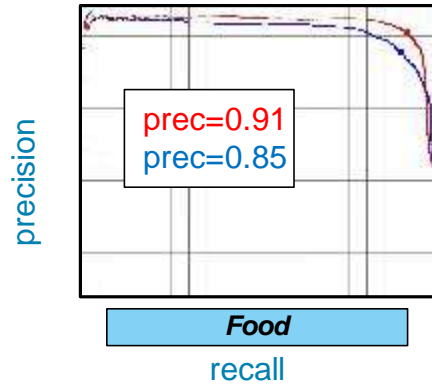
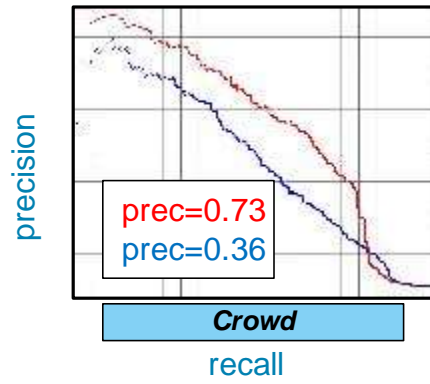
## [ACTIVITY] Cricket



Scene for *canoe* can be a *beach* or *park*,  
and *canoe* can have *people* or *no people*  
→ **not mutually exclusive across facets**

Scene for *cricket* can be a *park*,  
and *cricket* can have *people* or *no people*  
→ **not mutually exclusive across facets**

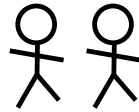
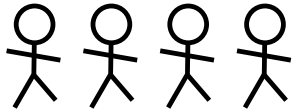
Results (582 labels): **Faceted Hierarchy (acc=0.95)** vs. **Taxonomy (acc=0.89)**



# Putting the User in-the-Loop in Visual Semantic Learning

Lots of supervision

Little supervision



Crowd Truth

Transfer Learning

Active Learning

Deep Learning

Unknown Data

Recognition

Models

Labels

Massive Crowd

Training Data

Unknown Data

Recognition

Models

Models

Labels (many)

Labels (few)

Training Data (lots)

Training Data (limited)

Unknown Data

Recognition

Models

Models

Models

Labels

Targeted Labels (few)

Training Data

Unknown Data

Recognition

Features

Models

Labels (few)

Training Data (unlabeled)

Training Data (limited)

Low cost – manual system

Low cost – semi-automatic

Low cost – semi-automatic

Minimizes cost – automatic

# Visual Learning Service – Simple Drag-and-Drop Training Example

## Upload Training Images

IMARS

HOME MODELS RESULTS CONTACT

Concept name: LIBERTY

48 images are shown in a grid, each with a file size (e.g., 0.4 MB, 7.4 KB, 87.7 KB, etc.).

## Learn Visual Classifier

IMARS

HOME MODELS RESULTS CONTACT

CONCEPT = STATUE\_OF\_LIBERTY ACCURACY = 0.992 ADD TRAINING IMAGES APPLY

Learned as positive for Statue\_of\_Liberty

80 images are shown in a grid, each with a confidence score (e.g., 0.810, 0.830, 0.852, etc.).

Demo

# Visual Learning can use Active Learning to Iteratively Train Models

## Unknown Images



IN



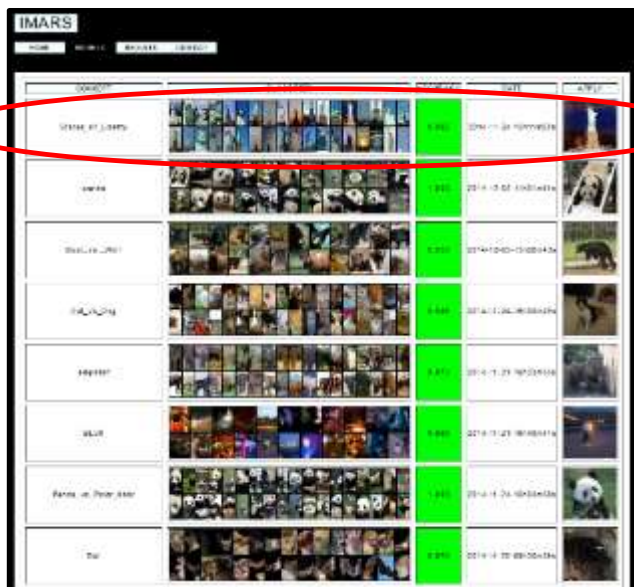
## Labels and Scores

Refine

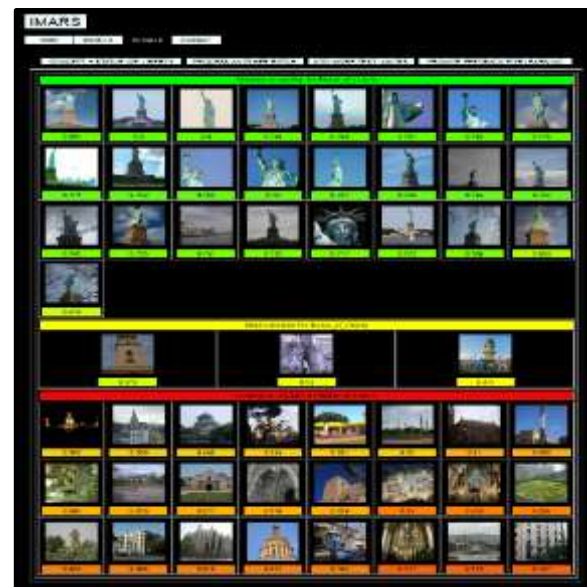


|    |                   |          |
|----|-------------------|----------|
| 14 | Statue_of_Liberty | 0.754721 |
| 15 | Statue_of_Liberty | 0.793386 |
| 16 | Statue_of_Liberty | 0.774210 |
| 17 | Statue_of_Liberty | 0.775921 |
| 18 | Statue_of_Liberty | 0.688645 |
| 19 | Statue_of_Liberty | 0.756334 |
| 20 | Statue_of_Liberty | 0.611773 |
| 21 | Statue_of_Liberty | 0.800345 |
| 22 | Statue_of_Liberty | 0.745605 |

OUT



Demo



## Key Messages:

**Image and video data is growing in volume and importance**

**Manual analysis is not scalable ... cannot keep up  
Vision system development has historically required deep expertise**

**Expert Analysis → Data-driven Visual Learning**

**Increasing availability of labeled training data  
Emergence of increasingly sophisticated visual learning algorithms**

**Focus needed on Visual Semantic Modeling**

**How to best model visual world (concepts and relationships)  
How to combine visual semantic modeling with learning systems**