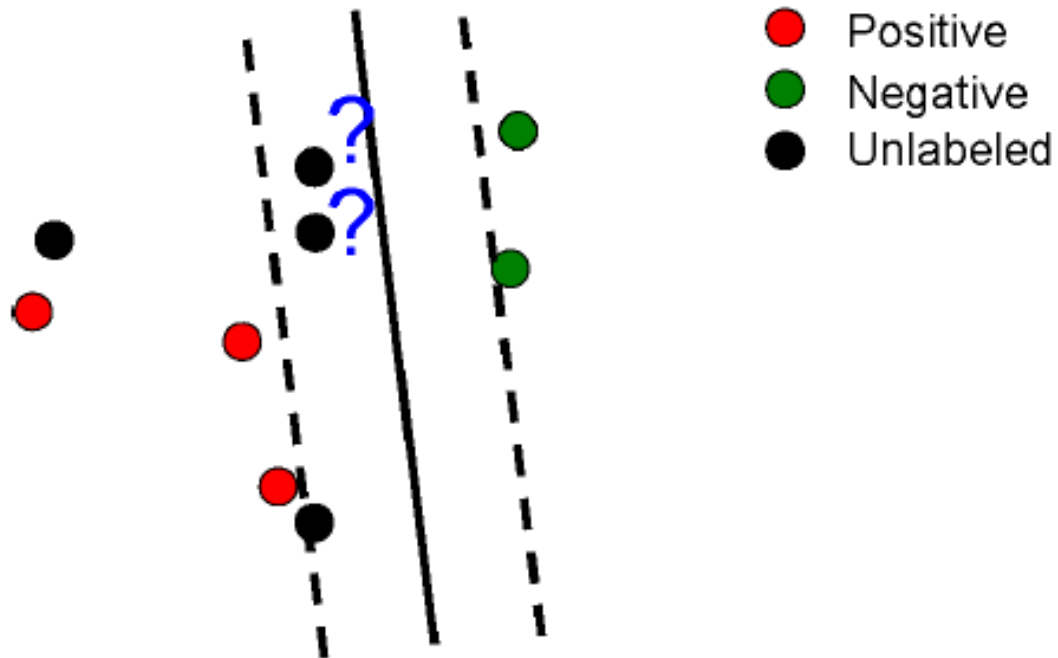


# What's It Going to Cost You?: Predicting Effort vs. Informativeness for Multi-Label Image Annotations

*Sudheendra Vijayanarasimhan and Kristen Grauman*

Deepti Ghadiyaram

# Active Learning



**Active selection is particularly complex for visual category learning.**

1) Real world images contain multiple objects



VS.



*Active learner must assess the value of an image containing some Unknown combination of categories.*

## 2) Different levels of information

Less  
expensive  
to obtain



More  
expensive  
to obtain

- ▶ **Weak labels:** informing about presence of an object



Phone



Phone



Not Phone

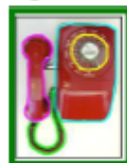


Not Phone

- ▶ **Strong labels:** outlines demarking the object



- ▶ **Stronger labels:** informing about labels of parts of objects

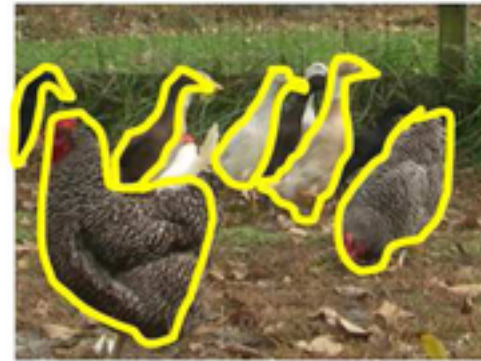


*Active learner must specify what type of annotation is currently most Helpful.*

3) Manual Effort dependent on annotation type and image content.



**Low effort**



**High effort**

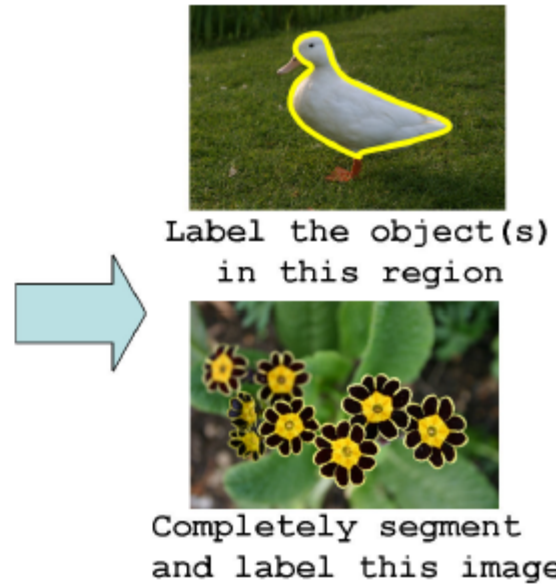
*Active learner should take into account the actual manual effort required to label the images.*

## PROBLEM STATEMENT

*How do we effectively learn from a mixture of strong and weak labels and select the most promising **{image + annotation type}** by balancing the **value of a new annotation** against the **time taken to receive it**.*

# PROPOSED APPROACH





(b) Unlabeled and partially labeled examples to survey

(c) Actively chosen queries sent to annotators

# 3 types of annotations



Name an object



What class is this region?

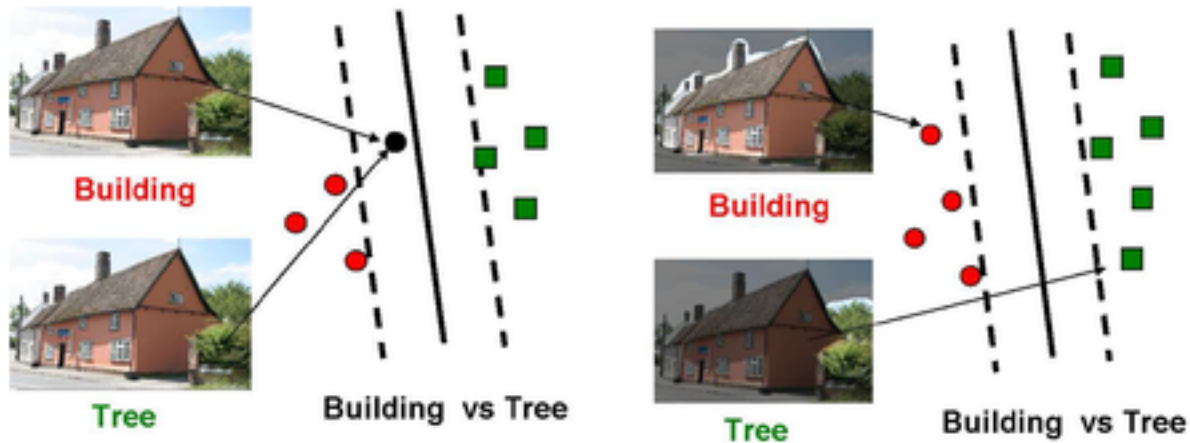


Segment this image

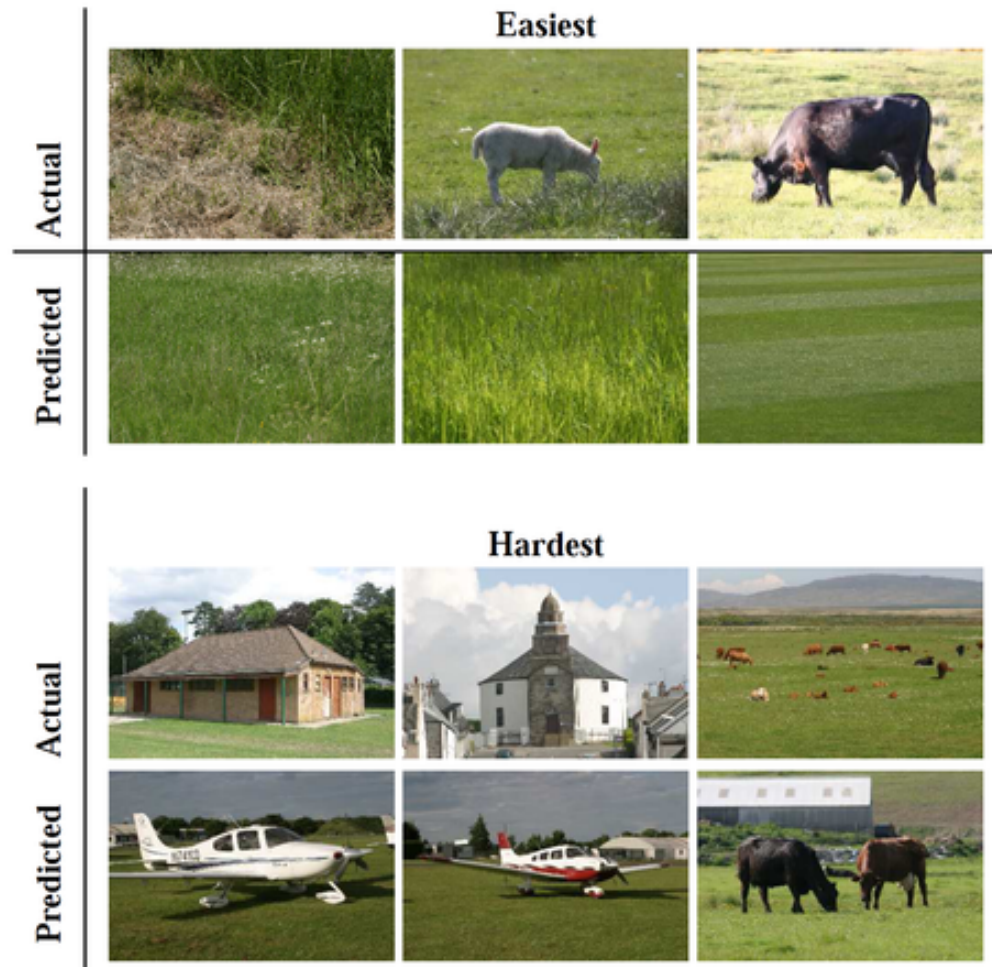
# **Key Ideas of the approach**

# 1. Multi-Label multiple-instance learning.

Multi-label Set Kernel based classifier



## 2. Predicting the cost of annotation based on image content.



### 3. Predicting the informativeness of an annotation ( $\mathbf{z}$ )

Change in the *Total Misclassification Risk* resulted from  $\mathbf{z}$ - Cost for obtaining that annotation

$$\begin{aligned} VOI(\mathbf{z}) &= T(\mathcal{X}_L, \mathcal{X}_U, \mathcal{X}_P) - T(\hat{\mathcal{X}}_L, \hat{\mathcal{X}}_U, \hat{\mathcal{X}}_P) \quad (6) \\ &= \mathcal{R}(\mathcal{X}_L) + \mathcal{R}(\mathcal{X}_U) + \mathcal{R}(\mathcal{X}_P) \\ &\quad - \left( \mathcal{R}(\hat{\mathcal{X}}_L) + \mathcal{R}(\hat{\mathcal{X}}_U) + \mathcal{R}(\hat{\mathcal{X}}_P) \right) - \mathcal{C}(\mathbf{z}), \end{aligned}$$

# EXPERIMENTS

Label 3: Building, Sky



Label 4: Aeroplane, Grass, Sky



Label 5: Cow, Grass, Mountain





# **1. ACTIVE V/S RANDOM SELECTION**

# Experimental Setup

## Active Learner

<b>Task</b>	<b>Actively</b> choose between tags, regions and complete segmentation
<b>Training &amp; Validation Data</b>	Super pixel segments
<b>Test Data</b>	Ground Truth Segments
<b>Number of iterations(number of samples added to classifier)</b>	30
<b>Initial Training set size</b>	3 ( 1 bag/image per class)
<b>Unlabeled Data</b>	87

## Random Selection

<b>Task</b>	<b>Randomly</b> choose between tags, regions and complete segmentation
<b>Training &amp; Validation Data</b>	Super pixel segments
<b>Test Data</b>	Ground Truth Segments
<b>Number of iterations(number of samples added to classifier)</b>	30
<b>Initial Training set size</b>	3 ( 1 bag/image per class)
<b>Unlabeled Data</b>	87

Risk /Cost prediction is not a part of random selection

# **ACTIVE LEARNER IN ACTION**



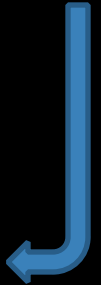
Learner: What class is this region?  
Oracle: Sky



Learner: What class is this region?  
Oracle: Grass



Learner: Name an object  
Oracle: Grass



Learner: What class is this region?  
Oracle: Sky



Learner: What class is this region?  
Oracle: Sky

# **RANDOM SELECTION**



Random: What class is this region?  
Oracle: Sky



Random: Segment this image



Random: Name an object  
Oracle: Sky



Random: Segment this image



Random: What class is this region?  
Oracle: Grass

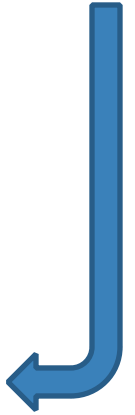
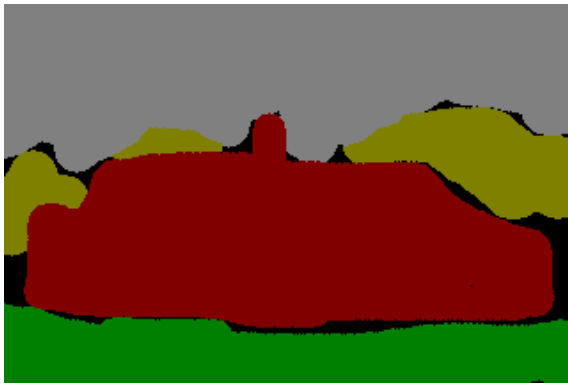




Random: What class is this region?  
Oracle: Sky

Random: Segment this image

Random: Name an object  
Oracle: Sky



Random: Segment this image

Random: What class is this region?  
Oracle: Grass

(At the end of 30 iterations)



**Observation:** Requesting for complete segmentation of few images doesn't *necessarily* yield to better classification performance or reduction in risk.

# Active Learning v/s Random sampling

At the end of 30 iterations...

	Active	Random ( average of 5 runs).
Execution time ( in secs)	261	<b>0</b>
Cost	<b>288</b>	326.967
Avg AUROI	<b>0.979867</b>	0.966164
Risk	<b>42.4003</b>	44.07574

# Information as a function of iterations

Steep gain from the first few picks.

The most informative selections are made in the first few iterations.

Consistent gain in information in active selection.

For a given number of instances, active learner ensures the best possible system.

2) WHAT TYPE OF ANNOTATION TO REQUEST?



1) Name an object



2) What class is this region?



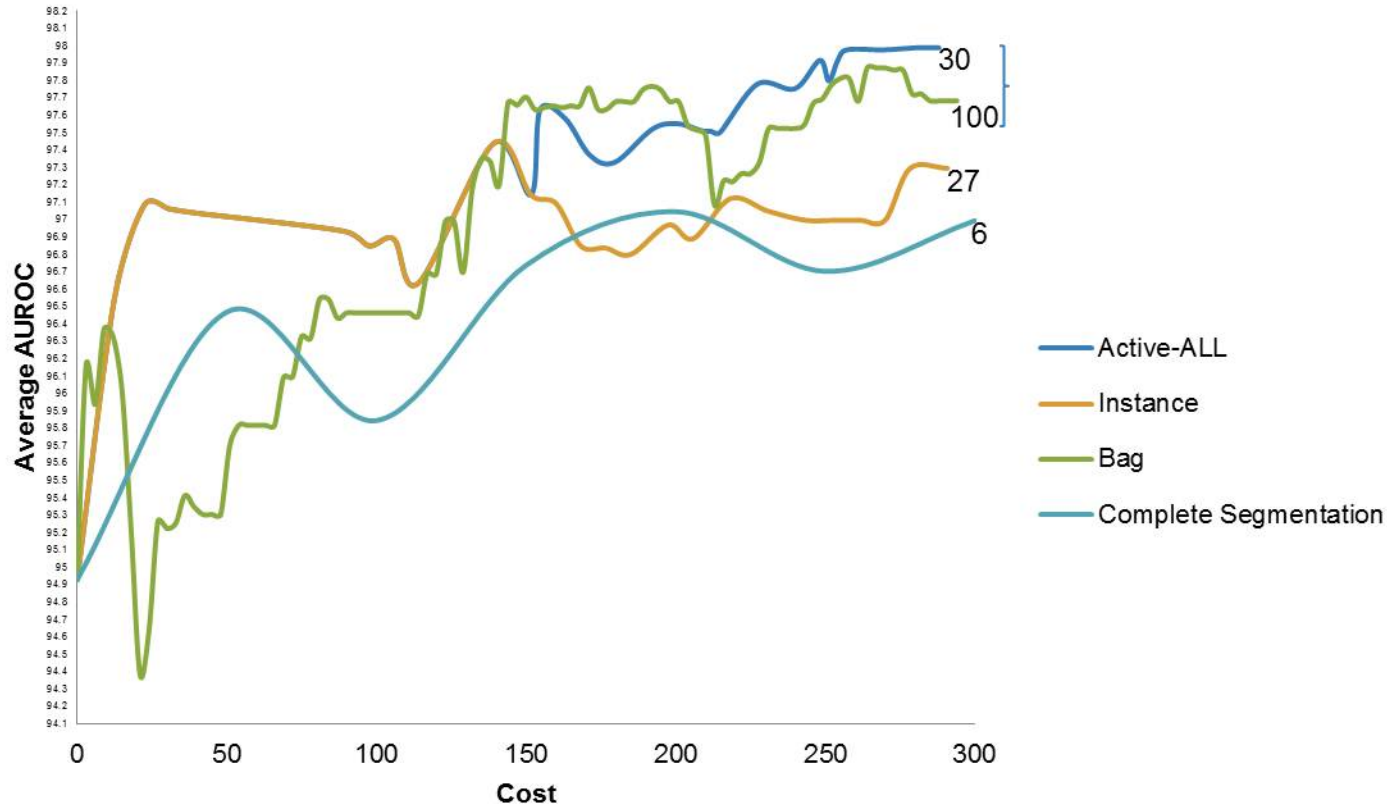
3) Segment this image

4) Any of the above

# Experimental Setup

<b>Task</b>	<b>Actively</b> choose between Case#1 Only tags Case#2 : Only regions Case#3 Only Complete segmentation Case#4 Any of the above
<b>Training &amp; Validation Data</b>	Super pixel segments
<b>Test Data</b>	Ground Truth Segments
<b>Number of iterations ( number of samples added to the training data)</b>	Variable ( Ranges between 6-100 for each of the above cases)
<b>Initial Training set size</b>	3 ( 1 bag per class)
<b>Unlabeled Data</b>	87

## Comparing different annotation types



### Observation

- A combination of annotation types is more beneficial than a fixed annotation type.

At the end of 30 iterations..

	Active-ALL	Active-bag	Active-instance	Complete segmentation
Execution time (in secs)	261	121	<b>23</b>	49
Cost	288	<b>84</b>	331.667	789
Global Mean Accuracy	<b>89.6774</b>	88.3871	87.7419	87.0968
Avg AUROI	<b>0.979867</b>	0.965411	0.973077	0.9723
Risk	<b>42.4003</b>	39.7979	45.003	46.8725

Execution time is **proportional** to the number of instances/bags to be considered.  
Better accuracy can be achieved by combining different types of annotations.

# 3. Ground Truth Segments v/s Super pixel Segments

**Key Idea:** Use ground truth segments instead of super pixel segments while training and testing MIML classifier.

**Aim:** To understand the upper bound of the active learner ( limitations from using super pixel segments)



# How noisy are the super pixel segments?

<b>For 90 images..</b>		Number of instances.
	Ground truth Segments	152
	Super pixel segments	272

Segment Type	Image	Instance Labels					
<b>Super pixel</b>	<b>I1</b>	5	5	3	5	5	5
<b>True</b>	<b>I1</b>	3	5				
<b>Super pixel</b>	<b>I2</b>	3	5	5	5	4	5
<b>True</b>	<b>I2</b>	3	4	5			

### True Learner

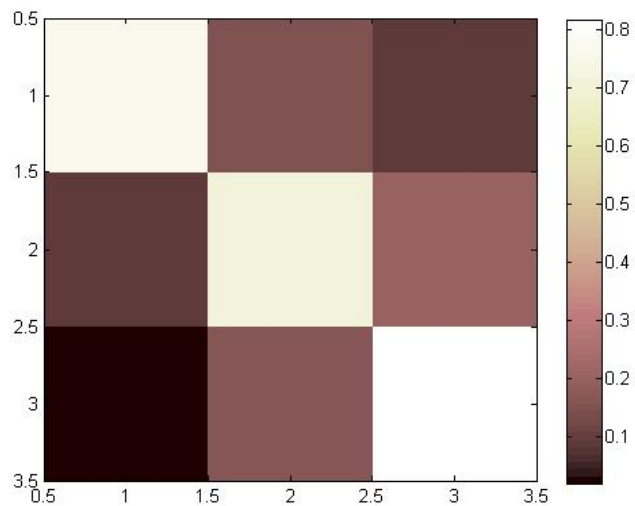
Task	<b>Actively</b> choose between tags, regions and complete segmentation
Training & Validation Data	<b>Ground Truth segments</b>
Test Data	<b>Super pixel Segments</b>
Number of iterations (number of samples added to classifier)	30
Initial Training set size	3 ( 1 bag per class)
Unlabeled Data	87

### Noisy Learner

Task	<b>Actively</b> choose between tags, regions and complete segmentation
Training & Validation Data	<b>Super pixel segments</b>
Test Data	<b>Super pixel Segments</b>
Number of iterations (number of samples added to classifier)	30
Initial Training set size	3 ( 1 bag per class)
Unlabeled Data	87

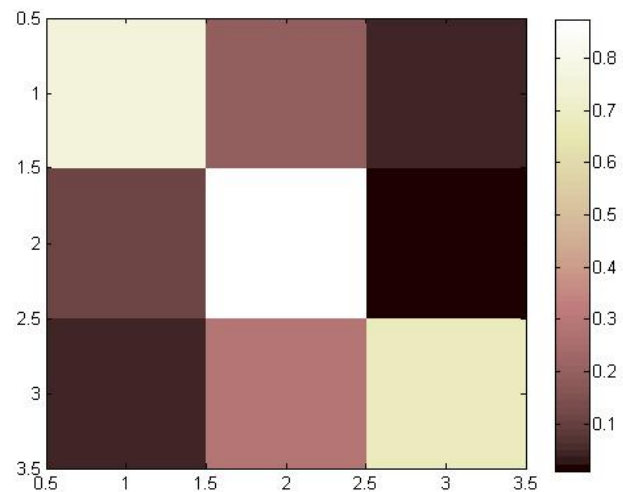
# Confusion matrix

## Noisy Learner



Per class mean accuracy: 76.2958

## True Learner



Per class mean accuracy: 77.0973

(At the end of 30 iterations)



# True Learner v/s Noisy Learner : Case1

At the end of 30 iterations...

	True Learner	Noisy Learner
Total Execution time ( in secs)	<b>53</b>	195
Cost	<b>128.000000</b>	194.166667
Global Mean Accuracy	<b>78.8235</b>	75.2941
Avg AUROI	<b>0.960086</b>	0.936850
Risk	<b>84.6304</b>	98.0278

- Lesser number of instances to process during training and validation for the true learner

### True Learner

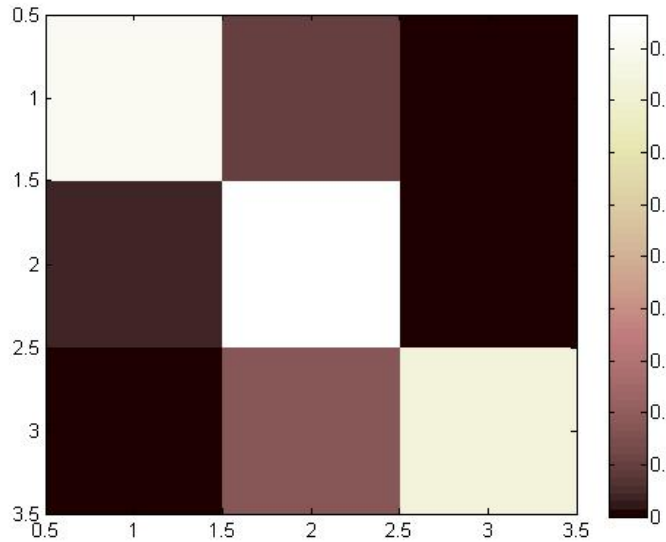
Task	<b>Actively</b> choose between tags, regions and complete segmentation
Training & Validation Data	<b>Ground Truth segments</b>
Test Data	<b>Ground Truth segments</b>
Number of iterations (number of samples added to classifier)	30
Initial Training set size	3 ( 1 bag per class)
Unlabeled Data	87

### Noisy Learner

Task	<b>Actively</b> choose between tags, regions and complete segmentation
Training & Validation Data	<b>Super pixel segments</b>
Test Data	<b>Ground Truth segments</b>
Number of iterations (number of samples added to classifier)	30
Initial Training set size	3 ( 1 bag per class)
Unlabeled Data	87

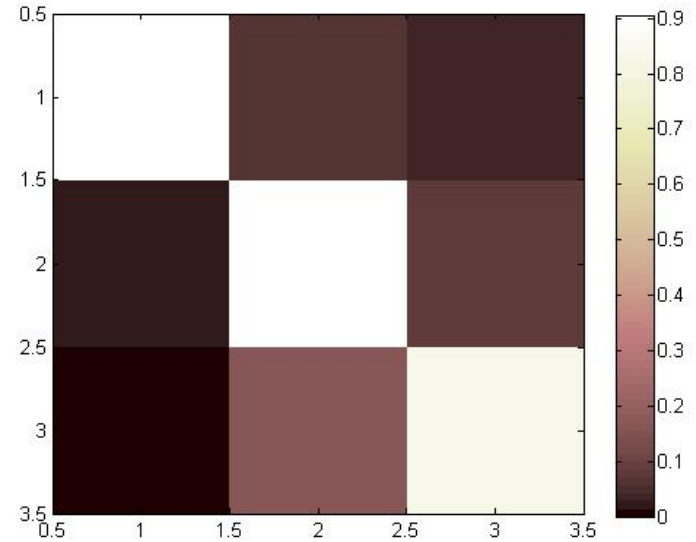
# True Learner v/s Noisy Learner : Case2

True Learner



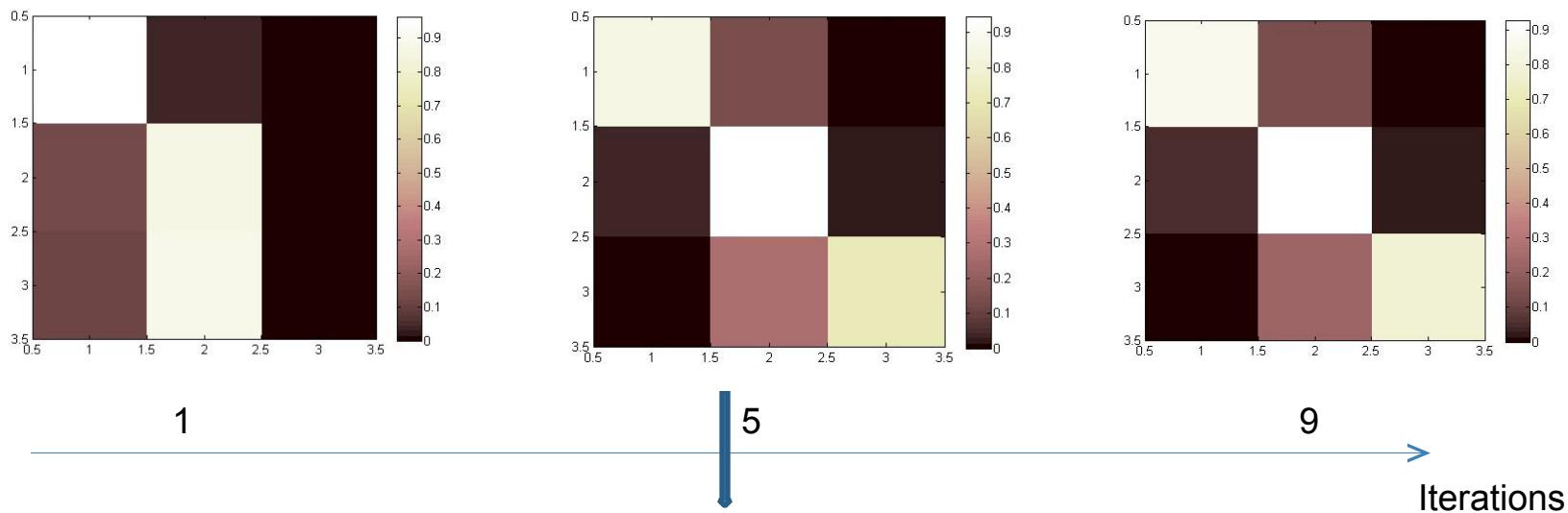
Per class mean accuracy: 90.012

Noisy Learner



Per class mean accuracy: 88.1252

# Initial misclassification with true segments



Learner: What class is this region?  
Oracle: Grass



The first four iterations have always misclassified the images with class label=5, this is changed when the above image is added.



# **True Learner v/s Noisy Learner : Case2**

True Learner v/s Noisy Learner : Case2

At the end of 30 iterations...

	True Learner	Noisy Learner
Execution time ( secs)	136	261
Cost	128.	288
Total Mean Accuracy	<b>91.6129</b>	89.6774
Avg AUROI	0.987389	0.979867
Risk	38.1211	42.4003

### True Learner

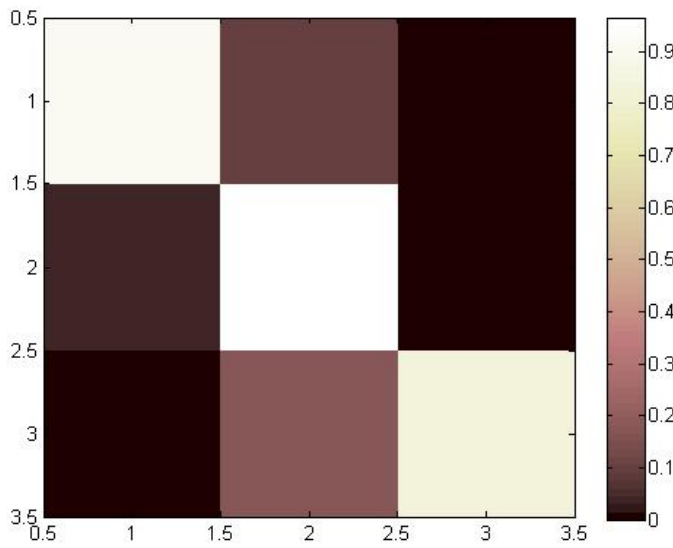
Task	<b>Actively</b> choose between tags, regions and complete segmentation
Training & Validation Data	Ground Truth segments
Test Data	Ground Truth segments
Number of iterations (number of samples added to classifier)	30
Initial Training set size	3 ( 1 bag per class)
Unlabeled Data	87

### Noisy Learner

Task	<b>Actively</b> choose between tags, regions and complete segmentation
Training & Validation Data	Super pixel segments
Test Data	Super pixel segments
Number of iterations (number of samples added to classifier)	30
Initial Training set size	3 ( 1 bag per class)
Unlabeled Data	87

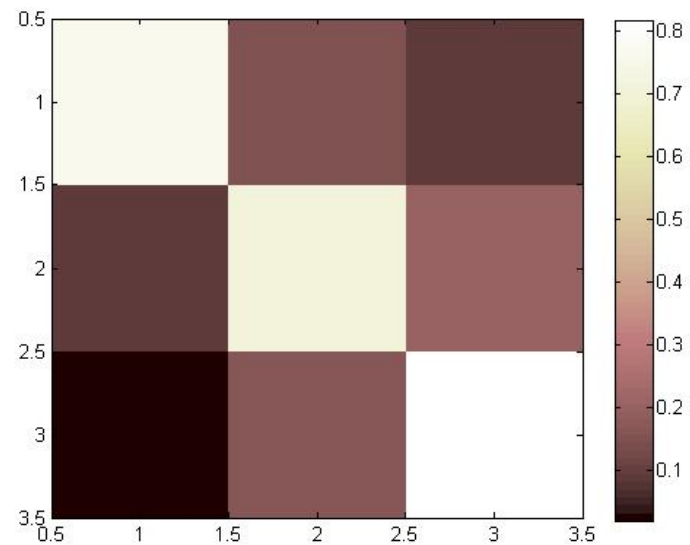
- Noisy Learner reaches the cost of that of a true learner is in 30 iterations, in only 18 iterations of learning.
- True learner outperforms noisy learner by a very high margin.
- ***This test case demonstrates the upper bound of the active learner***

# Confusion Matrix



Per class mean accuracy: 90.012

(At the end of 30 iterations)



Per class mean accuracy: 76.2958

# Active Learning: True v/s Noisy Learner

At the end of 30 iterations...

	True Learner	Noisy Learner
Execution time ( secs)	<b>53</b>	195
Cost	<b>128</b>	194.167
Total Mean Accuracy	<b>91.6129</b>	75.2941
Avg AUROI	<b>0.987389</b>	0.936850
Risk	<b>38.1211</b>	98.0278

True learner is faster, more accurate and has lesser total risk.

Test Cases	Risk	
Train: Ground Truth Segments Test: Ground Truth Segments	38.1211	→ Best case
Train: Super Pixel Segments Test: Super Pixel Segments	98.0278	→ Real System
Train: Ground Truth Segments Test: Super Pixel Segments	84.6304	
Train: Super Pixel Segments Test: Ground Truth Segments	42.4003	

- Super pixel segmentation imposes limitations on the performance of the active learner.
- The total risk of the system is lesser when ground truth segments are used.
- The computational cost of the system also varies with the correctness of the segments.

## 4) The contribution of each variable to the VOI for an annotation

$$\begin{aligned} VOI(\mathbf{z}) &= T(\mathcal{X}_L, \mathcal{X}_U, \mathcal{X}_P) - T(\hat{\mathcal{X}}_L, \hat{\mathcal{X}}_U, \hat{\mathcal{X}}_P) \quad (6) \\ &= \mathcal{R}(\mathcal{X}_L) + \mathcal{R}(\mathcal{X}_U) + \mathcal{R}(\mathcal{X}_P) \\ &\quad - \left( \mathcal{R}(\hat{\mathcal{X}}_L) + \mathcal{R}(\hat{\mathcal{X}}_U) + \mathcal{R}(\hat{\mathcal{X}}_P) \right) - \mathcal{C}(\mathbf{z}), \end{aligned}$$

- Importance of cost prediction : **C(z)**
- Effect of the risk parameter :  $r_L$
- VOI from Labeled data
- VOI from Partially labeled data.
- VOI from Unlabeled data.



## Importance of cost prediction.

Effect of the risk parameter.

VOI from Labeled data.

VOI from Partially labeled data.

VOI from Unlabeled data.

$$\begin{aligned} VOI(\mathbf{z}) &= T(\mathcal{X}_L, \mathcal{X}_U, \mathcal{X}_P) - T(\hat{\mathcal{X}}_L, \hat{\mathcal{X}}_U, \hat{\mathcal{X}}_P) \quad (6) \\ &= \mathcal{R}(\mathcal{X}_L) + \mathcal{R}(\mathcal{X}_U) + \mathcal{R}(\mathcal{X}_P) \\ &\quad - \left( \mathcal{R}(\hat{\mathcal{X}}_L) + \mathcal{R}(\hat{\mathcal{X}}_U) + \mathcal{R}(\hat{\mathcal{X}}_P) \right) - \boxed{\mathcal{C}(\mathbf{z})} \end{aligned}$$

# Experimental Setup

<b>Task</b>	<b>Actively</b> choose between tags, regions and complete segmentations.  Case#1 – Without Annotation Cost. Case#2 - With Annotation Cost.
<b>Training &amp; Validation Data</b>	Super pixel segments
<b>Test Data</b>	Ground Truth Segments
<b>Number of iterations ( number of samples added to the training data)</b>	Variable ( Ranges between 6-100 for each of the above cases)
<b>Initial Training set size</b>	3 ( 1 bag per class)
<b>Unlabeled Data</b>	87

Iterations=30	Without Annotation Cost	<b>With Annotation Cost</b>
Total Cost	<b>756</b>	288
Global Mean Accuracy	86.4516	89.6774

- Without  $C(z)$ , VOI is measured only in terms of estimate risk of misclassification.
- Having the penalty on cost is useful in making better choices.

Importance of cost prediction.

**Effect of the risk parameter.**

VOI from Labeled data.

VOI from Partially labeled data.

VOI from Unlabeled data.

# The effect of the risk parameter ( $r_L$ )

$$\mathcal{R}(\mathcal{X}_L) = \sum_{X_i \in \mathcal{X}_L} \sum_{l \in L_i} r_l (1 - p(l|X_i))$$

$$\mathcal{R}(\mathcal{X}_U) = \sum_{X_i \in \mathcal{X}_U} \sum_{l=1}^C r_l (1 - p(l|X_i)) \Pr(l|X_i),$$

$$\begin{aligned} \mathcal{R}(\mathcal{X}_P) &= \sum_{X_i \in \mathcal{X}_P} \sum_{l \in L_i} r_l (1 - p(l|X_i)) \\ &+ \sum_{l \in U_i} r_l (1 - p(l|X_i)) p(l|X_i), \end{aligned}$$

30

100

### **Observations:**

- Without  $r_L$ , the effect of the risk estimations is negligible and choice of instances is dominated by  $C(z)$
- We get better accuracy with lesser number of instances when risk estimation is also included.
- ***Thus, an equal contribution of both cost estimation and risk estimation leads to more informative learning.***

Importance of cost prediction.

Effect of the risk parameter.

VOI from Labeled data  $R(\mathcal{X}_L)$

VOI from Partially labeled data  $R(\mathcal{X}_P)$

VOI from Unlabeled data  $R(\mathcal{X}_U)$

$$\begin{aligned} VOI(\mathbf{z}) &= T(\mathcal{X}_L, \mathcal{X}_U, \mathcal{X}_P) - T(\hat{\mathcal{X}}_L, \hat{\mathcal{X}}_U, \hat{\mathcal{X}}_P) \quad (6) \\ &= \mathcal{R}(\mathcal{X}_L) + \mathcal{R}(\mathcal{X}_U) + \mathcal{R}(\mathcal{X}_P) \\ &\quad - \left( \mathcal{R}(\hat{\mathcal{X}}_L) + \mathcal{R}(\hat{\mathcal{X}}_U) + \mathcal{R}(\hat{\mathcal{X}}_P) \right) - \mathcal{C}(\mathbf{z}), \end{aligned}$$

Iterations=30	Without R(L)	Without R(P)	Without R(U)
Total Cost	<b>394</b>	84	288

- Exclusion of R(L) leads to high risk and high total cost.
- This result shows the real contribution of each pool to the decision making.
- Since most changes are happening to the labeled pool of data, with every iteration, it has the highest contribution to the VOI.



## 5) ANNOTATION DATA



Interface on  
Mechanical Turk



<b>Images</b>	240
<b>Users</b>	~70

# Most picked v/s Least picked Images (Avg=22)



31



13

# Least v/s most agreed upon Images (Avg= ~27)



3.3545



149.32