

## Dear Author

Here are the proofs of your article.

- You can submit your corrections **online** or by **fax**.
- For **online** submission please insert your corrections in the online correction form. Always indicate the line number to which the correction refers.
- Please return your proof together with the permission to publish confirmation.
- For **fax** submission, please ensure that your corrections are clearly legible. Use a fine black pen and write the correction in the margin, not too close to the edge of the page.
- Remember to note the journal title, article number, and your name when sending your response via e-mail, fax or regular mail.
- **Check** the metadata sheet to make sure that the header information, especially author names and the corresponding affiliations are correctly shown.
- **Check** the questions that may have arisen during copy editing and insert your answers/corrections.
- **Check** that the text is complete and that all figures, tables and their legends are included. Also check the accuracy of special characters, equations, and electronic supplementary material if applicable. If necessary refer to the *Edited manuscript*.
- The publication of inaccurate data such as dosages and units can have serious consequences. Please take particular care that all such details are correct.
- Please **do not** make changes that involve only matters of style. We have generally introduced forms that follow the journal's style.  
Substantial changes in content, e.g., new results, corrected values, title and authorship are not allowed without the approval of the responsible editor. In such a case, please contact the Editorial Office and return his/her consent together with the proof.
- If we do not receive your corrections **within 48 hours**, we will send you a reminder.

### Please note

Your article will be published **Online First** approximately one week after receipt of your corrected proofs. This is the **official first publication** citable with the DOI.  
**Further changes are, therefore, not possible.**

After online publication, subscribers (personal/institutional) to this journal will have access to the complete article via the DOI using the URL:

<http://dx.doi.org/10.1007/s10626-009-0071-x>

If you would like to know when your article has been published online, take advantage of our free alert service. For registration and further information, go to:  
<http://www.springerlink.com>.

Due to the electronic nature of the procedure, the manuscript and the original figures will only be returned to you on special request. When you return your corrections, please inform us, if you would like to have these documents returned.

The **printed version** will follow in a forthcoming issue.

To: **Springer Customer Support 2**  
E-mail: CorrAdmin2@spi-bpo.com  
Fax: +1-703-5621873

SPi  
Re: **SPi Building, Sacsac Bacong**  
**Oriental Negros 6216**  
**Philippines**

**Discrete Event Dynamic Systems DOI 10.1007/s10626-009-0071-x**

Partially Observable Markov Decision Process Approximations for Adaptive  
Sensing  
**Chong · Kreucher · Hero**

## **Permission to publish**

I have checked the proofs of my article and

- I have **no corrections**. The article is ready to be published without changes.
- I have **a few corrections**. I am enclosing the following pages:
- I have made **many corrections**. Enclosed is the **complete article**.

Date / signature: \_\_\_\_\_

## Metadata of the article that will be visualized in Online

---

**Please note: Image will appear in color online but will be printed in black and white.**

---

1	Article Title	<b>Partially Observable Markov Decision Process Approximations for Adaptive Sensing</b>
2	Article Sub- Title	
3	Article Copyright - Year	<b>Springer Science + Business Media, LLC 2009 (This will be the copyright line in the final PDF)</b>
4	Journal Name	Discrete Event Dynamic Systems
5	Family Name	<b>Chong</b>
6	Particle	
7	Given Name	<b>Edwin K. P.</b>
8	Corresponding Suffix	
9	Author Organization	Colorado State University
10	Division	
11	Address	Fort Collins , CO, USA
12	e-mail	edwin.chong@colostate.edu
13	Family Name	<b>Kreucher</b>
14	Particle	
15	Given Name	<b>Christopher M.</b>
16	Suffix	
17	Author Organization	Integrity Applications Incorporated
18	Division	
19	Address	Ann Arbor , MI, USA
20	e-mail	ckreuche@umich.edu
21	Family Name	<b>Hero</b>
22	Particle	
23	Given Name	<b>Alfred O.</b>
24	Author Suffix	<b>III</b>
25	Organization	University of Michigan
26	Division	
27	Address	Ann Arbor , MI, USA
28	e-mail	hero@umich.edu
29	Received	18 July 2008
30	Schedule Revised	
31	Accepted	14 May 2009
32	Abstract	Adaptive sensing involves actively managing sensor resources to achieve a sensing task, such as object detection, classification, and tracking, and represents a promising direction for new applications of discrete event system methods. We describe an approach to adaptive sensing based on

approximately solving a partially observable Markov decision process (POMDP) formulation of the problem. Such approximations are necessary because of the very large state space involved in practical adaptive sensing problems, precluding exact computation of optimal solutions. We review the theory of POMDPs and show how the theory applies to adaptive sensing problems. We then describe a variety of approximation methods, with examples to illustrate their application in adaptive sensing. The examples also demonstrate the gains that are possible from nonmyopic methods relative to myopic methods, and highlight some insights into the dependence of such gains on the sensing resources and environment.

---

33	Keywords separated by ' - '	Markov decision process - POMDP - Sensing - Tracking - Scheduling
34	Foot note information	This material is based in part upon work supported by the Air Force Office of Scientific Research under Award FA9550-06-1-0324 and by DARPA under Award FA8750-05-2-0285. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the Air Force or of DARPA. Approved for Public Release, Distribution Unlimited.

---

## Partially Observable Markov Decision Process Approximations for Adaptive Sensing

**Edwin K. P. Chong · Christopher M. Kreucher ·  
Alfred O. Hero III**

Received: 18 July 2008 / Accepted: 14 May 2009  
© Springer Science + Business Media, LLC 2009

**Abstract** Adaptive sensing involves actively managing sensor resources to achieve 1 a sensing task, such as object detection, classification, and tracking, and represents 2 a promising direction for new applications of discrete event system methods. We 3 describe an approach to adaptive sensing based on approximately solving a partially 4 observable Markov decision process (POMDP) formulation of the problem. Such ap- 5 proximations are necessary because of the very large state space involved in practical 6 adaptive sensing problems, precluding exact computation of optimal solutions. We 7 review the theory of POMDPs and show how the theory applies to adaptive sensing 8 problems. We then describe a variety of approximation methods, with examples to 9 illustrate their application in adaptive sensing. The examples also demonstrate the 10 gains that are possible from nonmyopic methods relative to myopic methods, and 11 highlight some insights into the dependence of such gains on the sensing resources 12 and environment. 13

**Keywords** Markov decision process · POMDP · Sensing · Tracking · Scheduling

14

---

This material is based in part upon work supported by the Air Force Office of Scientific Research under Award FA9550-06-1-0324 and by DARPA under Award FA8750-05-2-0285. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the Air Force or of DARPA. Approved for Public Release, Distribution Unlimited.

E. K. P. Chong (✉)  
Colorado State University, Fort Collins, CO, USA  
e-mail: edwin.chong@colostate.edu

C. M. Kreucher  
Integrity Applications Incorporated, Ann Arbor, MI, USA  
e-mail: ckreuche@umich.edu

A. O. Hero III  
University of Michigan, Ann Arbor, MI, USA  
e-mail: hero@umich.edu

## 15 1 Introduction

### 16 1.1 What is adaptive sensing?

17 In its broadest sense, *adaptive sensing* has to do with actively managing sensor  
18 resources to achieve a sensing task. As an example, suppose our goal is to determine  
19 the presence or absence of an object, and we have at our disposal a single sensor  
20 that can interrogate the scene with any one of  $K$  waveforms. Depending on which  
21 waveform is used to irradiate the scene, the response may vary greatly. After each  
22 measurement, we can decide whether to continue taking measurements using that  
23 waveform, change waveforms and take further measurements, or stop and declare  
24 whether or not the object is present. In adaptive sensing, this decision making is  
25 allowed to take advantage of the knowledge gained from the measurements so far.  
26 In this sense, the act of sensing “adapts” to what we know so far. What guides  
27 this adaptation is a performance objective that is determined beforehand—in our  
28 example above, this might be the average number of interrogations needed so that  
29 we can declare the presence or absence of the object with a confidence that exceeds  
30 some threshold (say, 90%).

31 Adaptive sensing problems arise in a variety of application areas, and represent a  
32 promising direction for new applications of discrete event system methods. Here, we  
33 outline only a few.

34 *Medical diagnostics* Perhaps the most familiar example of adaptive sensing takes  
35 place between a doctor and a patient. The task here is to diagnose an illness from  
36 a set of symptoms, using a variety of medical tests at the doctor’s disposal. These  
37 include physical examinations, blood tests, radiographs (X-ray images), computer-  
38 ized tomography (CT) scans, and magnetic resonance imaging (MRI). Doctors use  
39 results from tests so far to determine what test to perform next, if any, before making  
40 a diagnosis.

41 *Nondestructive testing* In nondestructive testing, the goal is to use noninvasive  
42 methods to determine the integrity of a material or to measure some characteristic  
43 of an object. A wide variety of methods are used in nondestructive testing, ranging  
44 from optical to microwave to acoustic. Often, several methods must be used before  
45 a determination can be made. The test results obtained so far inform what method  
46 to use next (including what waveform to select), thus giving rise to an instance of  
47 adaptive sensing.

48 *Sensor scheduling for target detection, identification, and tracking* Imagine a group  
49 of airborne sensors—say, radars on unmanned aerial vehicles (UAVs)—with the  
50 task of detecting, identifying, and tracking one or more targets on the ground. For a  
51 variety of reasons, we can use at most one sensor at any given time. These reasons  
52 include limitations in communication resources needed to transmit data from the  
53 sensors, and the desire to minimize radar usage to maintain covertness. The selection  
54 of which sensor to use over time is called sensor scheduling, and is an adaptive sensing  
55 problem.

56 *Waveform selection for radar imaging* Radar systems have become sufficiently agile  
57 that they can be programmed to use waveform pulses from a library of waveforms.

The response of a target in the scene can vary greatly depending on what waveform 58  
is used to radiate the area due to intrapulse characteristics (e.g., frequency and 59  
bandwidth) or interpulse characteristics (e.g., pulse repetition interval). The main 60  
issue in the operation of such agile radar systems is the selection of waveforms to 61  
use in a particular scenario. If past responses can be used to guide the selection of 62  
waveforms, then this issue is an instance of adaptive sensing. 63

*Laser pulse shaping* Similar to the last example, optical waveforms can also be 64  
designed to generate a variety of responses, only at much smaller wavelengths. By 65  
carefully tailoring the shape of intense light pulses, the interaction of light with even 66  
a single atom can be controlled (Bartels et al. 2000). The possibility of such controlled 67  
interactions of light with atoms has many promising applications. As in the previous 68  
example, these applications give rise to adaptive sensing problems. 69

## 1.2 Nonmyopic adaptive sensing 70

In our view, adaptive sensing is fundamentally a *resource management* problem, in 71  
the sense that the main task is to make decisions over time on the use of sensor 72  
resources to maximize sensing performance. It is informative to distinguish between 73  
*myopic* and *nonmyopic* (also known as *dynamic* or *multistage*) resource management, 74  
a topic of much current interest (see, e.g., Kreucher et al. 2004; He and Chong 2004, 75  
2006; Bertsekas 2005; Krakow et al. 2006; Li et al. 2006, 2007; Ji et al. 2007). In 76  
myopic resource management, the objective is to optimize performance on a per- 77  
decision basis. For example, consider the problem of *sensor scheduling* for tracking 78  
a single target, where the problem is to select, at each decision epoch, a single sensor 79  
to activate. An example sensor-scheduling scheme is *closest point of approach*, which 80  
selects the sensor that is perceived to be the closest to the target. Another (more 81  
sophisticated) example is the method described in Kreucher et al. (2005b), where 82  
the authors present a sensor scheduling method using alpha-divergence (or Rényi 83  
divergence) measures. Their approach is to make the decision that maximizes the 84  
expected information gain (in terms of the alpha-divergence). 85

Myopic adaptive sensing may not be ideal when the performance is measured over 86  
a horizon of time. In such situations, we need to consider schemes that trade off short- 87  
term for long-term performance. We call such schemes *nonmyopic*. Several factors 88  
motivate the consideration of nonmyopic schemes, easily illustrated in the context of 89  
sensor scheduling for target tracking: 90

*Heterogeneous sensors* If we have sensors with different locations, waveform char- 91  
acteristics, usage costs, and/or lifetimes, the decision of whether or not to use a 92  
sensor, and with what waveform, should consider the overall performance, not 93  
whether or not its use maximizes the current performance. 94

*Sensor motion* The future location of a sensor affects how we should act now. 95  
To optimize a long-term performance measure, we need to be opportunistic in our 96  
choice of sensor decisions. 97

*Target motion* If a target is moving, there is potential benefit in sensing the target 98  
before it becomes unresolvable (e.g., too close to other targets or to clutter, or 99

100 shadowed by large objects). In some scenarios, we may need to identify multiple  
101 targets before they cross, to aid in data association.

102 *Environmental variation* Time-varying weather patterns affect target visibility in  
103 a way that potentially benefits from nonmyopic decision making. In particular,  
104 by exploiting models of target visibility maps, we can achieve improved sensing  
105 performance by careful selection of waveforms and beam directions over time. We  
106 show an example along these lines in Section 8.

107 The main focus of this paper is on nonmyopic adaptive sensing. The basic  
108 methodology presented here consists of two steps:

- 109 1) Formulating the adaptive sensing problem as a partially observable Markov  
110 decision process (POMDP); and
- 111 2) Applying an approximation to the optimal policy for the POMDP, because  
112 computing the exact solution is intractable.

113 Our contribution is severalfold. First, we show in detail how to formulate adaptive  
114 sensing problems in the framework of POMDPs. Second, we survey a number of  
115 approximation methods for such POMDPs. Our treatment of these methods includes  
116 their underlying foundations and practical considerations in their implementation.  
117 Third, we illustrate the performance gains that can be achieved via examples.  
118 Fourth, in our illustrative examples, we highlight some insights that are relevant to  
119 adaptive sensing problems: (1) with very limited sensing resources, nonmyopic sensor  
120 and waveform scheduling can significantly outperform myopic methods with only  
121 moderate increase in computational complexity; and (2) as the number of available  
122 resources increases, the nonmyopic advantage decreases.

123 Significant interest in nonmyopic adaptive sensing has arisen in the recent robotics  
124 literature. For example, the recent book by Thrun et al. (2005) describes examples of  
125 such approaches, under the rubric of *probabilistic robotics*. Our paper aims to address  
126 increasing interest in the subject in the signal processing area as well. Our aim is to  
127 provide an accessible and expository treatment of the subject, introducing a class of  
128 new solutions to what is increasingly recognized to be an important new problem.

### 129 1.3 Paper organization

130 This paper is organized as follows. In Section 2, we give a concrete motivating  
131 example that advocates the use of nonmyopic methods. We then describe, in  
132 Section 3, a formulation of the adaptive sensing problem as a partially observable  
133 Markov decision process (POMDP). We provide three examples to illustrate how to  
134 formulate adaptive sensing problems in the POMDP framework. Next, in Section 4,  
135 we review the basic principles behind  $Q$ -value approximation, the key idea in our  
136 approach. Then, in Section 5, we illustrate the basic lookahead control framework  
137 and describe the constituent components. In Section 6, we describe a host of  $Q$ -  
138 value approximation methods. Among others, this section includes descriptions of  
139 Monte Carlo sampling methods, heuristic approximations, rollout methods, and  
140 the traditional reinforcement learning approach. In Sections 7 and 8, we provide  
141 simulation results on model problems that illustrate several of the approximate  
142 nonmyopic methods described in this paper. We conclude in Section 9 with some  
143 summary remarks.

In addition to providing an expository treatment on the application of POMDPs to the adaptive sensing problem, this paper includes several new and important contributions. First, we introduce a model problem that includes time-varying intervisibility which has all of the desirable properties to completely explore the trade between nonmyopic and myopic scheduling. Second, we introduce several potentially tractable and general numerical methods for generating approximately optimal nonmyopic policies, and show explicitly how they relate to the optimal solution. These include belief-state simplification, completely observable rollout, and reward surrogation, as well as a heuristic based on an information theoretic approximation to the value-to-go function which is applicable in a broad array of scenarios (these contributions have never appeared in journal publications). Finally, these new techniques are compared on a model problem, followed by an in-depth illustration of the value of nonmyopic scheduling on the model problem.

144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156

## 2 Motivating example

157

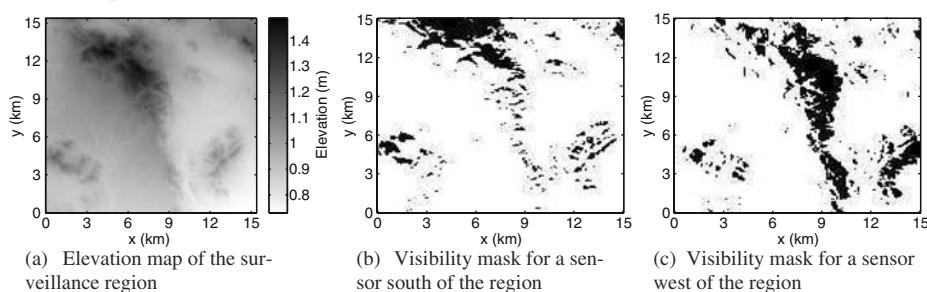
We now present a concrete motivating example that will be used to explain and justify the heuristics and approximations used in this paper. This example involves a remote sensing application where the goal is to learn the contents of a surveillance region via repeated interrogation. (See Hero et al. 2008 for a more complete exposition of adaptive sensing applied to such problems.)

158  
159  
160  
161  
162

Consider a single airborne sensor which is able to image a portion of a ground surveillance region to determine the presence or absence of moving ground targets. At each time epoch, the sensor is able to direct an electrically scanned array so as to interrogate a small area on the ground. Each interrogation yields some (imperfect) information about the small area. The objective is to choose the sequence of pointing directions that lead to the best ability to estimate the entire contents of the surveillance region.

163  
164  
165  
166  
167  
168  
169

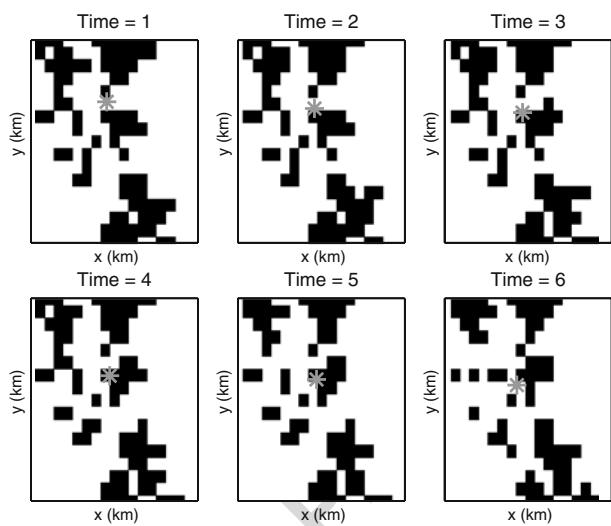
Further complicating matters is the fact that at each time epoch the sensor position causes portions of the ground to be unobservable due to the terrain elevation

170  
171

**Fig. 1** **a** A digital terrain elevation map for a surveillance region, indicating the height of the terrain in the region. **b, c** Visibility masks for a sensor positioned to the south and to the west, respectively, of the surveillance region. We show binary visibility masks (nonvisible areas are black and visible areas are white). In general, visibility may be between 0 and 1 indicating areas of reduced visibility, e.g., regions that are partially obscured by foliage

Q1

**Fig. 2** A six time step vignette where a target moves through an obscured area. Other targets are present elsewhere in the surveillance region. The target is depicted by an asterisk. Areas obscured to the sensor are black and areas that are visible are white. Extra dwells just before becoming obscured (time = 1) aid in localization after the target emerges (time = 6)



172 between the sensor and the ground. Given its position and the terrain elevation, the  
173 sensor can compute a visibility mask which determines how well a particular spot  
174 on the ground can be seen. As an example, in Fig. 1 we give binary visibility masks  
175 computed from a sensor positioned (a) south and (b) to the west of the topologically  
176 nonhomogeneous surveillance region (these plots come from real digital terrain  
177 elevation maps). As can be seen from the figures, sensor position causes “shadowing”  
178 of certain regions. These regions, if measured, would provide no information to  
179 the sensor. A similar target masking effect occurs with atmospheric propagation  
180 attenuation from disturbances such as fog, rain, sleet, or dust, as illustrated in  
181 Section 8. This example illustrates a situation where nonmyopic adaptive sensing is  
182 highly beneficial. Using a known sensor trajectory and known topological map, the  
183 sensor can predict locations that will be obscured in the future. This information can  
184 be used to prioritize resources so that they are used on targets that are predicted to  
185 become obscured in the future. Extra sensor dwells immediately before obscuration  
186 (at the expense of not interrogating other targets) will sharpen the estimate of target  
187 location. This sharpened estimate will allow better prediction of where and when the  
188 target will emerge from the obscured area. This is illustrated graphically with a six  
189 time-step vignette in Fig. 2.

### 190 3 Formulating adaptive sensing problems

#### 191 3.1 Partially observable Markov decision processes

192 An adaptive sensing problem can be posed formally as a *partially observable Markov*  
193 *decision process* (POMDP). Before discussing exactly how this is done, we first need  
194 to introduce POMDPs. Our level of treatment will not be as formal and rigorous as  
195 one would expect from a fullblown course on this topic. Instead, we seek to describe  
196 POMDPs in sufficient detail to allow the reader to see how an adaptive sensing

problem can be posed as a POMDP, and to explore methods to approximate optimal 197  
solutions. Our exposition assumes knowledge of probability, stochastic processes, 198  
and optimization. In particular, we assume some knowledge of Markov processes, 199  
including Markov decision processes, a model that should be familiar to the discrete 200  
event system community. For completeness, we will introduce POMDPs in sufficient 201  
detail to allow the reader to see how an adaptive sensing problem can be posed 202  
as a POMDP, and to explore methods to approximate optimal solutions. For a full 203  
treatment of POMDPs and related background, see Bertsekas (2007). 204

A POMDP is specified by the following ingredients: 205

- A set of states (the state space) and a distribution specifying the random initial 206  
state. 207
- A set of possible actions (the action space). 208
- A state-transition law specifying the next-state distribution given an action taken 209  
at a current state. 210
- A reward function specifying the reward (real number) received given an action 211  
taken at a state. 212
- A set of possible observations (the observation space). 213
- An observation law specifying the distribution of observations given an action 214  
taken at a state. 215

A POMDP is a controlled dynamical process in discrete time. The process begins 216  
at time  $k = 0$  with a (random) initial state. At this state, we perform an action and 217  
receive a reward, which depends on the action and the state. At the same time, we 218  
receive an observation, which again depends on the action and the state. The state 219  
then transitions to some random next state, whose distribution is specified by the 220  
state-transition law. The process then repeats in the same way—at each time, the 221  
process is at some state, and the action taken at that state determines the reward, 222  
observation, and next state. As a result, the state evolves randomly over time in 223  
response to actions, generating observations along the way. 224

We have not said anything so far about the finiteness of the state space or the 225  
sets of actions and observations. The advantage to leaving this issue open is that 226  
it frees us to construct models in the most natural way. Of course, if we are to 227  
represent any such model in a computer, we can only do so in a finite way (though 228  
the finite numbers that can be represented in a computer are typically sufficiently 229  
large to meet practical needs). For example, if we model the motion of a target 230  
on the ground in terms of its Cartesian coordinates, we can deal with this model 231  
in a computer only in a finite sense—specifically, there are only a finite number of 232  
possible locations that can be captured on a standard digital computer. Moreover, 233  
the theory of POMDPs becomes much more technically involved if we are to deal 234  
rigorously with infinite sets. For the sake of technical formality, we will assume 235  
henceforth that the state space, the action space, and the observation space are 236  
all finite (though not necessarily “small”—we stress that this assumption is merely 237  
for technical reasons). However, when thinking about models, we will not explicitly 238  
restrict ourselves to finite sets. For example, it is convenient to use a motion model 239  
for targets in which we view the Cartesian coordinates as real numbers. There is no 240  
harm in this dichotomous approach as long as we understand that ultimately we are 241  
computing only with finite sets. 242

243 3.2 Belief state

244 As a POMDP evolves over time, we do not have direct access to the states that occur.  
245 Instead, all we have are the observations generated over time, providing us with  
246 clues of the actual underlying states (hence the term *partially observable*). These  
247 observations might, in some cases, allow us to infer exactly what states actually  
248 occurred. However, in general, there will be some uncertainty in our knowledge  
249 of the states that actually occurred. This uncertainty is represented by the *belief*  
250 *state* (or *information state*), which is the *a posteriori* (or *posterior*) distribution of the  
251 underlying state given the history of observations.

252 Let  $\mathcal{X}$  denote the state space (the set of all possible states in our POMDP), and  
253 let  $\mathcal{B}$  be the set of distributions over  $\mathcal{X}$ . Then a belief state is simply an element  
254 of  $\mathcal{B}$ . Just as the underlying state changes over time, the belief state also changes  
255 over time. At time  $k = 0$ , the (initial) belief state is equal to the given initial state  
256 distribution. Then, once an action is taken and an observation is received, the belief  
257 state changes to a new belief state, in a way that depends on the observation received  
258 and the state-transition and observation laws. This change in the belief state can be  
259 computed explicitly using Bayes' rule.

260 To elaborate, suppose that the current time is  $k$ , and the current belief state is  
261  $b_k \in \mathcal{B}$ . Note that  $b_k$  is a probability distribution over  $\mathcal{X}$ —we use the notation  $b_k(x)$   
262 for the probability that  $b_k$  assigns to state  $x \in \mathcal{X}$ . Let  $\mathcal{A}$  represent the action space.  
263 Suppose that at time  $k$  we take action  $a_k \in \mathcal{A}$  and, as a result, we receive observation  
264  $y_k$ . Denote the state-transition law by  $P_{\text{trans}}$ , so that the probability of transitioning  
265 to state  $x'$  given that action  $a$  is taken at state  $x$  is  $P_{\text{trans}}(x'|x, a)$ . Similarly, denote the  
266 observation law by  $P_{\text{obs}}$ , so that the probability of receiving observation  $y$  given that  
267 action  $a$  is taken at state  $x$  is  $P_{\text{obs}}(y|x, a)$ . Then, the next belief state given action  $a_k$   
268 is computed using the following two-step update procedure:

269 1. Compute the “updated” belief state  $\hat{b}_k$  based on the observation  $y_k$  of the state  
270  $x_k$  at time  $k$ , using Bayes' rule:

$$\hat{b}_k(x) = \frac{P_{\text{obs}}(y_k|x, a_k)b_k(x)}{\sum_{s \in \mathcal{X}} P_{\text{obs}}(y_k|s, a_k)b_k(s)}, \quad x \in \mathcal{X}.$$

271 2. Compute the belief state  $b_{k+1}$  using the state-transition law:

$$b_{k+1}(x) = \sum_{s \in \mathcal{X}} \hat{b}_k(s) P_{\text{trans}}(x|s, a_k), \quad x \in \mathcal{X}.$$

272 This two-step procedure is commonly realized in terms of a Kalman filter or a particle  
273 filter (Ristic et al. 2004).

274 It is useful to think of a POMDP as a random process of evolving belief states. Just  
275 as the underlying state transitions to some random new state with the performance  
276 of an action at each time, the belief state also transitions to some random new  
277 belief state. So the belief state process also has some “belief-state-transition” law  
278 associated with it, which depends intimately on the underlying state-transition and  
279 the observation laws. But, unlike the underlying state, the belief state is fully  
280 accessible.

281 Indeed, any POMDP may be viewed as a *fully observable* Markov decision process  
282 (MDP) with state space  $\mathcal{B}$ , called the *belief-state MDP* or *information-state MDP*.

(see Bertsekas 2007). To complete the description of this MDP, we will show how 283 to write its reward function, which specifies the reward received when action  $a$  is 284 taken at belief-state  $b$ . Suppose  $b \in \mathcal{B}$  is some belief state and  $a$  is an action. Let 285  $R(x, a)$  be the reward received if action  $a$  is taken at underlying state  $x$ . Then let 286  $r(b, a) = \sum_{x \in \mathcal{X}} b(x)R(x, a)$  be the expected reward with respect to belief-state 287  $b$ , given action  $a$ . This reward  $r(b, a)$  then represents the reward function of the belief- 288 state MDP. 289

### 3.3 Optimization objective

290

Given a POMDP, our goal is to select actions over time to maximize the expected 291 cumulative reward (we take expectation here because the cumulative reward is 292 a random variable). To be specific, suppose we are interested in the expected 293 cumulative reward over a time horizon of length  $H$ :  $k = 0, 1, \dots, H - 1$ . Let  $x_k$  and 294  $a_k$  be the state and action at time  $k$ , and let  $R(x_k, a_k)$  be the resulting reward received. 295 Then, the cumulative reward over horizon  $H$  is given by 296

$$V_H = E \left[ \sum_{k=0}^{H-1} R(x_k, a_k) \right],$$

where  $E$  represents expectation. It is important to realize that this expectation is with 297 respect to  $x_0, x_1, \dots$ ; i.e., the random initial state and all the subsequent states in the 298 evolution of the process, given the actions  $a_0, a_1, a_2, \dots$  taken over time. The goal is 299 to pick these actions so that the objective function is maximized. 300

We have assumed without loss of generality that the reward is a function only 301 of the current state and the action. Indeed, suppose we write the reward such 302 that it depends on the current state, the next state, and the action. We can then 303 take the conditional mean of this reward with respect to the next state, given the 304 current state and action (the conditional distribution of the next state is given by 305 the state-transition law). Because the overall objective function involves expectation, 306 replacing the original reward with its conditional mean in the way described above 307 results in no loss of generality. Finally, notice that the conditional mean of the 308 original reward is a function of the current state and the action, but not the next 309 state. 310

Note that we can also represent the objective function in terms of  $r$  (the reward 311 function of the belief-state MDP) instead of  $R$ : 312

$$V_H(b_0) = E \left[ \sum_{k=0}^{H-1} r(b_k, a_k) \middle| b_0 \right].$$

where  $E[\cdot | b_0]$  represents conditional expectation given  $b_0$ . The expectation now is 313 with respect to  $b_0, b_1, \dots$ ; i.e., the initial belief state and all the subsequent belief 314 states in the evolution of the process. We leave it to the reader to verify this 315 expression involving belief states indeed gives rise to the same objective function 316 value as the earlier expression involving states. In Section 4 we will discuss an 317 equation, due to Bellman, that characterizes this conditional form of the objective 318 function. 319

It is often the case that the horizon  $H$  is very large. In such cases, for technical 320 reasons relevant to the analysis of POMDPs, the objective function is often expressed 321

322 as a limit. A sensible limiting objective function is the *infinite-horizon* (or *long-term*)  
323 *average* reward:

$$\lim_{H \rightarrow \infty} E \left[ \frac{1}{H} \sum_{k=0}^{H-1} R(x_k, a_k) \right].$$

324 Another common limiting objective function is the *infinite-horizon cumulative dis-*  
325 *counted* reward:

$$\lim_{H \rightarrow \infty} E \left[ \sum_{k=0}^{H-1} \gamma^k R(x_k, a_k) \right],$$

326 where  $\gamma \in (0, 1)$  is called the *discount factor*. In this paper, our focus is not on  
327 analytical approaches to solving POMDPs. Therefore, even when dealing with large  
328 horizons, we will not be concerned with the technical considerations involved in  
329 taking the kinds of limits in the above infinite-horizon objective functions (Bertsekas  
330 2007). Instead, we will often imagine that  $H$  is very large but still use the nonlimiting  
331 form.

### 332 3.4 Optimal policy

333 In general, the action chosen at each time should be allowed to depend on the entire  
334 history up to that time (i.e., the action at time  $k$  is a random variable that is a function  
335 of all observable quantities up to time  $k$ ). However, it turns out that if an optimal  
336 choice of such a sequence of actions exists, then there is an optimal choice of actions  
337 that depends only on “belief-state feedback” (see Smallwood and Sondik 1973 and  
338 references therein for the origins of this result). In other words, it suffices for the  
339 action at time  $k$  to depend only on the belief-state  $b_k$  at time  $k$ . So what we seek is,  
340 at each time  $k$ , a mapping  $\pi_k^* : \mathcal{B} \rightarrow \mathcal{A}$  such that if we perform action  $a_k = \pi_k^*(b_k)$ ,  
341 then the resulting objective function is maximized. As usual, we call such a mapping  
342 a *policy*. So, what we seek is an *optimal policy*.

### 343 3.5 POMDPs for adaptive sensing

344 POMDPs form a very general framework based on which many different stochastic  
345 control problems can be posed. Thus, it is no surprise that adaptive sensing problems  
346 can be posed as POMDPs.

347 To formulate an adaptive sensing problem as a POMDP, we need to specify  
348 the POMDP ingredients in terms of the given adaptive sensing problem. This  
349 specification is problem specific. To show the reader how this is done, here we  
350 provide some examples of what aspects of adaptive sensing problems influence how  
351 the POMDP ingredients are specified. As a further illustration, in the next three  
352 sections we specify POMDP models for three example problems, including the  
353 motivating example in Section 2 and the simulations.

354 *States* The POMDP state represents those features in the system (directly observ-  
355 able or not) that possibly evolve over time. Typically, the state is composed of several  
356 parts. These include target positions and velocities, sensor modes of operation,

sensor parameter settings, battery status, data quality, which sensors are active, states 357  
that are internal to tracking algorithms, the position and connectivity of sensors, and 358  
communication resource allocation. 359

*Actions* To specify the actions, we need to identify all the controllable aspects 360  
of the sensing system (those aspects that we wish to control over time in our 361  
adaptive sensing problem). These include sensor mode switching (e.g., waveform 362  
selection or carrier frequencies), pointing directions, sensor tunable parameters, sen- 363  
sor activation status (on/off), sensor position changes, and communication resource 364  
reallocation. 365

*State-transition law* The state-transition law is derived from models representing 366  
how states change over time. Some of these changes are autonomous, while some 367  
are in response to actions. Examples of such changes include target motion, which 368  
sensors were most recently activated, changes in sensor parameter settings, sensor 369  
failures over time, battery status changes based on usage, and changes in the position 370  
and connectivity of sensors. 371

*Reward function* To determine the reward function, we need to first decide on 372  
our overall objective function. To be amenable to POMDP methods, this objective 373  
function must be of the form shown before, namely the mean sum of per-time-step 374  
rewards. Writing the objective function this way automatically specifies the reward 375  
function. For example, if the objective function is the mean cumulative tracking 376  
error, then the reward function simply maps the state at each time to the mean 377  
tracking error at that time. 378

*Observations* The observation at each time represents those features of the system 379  
that depend on the state and are accessible to the controlling agent (i.e., can be used 380  
to inform control decisions). These include sensor outputs (e.g., measurements of 381  
target locations and velocities), and those parts of state that are directly observable 382  
(e.g., battery status), including prior actions. 383

*Observation law* The observation law is derived from models of how the observa- 384  
tions are related to the underlying states. In particular, we will need to use models 385  
of sensors (i.e., the relationship between the sensor outputs and the quantities being 386  
measured), and also models of the sensor network configuration. 387

In the next three sections, we provide examples to illustrate how to formulate 388  
adaptive sensing problems as POMDPs. In the next section, we show how to 389  
formulate an adaptive *classification* problem as a POMDP (with detection problems 390  
being special cases). Then, in the section that follows, we show how to formulate an 391  
adaptive *tracking* problem as a POMDP. Finally, we consider the airborne sensing 392  
problem in Section 2 and describe a POMDP formulation for it. (which also applies 393  
to the simulation example in Section 7). 394

### 3.6 POMDP for an adaptive classification problem 395

We now consider a simple classification problem and show how the POMDP frame- 396  
work can be used to formulate this problem. In particular, we will give specific forms 397

398 for each of the ingredients described in Section 3.5. This simple classification problem  
 399 statement can be used to model problems such as medical diagnostics, nondestructive  
 400 testing, and sensor scheduling for target detection.

401 Our problem is illustrated in Fig. 3. Suppose an object belongs to a particular un-  
 402 known class  $c$ , taking values in a set  $\mathcal{C}$  of possible classes. We can take measurements  
 403 on the object that provide us with information from which we will infer the unknown  
 404 class. These measurements come from a “controlled sensor” at our disposal, which  
 405 we can use at will. Each time we use the sensor, we first have to choose a control  
 406  $u \in \mathcal{U}$ . For each chosen control  $u$ , we get a measurement whose distribution depends  
 407 on  $c$  and  $u$ . Call this distribution  $P_{\text{sensor}}(\cdot|c, u)$  (repeated uses of the sensor generate  
 408 independent measurements). Each time we apply control  $u$ , we incur a cost of  $\kappa(u)$   
 409 (i.e., the cost of using the controlled sensor depends on the control applied). The  
 410 controlled sensor may represent a particular measurement instrument that can be  
 411 controlled (e.g., with different configurations or settings) or may represent a set  
 412 of fixed sensors from which to choose (e.g., a seismic, radar, and induction sensor  
 413 for landmine detection, as discussed in Scott et al. 2004). Notice that detection (i.e.,  
 414 hypothesis testing) is a special case of our problem because it reduces the case where  
 415 there are two classes: present and absent.

416 After each measurement is taken, we have to choose whether or not to produce  
 417 a classification (i.e., an estimate  $\hat{c} \in \mathcal{C}$ ). If we choose to produce such a classification,  
 418 the scenario terminates. If not, we can continue to take another measurement by  
 419 selecting a sensor control. The performance metric of interest here (to be maximized)  
 420 is the probability of correct classification minus the total cost of sensors used.

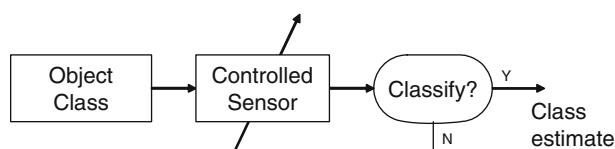
421 To formulate this problem as a POMDP, we must specify the ingredients described  
 422 in Section 3.5: states, actions, state-transition law, reward function, observations, and  
 423 observation law.

424 *States* The possible states in our POMDP formulation of this classification problem  
 425 are the possible classes, together with an extra state to represent that the scenario has  
 426 terminated, which we will denote by  $\tau$ . Therefore, the state space is given by  $\mathcal{C} \cup \{\tau\}$ .  
 427 Note that the state changes only when we choose to produce a classification, as we  
 428 will specify in the state-transition law below.

429 *Actions* The actions here are of two kinds: we can either choose to take a mea-  
 430 surement, in which case the action is the sensor control  $u \in \mathcal{U}$ , or we can choose to  
 431 produce a classification, in which case the action is the class  $\hat{c} \in \mathcal{C}$ . Hence, the action  
 432 space is given by  $\mathcal{U} \cup \mathcal{C}$ .

433 *State-transition law* The state-transition law represents how the state evolves at  
 434 each time step as a function of the action. As pointed out before, as long as we are  
 435 taking measurements, the state does not change (because it represents the unknown

**Fig. 3** An adaptive classification system



object class). As soon as we choose to produce a classification, the state changes to 436  
the terminal state  $\tau$ . Therefore, the state-transition law  $P_{\text{trans}}$  is given by 437

$$P_{\text{trans}}(x'|x, a) = \begin{cases} 1 & \text{if } a \in \mathcal{U} \text{ and } x' = x \\ 1 & \text{if } a \in \mathcal{C} \text{ and } x' = \tau \\ 0 & \text{otherwise.} \end{cases}$$

*Reward function* The reward function  $R$  here is given by 438

$$R(x, a) = \begin{cases} -\kappa(a) & \text{if } a \in \mathcal{U} \text{ and } x \neq \tau \\ 1 & \text{if } a \in \mathcal{C} \text{ and } x = a \\ 0 & \text{otherwise.} \end{cases}$$

If we produce a classification, then the reward is 1 if the classification is correct, and 439  
otherwise it is 0. Hence, the mean of the reward when producing a classification is 440  
the probability that the classification is correct. If we use the finite-horizon objective 441  
function with horizon  $H$ , then the objective function represents the probability of 442  
producing a correct classification within the time horizon of  $H$  (e.g., representing 443  
some maximum time limit for producing a classification) minus the total sensing cost. 444

*Observations* The observations in this problem represent the sensor outputs (mea- 445  
surements). The observation space is therefore the set of possible measurements. 446

*Observation law* The observation law specifies the distribution of the observations 447  
given the state and action. So, if  $x \in \mathcal{C}$  and  $a \in \mathcal{U}$ , then the observation law is given by 448  
 $P_{\text{sensor}}(\cdot|x, a)$ . If  $x = \tau$ , then we can define the observation law arbitrarily, because it 449  
does not affect the solution to the problem (recall that after the scenario terminates, 450  
represented by being in state  $\tau$ , we no longer take any measurements). 451

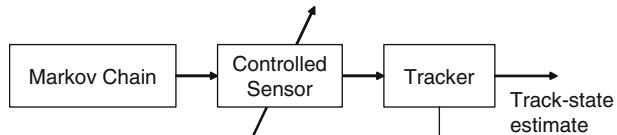
Note that as long as we are still taking measurements and have not yet produced a 452  
classification, the belief state for this problem represents the *a posteriori* distribution 453  
of the unknown class being estimated. It is straightforward to show that the optimal 454  
policy for this problem will always produce a classification that maximizes the *a* 455  
*posteriori* probability (i.e., is a “MAP” classifier). However, it is not straightforward 456  
to deduce exactly when we should continue to take measurements and when we 457  
should produce a classification. Determining such an optimal policy requires solving 458  
the POMDP. 459

### 3.7 POMDP for an adaptive tracking problem

We now consider a simple tracking problem and show how to formulate it using a 461  
POMDP framework. Our problem is illustrated in Fig. 4. We have a Markov chain 462  
with state space  $\mathcal{S}$  evolving according to a state-transition law given by  $T$  (i.e., for 463  
 $s, s' \in \mathcal{S}$ ,  $T(s'|s)$  is the probability of transitioning to state  $s'$  given that the state is 464  
 $s$ ). We assume that  $\mathcal{S}$  is a metric space—there is a function  $d : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$  such that 465  
 $d(s, s')$  represents a “distance” measure between  $s$  and  $s'$ .<sup>1</sup> The states of this Markov 466

<sup>1</sup>For the case where  $\mathcal{S}$  represents target kinematic states in Cartesian coordinates, we typically use the Euclidean norm for this metric.

**Fig. 4** An adaptive tracking system



467 chain are not directly accessible—they represent quantities to be tracked over time  
468 (e.g., the coordinates and velocities of targets).

469 To do the tracking, as in the last section, we exploit measurements from a  
470 “controlled sensor” over time. At each time step, we first have to choose a control  
471  $u \in \mathcal{U}$ . For each chosen control  $u$ , we get a measurement whose distribution depends  
472 on the Markov chain state  $s$  and control  $u$ , denoted  $P_{\text{sensor}}(\cdot|s, u)$  as before (again,  
473 we assume that sensor measurements over time are independent). Each time we  
474 apply control  $u$ , we incur a cost of  $\kappa(u)$  (i.e., as in the last example, the cost of using  
475 the controlled sensor depends on the control applied). As in the last example, the  
476 controlled sensor may represent a particular measurement instrument that can be  
477 controlled (e.g., with different configurations or settings) or may represent a set of  
478 fixed sensor assets from which to choose (e.g., multiple sensors distributed over a  
479 geographical region, where the control here is which subset of sensors to activate, as  
480 in He and Chong (2004, 2006), Krakow et al. (2006), Li et al. (2006, 2007)).

481 Each measurement is fed to a tracker, which is an algorithm that produces an  
482 estimate  $\hat{s}_k \in \mathcal{S}$  of the state at each time  $k$ . For example, the tracker could be a  
483 Kalman filter or a particle filter (Ristic et al. 2004). The tracker has an internal  
484 state, which we will denote  $z_k \in \mathcal{Z}$ . The internal state is updated as a function of  
485 measurements:

$$z_{k+1} = f_{\text{tracker}}(z_k, y_k),$$

486 where  $y_k$  is the measurement generated at time  $k$  as a result of control  $u_k$  (i.e., if  
487 the Markov chain state at time  $k$  is  $s_k$ , then  $y_k$  has distribution  $P_{\text{sensor}}(\cdot|s_k, u_k)$ ). The  
488 estimate  $\hat{s}_k$  is a function of this internal state  $z_k$ . For example, in the case of a Kalman  
489 filter, the internal state represents a mean vector together with a covariance matrix.  
490 The output  $\hat{s}_k$  is usually simply the mean vector. In the case of a particle filter,  
491 the internal state represents a set of particles. See Ristic et al. (2004) for explicit  
492 equations to represent  $f_{\text{tracker}}$ .

493 The performance metric of interest here (to be maximized) is the negative mean  
494 of the sum of the cumulative tracking error and the sensor usage cost over a horizon  
495 of  $H$  time steps. To be precise, the tracking error at time  $k$  is the “distance” between  
496 the output of the tracker,  $\hat{s}_k$ , and the true Markov chain state,  $s_k$ . Recall that the  
497 “distance” here is well-defined because we have assumed that  $\mathcal{S}$  is a metric space. So  
498 the tracking error at time  $k$  is  $d(\hat{s}_k, s_k)$ .

499 As in the last section, to formulate this adaptive tracking problem as a POMDP,  
500 we must specify the ingredients described in Section 3.5: states, actions, state-  
501 transition law, reward function, observations, and observation law.

502 *States* It might be tempting to define the state space for this problem simply to be  
503 the state space for the Markov chain,  $\mathcal{S}$ . However, it is important to point out that  
504 the tracker also contains an internal state, and the POMDP state should take both  
505 into account. Accordingly, for this problem we will take the state at time  $k$  to be the

pair  $[s_k, z_k]$ , where  $s_k$  is the state of the Markov chain to be tracked, and  $z_k$  is the tracker state. Hence, the state space is  $\mathcal{S} \times \mathcal{Z}$ . 506  
507

*Actions* The actions here are the controls applied to the controlled sensor. Hence, 508  
the action space is simply  $\mathcal{U}$ . 509

*State-transition law* The state-transition law specifies how the state changes at 510 each time  $k$ , given the action  $a_k$  at that time. Recall that the state at time  $k$  is 511 the pair  $[s_k, z_k]$ . The Markov chain state  $s_k$  makes a transition according to the 512 transition probability  $T(\cdot|s_k)$ . The tracker state  $z_k$  makes a transition depending on 513 the observation  $y_k$ . In other words, the transition distribution for the next tracker 514 state given  $z_k$  is the distribution of  $f_{\text{tracker}}(z_k, y_k)$  (which in turn depends on the 515 measurement distribution  $P_{\text{sensor}}(\cdot|s_k, a_k)$ ). This completely specifies the distribution 516 of  $[s_{k+1}, z_{k+1}]$  as a function of  $[s_k, z_k]$  and  $a_k$ . 517

*Reward function* The reward function is given by 518

$$R([s_k, z_k], a_k) = -(d(\hat{s}_k, s_k) + \kappa(a_k)),$$

where the reader should recall that the tracker output  $\hat{s}_k$  is a function of  $z_k$ . Notice 519 that the first term in the (per-time-step) reward, which represents tracking error, is 520 not a function of  $a_k$ . Instead, the tracking errors depend on the actions applied over 521 time through the track estimates  $\hat{s}_k$  (which in turn depend on the actions through the 522 distributions of the measurements). 523

*Observations* As in the previous example, the observations here represent the sensor 524 outputs (measurements). The observation space is therefore the set of possible 525 measurements. 526

*Observation law* The observation law is given by the measurement distribution 527  $P_{\text{sensor}}(\cdot|s_k, a_k)$ . Note that the observation law does not depend on  $z_k$ , the tracker 528 state, even though  $z_k$  is part of the POMDP state. 529

### 3.8 POMDP for motivating example 530

In this section, we give mathematical forms for each of the ingredients listed in 531 Section 3.5 for the motivating example described in Section 2 (these also apply to 532 the simulation example in Section 7). To review, the motivating example dealt with 533 an airborne sensor charged with detecting and tracking multiple moving targets. The 534 airborne sensor is agile in that it can steer its beam to different ground locations. Each 535 interrogation of the ground results in an observation as to the absence or presence 536 of targets in the vicinity. The adaptive sensing problem is to use the collection of 537 measurements made up to the current time to determine the best place to point next. 538

*States* In this motivating problem, we are detecting and tracking  $N$  moving ground 539 targets. For the purposes of this discussion we assume that  $N$  is known and fixed, and 540 that the targets are moving in 2 dimensions (a more general treatment, where the 541 number of targets is both unknown and time varying, is given elsewhere (Kreucher 542 et al. 2005c)). We denote these positions as  $x_1, \dots, x_N$  where  $x_i$  is a 2-dimensional 543

544 vector corresponding to target  $i$ . Furthermore, because of the terrain, the position  
 545 of the sensor influences the visibility of certain locations on the ground, so sensor  
 546 position is an important component of the state. Denote the (directly observable)  
 547 3-dimensional sensor position by  $\sigma$ . Then the state space  $\mathcal{X}$  consists of real-valued  
 548 vectors in  $\mathbb{R}^{2N+3}$ , i.e., each state takes the form

$$x = [x_1, x_2, \dots, x_{N-1}, x_N, \sigma].$$

549 Although not explicitly shown here, the surveillance region topology is assumed  
 550 known and considered part of the problem specification. This specification affects the  
 551 observation law, as we shall see below.

552 *Actions* The airborne sensor is able to measure a single detection cell and make  
 553 an imperfect measurement as to the presence or absence of a target in that cell.  
 554 Therefore, the action  $a \in \{1, \dots, C\}$  is an integer specifying which of the  $C$  discrete  
 555 cells is measured.

556 *State-transition law* The state-transition law describes the distribution of the next  
 557 state vector  $x' = [x'_1, x'_2, \dots, x'_N, \sigma']$  conditioned on the current state vector  $x =$   
 558  $[x_1, x_2, \dots, x_N, \sigma]$  and the action  $a$ . Because our states are vectors in  $\mathbb{R}^{2N+3}$ , we will  
 559 specify the state-transition law as a conditional density function. For simplicity, we  
 560 have chosen to model the evolution of each of the  $N$  targets as independent and  
 561 following a Gaussian law, i.e.,

$$T_{\text{single target}}(x'_i|x_i) = \frac{1}{2\pi|\Sigma|^{1/2}} \exp^{-\frac{1}{2}(x_i - x'_i)^\top \Sigma^{-1}(x_i - x'_i)}, \quad i = 1, \dots, N$$

562 (where  $x_i$  and  $x'_i$  are treated here as column vectors). In other words, each target  
 563 moves according to a random walk (purely diffusive). Because of our independence  
 564 assumption, we can write the joint target-motion law as

$$T_{\text{target}}(x'_1, \dots, x'_N|x_1, \dots, x_N) = \prod_{i=1}^N T_{\text{single target}}(x'_i|x_i).$$

565 The temporal evolution of the sensor position is assumed deterministic and known  
 566 precisely (i.e., the aircraft is flying a pre-planned pattern). We use  $f(\sigma)$  to denote  
 567 the sensor trajectory function, which specifies the next position of the sensor given  
 568 current sensor position  $\sigma$ ; i.e., if the current sensor position is  $\sigma$ , then  $f(\sigma)$  is exactly  
 569 the next sensor position. Then, the motion law for the sensor is

$$T_{\text{sensor}}(\sigma'|\sigma) = \delta(\sigma' - f(\sigma)).$$

570 With these assumptions, the state-transition law is completely specified by

$$P_{\text{trans}}(x'|x, a) = T_{\text{target}}(x'_1, \dots, x'_N|x_1, \dots, x_N) T_{\text{sensor}}(\sigma'|\sigma).$$

571 Note that according to our assumptions, the actions taken do not affect the state  
 572 evolution. In particular, we assume that the targets do not know they are under  
 573 surveillance and consequently they do not take evasive action (see Kreucher et al.  
 574 2006 for a model that includes evasion).

*Reward function* In previous work (Kreucher et al. 2005b), we have found that 575  
*information gain* provides a useful metric that captures a wide variety of goals. 576  
Information gain is a metric that measures the relative information increase between 577  
a prior belief state and a posterior belief state, i.e., it measures the benefit a particular 578  
observation has yielded. An information theoretic metric is intuitively pleasing as it 579  
measures different types of benefits (e.g., information about the number of targets 580  
present versus information about the positions of individual targets) on an equal 581  
footing, that of information gain. Furthermore, it has been shown that information 582  
gain can be viewed as a near universal proxy for any risk function (Kreucher et al. 583  
2005a). Therefore, the reward used in this application is the gain in information 584  
between the belief state before a measurement  $b_k$  and the (measurement updated) 585  
belief state after a measurement is made  $\hat{b}_k$ . We use a particular information metric 586  
called the Renyi divergence, defined as follows. The Renyi divergence of two belief 587  
states  $p$  and  $q$  is given by 588

$$D_\alpha(p||q) = \frac{1}{\alpha - 1} \ln \sum_{x \in \mathcal{X}} p(x)^\alpha q(x)^{1-\alpha}$$

where  $\alpha > 0$ . To define the reward  $r(b, a)$  in our context, given a belief state  $b$  and 589  
an action  $a$ , we first write, 590

$$\Delta_\alpha(b, a, y) = D_\alpha(\hat{b} || b),$$

where  $y$  is an observation with distribution given by the observation law  $P_{\text{obs}}(\cdot | b, a)$  591  
and  $\hat{b}$  is the “updated” belief state computed as described earlier in Section 3.2 using 592  
Bayes’ rule and knowledge of  $b, a$ , and  $y$ . Note that  $\Delta_\alpha(b, a, y)$  is a random variable 593  
because it is a function of the random observation  $y$ , and hence its distribution 594  
depends on  $a$ . We will call this random variable the *myopic information gain*. 595  
The reward function is defined in terms of the myopic information gain by taking 596  
expectation:  $r(b, a) = E[\Delta_\alpha(b, a, y) | b, a]$ . 597

*Observations* When a cell is interrogated, the sensor receives return energy and 598  
thresholds this energy to determine whether it is to be declared a detection or a 599  
nondetection. This imperfect measurement gives evidence as to the presence or 600  
absence of targets in the cell. Additionally, the current sensor position is directly 601  
observable. Therefore, the observation is given by  $[z, \sigma]$ , where  $z \in \{0, 1\}$  is the one- 602  
bit observation representing detection or nondetection, and  $\sigma$  is the position of the 603  
sensor. 604

*Observation law* Detection/nondetection is assumed to result from thresholding a 605  
Rayleigh-distributed random variable that characterizes the energy returned from an 606  
interrogation of the ground. The performance is completely specified by a probability 607  
of detection  $P_d$  and a false alarm rate  $P_f$ , which under the Rayleigh assumption are 608  
linked by a signal-to-noise-plus-clutter ratio,  $SNCR$ , by  $P_d = P_f^{1/(1+SNCR)}$ . 609

To precisely specify the observation model, we make the following notational 610  
definitions. First, let  $o_a(x_1, \dots, x_N)$  denote the occupation indicator function for cell 611  
 $a$ , defined as  $o_a(x_1, \dots, x_N) = 1$  when at least one of the targets projects into sensor 612  
cell  $a$  (i.e., at least one of the  $x_i$  locations are within cell  $a$ ), and  $o_a(x_1, \dots, x_N) = 0$  613  
otherwise. Furthermore, let  $v_a(\sigma)$  denote the visibility indicator function for cell  $a$ , 614

615 defined as  $v_a(\sigma) = 1$  when cell  $a$  is visible from a sensor positioned at  $\sigma$  (i.e., there is  
616 no line of sight obstruction between the sensor and the cell), and  $v_a(\sigma) = 0$  otherwise.  
617 Then the probability of receiving a detection given state  $x = [x_1, \dots, x_N, \sigma]$  and  
618 action  $a$  is

$$P_{\text{det}}(x, a) = \begin{cases} P_d & \text{if } o_a(x_1, \dots, x_N)v_a(\sigma) = 1 \\ P_f & \text{if } o_a(x_1, \dots, x_N)v_a(\sigma) = 0. \end{cases}$$

619 Therefore, the observation law is specified completely by

$$P_{\text{obs}}(z|x, a) = \begin{cases} P_{\text{det}}(x, a) & \text{if } z = 1 \\ 1 - P_{\text{det}}(x, a) & \text{if } z = 0. \end{cases}$$

## 620 4 Basic principle: *Q*-value approximation

### 621 4.1 Overview and history

622 In this section, we describe the basic principle underlying approximate methods to  
623 solve adaptive sensing problems that are posed as POMDPs. This basic principle is  
624 due to Bellman, and gives rise to a natural framework in which to discuss a variety of  
625 approximation approaches. Specifically, these approximation methods all boil down  
626 to the problem of approximating *Q*-values.

627 Methods for solving POMDPs have their roots in the field of optimal control,  
628 which dates back to the end of the seventeenth century with the work of Johann  
629 Bernoulli (Willems 1996). This field received significant interest in the middle of  
630 the twentieth century, when much of the modern methodology was developed, most  
631 notably by Bellman (1957), who applied *dynamic programming* to bear on optimal  
632 control, and Pontryagin et al. (1962), who introduced his celebrated *maximum*  
633 *principle* based on calculus of variations. Since then, the field of optimal control has  
634 enjoyed much fruit in its application to control problems arising in engineering and  
635 economics.

636 The recent history of methods to solve optimal stochastic decision problems took  
637 an interesting turn in the second half of the twentieth century with the work of  
638 computer scientists in the field of artificial intelligence seeking to solve “planning”  
639 problems (roughly analogous to what engineers and economists call optimal control  
640 problems). The results of their work most relevant to the POMDP methods discussed  
641 here are reported in a number of treatises from the 80s and 90s (Cheng 1988;  
642 Kaelbling et al. 1996, 1998; Zhang and Liu 1996). The methods developed in the  
643 artificial intelligence (machine learning) community aim to provide computationally  
644 feasible approximations to optimal solutions for complex planning problems under  
645 uncertainty. The operations research literature has also continued to reflect ongoing  
646 interest in computationally feasible methods for optimal decision problems (Lovejoy  
647 1991b; Chang et al. 2007; Powell 2007).

648 The connection between the significant work done in the artificial intelligence  
649 community and those of the earlier work on optimal control is noted by Bertsekas  
650 and Tsitsiklis in their 1996 book (Bertsekas and Tsitsiklis 1996). In particular, they  
651 note that the developments in *reinforcement learning*—the approach taken by arti-  
652 ficial intelligence researchers for solving planning problems—is most appropriately

understood in the framework of Markov decision theory and dynamic programming. 653  
 This framework is now widely reflected in the artificial intelligence literature (Kael- 654  
 bling et al. 1996, 1998; Zhang and Liu 1996; Thrun et al. 2005). Our treatment in this 655  
 paper rests on this firm and rich foundation (though our focus is not on reinforcement 656  
 learning methods). 657

#### 4.2 Bellman's principle and $Q$ -values 658

The key result in Markov decision theory relevant here is Bellman's principle. Let 659  
 $V_H^*(b_0)$  be the optimal objective function value (over horizon  $H$ ) with  $b_0$  as the initial 660  
 belief state. Then, *Bellman's principle* states that 661

$$V_H^*(b_0) = \max_a (r(b_0, a) + E[V_{H-1}^*(b_1)|b_0, a])$$

where  $b_1$  is the random next belief state (with distribution depending on  $a$ ), and 662  
 $E[\cdot|b_0, a]$  represents conditional expectation with respect to the random next state 663  
 $b_1$ , whose distribution depends on  $b_0$  and  $a$ . Moreover, 664

$$\pi_0^*(b_0) = \arg \max_a (r(b_0, a) + E[V_{H-1}^*(b_1)|b_0, a])$$

is an optimal policy. 665

Define the *Q-value* of taking action  $a$  at state  $b_k$  as 666

$$Q_{H-k}(b_k, a) = r(b_k, a) + E[V_{H-k-1}^*(b_{k+1})|b_k, a],$$

where  $b_{k+1}$  is the random next belief state (which depends on the observation  $y_k$  at 667  
 time  $k$ , as described in Section 3.2). Then, Bellman's principle can be rewritten as 668

$$\pi_k^*(b_k) = \arg \max_a Q_{H-k}(b_k, a)$$

i.e., the optimal action at belief-state  $b_k$  (at time  $k$ , with a horizon-to-go of  $H - k$ ) is 669  
 the one with largest *Q*-value at that belief state. This principle, called *lookahead*, is 670  
 the heart of POMDP solution approaches. 671

#### 4.3 Stationary policies 672

In general, an optimal policy is a function of time  $k$ . If  $H$  is sufficiently large, then 673  
 the optimal policy is approximately *stationary* (independent of  $k$ ). This is intuitively 674  
 clear: if the end of the time horizon is a million years away, then how we should act 675  
 today given a belief-state is the same as how we should act tomorrow with the same 676  
 belief state. Said differently, if  $H$  is sufficiently large, the difference between  $Q_H$  and 677  
 $Q_{H-1}$  is negligible. Moreover, if needed we can always incorporate time itself into the 678  
 definition of the state, so that dependence on time is captured simply as dependence 679  
 on state. 680

Henceforth we will assume for convenience there is a stationary optimal policy, 681  
 and this is what we seek. We will use the notation  $\pi$  for stationary policies (with 682  
 no subscript  $k$ )—this significantly simplifies the notation. Our approach is equally 683

684 applicable to the short-horizon, nonstationary case, with appropriate notational  
685 modification (to account for the time dependence of decisions).

#### 686 4.4 Receding horizon

687 Assuming  $H$  is sufficiently large and that we seek a stationary optimal policy, at any  
688 time  $k$  we write:

$$\pi^*(b) = \arg \max_a Q_H(b, a).$$

689 Notice that the horizon is taken to be fixed at  $H$ , regardless of the current time  $k$ . This  
690 is justified by our assumption that  $H$  is so large that at any time  $k$ , the horizon is still  
691 approximately  $H$  time steps away. This approach of taking the horizon to be fixed at  
692  $H$  is called *receding horizon control*. For convenience, we will also henceforth drop  
693 the subscript  $H$  from our notation (unless the subscript is explicitly needed).

#### 694 4.5 Approximating $Q$ -values

695 Recall  $Q(b, a)$  is the reward  $r(b, a)$  of taken action  $a$  at belief-state  $b$  plus the  
696 expected cumulative reward of applying the optimal policy for all future actions.  
697 This second term in the  $Q$ -value is in general difficult to obtain, especially when the  
698 belief-state is large. For this reason, approximation methods are necessary to obtain  
699  $Q$ -values. Note that the quality of an approximation is not so much in the accuracy  
700 of the actual  $Q$ -values obtained, but in the *ranking* of the actions reflected by their  
701 *relative* values.

702 In Section 6, we describe a variety of methods to approximate  $Q$ -values. But  
703 before discussing such methods, we first describe the basic control framework for  
704 using  $Q$ -values to inform control decisions.

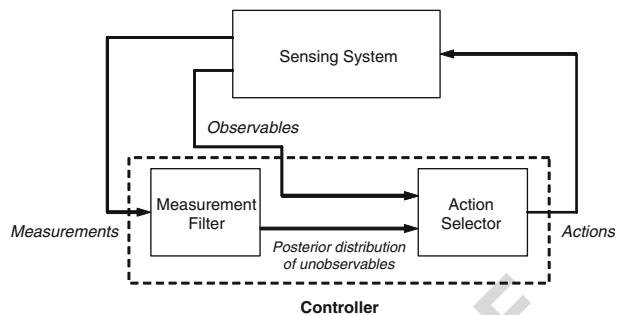
### 705 5 Basic control architecture

706 By Bellman's principle, knowing the  $Q$ -values allows us to make optimal control  
707 decisions. In particular, if we are currently at belief-state  $b$ , we need only find the  
708 action  $a$  with the largest  $Q(b, a)$ . This principle yields a basic control framework  
709 that is illustrated in Fig. 5. The top-most block represents the sensing system, which  
710 we treat as having an input and two forms of output. The input represents actions  
711 (external control commands) we can apply to control the sensing system. Actions  
712 usually include sensor-resource controls, such as which sensor(s) to activate, at what  
713 power level, where to point, what waveforms to use, and what sensing modes to  
714 activate. Actions may also include communication-resource controls, such as the data  
715 rate for transmission from each sensor.

716 The two forms of outputs from the sensing system represent:

- 717 1) Fully observable aspects of the internal state of the sensing system (called  
718 *observables*), and
- 719 2) Measurements (observations) of those aspects of the internal state that are not  
720 directly observable (which we refer to simply as *measurements*).

**Fig. 5** Basic lookahead framework



We assume that the underlying state-space is the Cartesian product of two sets, one representing unobservables and the other representing observables. Target states are prime examples of unobservables. So, measurements are typically the outputs of sensors, representing observations of target states. Observables include things like sensor locations and orientations, which sensors are activated, battery status readings, etc. In the remainder of this section, we describe the components of our control framework. Our description starts from the architecture of Fig. 5 and progressively fills in the details.

## 5.1 Controller

729

At each decision epoch, the *controller* takes the outputs (measurements and observables) from the sensing system and, in return, generates an action that is fed back to the sensing system. This basic closed-loop architecture is familiar to mainstream control system design approaches.

730  
731  
732  
733

The controller has two main components. The first is the *measurement filter*, which takes as input the measurements, and provides as output the *a posteriori* (posterior) distribution of unobservable internal states (henceforth called *unobservables*). In the typical situation where the unobservables are target states, the measurement filter outputs a posterior distribution on target states given the measurement history. The measurement filter is discussed further below. The posterior distribution of the unobservables in addition to the observables form the belief state, the posterior distribution of the underlying state. The second component is the *action selector*, which takes the belief state and computes an action (the output of the controller). The basis for action selection is Bellman's principle, using *Q*-values. This is discussed below.

734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744

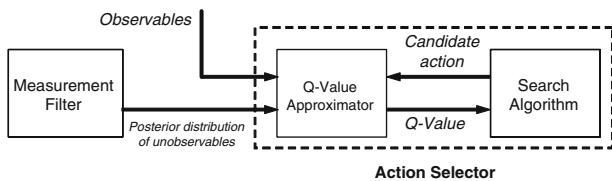
## 5.2 Measurement filter

745

The measurement filter computes the posterior distribution given measurements. This component is present in virtually all target-tracking systems. It turns out that the posterior distribution can be computed iteratively: each time we obtain a new measurement, the posterior distribution can be obtained by updating the previous posterior distribution based on knowing the current action, the transition law, and the observation law. This update is based on Bayes' rule, described earlier in Section 3.2.

746  
747  
748  
749  
750  
751

**Fig. 6** Basic components of the action selector



752 The measurement filter can be constructed in a number of ways. If the posterior  
 753 distribution always resides within a family of distributions that is conveniently para-  
 754 meterized, then all we need to do is keep track of the belief-state parameters. This  
 755 is the case, for example, if the belief state is Gaussian. Indeed, if the unobservables  
 756 evolve in a linear fashion, then these Gaussian parameters can be updated using a  
 757 Kalman filter. In general, however, it is not practical to keep track of the exact belief  
 758 state. Indeed, a variety of options have been explored for belief-state representation  
 759 and simplification (e.g., Rust 1997; Roy et al. 2005; Yu and Bertsekas 2004). We will  
 760 have more to say about belief-state simplification in Section 6.11.

761 Particle filtering is a Monte Carlo sampling method for updating posterior distri-  
 762 butions. Instead of maintaining the exact posterior distribution, we maintain a set of  
 763 representative samples from that distribution. It turns out that this method dovetails  
 764 naturally with Monte Carlo sampling-based methods for  $Q$ -value approximation, as  
 765 we will describe later in Section 6.8.

### 766 5.3 Action selector

767 As shown in Fig. 6, the action selector consists of a search (optimization) algorithm  
 768 that optimizes an objective function, the  $Q$ -function, with respect to an action. In  
 769 other words, the  $Q$ -function is a function of the action—it maps each action, at a  
 770 given belief state, to its  $Q$ -value. The action that we seek is one that maximizes the  
 771  $Q$ -function. So, we can think of the  $Q$ -function as a kind of “action-utility” function  
 772 that we wish to maximize. The search algorithm iteratively generates a candidate  
 773 action and evaluates the  $Q$ -function at this action (this numerical quantity is the  $Q$ -  
 774 value), searching over the space of candidate actions for one with the largest  $Q$ -value.  
 775 Methods for obtaining (approximating) the  $Q$ -values is described in the next section.

## 776 6 $Q$ -value approximation methods

### 777 6.1 Basic approach

778 Recall the definition of the  $Q$ -value,

$$Q(b, a) = r(b, a) + \mathbb{E} [V^*(b') | b, a], \quad (1)$$

779 where  $b'$  is the random next belief state (with distribution depending on  $a$ ). In all but  
 780 very special problems, it is impossible to compute the  $Q$ -value exactly. In this section,  
 781 we describe a variety of methods to approximate the  $Q$ -value. Because the first term  
 782 on the right-hand side of (1) is usually easy to compute, most approximation methods  
 783 focus on the second term. As pointed out before, it is important to realize that the

quality of an approximation to the  $Q$ -value is not so much in the accuracy of the 784  
actual values obtained, but in the *ranking* of the actions reflected by their *relative* 785  
values. 786

We should point out that each of the approximation methods presented in this 787  
section has its own domain of applicability. Traditional reinforcement learning 788  
approaches (Section 6.6), predicated on running a large number of simulations to 789  
“train,” are broadly applicable as they only require a generative model. However, 790  
these methods often have infeasible computational burden owing to the long training 791  
time required for some problems. Furthermore, there is an extensibility problem, 792  
where a trained function may perform very poorly if the problem changes slightly 793  
between the training stage and the application stage. To address these concerns, 794  
we present several sampling techniques (Sections 6.2, 6.8, 6.9, 6.11) which are also 795  
very broadly applicable as they only require a generative model. These methods 796  
do not require a training phase, *per se*, but do on-line estimation. However, in 797  
some instances, these too may require more computations than desirable. Simi- 798  
larly, parametric approximations (Section 6.5) and action-sequence approximations 799  
(Section 6.7) are general in applicability but may entail excessive computational 800  
requirements. Relaxation methods (Section 6.3) and heuristics (Section 6.4) may 801  
provide reduced computation but require advanced domain knowledge. 802

## 6.2 Monte Carlo sampling 803

In general, we can think of Monte Carlo methods simply as the use of computer 804  
generated random numbers in computing expectations of random variables through 805  
averaging over many samples. With this in mind, it seems natural to consider using 806  
Monte Carlo methods to compute the value function directly based on Bellman’s 807  
equation: 808

$$V_H^*(b_0) = \max_{a_0} (r(b_0, a_0) + E[V_{H-1}^*(b_1)|b_0, a_0]).$$

Notice that the second term on the right-hand side involves expectations (one per 809  
action candidate  $a_0$ ), which can be computed using Monte Carlo sampling. However, 810  
the random variable inside each expectation is itself an objective function value 811  
(with horizon  $H - 1$ ), and so it too involves a max of an expectation via Bellman’s 812  
equation: 813

$$V_H^*(b_0) = \max_{a_0} \left( r(b_0, a_0) + E \left[ \max_{a_1} (r(b_1, a_1) + E[V_{H-2}^*(b_2)|b_1, a_1]) \middle| b_0, a_0 \right] \right).$$

Notice we now have two “layers” of max and expectation, one “nested” within 814  
the other. Again, we see the inside expectation involves the value function (with 815  
horizon  $H - 2$ ), which again can be written as a max of expectations. Proceeding 816  
this way, we can write  $V_H^*(b_0)$  in terms of  $H$  layers of max and expectations. Each 817  
expectation can be computed using Monte Carlo sampling. The remaining question 818  
is how computationally burdensome is this task? 819

Kearns et al. (1999) have provided a method to calculate the computational 820  
burden of approximating the value function using Monte Carlo sampling as described 821  
above, given some prescribed accuracy in the approximation of the value function. 822

823 Unfortunately, it turns out that for practical POMDP problems this computational  
824 burden is prohibitive, even for modest degrees of accuracy. So, while Bellman's  
825 equation suggests a natural Monte Carlo method for approximating the value  
826 function, the method is not useful in practice. For this reason, we seek alternative  
827 approximation methods. In the next few subsections, we explore some of these  
828 methods.

### 829 6.3 Relaxation of optimization problem

830 Some problems that are difficult to solve become drastically easier if we *relax* certain  
831 aspects of the problem. For example, by removing a constraint in the problem,  
832 the "relaxed" problem may yield to well-known solution methods. This constraint  
833 relaxation enlarges the constraint set, and so the solution obtained may no longer  
834 be feasible in the original problem. However, the objective function value of the  
835 solution *bounds* the optimal objective function value of the original problem.

836 The  $Q$ -value involves the quantity  $V^*(b')$ , which can be viewed as the optimal  
837 objective function value corresponding to some optimization problem. The method  
838 of relaxation, if applicable, gives rise to a bound on  $V^*(b')$ , which then provides an  
839 approximation to the  $Q$ -value. For example, a relaxation of the original POMDP  
840 may result in a bandit problem (see Krishnamurthy and Evans 2001; Krishnamurthy  
841 2005); or may be solvable via linear programming (see de Farias and Van Roy  
842 2003, 2004). (See also specific applications to sensor management Castanon 1997;  
843 Washburn et al. 2002.) In general, the quality of this approximation is a function of  
844 the specific relaxation and is very problem specific. For example, Castanon (1997)  
845 suggests that in his setting his relaxation approach is feasible for generating near-  
846 optimal solutions. Additionally, Washburn et al. (2002) show that the performance of  
847 their index rule is eclipsed by that of multi-step lookahead under certain conditions  
848 of the process noise, while being much closer in the low-noise situation. While it  
849 is sometimes possible to apply analytical approaches to a relaxed version of the  
850 problem, it is generally accepted that problems that can be posed as POMDPs are  
851 unlikely to be amenable to analytical solution approaches.

852 Bounds on the optimal objective function value can also be obtained by approx-  
853 imating the state space. Lovejoy (1991a) shows how to approximate the state space  
854 by a finite grid of points, and use that grid to construct upper and lower bounds on  
855 the optimal objective function.

### 856 6.4 Heuristic approximation

857 In some applications we are unable to compute  $Q$ -values directly, but can use domain  
858 knowledge to develop an idea of its behavior. If so, we can heuristically construct a  
859  $Q$ -function based on this knowledge.

860 Recall from (1) that the  $Q$ -value is the sum of two terms, where the first term  
861 (the immediate reward) is usually easy to compute. Therefore, it often suffices  
862 to approximate only the second term in (1), which is the mean optimal objective  
863 function value starting at the next belief state, which we call the *expected value-to-go*  
864 (EVTG). (Note the EVTG is a function of both  $b$  and  $a$ , because the distribution  
865 of the next belief state is a function of  $b$  and  $a$ .) In some problems, it is possible to  
866 construct a heuristic EVTG based on domain knowledge. If the constructed EVTG

properly reflects tradeoffs in the selection of alternative actions, then the ranking of these actions via their  $Q$ -values will result in the desired “lookahead.”

For example, consider the motivating example of tracking multiple targets with a single sensor. Suppose we can only measure the location of one target per decision epoch. The problem then is to decide which location to measure and the objective function is the aggregate (multi-target) tracking error. The terrain over which the targets are moving is such that the measurement errors are highly location dependent, for example because of the presence of topological features which cause some areas to be invisible from a future sensor position. In this setting, it is intuitively clear that if we can predict sensor and target motion so that we expect a target is about to be obscured, then we should focus our measurements on that target immediately before the obscuration so that its track accuracy is improved and the overall tracking performance maximized in light of the impending obscuration.

The same reasoning applies in a variety of other situations, including those where targets are predicted to become unresolvable to the sensor (e.g., two targets that cross) or where the target and sensor motion is such that future measurements are predicted to be less reliable (e.g., a bearings-only sensor that is moving away from a target). In these situations, we advocate a heuristic method that replaces the EVTG by a function that captures the long-term benefit of an action in terms of an “opportunity cost” or “regret.” That is, we approximate the  $Q$ -value as

$$Q(b, a) \approx r(b, a) + wN(b, a)$$

where  $N(b, a)$  is an easily computed heuristic approximation of the long-term value, and  $w$  is a weighting term that allows us to trade the influence of the immediate value and the long-term value. As a concrete example of a useful heuristic, we have used the “gain in information for waiting” as a choice of  $N(b, a)$  (Kreucher et al. 2004). Specifically, let  $\bar{g}_a^k$  denote the expected value of the Rényi divergence between the belief state at time  $k$  and the updated belief state at time  $k$  after taking action  $a$ , as defined in Section 3.8 (i.e., the *myopic information gain*). Note that this myopic information gain is a random variable whose distribution depends on  $a$ , as explained in Section 3.8. Let  $p_a^k(\cdot)$  denote the distribution of this random variable. Then a useful approximation of the long-term value of taking action  $a$  is the gain (loss) in information received by waiting until a future time step to take the action,

$$N(b, a) \approx \sum_{m=1}^M \gamma^m \operatorname{sgn}(\bar{g}_a^k - \bar{g}_a^{k+m}) D_\alpha(p_a^k(\cdot) || p_a^{k+m}(\cdot))$$

where  $M$  is the number of time steps in the future that are considered.

Each term in the summand of  $N(b, a)$  has two components. First,  $\operatorname{sgn}(\bar{g}_a^k - \bar{g}_a^{k+m})$  signifies if the expected reward for taking action  $a$  in the future is more or less than the present. A negative value implies that the future is better and that the action ought to be discouraged at present. A positive value implies that the future is worse and that the action ought to be encouraged at present. This may happen, for example, when the visibility of a given target is getting worse with time. The second term,  $D_\alpha(p_a^k(\cdot) || p_a^{k+m}(\cdot))$ , reflects the magnitude of the change in reward using the divergence between the density on myopic rewards at the current time step and at a future time step. A small number implies the present and future rewards are very similar, and therefore the nonmyopic term should have little impact on the decision making.

910 Therefore,  $N(b, a)$  is positive if an action is less favorable in the future (e.g.,  
911 the target is about to become obscured). This encourages taking actions that are  
912 beneficial in the long term, and not just taking actions based on their immediate  
913 reward. Likewise, the term is negative if the action is more favorable in the future  
914 (e.g., the target is about to emerge from an obscuration). This discourages taking  
915 actions now that will have more value in the future.

## 916 6.5 Parametric approximation

917 In situations where a heuristic  $Q$ -function is difficult to construct, we may consider  
918 methods where the  $Q$ -function is approximated by a parametric function (by this  
919 we mean that we have a function approximator parameterized by one or more  
920 parameters). Let us denote this approximation by  $\tilde{Q}(b, \theta)$ , where  $\theta$  is a parameter  
921 (to be tuned appropriately). For this approach to be useful, the computation of  
922  $\tilde{Q}(b, \theta)$  has to be relatively simple, given  $b$  and  $\theta$ . Typically, we seek approximations  
923 for which it is easy to set the value of the parameter  $\theta$  appropriately, given some  
924 information of how the  $Q$ -values "should" behave (e.g., from expert knowledge,  
925 empirical results, simulation, or on-line observation). This adjustment or tuning of  
926 the parameter  $\theta$  is called *training*. In contrast to on-line approximation methods  
927 discussed in this section, the training process in parametric approximation is often  
928 done off-line.

929 As in the heuristic approximation approach, the approximation of the  $Q$ -function  
930 by the parametric function approximator is usually accomplished by approximating  
931 the EVTG, or even directly approximating the objective function  $V^*$ .<sup>2</sup> In the usual  
932 parametric approximation approach, the belief state  $b$  is first mapped to a set of  
933 *features*. The features are then passed through a parametric function to approximate  
934  $V^*(b)$ . For example, in the problem of tracking multiple targets with a single sensor,  
935 we may extract from the belief state some information on the location of each target  
936 relative to the sensor, taking into account the topology. These constitute features.  
937 For each target, we then assign a numerical value to these features, reflecting the  
938 measurement accuracy. Finally, we take a linear combination of these numerical  
939 values, where the coefficients of this linear combination serve the role of the  
940 parameters to be tuned.

941 The parametric approximation method has some advantages over methods based  
942 only on heuristic construction. First, the training process usually involves numerical  
943 optimization algorithms, and thus well-established methodology can be brought to  
944 bear on the problem. Second, even if we lack immediate expert knowledge on our  
945 problem, we may be able to experiment with the system (e.g., by using a simulation  
946 model). Such empirical output is useful for training the function approximator.  
947 Common training methods found in the literature go by the names of reinforcement  
948 learning,  $Q$ -learning, neurodynamic programming, and approximate dynamic pro-  
949 gramming. We have more to say about reinforcement learning in the next section.

<sup>2</sup>In fact, given a POMDP, the  $Q$ -value can be viewed as the objective function value for a related problem; see Bertsekas and Tsitsiklis (1996).

The parametric approximation approach may be viewed as a systematic method 950  
to implement the heuristic approach. But note that even in the parametric approach, 951  
some heuristics are still needed in the choice of features and in the form of the 952  
function approximator. For further reading, see Bertsekas and Tsitsiklis (1996). 953

## 6.6 Reinforcement learning 954

A popular method for approximating the  $Q$ -function based on the parametric 955  
approximation approach is *reinforcement learning* or  *$Q$ -learning* (Watkins 1989). 956  
Recall that the  $Q$ -function satisfies the equation 957

$$Q(b, a) = r(b, a) + E \left[ \max_{\alpha} Q(b', \alpha) \mid b, a \right]. \quad (2)$$

In  $Q$ -learning, the  $Q$ -function is estimated from multiple trajectories of the process. 958  
Assuming as usual that the number of states and actions are finite, we can represent 959  
 $Q(b, a)$  as a lookup table. In this case, given an arbitrary initial value of  $Q(b, a)$ , 960  
the one-step  $Q$ -learning algorithm (Sutton and Barto 1998) is given by the repeated 961  
application of the update equation: 962

$$Q(b, a) \leftarrow (1 - \beta) Q(b, a) + \beta \left( r(b, a) + \max_{\alpha} Q(b', \alpha) \right), \quad (3)$$

where  $\beta$  is a parameter in  $(0, 1)$  representing a “learning rate,” and each of the 4- 963  
tuples  $\{b, a, b', r\}$  are examples of states, actions, next states, and rewards incurred 964  
during the training phase. With enough examples of belief states and actions, the 965  
 $Q$ -function can be “learned” via simulation or on-line. 966

Unfortunately, in most realistic problems (the problems considered in this paper 967  
included) it is infeasible to represent the  $Q$ -function as a lookup table. This is 968  
either due to the large number of possible belief states (our case), actions, or both. 969  
Therefore, as pointed out in the last section, function approximation is required. A 970  
standard and simplest class of  $Q$ -function approximators are linear combinations of 971  
basis functions (also called features): 972

$$Q(b, a) = \theta(a)^T \phi(b), \quad (4)$$

where  $\phi(b)$  is a feature vector (often constructed by a domain expert) associated 973  
with state  $b$  and the coefficients of  $\theta(a)$  are to be estimated, i.e., the training data 974  
is used to learn the best approximation to  $Q(b, a)$  among all linear combinations of 975  
the features. Gradient descent is used with the training data to update the estimate 976  
of  $\theta(a)$ : 977

$$\begin{aligned} \theta(a) &\leftarrow \theta(a) + \beta \left( r(b, a) + \max_{a'} Q(b', a') - Q(b, a) \right) \nabla_{\theta} Q(b, a) \\ &= \theta(a) + \beta \left( r(b, a) + \max_{a'} \theta(a')^T \phi(b') - \theta(a)^T \phi(b) \right) \phi(b). \end{aligned}$$

Note that we have taken advantage of the fact that for the case of a linear function 978  
approximator, the gradient is given by  $\nabla Q(b, a) = \phi(b)$ . Hence, at every iteration, 979

980  $\theta(a)$  is updated in the direction that minimizes the empirical error in (2). When  
 981 a lookup table is used in (4), this algorithm reduces to (3). Once the learning  
 982 of the vector  $\theta(a)$  is completed, optimal actions can be computed according to  
 983  $\arg \max_a \theta(a)^\top \phi(b)$ . Determining the learning rate ( $\beta$ ) and the number of training  
 984 episodes required is a matter of active research.

985 Selecting a set of features that simultaneously provide both an adequate descrip-  
 986 tion of the belief state and a parsimonious representation of the state space requires  
 987 domain knowledge. For the illustrative example that we use in this paper (see  
 988 Section 3.8), the feature vector  $\phi(b)$  should completely characterize the surveillance  
 989 region and capture its nonstationary nature. For consistency in comparison to other  
 990 approaches, we appeal to features that are based on information theory, although  
 991 this is simply one possible design choice. In particular, we use the expected myopic  
 992 information gain at the current time step and the expected myopic information  
 993 gain at the next time step as features which characterize the state. Specifically, let  
 994  $r(b, a) = E[\Delta_\alpha(b, a, y)|b, a]$  be defined as in Section 3.8. Next, define  $b'$  to be the  
 995 belief state at the hypothetical “next” time step starting at the current belief state  $b$ ,  
 996 computed using the second of the two-step update procedure in Section 3.2. In other  
 997 words,  $b'$  is what results in the next step if only a state transition takes place, without  
 998 an update based on incorporating a measurement. Then, the feature vector is

$$\phi(b) = [r(b, 1), \dots, r(b, C), r(b', 1), \dots, r(b', C)]$$

999 where  $C$  is the number of cells (and also the number of actions). In the situation  
 1000 of time-varying visibility, these features capture the immediate value of various  
 1001 actions and allow the system to learn the long-term value by looking at the change in  
 1002 immediate value of the actions over time. In a more general version of this problem,  
 1003 actions might include more than just which cell to measure—for example, actions  
 1004 might also involve which waveform to transmit. In these more general cases, the  
 1005 feature vector will have more components to account for the larger set of possible  
 1006 actions.

## 1007 6.7 Action-sequence approximations

1008 Let us write the value function (optimal objective function value as a function of  
 1009 belief state) as

$$\begin{aligned} V^*(b) &= \max_{\pi} E \left[ \sum_{k=0}^{H-1} r(b_k, \pi(b_k)) \middle| b, \pi(b) \right] \\ &= E \left[ \max_{a_0, \dots, a_{H-1}: a_k = \pi(b_k)} \sum_{k=0}^{H-1} r(b_k, a_k) \middle| b \right], \end{aligned} \quad (5)$$

1010 where the notation  $\max_{a_0, \dots, a_{H-1}: a_k = \pi(b_k)}$  means maximization subject to the constraint  
 1011 that each action  $a_k$  is a (fixed) function of the belief state  $b_k$ . If we relax this constraint  
 1012 on the actions and allow them to be arbitrary random variables, then we have an  
 1013 upper bound on the value function:

$$\hat{V}_{HO}(b) = E \left[ \max_{a_0, \dots, a_{H-1}} \sum_{k=0}^{H-1} r(b_k, a_k) \middle| b \right].$$

In some applications, this upper bound provides a suitable approximation to the value function. The advantage of this method is that in certain situations the computation of the “max” above involves solving a relatively easy optimization problem. This method is called *hindsight optimization* (Chong et al. 2000; Wu et al. 2002).

One implementation involves averaging over many Monte Carlo simulation runs to compute the expectation above. In this case, the “max” is computed for each simulation run by first generating all the random numbers for that run, and then applying a static optimization algorithm to compute optimal actions  $a_0, \dots, a_{H-1}$ . It is easy now to see why we call the method “hindsight” optimization: the optimization of the action sequence is done after knowing all uncertainties over time, as if making decisions in hindsight.

As an alternative to relaxing the constraint in (5) (that each action  $a_k$  is a fixed function of the belief state  $b_k$ ), suppose we further *restrict* each action to be simply fixed (not random). This restriction gives rise to a lower bound on the value function:

$$\hat{V}_{FO}(b) = \max_{a_0, \dots, a_{H-1}} E[r(b_0, a_0) + \dots + r(b_{H-1}, a_{H-1}) | b, a_0, \dots, a_{H-1}].$$

To use analogous terminology to “hindsight optimization,” we call this method *foresight optimization*—we make decisions before seeing what actually happens, based on our expectation of what will happen. The method is also called *open loop feedback control* (Bertsekas 2007). For a tracking application of this, see Chhetri et al. (2004).

We should also point out some alternatives to the simple hindsight or foresight approaches above. In Yu and Bertsekas (2004), more sophisticated bounds are described that do not involve simulation, but instead rely on convexity. The method in Miller et al. (2009) also does not involve simulation, but approximates the future belief-state evolution using a single sample path.

## 6.8 Rollout

In this section, we describe the method of *policy rollout* (or simply *rollout*) (Bertsekas and Castanon 1999). The basic idea is simple. First let  $V^\pi(b_0)$  be the objective function value corresponding to policy  $\pi$ . Recall that  $V^* = \max_\pi V^\pi$ . In the method of rollout, we assume that we have a candidate policy  $\pi_{\text{base}}$  (called the *base policy*), and we simply replace  $V^*$  in (1) by  $V^{\pi_{\text{base}}}$ . In other words, we use the following approximation to the  $Q$ -value:

$$Q^{\pi_{\text{base}}}(b, a) = r(b, a) + E[V^{\pi_{\text{base}}}(b') | b, a].$$

We can think of  $V^{\pi_{\text{base}}}$  as the performance of applying  $\pi_{\text{base}}$  in our system. In many situations of interest,  $V^{\pi_{\text{base}}}$  is relatively easy to compute, either analytically, numerically, or via Monte Carlo simulation.

It turns out that the policy  $\pi$  defined by

$$\pi(b) = \arg \max_a Q^{\pi_{\text{base}}}(b, a) \quad (6)$$

1050 is at least as good as  $\pi_{\text{base}}$  (in terms of the objective function); in other words,  
1051 this step of using one policy to define another policy has the property of *policy*  
1052 *improvement*. This result is the basis for a method known as *policy iteration*, where  
1053 we iteratively apply the above policy-improvement step to generate a sequence  
1054 of policies converging to the optimal policy. However, policy iteration is difficult  
1055 to apply in problems with large belief-state spaces, because the approach entails  
1056 explicitly representing a policy and iterating on it (remember that a policy is a  
1057 mapping with the belief-state space  $\mathcal{B}$  as its domain).

1058 In the method of policy rollout, we do not explicitly construct the policy  $\pi$  in (6).  
1059 Instead, at each time step, we use (6) to compute the output of the policy at the  
1060 current belief-state. For example, the term  $E[V^{\pi_{\text{base}}}(b')|b, a]$  can be computed using  
1061 Monte Carlo sampling. To see how this is done, observe that  $V^{\pi_{\text{base}}}(b')$  is simply the  
1062 mean cumulative reward of applying policy  $\pi_{\text{base}}$ , a quantity that can be obtained  
1063 by Monte Carlo simulation. The term  $E[V^{\pi_{\text{base}}}(b')|b, a]$  is the mean with respect to  
1064 the random next belief-state  $b'$  (with distribution that depends on  $b$  and  $a$ ), again  
1065 obtainable via Monte Carlo simulation. We provide more details in Section 6.10. In  
1066 our subsequent discussion of rollout, we will focus on its implementation using Monte  
1067 Carlo simulation. For an application of the rollout method to sensor scheduling for  
1068 target tracking, see He and Chong (2004, 2006), Krakow et al. (2006), Li et al. (2006,  
1069 2007).

## 1070 6.9 Parallel rollout

1071 An immediate extension to the method of rollout is to use multiple base policies. So  
1072 suppose that  $\Pi_B = \{\pi^1, \dots, \pi^n\}$  is a set of base policies. Then replace  $V^*$  in (1) by

$$\hat{V}(b) = \max_{\pi \in \Pi_B} V^\pi(b).$$

1073 We call this method *parallel rollout* (Chang et al. 2004). Notice that the larger the set  
1074  $\Pi_B$ , the tighter  $\hat{V}(b)$  becomes as a bound on  $V^*(b)$ . Of course, if  $\Pi_B$  contains the  
1075 optimal policy, then  $\hat{V} = V^*$ . It follows from our discussion of rollout that the policy  
1076 improvement property also holds here. As with the rollout method, parallel rollout  
1077 can be implemented using Monte Carlo sampling.

## 1078 6.10 Control architecture in the Monte Carlo case

1079 The method of rollout provides a convenient turnkey (systematic) procedure for  
1080 Monte-Carlo-based decision making and control. Here, we specialize the general  
1081 control architecture of Section 5 to the use of particle filtering for belief-state  
1082 updating and a Monte Carlo method for  $Q$ -value approximation (e.g., rollout). We  
1083 note that there is increasing interest in Monte Carlo methods for solving Markov  
1084 decision processes (Thrun et al. 2005; Chang et al. 2007). Particle filtering, which  
1085 is a Monte Carlo sampling method for updating posterior distributions, dovetails  
1086 naturally with Monte Carlo methods for  $Q$ -value approximation. An advantage  
1087 of the Monte Carlo approach is that it does not rely on analytical tractability—it  
1088 is straightforward in this approach to incorporate sophisticated models for sensor  
1089 characteristics and target dynamics.

**Fig. 7** Basic control architecture with particle filtering

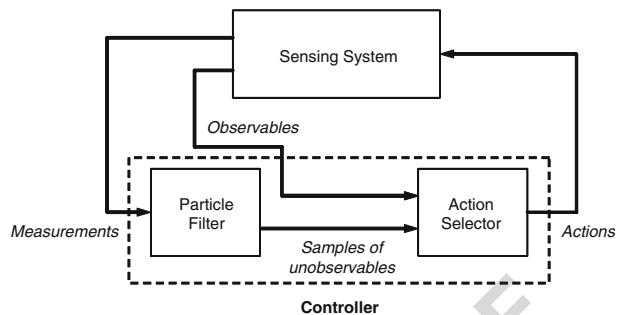
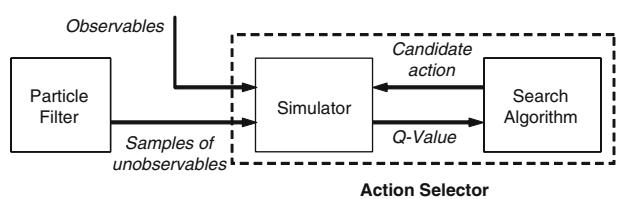


Figure 7 shows the control architecture specialized to the Monte Carlo setting. In contrast to Fig. 5, a particle filter plays the role of the measurement filter, and its output consists of samples of the unobservables. Figure 8 shows the action selector in this setting. Contrasting this with Fig. 6, we see that a Monte Carlo simulator plays the role of the  $Q$ -value approximator (e.g., via rollout). Search algorithms that are suitable here include the method of Shi and Chen (2000), which is designed for such problems, dovetails well with a simulation-based approach, and accommodates heuristics to guide the search within a rigorous framework.

As a specific example, consider applying the method of rollout. In this case, the evaluation of the  $Q$ -value for any given candidate action relies on a simulation model of the sensing system with some base policy. This simulation model is a “dynamic” model in that it evaluates the behavior of the sensing system over some horizon of time (specified beforehand). The simulator requires as inputs the current observables and samples of unobservables from the particle filter (to specify initial conditions) and a candidate action. The output of the simulator is a  $Q$ -value corresponding to the current measurements and observables, for the given candidate action. The output of the simulator represents the mean performance of applying the base policy, depending on the nature of the objective function. For example, the performance measure of the system may be the negative mean of the sum of the cumulative tracking error and the sensor usage cost over a horizon of  $H$  time steps, given the current system state and candidate action.

To elaborate on exactly how the  $Q$ -value approximation using rollout is implemented, suppose we are given the current observables and a set of samples of the unobservables (from the particle filter). The current observables together with a single sample of unobservables represent a candidate current underlying state of the sensing system. Starting from this candidate current state, we simulate the application of the given candidate action (which then leads to a random next state), followed by application of the base policy for the remainder of the time horizon—during this time

**Fig. 8** Components of the action selector



1118 horizon, the system state evolves according to the dynamics of the sensing system as  
1119 encoded within the simulation model. For this single simulation run, we compute  
1120 the “action utility” of the system (e.g., the negative of the sum of the cumulative  
1121 tracking error and sensor usage cost over that simulation run). We do this for each  
1122 sample of the unobservables, and then average over the performance values from  
1123 these multiple simulation runs. This average is what we output as the  $Q$ -value.

1124 The samples of the unobservables from the particle filter that are fed to the  
1125 simulator (as candidate initial conditions for unobservables) may include all the  
1126 particles in the particle filter (so that there is one simulation run per particle), or  
1127 may constitute only a subset of the particles. In principle, we may even run multiple  
1128 simulation runs per particle.

1129 The above Monte Carlo method for approximating POMDP solutions has some  
1130 beneficial features. First, it is flexible in that a variety of adaptive sensing scenarios  
1131 can be tackled using the same framework. This is important because of the wide  
1132 variety of sensors encountered in practice. Second, the method does not require  
1133 analytical tractability; in principle, it is sufficient to simulate a system component,  
1134 whether or not its characteristics are amenable to analysis. Third, the framework  
1135 is modular in the sense that models of individual system components (e.g., sensor  
1136 types, target motion) may be treated as “plug-in” modules. Fourth, the approach  
1137 integrates naturally with existing simulators (e.g., Umbra (Gottlieb and Harrigan  
1138 2001)). Finally, the approach is inherently nonmyopic, allowing the tradeoff of short-  
1139 term gains for long-term rewards.

#### 1140 6.11 Belief-state simplification

1141 If we apply the method of rollout to a POMDP, we need a base policy that maps  
1142 belief states to actions. Moreover, we need to simulate the performance of this  
1143 policy—in particular, we have to sample future belief states as the system evolves  
1144 in response to actions resulting from this policy. Because belief states are probability  
1145 distributions, keeping track of them in a simulation is burdensome.

1146 A variety of methods are available to approximate the belief state. For example,  
1147 we could simulate a particle filter to approximate the evolution of the belief state  
1148 (as described previously), but even this may be unduly burdensome. As a further  
1149 simplification, we could use a Gaussian approximation and keep track only of the  
1150 mean and covariance of the belief state using a Kalman filter or any of its extensions,  
1151 including *extended Kalman filters* and *unscented Kalman filters* (Julier and Uhlmann  
1152 2004). Naturally, we would expect that the more accurate the approximation of the  
1153 belief state, the more burdensome the computation.

1154 An extreme special case of the above tradeoff is to use a Dirac delta distribution  
1155 for belief states in our simulation of the future. In other words, in our lookahead  
1156 simulation, we do away with keeping track of belief states altogether and instead  
1157 simulate only a *completely observable* version of the system. In this case, we need only  
1158 consider a base policy that maps underlying states to actions—we could simply apply  
1159 rollout to this policy, and not have to maintain any belief states in our simulation.  
1160 Call this method *completely observable (CO) rollout*. It turns out that in certain  
1161 applications, such as in sensor scheduling for target tracking, a CO-rollout base policy  
1162 is naturally available (see He and Chong 2004, 2006; Krakow et al. 2006; Li et al. 2006,  
1163 2007). Note that we will still need to keep track of (or estimate) the actual belief state

of the system, even if we use CO rollout. The benefit of CO rollout is that it allows 1164 us to avoid keeping track of (simulated) belief states in our *simulation* of the future 1165 evolution of the system. 1166

In designing lookahead methods with a simplified belief state, we must ensure the 1167 simplification does not hide the good or bad effects of actions. The resulting  $Q$ -value 1168 approximation must properly rank current actions. This requires a carefully designed 1169 simplification of the belief state together with a base policy that appropriately reflects 1170 the effects of taking specific current actions. 1171

For example, suppose that a particular current action results in poor future 1172 rewards because it leads to belief states with large variances. Then, if we use the 1173 method of CO rollout, we have to be careful to ensure that this detrimental effect of 1174 the particular current action be reflected as a cost in the lookahead. (Otherwise, the 1175 effect would not be accounted for properly, because in CO rollout we do not keep 1176 track of belief states in our simulation of the future effect of current actions.) 1177

Another caveat in the use of simplified belief states in our lookahead is that the 1178 resulting rewards in the lookahead may also be affected (and this may have to be 1179 taken into account). For example, consider again the problem of sensor scheduling 1180 for target tracking, where the per-step reward is the negative mean of the sum of 1181 the tracking error and the sensor usage cost. Suppose that we use a particle filter 1182 for tracking (i.e., for keeping track of the actual belief state). However, for our 1183 lookahead, we use a Kalman filter to keep track of future belief states in our rollout 1184 simulation. In general, the tracking error associated with the Kalman filter is different 1185 from that of the particle filter. Therefore, when summed with the sensor usage cost, 1186 the relative contribution of the tracking error to the overall reward will be different 1187 for the Kalman filter compared to the particle filter. To account for this, we will need 1188 to scale the tracking error (or sensor usage cost) in our simulation so that the effect of 1189 current actions are properly reflected in the  $Q$ -value approximations from the rollout 1190 with the simplified belief state calculation. 1191

## 6.12 Reward surrogation

1192

In applying a POMDP approximation method, it is often useful to substitute the 1193 reward function for an alternative (*surrogate*), for a number of reasons. First, we 1194 may have a surrogate reward that is much simpler (or more reliable) to calculate 1195 than the actual reward (e.g., the method of reduction to classification (Blatt and 1196 Hero 2006a, b)). Second, it may be desirable to have a single surrogate reward for 1197 a range of different actual rewards. For example, Kreucher et al. (2005b), Hero 1198 et al. (2008) shows that average Rényi information gain can be interpreted as a near 1199 universal proxy for any bounded performance metric. Third, reward surrogation may 1200 be necessitated by the use of a belief-state simplification technique. For example, if 1201 we use a Kalman filter to update the mean and covariance of the belief state, then 1202 the reward can only be calculated using these entities. 1203

The use of a surrogate reward can lead to many benefits. But some care must 1204 be taken in the design of a suitable surrogate reward. Most important is that the 1205 surrogate reward be sufficiently reflective of the true reward that the ranking of 1206 actions with respect to the approximate  $Q$ -values be preserved. A superficially 1207 benign substitution may in fact have unanticipated but significant impact on the 1208 ranking of actions. For example, recall the example raised in the previous section on 1209

1210 belief-state simplification, where we substitute the tracking error of a particle filter  
 1211 for the tracking error of a Kalman filter. Superficially, this substitute appears to be  
 1212 hardly a "surrogate" at all. However, as pointed out before, the tracking error of the  
 1213 Kalman filter may be significantly different in magnitude from that of a particle filter.

## 1214 7 Illustration: spatially adaptive airborne sensing

1215 In this section, we illustrate the performance of several of the strategies discussed  
 1216 in this paper on a common model problem. The model problem has been chosen  
 1217 to have the characteristics of the motivating example given earlier, while remaining  
 1218 simple enough so that the workings of each method are transparent.

1219 In the model problem, there are two targets, each of which is described by a  
 1220 one-dimensional position (see Fig. 9). The state is therefore a 2-dimensional real  
 1221 number describing the target locations plus the sensor position, as described in  
 1222 Section 3.8. Targets move according to a pure diffusion model (given explicitly  
 1223 in Section 3.8 as  $T_{\text{single target}}(y|x)$ ), and the belief state is propagated using this  
 1224 model. Computationally, the belief state is estimated by a multi-target particle filter,  
 1225 according to the algorithm given in Kreucher et al. (2005c).

1226 The sensor may measure any one of 16 cells, which span the possible target  
 1227 locations (again, see Fig. 9). The sensor is capable of making three (not necessarily  
 1228 distinct) measurements per time step, receiving binary returns independent from  
 1229 dwell to dwell. The three measurements are fused sequentially: after each measure-  
 1230 ment, we update the belief state by incorporating the measurement using Bayes' rule,  
 1231 as discussed in Section 3.2. In occupied cells, a detection is received with probability  
 1232  $P_d = 0.9$ . In cells that are unoccupied a detection is received with probability  $P_f$  (set  
 1233 here at 0.01). This sensor model is given explicitly in Section 3.8 by  $P_{\text{obs}}(z|x, a)$ .

1234 At the onset, positions of the targets are known only probabilistically. The belief  
 1235 state for the first target is uniform across sensor cells  $\{2, \dots, 6\}$  and for the second  
 1236 target is uniform across sensor cells  $\{11, \dots, 15\}$ . The particle filter used to estimate  
 1237 the belief state is initialized with this uncertainty.

1238 Visibility of the cells changes with time as in the motivating example of Section 3.8.  
 1239 At time 1, all cells are visible. At times 2, 3, and 4, cells  $\{11, \dots, 15\}$  become obscured.  
 1240 At time 5, all cells are visible again. This time varying visibility map is known to  
 1241 the sensor management algorithm and should be exploited to best choose sensing  
 1242 actions.

	0-1	1-2	2-3	3-4	4-5	5-6	6-7	7-8	8-9	9-10	10-11	11-12	12-13	13-14	14-15	15-16
Cell 1	X															
Time 1	X															
Time 2																
Time 3																
Time 4																
Time 5																

**Fig. 9** The model problem. At the onset, the belief state for target 1 is uniformly distributed across cells  $\{2, \dots, 6\}$  and the belief state for target 2 is uniformly distributed across cells  $\{11, \dots, 15\}$ . At time 1 all cells are visible. At times 2, 3, and 4, cells  $\{11, \dots, 15\}$  are obscured. This is a simple case where a target is initially visible, becomes obscured, and then reemerges

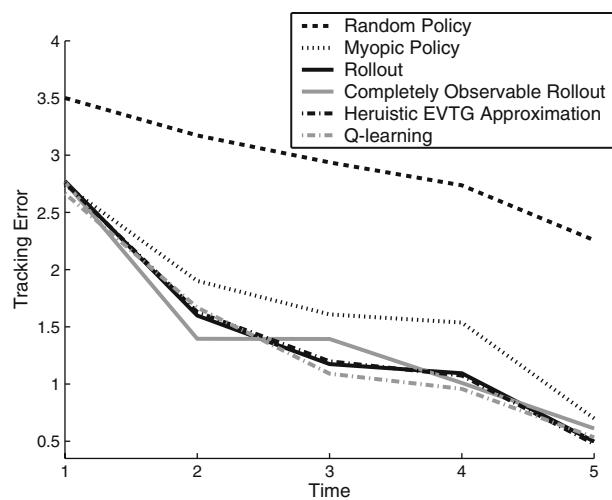
Sensor management decisions are made by using the belief state to predict which actions are most valuable. In the following paragraphs, we contrast the decisions made by a number of different strategies that have been described earlier.

At time 1 a myopic strategy, using no information about the future visibility, will choose to measure cells uniformly from the set  $\{2, \dots, 6\} \cup \{11, \dots, 15\}$  as they all have the same expected immediate reward. As a result, target 1 and target 2 will on the average be given equal attention. A nonmyopic strategy, on the other hand, will choose to measure cells from  $\{11, \dots, 15\}$  as they are soon to become obscured. That is, the policy of looking for target 2 at time 1 followed by looking for target 1 is best.

Figure 10 shows the performance of several of the on-line strategies discussed in this paper on this common model problem. The performance of each scheduling strategy is measured in terms of the mean squared tracking error at each time step. The curves represent averages over 10,000 realizations of the model problem. Each realization has randomly chosen initial positions of the targets and measurements corrupted by random mistakes as discussed above. The five policies are as follows.

- A **random** policy that simply chooses one of the 16 cells randomly for interrogation. This policy provides a worst-case performance and will bound the performance of the other policies.
- A **myopic** policy that takes the action expected to maximize immediate reward. Here the surrogate reward is myopic information gain as defined in Section 6.4, measured in terms of the expected Rényi divergence with  $\alpha = 0.5$  (see Kreucher et al. 2005b). So the value of an action is estimated by the amount of information it gains. The myopic policy is sub-optimal because it does not consider the long term ramifications of its choices. In particular, at time 1 the myopic strategy has no preference as to which target to measure because both are unobsured and have uncertain position. Therefore, half of the time, target 1 is measured, resulting in an opportunity cost because target 2 is about to disappear.
- The **reinforcement learning** approach described in Section 6.6. The  $Q$ -function was learned using a linear function approximator, as described in detail in Section 6.6, by running a large number ( $10^5$ ) of sample vignettes. Each sample

**Fig. 10** The performance of the five policies discussed above. Performance is measured in terms of mean squared tracking error at each time step, averaged over a  $10^4$  Monte Carlo trials



vignette proceeds as follows. An action is taken randomly. The resulting immediate gain (as measured by the expected information gain) is recorded and the resulting next-state computed. This next-state is used to predict the long-term gain using the currently available  $Q$ -function. The  $Q$ -function is then refined given this information (in practice this is done in blocks of many vignettes, but the principle is the same). Training the  $Q$ -function is a very time consuming process. In this case, for each of the  $10^5$  sample vignettes, the problem was simulated from beginning to end, and the state and reward variables were saved along the way. It is also unclear as to how the performance of the trained  $Q$ -function will change if the problem is perturbed. However, with these caveats in mind, once the  $Q$ -function has been learned, decision making is very quick and the resulting policy in this case is very good.

- The **heuristic EVTG approximation** described in Section 6.4 favors actions expected to be more valuable now than in the future. In particular, actions corresponding to measuring target 2 have additional value because target 2 is predicted to be obscured in the future. This makes the ranking of actions that measure target 2 higher than those that measure target 1. Therefore, this policy (like the other nonmyopic approximations described here) outperforms the myopic policy. The computational burden is on the order of  $H$  times the myopic policy, where  $H$  is the horizon length.
- The **rollout** policy described in Section 6.8. The base policy used here is to take each of the three measurements sequentially at the location where the target is expected to be, which is a function of the belief state that is current to the particular measurement. This expectation is computed using the predicted future belief state, which requires the belief state to be propagated in time. This is done using a particle filter. We again use information gain as the surrogate reward to approximate  $Q$ -values. The computational burden of this method is on the order of  $NH$  times that of the myopic policy, where  $H$  is the horizon length and  $N$  is the number of Monte Carlo trials used in the approximation (here  $H = 5$  and  $N = 25$ ).
- The **completely observable rollout** policy described in Section 6.11. As in the rollout policy above, the base policy here is to take measurements sequentially at locations where the target is expected to be, but enforces the criterion that the sensor should alternate looking at the two targets. This slight modification is necessary due to the delta-function representation of future belief states. Since the completely observable policy does not predict the posterior into the future, it is significantly faster than standard rollout (an order of magnitude faster in these simulations). However, it requires a different surrogate reward (one that does not require the posterior like the information gain surrogate metric). Here we have chosen as a surrogate reward to count the number of detections received, discounting multiple detections of the same target.

Our main intent here is simply to convey that, from Fig. 10, the nonmyopic policies perform similarly, and are better than the myopic and random policies, though at the cost of additional computational burden. The nonmyopic techniques perform similarly since they ultimately choose similar policies. Each one prioritizes measuring the target that is about to disappear over the target that is in the clear. On the other hand, the myopic policy is “losing” the target more often, resulting in higher mean error as there are more catastrophic events.

**8 Illustration: multi-mode adaptive airborne sensing**

1321

In this section, we turn our attention to adaptive sensing with a waveform-agile sensor. In particular, we investigate how the availability of multiple waveform choices effects the myopic/nonmyopic trade. The model problem considered here again focuses on detection and tracking in a visibility impaired environment. The target dynamics, belief-state update, and observation law are identical to that described in the first simulation. However, in this section we look at a sensor that is agile over waveform as well as pointing direction (i.e., can choose both where to interrogate as well as what waveform to use). Furthermore, the different waveforms are subject to different (time-varying) visibility maps. Simulations show that the addition of waveform agility (and corresponding visibility differences) changes the picture. In this section, we restrict our attention to the EVTG heuristic for approximate nonmyopic planning. Earlier simulations have shown that in model problems of this type, the various approaches presented here perform similarly.

1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334**8.1 A study with a single waveform**

1335

We first present a baseline result comparing random, myopic, and heuristic EVTG (HECTG) approximation based performance in the (modified) model problem. The model problem again covers a surveillance area broken into 16 regions with a target that is to be detected and tracked. The single target moves according to a purely diffusive model, and the belief state is propagated using this model. However, in this simulation the model problem is modified in that there is only one sensor allocation per time step and the detection characteristics are severely degraded. The region is occluded by a time-varying visibility map that obscures certain sub-regions at each time step, degrading sensor effectiveness in those regions at that time step. The visibility map is known exactly *a priori* and can be used both to predict which portions of the region are useless to interrogate at the present time (because of current occlusion) and to predict which regions will be occluded in the future. The sensor management choice in the case of a single waveform is to select the pointing direction (one of the 16 sub-regions) to interrogate. If a target is present and the sub-region is not occluded, the sensor reports a detection with  $p_d = 0.5$ . If the target is not present or the sub-region is occluded the sensor reports a detection with  $p_f = .01$ .

1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351

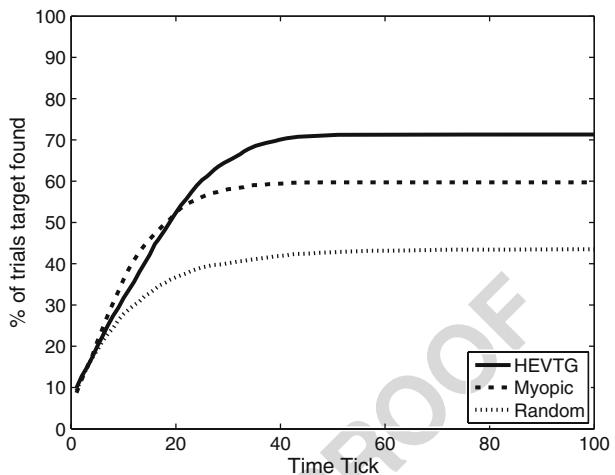
Both the myopic and nonmyopic information based methods discount the value of looking at occluded sub-regions. Prediction of myopic information gain uses visibility maps to determine that interrogating an occluded cell provides no information because the outcome is certain (it follows the false alarm distribution). However, the nonmyopic strategy goes further: It uses future visibility maps to predict which sub-regions will be occluded in the future and gives higher priority to their interrogation at present.

1352  
1353  
1354  
1355  
1356  
1357  
1358

The simulation results shown in Fig. 11 indicate that the HEVTG approximation to the nonmyopic scheduler provides substantial performance improvement with respect to a myopic policy in the single waveform model problem. The gain in performance for the policy that looks ahead is primarily ascribable to the following. It is important to promote interrogation of sub-regions that are about to become occluded over those that will remain visible. If a sub-region is not measured and then becomes occluded, the opportunity to determine target presence in that region

1359  
1360  
1361  
1362  
1363  
1364  
1365

**Fig. 11** Performance of the scheduling policies with a pointing-agile single waveform sensor



1366 is lost until the region becomes visible again. This opportunity cost is captured in  
 1367 the HEVTG approximation as it predicts which actions will have less value in the  
 1368 future and promotes them at the present. The myopic policy merely looks at the  
 1369 current situation and takes the action with maximal immediate gain. As a result of  
 1370 this greediness, it misses opportunities that have long term benefit. As a result of this  
 1371 greediness, the myopic policy may outperform the HEVTG in the short term but  
 1372 ultimately underperforms.

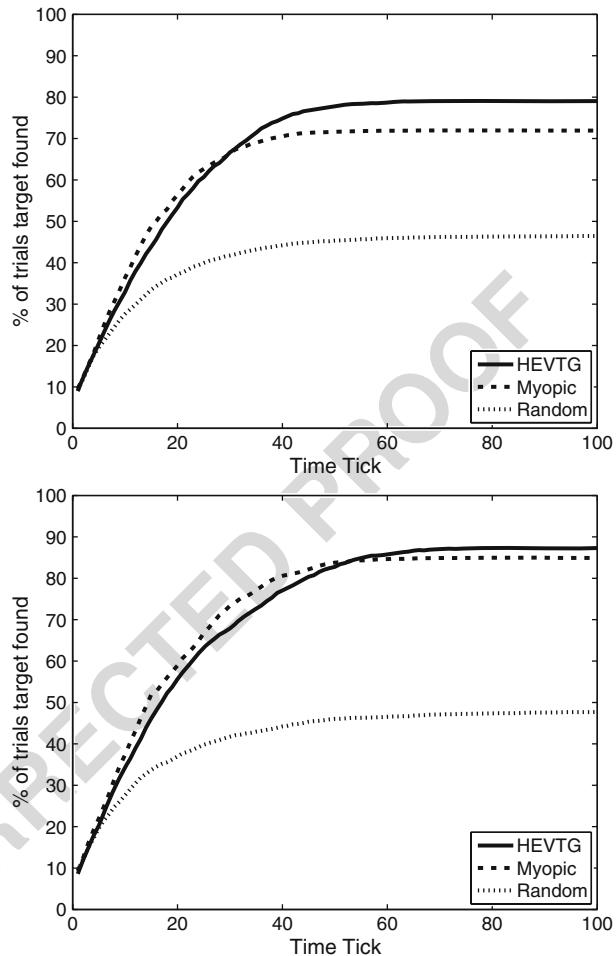
### 1373 8.2 A study with multiple independent waveforms

1374 This subsection explores the effect of multiple waveforms on the nonmyopic/myopic  
 1375 trade. We consider multiple *independent* waveforms, where independent means the  
 1376 time-varying visibility maps for the different waveforms are not coupled in any way.  
 1377 This assumption is relaxed in the following subsection.

1378 Each waveform has an associated time-varying visibility map drawn indepen-  
 1379 dently from the others. The sensor management problem is one of selecting both  
 1380 pointing direction and the waveform. All other simulation parameters are set identi-  
 1381 cally to the previous simulation (i.e., detection and false alarm probabilities, and  
 1382 target kinematics). Figure 12 shows performance curves for two and five independent  
 1383 waveforms. In comparison to the single waveform simulation, these simulations (a)  
 1384 have improved overall performance, and (b) have a narrowed gap in performance  
 1385 between nonmyopic and myopic schedulers.

1386 Figure 13 provides simulation results as the number of waveforms available is  
 1387 varied. These results indicate that as the number of independent waveforms available  
 1388 to the scheduler increase, the performance difference between a myopic policy and  
 1389 a nonmyopic policy narrows. This is largely due to the softened opportunity cost the  
 1390 myopic policy suffers. In the single waveform situation, if a region became occluded  
 1391 it could not be observed until the visibility for the single waveform changed. This puts  
 1392 a sharp penalty on a myopic policy. However, in the multiple independent waveform  
 1393 scenario, the penalty for myopic decision making is much less severe. In particular,  
 1394 if a region becomes occluded in waveform  $i$ , it is likely that some other waveform

**Fig. 12** Top: Performance of the strategies with a two-waveform sensor. Bottom: Performance curves with a five-waveform sensor



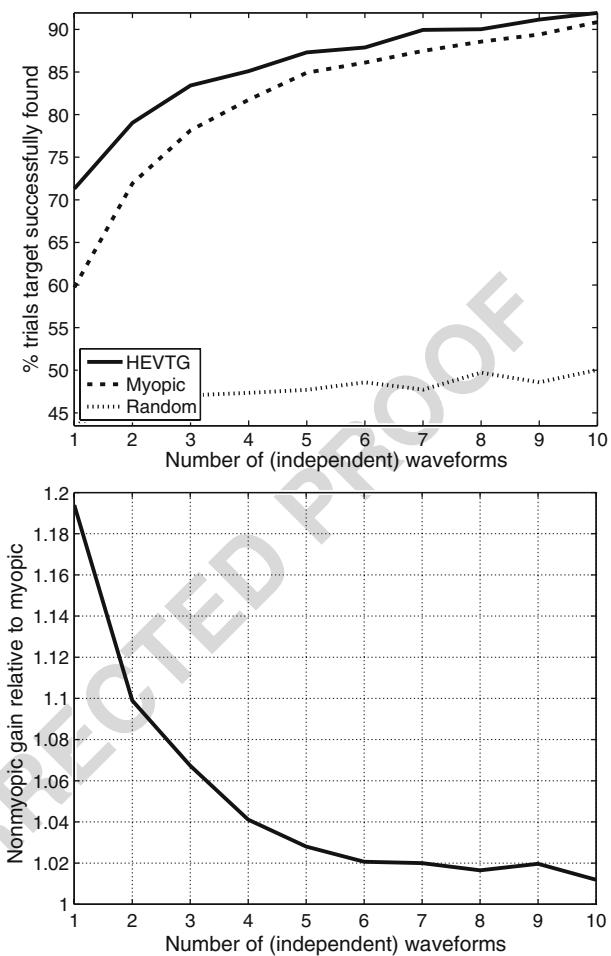
is still viable (i.e., the region is unoccluded to that waveform) and a myopic policy 1395 suffers little loss. As the number of independent waveforms available to the sensor 1396 increases, this effect is magnified until there is essentially no difference in the two 1397 policies. 1398

### 8.3 A study with multiple coupled waveforms

1399

A more realistic multiple waveform scenario is one in which the visibility occlusions between waveforms are highly coupled. Consider the case where a platform may 1400 choose between the following 5 waveforms (modalities) for interrogation of a region: 1401 electro-optical (EO), infra-red (IR), synthetic aperture radar (SAR), foliage pen- 1402 etrating radar (FOPEN), and moving target indication radar (MTI). In this situation, 1403 the visibility maps for the 5 waveforms are highly coupled through the environmental 1404 conditions (ECs) present in the region. For example, clouds effect the visibility of 1405

**Fig. 13** *Top:* The terminal performance of the scheduling algorithms versus number of waveforms. *Bottom:* The gain (performance improvement) of the nonmyopic policy with respect to the myopic policy



1407 both EO and IR. Similarly, tree cover effects the performance of all modes except  
 1408 FOPEN, and so on.

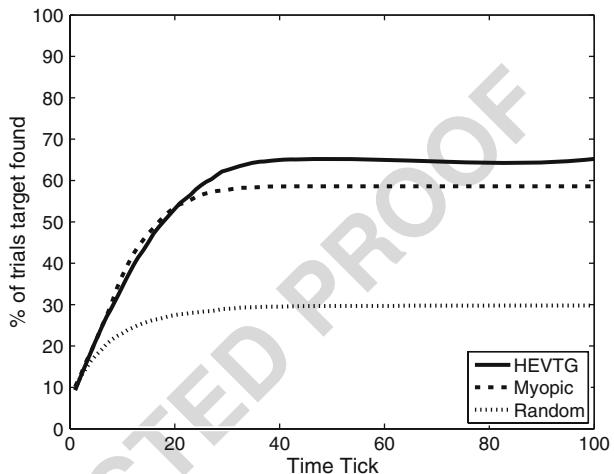
1409 Therefore, a more realistic study of multiple waveform performance is to model  
 1410 the time-varying nature of a collection of environmental conditions and generate the  
 1411 (now coupled) waveform visibility maps from the ECs. For this simulation study, we  
 Q3 1412 choose the nominal causation map shown in Fig. 14.

1413 The time-varying maps of each EC are chosen to resemble a passover, where for  
 1414 example the initial cloud map is chosen randomly and then it moves at a random ori-  
 1415 entation and random velocity through the region over the simulation time. The wave-  
 1416 form visibility maps are then formed by considering all obscuring ECs and choosing  
 1417 the maximum obscuration. This setup results in fewer than five independent wave-  
 1418 forms available to the sensor because the viability maps are coupled through the ECs.

1419 Figure 14 (right) shows a simulation result of the performance for a five waveform  
 1420 sensor. The simulation shows the gap between the myopic policy and the nonmyopic  
 1421 policy widens from where it was in the independent waveform simulation. In fact,

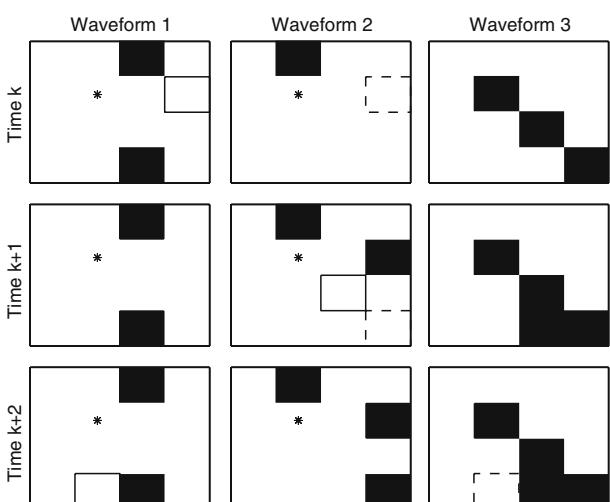
**Fig. 14** Top: EC Causation map. Bottom: Performance of the scheduling strategies with a pointing-agile five waveform sensor, where the visibility maps are coupled through the presence of environmental conditions

	Cloud	Rain	Wind	Fog	Foliage
EO	X	X	X	X	X
SAR		X			X
FOPEN		X	X		
IR	X	X		X	X
GMTI		X			X



in this scenario, the 5 dependent waveforms have performance characteristics that 1422  
are similar to 2 independent waveforms, as measured by the ratio of nonmyopic 1423  
scheduler performance to myopic scheduler performance. Figure 15 illustrates the 1424  
difference among the three policies being compared here, highlighting the “looka- 1425  
head” property of the nonmyopic scheme. 1426

**Fig. 15** Three time steps from a three waveform simulation. Obscured areas are shown with filled black squares and unobscured areas are white. The true target position is shown by an asterisk for reference. The decisions (waveform choice and pointing direction) are shown with solid-bordered squares (myopic policy) and dashed-bordered squares (nonmyopic policy). This illustrates “lookahead,” where regions that are about to be obscured are measured preferentially by the nonmyopic policy



1427 **9 Conclusions**

1428 This paper has presented methods for adaptive sensing based on approximations  
1429 for partially observable Markov decision processes, a special class of discrete event  
1430 system models. Though we have not specifically highlighted the event-driven nature  
1431 of these models, our framework is equally applicable to models that are more  
1432 appropriately viewed as event driven. The methods have been illustrated on the  
1433 problem of waveform-agile sensing, wherein it has been shown that intelligently  
1434 selecting waveforms based on past outcomes provides significant benefit over naive  
1435 methods. We have highlighted, via simulation, computationally approaches based  
1436 on rollout and a particular heuristic related to information gain. We have detailed  
1437 some of the design choices that go into finding appropriate approximations, including  
1438 choice of surrogate reward and belief-state representation.

1439 Throughout this paper we have taken special care to emphasize the limitations of  
1440 the methods. Broadly speaking, all tractable methods require domain knowledge in  
1441 the design process. Rollout methods require a base policy specially designed for the  
1442 problem at hand; relaxation methods require one to identify the proper constraint(s)  
1443 to remove; heuristic approximations require identification of appropriate value-to-  
1444 go approximations, and so on. That being said, when domain knowledge is available  
1445 it can often yield dramatic improvement in system performance over traditional  
1446 methods at a fixed computational cost. Formulating a problem as a POMDP itself  
1447 poses a number of challenges. For example, it might not be straightforward to cast  
1448 the optimization objective of the problem into an expected cumulative reward (with  
1449 stagewise additivity).

1450 A number of extensions to the basic POMDP framework are possible. First, of  
1451 particular interest to discrete event systems is the possibility of event-driven sensing,  
1452 where actions are taken only after some event occurs or some condition is met. In this  
1453 case, the state evolution is more appropriately modeled as a semi-Markov process  
1454 (though with some manipulation it can be converted into an equivalent standard  
1455 Markovian model) (Tijms 2003, Ch. 7). A second extension is to incorporate explicit  
1456 constraints into the decision-making framework (Altman 1998; Chen and Wagner  
1457 2007; Zhang et al. 2008).

1458 **References**

- 1459 Altman E (1998) Constrained Markov decision processes. Chapman and Hall/CRC, London  
1460 Bartels R, Backus S, Zeek E, Misoguti L, Vdovin G, Christov IP, Murnane MM, Kapteyn HC (2000)  
1461 Shaped-pulse optimization of coherent soft X-rays. *Nature* 406:164–166  
1462 Bellman R (1957) Dynamic programming. Princeton University Press, Princeton  
1463 Bertsekas DP (2005) Dynamic programming and suboptimal control: a survey from ADP to MPC.  
1464 In: Proc. joint 44th IEEE conf. on decision and control and European control conf., Seville, 12–15  
1465 December 2005  
1466 Bertsekas DP (2007) Dynamic programming and optimal control, vol I, 3rd edn, 2005; vol II, 3rd edn.  
1467 Athena Scientific, Belmont  
1468 Bertsekas DP, Castanon DA (1999) Rollout algorithms for stochastic scheduling problems. *Journal*  
1469 *of Heuristics* 5:89–108  
1470 Bertsekas DP, Tsitsiklis JN (1996) Neuro-dynamic programming. Athena Scientific, Belmont  
1471 Blatt D, Hero AO III (2006a) From weighted classification to policy search. In: Advances in neural  
1472 information processing systems (NIPS) vol 18, pp 139–146

Blatt D, Hero AO III (2006b) Optimal sensor scheduling via classification reduction of policy search (CROPS). In: Proc. int. conf. on automated planning and scheduling (ICAPS)	1473 1474
Castanon D (1997) Approximate dynamic programming for sensor management. In: Proc. 36th IEEE conf. on decision and control, San Diego, pp 1202–1207	1475 1476
Chang HS, Givan RL, Chong EKP (2004) Parallel rollout for online solution of partially observable Markov decision processes. <i>Discret Event Dyn Syst</i> 14(3):309–341	1477 1478
Chang HS, Fu MC, Hu J, Marcus SI (2007) Simulation-based algorithms for Markov decision processes. Springer series in communications and control engineering. Springer, Berlin Heidelberg New York	1479 1480 1481
Chen RC, Wagner K (2007) Constrained partially observed Markov decision processes for adaptive waveform scheduling. In: Proc. int. conf. on electromagnetics in advanced applications, Torino, 17–21 September 2007, pp 454–463	1482 1483 1484
Cheng HT (1988) Algorithms for partially observable Markov decision processes. PhD dissertation, University of British Columbia	1485 1486
Chhetri A, Morrell D, Papandreou-Suppappola A (2004) Efficient search strategies for non-myopic sensor scheduling in target tracking. In: Asilomar conf. on signals, systems, and computers	1487 1488
Cinlar E (1975) Introduction to stochastic processes. Prentice-Hall, Englewood Cliffs	1489
Chong EKP, Givan RL, Chang HS (2000) A framework for simulation-based network control via hindsight optimization. In: Proc. 39th IEEE conf. on decision and control, Sydney, 12–15 December 2000, pp 1433–1438	1490 1491 1492
de Farias DP, Van Roy B (2003) The linear programming approach to approximate dynamic programming. <i>Oper Res</i> 51(6):850–865	1493 1494
de Farias DP, Van Roy B (2004) On constraint sampling in the linear programming approach to approximate dynamic programming. <i>Math Oper Res</i> 29(3):462–478	1495 1496
Gottlieb E, Harrigan R (2001) The Umbra simulation framework. Sandia Tech Report SAND2001-1533 (Unlimited Release)	1497 1498
Q4 Gubner JA (2006) Probability and random processes for electrical and computer engineers. Cambridge University Press, New York	1499 1500
He Y, Chong EKP (2004) Sensor scheduling for target tracking in sensor networks. In: Proc. 43rd IEEE conf. on decision and control (CDC'04), 14–17 December 2004, pp 743–748	1501 1502
He Y, Chong EKP (2006) Sensor scheduling for target tracking: a Monte Carlo sampling approach. <i>Digit Signal Process</i> 16(5):533–545	1503 1504
Hero A, Castanon D, Cochran D, Kastella K (eds) (2008) Foundations and applications of sensor management. Springer, Berlin Heidelberg New York	1505 1506
Ji S, Parr R, Carin L (2007) Nonmyopic multiaspect sensing with partially observable Markov decision processes. <i>IEEE Trans Signal Process</i> 55(6):2720–2730 (Part 1)	1507 1508
Julier S, Uhlmann J (2004) Unscented filtering and nonlinear estimation. <i>Proc IEEE</i> 92(3):401–422	1509
Krakow LW, Li Y, Chong EKP, Groom KN, Harrington J, Rigdon B (2006) Control of perimeter surveillance wireless sensor networks via partially observable Markov decision process. In: Proc. 2006 IEEE int Carnahan conf on security technology (ICCST), Lexington, 17–20 October 2006	1510 1511 1512 1513
Kearns MJ, Mansour Y, Ng AY (1999) A sparse sampling algorithm for near-optimal planning in large Markov decision processes. In: Proc. 16th int. joint conf. on artificial intelligence, pp 1324–1331	1514 1515 1516
Kaelbling LP, Littman ML, Moore AW (1996) Reinforcement learning: a survey. <i>J Artif Intell Res</i> 4:237–285	1517 1518
Kaelbling LP, Littman ML, Cassandra AR (1998) Planning and acting in partially observable stochastic domains. <i>Artif Intell</i> 101:99–134	1519 1520
Kreucher CM, Hero A, Kastella K (2005a) A comparison of task driven and information driven sensor management for target tracking. In: Proc. 44th IEEE conf. on decision and control (CDC'05), 12–15 December 2005	1521 1522 1523
Kreucher CM, Kastella K, Hero AO III (2005b) Sensor management using an active sensing approach. <i>Signal Process</i> 85(3):607–624	1524 1525
Kreucher CM, Kastella K, Hero AO III (2005c) Multitarget tracking using the joint multitarget probability density. <i>IEEE Trans Aerosp Electron Syst</i> 41(4):1396–1414	1526 1527
Kreucher CM, Blatt D, Hero AO III, Kastella K (2006) Adaptive multi-modality sensor scheduling for detection and tracking of smart targets. <i>Digit Signal Process</i> 16:546–567	1528 1529
Kreucher CM, Hero AO, Kastella K, Chang D (2004) Efficient methods of non-myopic sensor management for multitarget tracking. In: Proc. 43rd IEEE conf. on decision and control (CDC'04), 14–17 December 2004	1530 1531 1532

- 1533 Krishnamurthy V (2005) Emission management for low probability intercept sensors in network  
1534 centric warfare. *IEEE Trans Aerosp Electron Syst* 41(1):133–151
- 1535 Krishnamurthy V, Evans RJ (2001) Hidden Markov model multiarm bandits: a methodology for  
1536 beam scheduling in multitarget tracking. *IEEE Trans Signal Process* 49(12):2893–2908
- 1537 Li Y, Krakow LW, Chong EKP, Groom KN (2006) Dynamic sensor management for multisensor  
1538 multitarget tracking. In: Proc. 40th annual conf. on information sciences and systems, Princeton,  
1539 22–24 March 2006, pp 1397–1402
- 1540 Li Y, Krakow LW, Chong EKP, Groom KN (2007) Approximate stochastic dynamic programming  
1541 for sensor scheduling to track multiple targets. *Digit Signal Process*. doi:[10.1016/j.dsp.2007.05.004](https://doi.org/10.1016/j.dsp.2007.05.004)
- 1542 Lovejoy WS (1991a) Computationally feasible bounds for partially observed Markov decision  
1543 processes. *Oper Res* 39:162–175
- 1544 Lovejoy WS (1991b) A survey of algorithmic methods for partially observed Markov decision  
1545 processes. *Ann Oper Res* 28(1):47–65
- Q4 1546 Meyn SP, Tweedie RL (1993) Markov chains and stochastic stability. Springer, London
- 1547 Miller SA, Harris ZA, Chong EKP (2009) A POMDP framework for coordinated guidance of  
1548 autonomous UAVs for multitarget tracking. *EURASIP J Appl Signal Process* (Special Issue  
1549 on Signal Processing Advances in Robots and Autonomy). doi:[10.1155/2009/724597](https://doi.org/10.1155/2009/724597)
- 1550 Pontryagin LS, Boltyansky VG, Gamkrelidze RV, Mishchenko EF (1962) The mathematical theory  
1551 of optimal processes. Wiley, New York
- 1552 Powell WB (2007) Approximate dynamic programming: solving the curses of dimensionality. Wiley-  
1553 Interscience, New York
- 1554 Ristic B, Arulampalam S, Gordon N (2004) Beyond the Kalman filter: particle filters for tracking  
1555 applications. Artech House, Norwood
- Q4 1556 Ross SM (1970) Applied probability models with optimization applications. Dover, New York
- 1557 Roy N, Gordon G, Thrun S (2005) Finding approximate POMDP solutions through belief compres-  
1558 sion. *J Artif Intell Res* 23:1–40
- 1559 Rust J (1997) Using randomization to break the curse of dimensionality. *Econometrica* 65(3):487–516
- 1560 Scott WR Jr, Kim K, Larson GD, Gurbuz AC, McClellan JH (2004) Combined seismic, radar, and  
1561 induction sensor for landmine detection. In: Proc. 2004 int. IEEE geoscience and remote sensing  
1562 symposium, Anchorage, 20–24 September 2004, pp 1613–1616
- 1563 Shi L, Chen C-H (2000) A new algorithm for stochastic discrete resource allocation optimization.  
1564 *Discret Event Dyn Syst* 10:271–294
- 1565 Smallwood RD, Sondik EJ (1973) The optimal control of partially observable Markov processes over  
1566 a finite horizon. *Oper Res* 21(5):1071–1088
- 1567 Sutton RS, Barto AG (1998) Reinforcement learning. MIT, Cambridge
- 1568 Thrun S, Burgard W, Fox D (2005) Probabilistic robotics. MIT, Cambridge
- 1569 Tijms HC (2003) A first course in stochastic models. Wiley, New York
- 1570 Washburn R, Schneider M, Fox J (2002) Stochastic dynamic programming based approaches to  
1571 sensor resource management. In: 5th int conf on information fusion
- 1572 Watkins CJCH (1989) Learning from delayed rewards. PhD dissertation, King's College, University  
1573 of Cambridge
- 1574 Willems JC (1996) 1969: the birth of optimal control. In: Proc. 35th IEEE conf. on decision and  
1575 control (CDC'96), pp 1586–1587
- 1576 Wu G, Chong EKP, Givan RL (2002) Burst-level congestion control using hindsight optimization.  
1577 *IEEE Trans Automat Control* (Special Issue on Systems and Control Methods for Communica-  
1578 tion Networks) 47(6):979–991
- 1579 Yu H, Bertsekas DP (2004) Discretized approximations for POMDP with average cost. In: Proc. 20th  
1580 conf. on uncertainty in artificial intelligence, Banff, pp 619–627
- 1581 Zhang NL, Liu W (1996) Planning in stochastic domains: problem characteristics and approximation.  
1582 Tech. report HKUST-CS96-31, Dept. of Computer Science, Hong Kong University of Science  
1583 and Technology
- 1584 Zhang Z, Moola S, Chong EKP (2008) Approximate stochastic dynamic programming for oppor-  
1585 tunistic fair scheduling in wireless networks. In: Proc. 47th IEEE conf. on decision and control,  
1586 Cancun, 9–11 December 2008, pp 1404–1409



**Edwin K. P. Chong** received the BE(Hons) degree with First Class Honors from the University of Adelaide, South Australia, in 1987; and the MA and PhD degrees in 1989 and 1991, respectively, both from Princeton University, where he held an IBM Fellowship. He joined the School of Electrical and Computer Engineering at Purdue University in 1991, where he was named a University Faculty Scholar in 1999, and was promoted to Professor in 2001. Since August 2001, he has been a Professor of Electrical and Computer Engineering and a Professor of Mathematics at Colorado State University. His research interests span the areas of communication and sensor networks, stochastic modeling and control, and optimization methods. He coauthored the recent best-selling book, An Introduction to Optimization, 3rd Edition, Wiley-Interscience, 2008. He is currently on the editorial board of the IEEE Transactions on Automatic Control, Computer Networks, Journal of Control Science and Engineering, and IEEE Expert Now. He is a Fellow of the IEEE, and served as an IEEE Control Systems Society Distinguished Lecturer. He received the NSF CAREER Award in 1995 and the ASEE Frederick Emmons Terman Award in 1998. He was a co-recipient of the 2004 Best Paper Award for a paper in the journal Computer Networks. He has served as Principal Investigator for numerous funded projects from NSF, DARPA, and other funding agencies. 1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601



**Christopher M. Kreucher** received the BS, MS, and PhD degrees in Electrical Engineering from the University of Michigan in 1997, 1998, and 2005, respectively. He is currently a Senior Systems Engineer at Integrity Applications Incorporated in Ann Arbor, Michigan. His current research interests include nonlinear filtering (specifically particle filtering), Bayesian methods of fusion and multitarget tracking, self localization, information theoretic sensor management, and distributed swarm management. 1602  
1603  
1604  
1605  
1606  
1607



1608 **Alfred O. Hero III** received the BS (summa cum laude) from Boston University (1980) and the PhD  
1609 from Princeton University (1984), both in Electrical Engineering. Since 1984 he has been with the  
1610 University of Michigan, Ann Arbor, where he is a Professor in the Department of Electrical Engi-  
1611 neering and Computer Science and, by courtesy, in the Department of Biomedical Engineering and  
1612 the Department of Statistics. He has held visiting positions at Massachusetts Institute of Technology  
1613 (2006), Boston University, I3S University of Nice, Sophia-Antipolis, France (2001), Ecole Normale  
1614 Supérieure de Lyon (1999), Ecole Nationale Supérieure des Télécommunications, Paris (1999),  
1615 Scientific Research Labs of the Ford Motor Company, Dearborn, Michigan (1993), Ecole Nationale  
1616 Supérieure des Techniques Avancées (ENSTA), Ecole Supérieure d'Électricité, Paris (1990), and  
1617 M.I.T. Lincoln Laboratory (1987–1989). His recent research interests have been in areas including:  
1618 inference for sensor networks, adaptive sensing, bioinformatics, inverse problems, and statistical  
1619 signal and image processing. He is a Fellow of the Institute of Electrical and Electronics Engineers  
1620 (IEEE), a member of Tau Beta Pi, the American Statistical Association (ASA), the Society for  
1621 Industrial and Applied Mathematics (SIAM), and the US National Commission (Commission C)  
1622 of the International Union of Radio Science (URSI). He has received a IEEE Signal Processing  
1623 Society Meritorious Service Award (1998), IEEE Signal Processing Society Best Paper Award  
1624 (1998), a IEEE Third Millennium Medal and a 2002 IEEE Signal Processing Society Distinguished  
1625 Lecturership. He was President of the IEEE Signal Processing Society (2006–2007) and during his  
1626 term served on the TAB Periodicals Committee (2006). He was a member of the IEEE TAB Society  
1627 Review Committee (2008) and is Director-elect of IEEE for Division IX (2009).

## AUTHOR QUERIES

### **AUTHOR PLEASE ANSWER ALL QUERIES**

- Q1. References in the caption of Fig. 1 to “Top” and “Bottom” were changed to “a” and “b and c”, respectively. Please check if this was appropriate.
- Q2. Please check Figure 9 if captured correctly.
- Q3. The letter “L” was deleted from a citation to “Fig. 14L”. Please check if this was appropriate.
- Q4. “Çinlar (1975)”, “Gubner (2006)”, “Meyn and Tweedie (1993)”, and “Ross (1970)” were listed in the reference list but were not cited in the text. Please provide citations to these reference items or, alternatively, delete them from the reference list.