

---

# Timely Object Recognition

---

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

In a large visual multi-class detection framework, the timeliness of results can be crucial. Our method for *timely* multi-class detection aims to give the best possible performance at any single point after a start time; it is terminated at a deadline time. Toward this goal, we formulate a dynamic, *closed-loop* policy that infers the contents of the image in order to decide which detector to deploy next. In contrast to previous work, our method significantly diverges from the predominant greedy strategies, and is able to learn to take actions with deferred values. We evaluate our method with a novel *timeliness* measure, computed as the area under an Average Precision vs. Time curve. Experiments are conducted on the eminent PASCAL VOC object detection dataset. If execution is stopped when only half the detectors have been run, our method obtains 66% better AP than a random ordering, and 14% better performance than an intelligent baseline. On the timeliness measure, our method obtains at least 11% better performance. Our code, to be made available upon publication, is easily extensible as it treats detectors and classifiers as black boxes and learns from execution traces using reinforcement learning.

## 1 Introduction

In most real-world applications of object recognition, performance is time-sensitive and inherently tied to the many-class nature of the world. In robotics, a small finite amount of processing power per unit time is all that is available for robust object detection if the robot is to usefully interact with humans. In large-scale detection systems (e.g. image search), results need to be obtained quickly per image as the number of items to process is constantly growing. In these cases, an acceptable answer at a reasonable time may be more valuable than the best answer given too late.

The conventional approach to evaluation of visual category recognition does not consider efficiency and evaluates performance independently across classes. We argue that the key to tackling problems of dynamic recognition resource allocation is to start asking a new question: *What is the best performance we can get on a budget?*

We propose a new *timeliness* measure of performance vs. time (shown in Figure 1), and present a method based on reinforcement learning that takes various detectors and classifiers as black boxes, and learns a dynamic policy for selecting detector and other actions to achieve the highest performance under this evaluation. We are able to obtain better performance than all baselines when there is less time available than is needed to exhaustively run all detectors.

A hypothetical system for vision-based advertising presents a case study. The system will have different values and accuracies for objects of different classes, and the queue of unprocessed images will vary in size. The detection strategy to maximize profit in such an environment has to be dynamic. We explore this scenario using the PASCAL VOC datasets and evaluation regimes, and show that our learned dynamic policy offers a better strategy than a random or optimal static ordering of detectors.

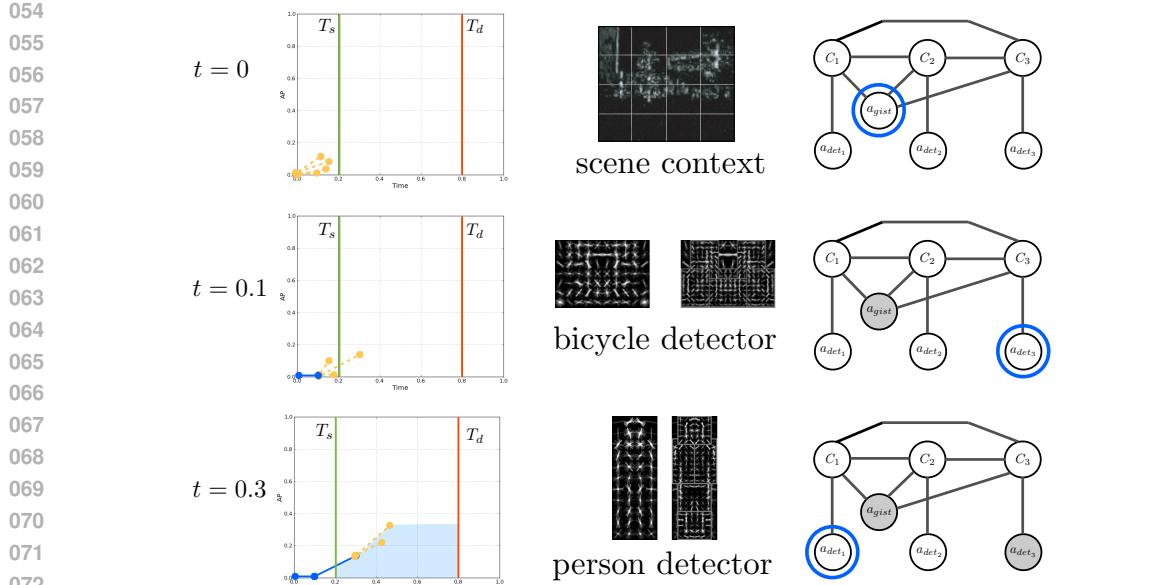


Figure 1: A sample trace of our method. At each time step beginning at  $t = 0$ , potential actions are considered according to their *value*, and the one with the highest value is picked. The selected action is performed and returns observations, which influence the selection of the next action. The final evaluation of a detection episode is the area of the AP vs. Time curve between the start and end times. The value of an action is the expected value of the final evaluation if the action is taken and the policy continues to be followed, which allows actions without an immediate benefit to be scheduled.

## 2 Recognition Problems and Related Work

We deal with a dataset of images  $\mathcal{D}$ , where each image  $\mathcal{I}$  contains zero or more objects. Each object is labeled with exactly one category label  $k \in \{1, \dots, K\}$ .

The multi-class, multi-label **classification** problem asks whether  $\mathcal{I}$  contains at least one object of class  $k$ . We write the ground truth for an image as  $\mathbf{C} = \{C_1, \dots, C_K\}$ , where  $C_k \in \mathbb{B} = \{0, 1\}$  is set to 1 if an object of class  $k$  is present.

The **detection** problem is to output a list of bounding boxes (sub-images defined by four coordinates), each with a real-valued confidence that it encloses a single instance of an object of class  $k$ , for each  $k$ . The answer for a single class is given by an algorithm  $detect(\mathcal{I}, k)$ , which outputs a list of sub-image bounding boxes  $B$  and their associated confidences.

The answer is evaluated by plotting precision vs. recall across dataset  $\mathcal{D}$  (by progressively lowering the confidence threshold for a positive detection). The area under the curve yields the Average Precision (AP) metric, which has become the standard evaluation for recognition performance on challenging datasets in vision [1]. A common measure of a correct detection is the PASCAL overlap: two bounding boxes are considered to match if they have the same label and the ratio of their intersection to their union is at least  $\frac{1}{2}$ .

To highlight the hierarchical structure of these problems, we note that (1) the confidences for each sub-image  $b \in B$  may be given by  $classify(b, k)$ ; (2) correct answer to the detection problem also answers the classification problem.

Multi-class performance is evaluated by averaging the individual per-class AP values. We generalize this metric to a weighted average, with the weights set by the *values* of the classes.

108  
109  
**2.1 Related Work**

110 The literature on object recognition is vast. Here we briefly summarize work that sufficiently con-  
111 textualizes our contribution.

112 **Single-class detection** The best recent performance has come from detectors that use gradient-based  
113 features to represent objects as either a collection of local patches or as object-sized windows [2, 3].  
114 Classifiers are then used to distinguish between featurizations of a given class and all other possible  
115 contents of an image window. Window proposal is most often done exhaustively over the image  
116 space, as a “sliding window”. Some approaches use “jump windows” (hypotheses voted on by local  
117 features) [4, 5], or a bounded search over the space of all possible windows [6].

118 None of the best-performing systems treat window proposal and evaluation as a closed-loop system,  
119 with feedback from evaluation to proposal. Some work has been done on this topic, mostly inspired  
120 by ideas from biological vision and attention research [7, 8].

121 **Context** Context has a foundational role in vision. One source of context is the scene or other  
122 non-detector cues; the most common scene-level feature is the GIST [9] of the image. For the  
123 commonly used PASCAL VOC dataset [1], GIST and other sources of context are quantitatively  
124 explored in [10]. Inter-object context has also been shown to be useful for improving detection [11].  
125 A critical summary of the main approaches to using context for object and scene recognition is given  
126 in [12].

127 **Multi-Class Detection** Work on inherently multi-class detection focuses largely on making detec-  
128 tion time sublinear in the number of classes through sharing features [13, 14]. A post-processing  
129 extension to detection systems uses structured prediction to incorporate multi-class context as a  
130 principled replacement for non-maximum suppression [15].

131 **Cascades** An early success in efficient object detection used simple, fast features to build up a *cascade*  
132 of classifiers, which then considered image regions in a sliding window regime [16]. Most  
133 recently, cyclic optimization has been applied to optimize cascades with respect to feature computa-  
134 tion cost as well as classifier performance [17]. However, cascades are not dynamic policies—they  
135 cannot change the order of execution based on observations obtained during execution, which is our  
136 goal.

137 **Anytime Algorithms and Active Classification** This surprisingly little-explored line of work in  
138 vision is closest to our approach. A recent application to the problem of visual detection picks  
139 features with maximum value of information in a Hough-voting framework [18]. There has also  
140 been work on active classification [19] and active sensing [20], in which intermediate results are  
141 considered in order to decide on the next classification step. Most commonly, the scheduling in  
142 these approaches is greedy with respect to some manual quantity such as expected information gain.  
143 In contrast, we learn policies that take actions without any immediate reward.

144  
145 **3 Multi-class Recognition Policy**  
146

147 Our goal is a multi-class recognition policy  $\pi$  that takes an image  $\mathcal{I}$  and outputs a list of multi-class  
148 detection results by running detector and global scene *actions* sequentially.

149 The policy repeatedly selects an action  $a_i \in \mathcal{A}$ , executes it, receiving observations  $o_i$ , and then  
150 selects the next action. The set of actions  $\mathcal{A}$  can include both classifiers and detectors: anything that  
151 would be useful for inferring the contents of the image.

152 Each action  $a_i$  has an expected cost  $c(a_i)$  of execution. Depending on the setting, the cost can be  
153 defined in terms of algorithmic runtime analysis, an idealized property such as number of *flops*, or  
154 simply the empirical runtime on specific hardware. We take the empirical approach: every executed  
155 action advances  $t$ , the *time into episode*, by its empirical runtime.

156 As shown in Figure 1, the system is given two times: the setup time  $T_s$  and deadline  $T_d$ . From the  
157 setup time to the deadline, we want to obtain the best possible answer if stopped at any given time. A  
158 single-number metric that corresponds to this objective is the area captured under the curve between  
159 the start and deadline bounds, normalized by the total area. We evaluate policies by this more robust  
160 metric and not simply by the final performance at deadline time for the same reason that Average  
161 Precision is used instead of a fixed Precision vs. Recall point in the conventional evaluations.

162    **3.1 Sequential Execution**  
 163

164 An *open-loop* policy takes actions in a sequence that does not depend on observations received from  
 165 previous actions. The common classifier cascade is an example [16]. In contrast, our goal is to learn  
 166 a dynamic, or *closed-loop*, policy, which would exploit the signal in scene and inter-object context  
 167 for a maximally efficient path through the actions.

168 We refer to the information available to the decision process as the *state*  $s$ . The state includes the  
 169 currently believed distribution over class presence variables  $P(\mathbf{C}) = \{P(C_0), \dots, P(C_K)\}$ , where  
 170 we write  $P(C_k)$  to mean  $P(C_k = 1)$ .

171 Additionally, the state records that an action  $a_i$  has been taken by adding it to the initially empty  
 172 set  $\mathcal{O}$  and recording the resulting observations  $o_i$ . We refer to the current set of observations as  
 173  $\mathbf{o} = \{o_i | a_i \in \mathcal{O}\}$ . The state also keeps track of the time into episode  $t$ , and the setup and deadline  
 174 times  $T_s, T_d$ .

175 A recognition *episode* takes an image  $\mathcal{I}$  and proceeds from the initial state  $s^0$  and action  $a^0$  to the  
 176 next pair  $(s^1, a^1)$ , and so on until  $(s^J, a^J)$ , where  $J$  is the last step of the process with  $t \leq T_d$ . At  
 177 that point, the policy is terminated, and a new episode can begin on a new image.

179 The specific actions we consider in the following exposition are detector actions  $a_{det_i}$ , where  $det_i$   
 180 is a detector class  $C_i$ , and a scene-level context action  $a_{gist}$ , which updates the probabilities of all  
 181 classes.

182    **3.2 Selecting actions**  
 183

184 As our goal is to pick actions dynamically, we want a function  $Q(s, a) : S \times \mathcal{A} \mapsto \mathbb{R}$ , where  $S$  is the  
 185 space of all possible states, to assign a value to a potential action  $a \in \mathcal{A}$  given the current state  $s$  of  
 186 the decision process. We can then define the desired policy  $\pi$  as simply taking the (untaken) action  
 187 with the maximum value:

$$\pi(s) = \operatorname{argmax}_{a_i \in \mathcal{A} \setminus \mathcal{O}} Q(s, a_i) \quad (1)$$

190 Although the action space  $\mathcal{A}$  is quite manageable, consisting of the detectors and global classifier  
 191 we would like to run on the image, the space of possible states  $S$  is intractable. Therefore we cannot  
 192 learn a tabular representation of  $Q(s, a)$ , and must use function approximation to represent it [21].  
 193 We featurize the state-action pair and assume linear structure:  $Q^\pi(s, a_i) = \theta_\pi^\top \phi(s, a_i)$ .

195 The policy's performance at time  $t$  is determined by the detections that are part of the set of obser-  
 196 vations  $\mathbf{o}^j$  at the last state  $s^j$  before  $t$ . Therefore, the final AP vs. Time evaluation of an episode is  
 197 a function  $eval(h, T_s, T_d)$  of the history of execution  $h = s^0, s^1, \dots, s^J$ . It is precisely the normal-  
 198 ized area under the AP vs. Time curve between  $T_s$  and  $T_d$ , as determined by the detections in  $\mathbf{o}^j$  for  
 199 all steps  $j$  in the episode.

200 As shown in 3b, this evaluation function is additive per action, as each action can generate detec-  
 201 tions that either raise or lower the mean AP of the results so far ( $\Delta ap$ ) and takes a certain time  
 202 ( $\Delta t$ ). From these and  $T_s$  and  $T_d$ , we can find the area under the curve that was contributed by the  
 203 action. So, we can then represent the final evaluation  $eval(h, T_s, T_d)$  in terms of individual rewards:  
 204  $\sum_{j=0}^J R(s^j, a^j)$ .

205 Specifically, as shown in Figure 3b, we define the *reward* of an action as

$$R(s^j, a_i) = \Delta ap_i(t_T^j - \frac{1}{2} \Delta t_i) \quad (2)$$

208 where  $t_T^j$  and  $ap^j$  are the time left until the deadline and the AP at state  $s^j$ , and  $\Delta t_i$  and  $\Delta ap_i$  are the  
 209 time taken and AP change produced by the action  $a_i$ . (For clarity of exposition, we do not account  
 210 for  $T_s$  here.)

212    **3.3 Learning the policy**  
 213

214 The expected value of the final evaluation can be written recursively in terms of the value function:

$$Q^\pi(s^j, a_i) = \mathbb{E}_{s^{j+1}}[R(s^j, a_i) + \gamma Q^\pi(s^{j+1}, \pi(s^{j+1}))] \quad (3)$$

216 where  $\gamma \in [0, 1]$  is a *discount* value that can mitigate the effects of increasing state-transition uncer-  
217 tainty over long episodes.

219 While we can't directly compute the expectation in (3), we can sample it by running actual episodes  
220 to gather  $\langle s, a, r, s' \rangle$  samples, where  $r$  is the reward obtained by taking action  $a$  in state  $s$ , and  $s'$   
221 is the following state.

222 Learning the policy is then a problem of repeatedly gathering samples with the current policy, mini-  
223 mizing the error between the discounted reward to the end of the episode as predicted by our current  
224  $Q(s^j, a_i)$  and the actual values gathered, and updating the policy with the resulting weights. This is  
225 fitted Q-iteration, a variant of generalized policy iteration [22, 21].

226 To ensure sufficient exploration of the state space, we implement  $\epsilon$ -greedy action selection during  
227 training: with a probability that decreases with each training iteration, a random action is selected  
228 instead of following the policy. During test time,  $\epsilon$  is set to 0.05.

229 We use  $L_2$ -regularized regression to minimize the error. We run 15 iterations of accumulating sam-  
230 ples by running 350 episodes, starting with a baseline policy which will be described in section 4,  
231 and cross-validating the regularization parameter at each iteration; samples are not thrown away  
232 between iterations.

233 A meta-parameter of the approach is the discount  $\gamma$ . With  $\gamma = 0$ , the value function is determined  
234 entirely by the immediate reward. Learning in this case can only result in completely greedy policies.  
235

236 With  $\gamma = 1$ , the value function is determined by the expected rewards to the end of the episode, and  
237 so should be a close approximation to the final evaluation metric. However, the highest value for  $\gamma$   
238 is not necessarily best if the action-value function is not expressive enough to represent the actual  
239 state transition behavior of the world. We experiment with several values of  $\gamma$ , and find a mid-level  
240 value (0.4) to work best.

### 241 3.4 Feature representation

242 Our policy is at its base determined by a linear function of the features of the state:  $\pi(s) =$   
243  $\text{argmax}_{a_i \in \mathcal{A} \setminus \mathcal{O}} \theta_\pi^\top \phi(s, a_i)$ . Since we want to be able to learn a dynamic policy, the observations  $\mathbf{o}$   
244 that are part of the state  $s$  should play a role in determining the value of a potential action.

245 We include the following quantities as features  $\phi(s, a)$ :

246  $P(C_a)$  The prior probability of the class that corresponds to the detector of action  $a$

247  $P(C_0|\mathbf{o}) \dots P(C_K|\mathbf{o})$  The probabilities of all classes, conditioned on the current set of observa-  
248 tions.

249  $H(C_0|\mathbf{o}) \dots H(C_K|\mathbf{o})$  The entropies of all classes, conditioned on the current set of observations.

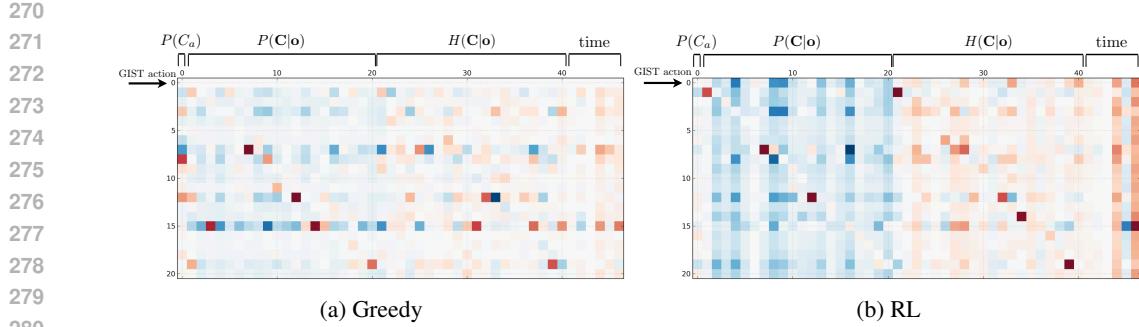
250 Additionally, we include the mean and maximum entropies of all classes, the expected time of the  
251 action  $c(a)$ , and time features that represent the times until start and deadline, for a total of  $F$   
252 features.

253 We note that this setup is commonly used to solve Markov Decision Processes [21]. There are two  
254 related limitations of MDPs when it comes to most systems of interesting complexity, however: the  
255 state has to be functionally approximated instead of exhaustively enumerated; and some aspects of  
256 the state are not observed, making the problem a Partially Observed MDP (POMDP), for which  
257 exact solution methods are intractable for all but rather small problems [23].

258 There isn't much to do about the necessity of approximating the state but design good features. Our  
259 initial solution to the partial observability problem is to include uncertainty features into the feature  
260 representation to *augment* the MDP [24].

261 To formulate learning the policy as a single regression problem, we represent the features in block  
262 form, where  $\phi(s, a_i)$  is a vector of size  $F|\mathcal{A}|$ , with all values set to 0 except for the block corre-  
263 sponding to  $a_i$ .

264 As an illustration, we visualize the learned weights on these features in Figure 2, reshaped such  
265 that each row shows the weights learned for an action, in order. The featurization for the  $a_{gist}$   
266 scene-context action, which concerns all classes, omits the first  $P(C_a)$  feature.



270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
Figure 2: Learned policy weights  $\theta_\pi$  (best viewed in color: red corresponds to positive, blue to negative values). The first row corresponds to the scene-level action, which does not generate detections itself but only helps reduce uncertainty about the contents of the image. Note that in the greedy learning case, this action is learned to never be taken—whereas it is shown to be useful in the reinforcement learning case.

### 3.5 Updating with observations

The bulk of our feature representation is formed by probability of individual class occurrence, conditioned on the observations so far:  $P(C_0|\mathbf{o}) \dots P(C_K|\mathbf{o})$ . This allows the action-value function to learn correlations between presence of different classes, and so the policy can look for the most probable classes given the observations.

However, higher-order co-occurrences are not well represented in this form. Additionally, updating  $P(C_i|\mathbf{o})$  presents choices regarding independence assumptions between the classes.

We evaluate two approaches for updating probabilities: *direct* and *MRF*. In the *direct* method,  $P(C_i|\mathbf{o}) = \text{score}(C_i)$  if  $\mathbf{o}$  includes the observations for class  $C_i$  and  $P(C_i|\mathbf{o}) = P(C_i)$  otherwise. This means that an observation of class  $i$  does not influence the estimated probability of any class but  $C_i$ .  $\text{score}(C_i)$  for  $a_{\text{det}_i}$  is obtained by training a probabilistic classifier on the detections output.  $\text{score}(C_i)$  for  $a_{\text{gist}}$  is obtained by training probabilistic classifiers on the GIST feature, for all classes.

The *MRF* approach employs a pairwise fully-connected Markov Random Field (MRF), as shown in Figure 1, with the observation nodes set to  $\text{score}(C_i)$  appropriately, or considered unobserved.

The graphical model structure is set as fully-connected, but some classes are overwhelmingly unlikely to co-occur in our dataset. Accordingly, the graph edge weights are learned with  $L_1$  regularization, which obtains the desired sparse structure [25]. All parameters of the model are trained on fully-observed data, and Loopy Belief Propagation inference is implemented with an open-source graphical model package [26].

## 4 Evaluation

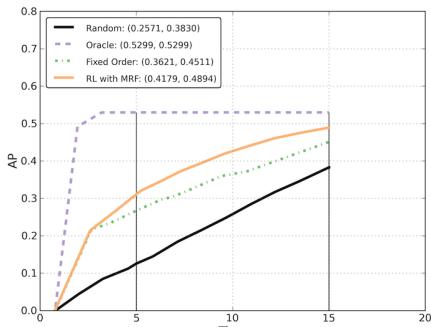
We evaluate our system on the multi-class, multi-label detection task as previously described. Each detection episode takes an image and outputs detections with associated times, based on the order of actions taken. We evaluate on a popular detection challenge task: the PASCAL VOC 2007 dataset [1], which exhibits only a rather modest amount of class co-occurrence: the “person” class is highly likely to occur, and less than 10% of the images have more than two classes.

We learn weights on the training and validation sets, and run our policy on all images in the testing set. The final evaluation pools all detections up to a certain time, and computes their multi-class AP per image, averaging over all images. This is done for different times to plot the AP vs. Time curve over the whole dataset. Our method of averaging per-image performance follows [15].

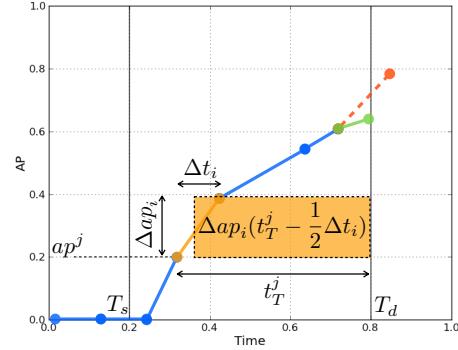
For the detector actions, we use one-vs-all cascaded deformable part-model detectors on a HOG featurization of the image [27], with associated linear classification of the detections. There are 20

324 classes in the PASCAL challenge task, so there are 20 detector actions. Running a detector on a  
 325 PASCAL image takes about 1 second.  
 326

327 We test three different settings of the start and deadline times. In the first one, the start time is  
 328 immediate and execution is cut off at 20 seconds, which is enough time to run basically all actions.  
 329 In the second one, execution is cut off after only 10 seconds. Lastly, we measure performance  
 330 between 5 seconds and 15 seconds. These operating points show how our method behaves when  
 331 deployed in different conditions. The results are given in rows of [Table 1](#).  
 332  
 333



334 (a)



335 (b)

336 Figure 3: (a) AP vs. Time curves for Random, Oracle, the Fixed Order baseline, and our best-  
 337 performing policy, with start time of 5 seconds and deadline of 15 seconds. (b) Reward function  
 338 explanation.

339 We establish the first baseline for our system by selecting actions randomly at each step. As shown  
 340 in Figure 3a, this **Random** policy results in a roughly linear gain of AP vs. time. This is expected:  
 341 the detectors are capable of obtaining a certain level of performance; if half the detectors are run,  
 342 the expected performance level is half of the maximum level.

343 To establish an upper bound on performance, we plot the **Oracle** policy, obtained by re-ordering the  
 344 actions with the hindsight of the results they produced.

345 We consider another baseline: selecting actions in a fixed order based on the value they bring to the  
 346 AP vs. Time evaluation, which is roughly proportional to their occurrence probability. We refer to  
 347 this as **Fixed Order**.

348 Then there are our proposed systems, described in the previous section : **RL w/ Direct** and **RL w/  
 349 MRF**. In experiments, the two were not exceedingly different in performance, with the MRF model  
 350 consistently performing slightly better.

351 In Figure 3a, we can see that due to the dataset bias, the fixed-order policy performs well at first, as  
 352 the person class is disproportionately likely to be in the image, but is significantly overtaken by our  
 353 model as execution goes on and more rare classes have to be detected.

354 Lastly, we include an additional scene-level GIST feature that updates the posterior probabilities of  
 355 all classes. This is considered one action, and takes about 0.3 seconds. This setting always uses the  
 356 MRF model to properly update the class probabilities with GIST observations. This brings another  
 357 boost in performance. All results are shown in [Table 1](#).

358 Additionally, we evaluated different settings of the discount parameter  $\gamma$ , and found that a mid-range  
 359 value works best. We only report results for that value on the test set. [Figure 2](#) shows how changing  
 360  $\gamma$  affects the learned weights: the GIST action is learned to never be taken in the greedy setting,  
 361 because it does not bring any immediate detections results. In the reinforcement learning setting,  
 362 however (higher  $\gamma$ ), GIST is learned to be taken for its value in quickly refining the class presence  
 363 probabilities.

364  
 365 [Figure 4](#) illustrates the action trajectories of different policies.

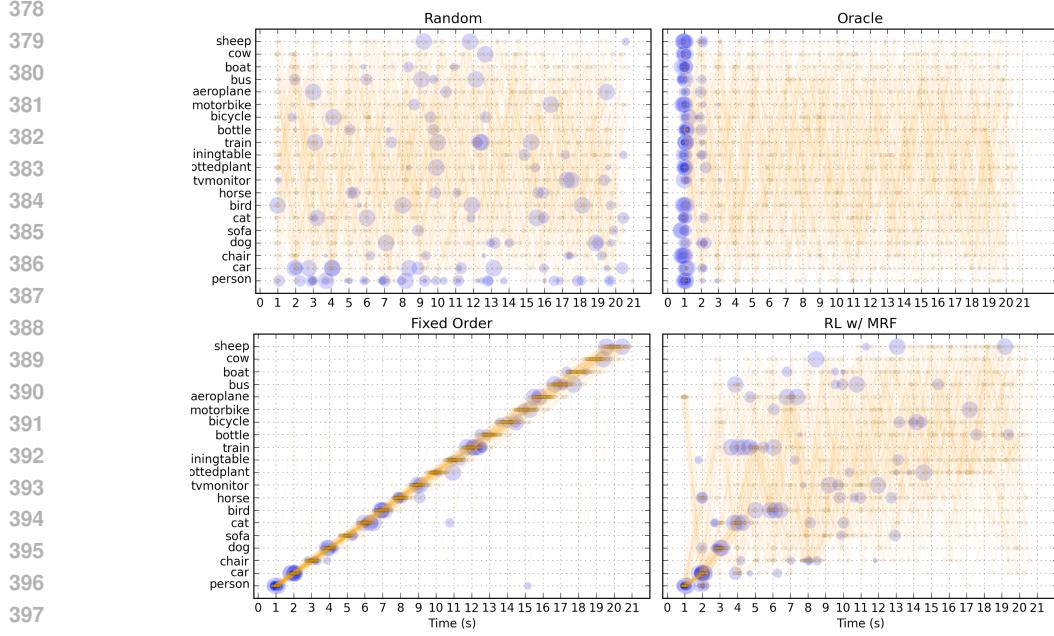


Figure 4: Visualizing the action trajectories of different policies. Action selection traces are plotted in orange over many episodes; the size of the blue circles correspond to the  $\Delta ap$  obtained by the action. We see that the **Random** policy selects actions and obtains rewards randomly, while the **Oracle** policy obtains all rewards in the first few actions. The **Fixed Order** policy selects actions in a static optimal order. Our policy does not stick a static order but selects actions dynamically to maximize the rewards obtained early on.

Table 1: The areas under the AP vs. Time curve for different experimental conditions. Our method reliably beats the intelligent static baseline; including the GIST action further boosts performance.

| Bounds | Random | Fixed Order | RL w/ MRF | w/ GIST      | Oracle |
|--------|--------|-------------|-----------|--------------|--------|
| (0,20) | 0.250  | 0.342       | 0.378     | <b>0.382</b> | 0.488  |
| (0,10) | 0.119  | 0.240       | 0.266     | <b>0.267</b> | 0.464  |
| (5,15) | 0.257  | 0.362       | 0.418     | <b>0.420</b> | 0.530  |

## 5 Conclusion

We presented a method for learning “closed-loop” policies for multi-class object recognition, given existing object detectors and classifiers and a metric to optimize. The method learns the optimal policy using reinforcement learning, by observing execution traces in training. If detection on an image is cut off after only half the detectors have been run, our method does 66% better than random, and 14% better than an intelligent baseline. In particular, our method learns to take action with no intermediate reward in order to improve the overall performance of the system.

As always with reinforcement learning problems, defining the reward function requires some manual work. Here, we derive it for the novel detection AP vs. Time evaluation that we suggest is most useful for evaluating efficiency in recognition. Although computation devoted to scheduling actions is much less significant than the computation due to running the actions, an interesting future research direction is to explicitly consider this decision-making cost. At time of publication, we will release the code of our modular framework, that allows for easy incorporation of other classifiers and detectors on different tasks and rewards.

432      **References**

- 433
- [1] M Everingham, L Van Gool, C K I Williams, J Winn, and A Zisserman. The PASCAL VOC Challenge  
434      2010 Results. <http://www.pascal-network.org/challenges/VOC/voc2010/workshop/index.html>, 2010. 2,  
435      3, 6
- [2] N Dalal and B Triggs. Histograms of Oriented Gradients for Human Detection. In *CVPR*, pages 886–893.  
436      Ieee, 2005. 3
- [3] David G Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of  
439      Computer Vision*, 60(2):91–110, November 2004. 3
- [4] Andrea Vedaldi, Varun Gulshan, Manik Varma, and Andrew Zisserman. Multiple kernels for object  
440      detection. *ICCV*, pages 606–613, September 2009. 3
- [5] Sudheendra Vijayanarasimhan and Kristen Grauman. Large-Scale Live Active Learning: Training Object  
442      Detectors with Crawled Data and Crowds. In *CVPR*, 2011. 3
- [6] Christoph H Lampert, Matthew B Blaschko, and Thomas Hofmann. Beyond sliding windows: Object  
444      localization by efficient subwindow search. *CVPR*, pages 1–8, June 2008. 3
- [7] N J Butko and J R Movellan. Optimal scanning for faster object detection. *CVPR*, (1):2751–2758, June  
446      2009. 3
- [8] Julia Vogel and Nando de Freitas. Target-directed attention: Sequential decision-making for gaze plan-  
448      ning. *ICRA*, pages 2372–2379, May 2008. 3
- [9] Aude Oliva and Antonio Torralba. Modeling the Shape of the Scene: A Holistic Representation of the  
450      Spatial Envelope. *IJCV*, 42(3):145–175, 2001. 3
- [10] S K Divvala, D Hoiem, J H Hays, A.a. Efros, and M Hebert. An empirical study of context in object  
452      detection. In *CVPR*, pages 1271–1278. Ieee, June 2009. 3
- [11] Antonio Torralba, Kevin P Murphy, and William T Freeman. Contextual Models for Object Detection  
454      Using Boosted Random Fields. *MIT Computer Science and Artificial Intelligence Laboratory Technical  
455      Report*, 2004. 3
- [12] Carolina Galleguillos and Serge Belongie. Context based object categorization: A critical survey. *Com-  
457      puter Vision and Image Understanding*, 114(6):712–722, June 2010. 3
- [13] Antonio Torralba, Kevin P Murphy, and William T Freeman. Sharing visual features for multiclass and  
459      multiview object detection. *PAMI*, 29(5):854–869, May 2007. 3
- [14] Xiaodong Fan. Efficient Multiclass Object Detection by a Hierarchy of Classifiers. In *CVPR*, pages  
461      716–723. Ieee, 2005. 3
- [15] Chaitanya Desai, Deva Ramanan, and Charless Fowlkes. Discriminative models for multi-class object  
462      layout. In *2009 IEEE 12th International Conference on Computer Vision*, pages 229–236. Ieee, September  
463      2009. 3, 6
- [16] Paul Viola and Michael Jones. Rapid object detection using a boosted cascade of simple features. In  
465      *CVPR*, 2001. 3, 4
- [17] Minmin Chen, Zhixiang Eddie, Xu Kilian, and Saint Louis. Classifier Cascade for Minimizing Feature  
467      Evaluation Cost. In *AISTATS*, 2012. 3
- [18] Sudheendra Vijayanarasimhan and Ashish Kapoor. Visual Recognition and Detection Under Bounded  
469      Computational Resources. In *CVPR*, pages 1006–1013, 2010. 3
- [19] Tianshi Gao and Daphne Koller. Active Classification based on Value of Classifier. In *NIPS*, pages 1–9,  
471      2011. 3
- [20] Shipeng Yu, Balaji Krishnapuram, Romer Rosales, and R Bharat Rao. Active Sensing. In *AISTATS*, pages  
473      639–646, 2009. 3
- [21] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998. 4, 5
- [22] Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-Based Batch Mode Reinforcement Learning.  
476      *Journal of Machine Learning Research*, 6:503–556, 2005. 5
- [23] Nicholas Roy and Geoffrey Gordon. Exponential Family PCA for Belief Compression in POMDPs. In  
478      *NIPS*, 2002. 5
- [24] Cody Kwok and Dieter Fox. Reinforcement Learning for Sensing Strategies. In *IROS*, 2004. 5
- [25] Su-In Lee, Varun Ganapathi, and Daphne Koller. Efficient Structure Learning of Markov Networks using  
481      L 1 -Regularization. In *NIPS*, 2006. 6
- [26] Ariel Jaimovich and Ian Mcgraw. FastInf : An Efficient Approximate Inference Library. *Journal of  
483      Machine Learning Research*, 11:1733–1736, 2010. 6
- [27] Pedro F Felzenszwalb, Ross B Girshick, and David McAllester. Cascade object detection with deformable  
485      part models. In *CVPR*, pages 2241–2248. IEEE, June 2010. 6