

# Sensor Scheduling for Energy-Efficient Target Tracking in Sensor Networks

George K. Atia, *Member, IEEE*, Venugopal V. Veeravalli, *Fellow, IEEE*, and Jason A. Fuemmeler, *Member, IEEE*

**Abstract**—In this paper, we study the problem of tracking an object moving randomly through a network of wireless sensors. Our objective is to devise strategies for scheduling the sensors to optimize the tradeoff between tracking performance and energy consumption. We cast the scheduling problem as a partially observable Markov decision process (POMDP), where the control actions correspond to the set of sensors to activate at each time step. Using a bottom-up approach, we consider different sensing, motion and cost models with increasing levels of difficulty. At the first level, the sensing regions of the different sensors do not overlap and the target is only observed within the sensing range of an active sensor. Then, we consider sensors with overlapping sensing range such that the tracking error, and hence the actions of the different sensors, are tightly coupled. Finally, we consider scenarios wherein the target locations and sensors' observations assume values on continuous spaces. Exact solutions are generally intractable even for the simplest models due to the dimensionality of the information and action spaces. Hence, we devise approximate solution techniques, and in some cases derive lower bounds on the optimal tradeoff curves. The generated scheduling policies, albeit suboptimal, often provide close-to-optimal energy-tracking tradeoffs.

**Index Terms**—Dynamic programming, Markov models, POMDP, sensor networks, target tracking.

## I. INTRODUCTION

IN large networks of inexpensive sensors with small batteries, the sensor nodes are required to operate on limited energy budgets. Sensor management can prolong the lifetime of a sensor network and conserve scarce energy resources. However, inefficient management could result in severe performance degradation. In this paper, we consider a network of  $n$  sensors tracking a single object. The sensors can be turned on or off at consecutive time steps and the goal is to select the subset of sensors to activate at each time step. This problem is challenging due to the inherent tradeoff between the value of information in

the measurements and the energy cost, combined with the combinatorial complexity of the decision space.

In previous work [1], two of the authors considered approximate strategies for *sensor sleeping*, where the sensors are put to sleep to save energy and decisions are made concerning their sleep duration (in time slots). Once in a sleep mode, a sensor would only wake up after its own sleep timer expires. Here, we consider a scheduling variant of the problem which can be thought of as a sleeping problem with an external wake-up mechanism, i.e., sensors can be woken up by external means (e.g., a low-power wake-up radio). At time  $k$ , the permissible control actions for an  $n$ -sensor scheduling problem are  $n$ -dimensional binary vectors, i.e., vectors in  $\{0, 1\}^n$  (corresponding to the set of sensor nodes to activate at each time step), in contrast to vectors in  $\mathbb{N}_0^{n_a(k)}$  for the sleeping problem (corresponding to the sleep durations of awake sensors), where  $\mathbb{N}_0$  is the set of non-negative integers and  $n_a(k)$  the number of awake sensors at time  $k$ . The simpler structure of the control space for the scheduling problem does not address the combinatorial nature of the control space, yet it enables efficient approximate solution methodologies for the more realistic models that we study in this paper.

A significant body of related research work considers sensor management for tasking sensors in dynamically evolving environments. Castanon [2] has developed an approximate dynamic programming approach for dynamic scheduling of multi-mode sensor resources for the classification of a large number of unknown objects. The goal is to achieve an accurate classification of each object at the end of a fixed finite horizon by assigning different sensor modes to different objects subject to periodic or total resource usage constraints. Mode allocation strategies are computed based on Lagrangian relaxation for an approximate optimization problem wherein sample-path resource constraints are replaced by expected value constraints. In the context of sensor scheduling for target tracking, information-based approaches [3]–[6] have been developed for optimizing tracking performance subject to an explicit constraint on communication costs in a decentralized setting. Williams *et al.* [3] also adopt a Lagrangian relaxation approach to solve a constrained dynamic program over a rolling horizon. There, the combinatorial complexity of the decision space is avoided by first selecting one leader node, followed by greedy sensor subset selection. Other related work on sensor scheduling include leader-based distributed tracking schemes [7]–[10], where at any time instant there is only one active sensor, namely, the leader sensor which changes dynamically as a function of the object state, while the rest of the network is idle. In [11] a scheduler chooses the least number of sensors necessary to reduce the covariance matrix

Manuscript received July 27, 2010; revised February 11, 2011 and May 24, 2011; accepted June 08, 2011. Date of publication June 20, 2011; date of current version September 14, 2011. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Mark Coates. This work was funded in part by a grant from the Motorola corporation, a U.S. Army Research Office MURI grant W911NF-06-1-0094 through a subcontract from Brown University at the University of Illinois, an NSF Graduate Research Fellowship, and by a Vodafone Fellowship. This work appeared in part at the Asilomar Conference on Signals, Systems, and Computers, Pacific Grove, CA, November 2010.

The authors are with the Coordinated Science Laboratory (CSL), University of Illinois at Urbana Champaign, Urbana IL 61801 USA (e-mail: atia1@illinois.edu; vvv@illinois.edu; fuemmele@illinois.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSP.2011.2160055

of the estimate error to a desired value. Other non-myopic policies including covariance-based and unscented transform-based schedulers were developed in [12] for scheduling the best sensor at each time step.

While previous work focused on developing distributed implementations of efficient sensor scheduling strategies, our goal here is to study the fundamental theory of sensor scheduling for tracking and surveillance applications. Specifically, to explicitly study the fundamental tradeoff between tracking performance and energy expenditure, we define a unified objective function combining tracking and energy costs trading-off the complexity of per-stage costs for tractability. We adopt a bottom-up approach where we consider a range of sensing, motion and cost models with increasing levels of difficulty and devise suboptimal scheduling policies to balance the tradeoff between energy expenditure and tracking performance. In some cases we are also able to derive lower bounds on the optimal energy-tracking tradeoff.

Due to noise and model uncertainties, natural limitations of the measurement devices, or incomplete data about the surroundings, we need to design scheduling policies when the system's state is only partially observable to the controller. Partially observable Markov decision processes (POMDPs) provide a natural framework for addressing sequential decision problems where the goal is to find a policy (strategy) for selecting actions based on the information available to the controller, while addressing both short-term and long-term benefits and costs. Solving POMDPs optimally is generally intractable. For example, the value function for a POMDP with a finite state space depends on information states consisting of conditional probability vectors of dimension equal to the number of states. This has led to a number of POMDP approximations and we refer the reader to Monahan [13] and Hauskrecht [14] for excellent surveys on these methods. Usually, no single approximation can be prescribed for all POMDPs, rather approximations can be judiciously used to exploit specific problem structures. In this paper, we use a subset of these approximate solution techniques, including reduced-uncertainty and point-based approximations [15]–[18]. The former approach assumes that more information would be available to the controller at future time steps, and the latter solves a reduced optimization problem based on a relatively small subset of sampled beliefs about the object's state. We devise different approaches to deal with the aforementioned computational complexity of the decision space. In one approach, instead of solving one large combinatorial problem, we solve a set of simpler subproblems based on the intuition gained from a simplistic sensing model. In another approach, we iteratively sample control actions from a reduced control space based on the sparsity of a reachable belief set combined with point-based value updates.

To summarize, previous work on sensor scheduling for tracking has either focused on scheduling a small number of sensors, or addressed the combinatorial complexity for larger networks using leader-based distributed algorithms and greedy sensor selection or non-myopic policies based on few-steps look ahead. The main contributions of our work are:

- **Developing approximate sensor scheduling policies:**

In this paper, we develop scheduling policies based on

an observable-after-control assumption ( $Q_{\text{MDP}}$ ) and point-based approximations for various sensing, transition and cost models with increasing level of difficulty. In particular,

- 1) We show that under a reduced-future-uncertainty assumption, the value function for a first level model (Section II-A) is significantly simplified as the problem decomposes into simpler subproblems (one per-sensor) leading to substantial complexity reduction. The resulting sensor scheduling policy is shown to be near optimal.
  - 2) We develop new sensor scheduling policies combining  $Q_{\text{MDP}}$  and reinforcement learning for more advanced models considered in Sections II-B and II-C wherein tracking errors—and hence actions of different sensors—are tightly coupled. We propose per-sensor surrogate value functions for artificially decoupled per-sensor subproblems and use reinforcement learning to learn the corresponding individual tracking costs.
  - 3) We develop point-based sensor scheduling policies which are generally shown to outperform their  $Q_{\text{MDP}}$  counterparts at the expense of an increase in computational complexity. We integrate our point-based schedulers with machinery to address the dimensionality of the control and observation spaces via action sampling and observation aggregation, respectively.
- **Lower bounds on optimal performance:** In some cases, we derive lower bounds on the optimal energy-tracking tradeoffs. Particularly, the  $Q_{\text{MDP}}$  surrogate value function is itself a lower bound on the optimal value function for our first level model in Section II-A. We also derive a lower bound on optimal performance for a continuous Gaussian observation model with Hamming distance tracking cost.

The remainder of this paper is organized as follows. In Section II, we describe the tracking problem and define the sensing, transition and cost models, as well as the optimization problem, for each of the considered models. In Section III we describe approximate strategies to generate suboptimal scheduling policies. In Section IV, we present some experimental results, and finally, in Section V, we provide some concluding remarks.

## II. SCHEDULING PROBLEM

In the following we consider different models with increasing level of difficulty. Depending on the structure of the model, we devise approximate methods to address the associated difficulties and generate efficient scheduling policies. For notation, vectors are denoted by bold lower-case letters. Superscript T denotes transposition and the indicator function is written as  $\mathbb{I}\{\cdot\}$ .

### A. Simple Sensing, Observation and Cost Models

In this model, the network is divided into  $n$  distinct cells, one for each sensor. In other words, each cell corresponds to the sensing range of one particular sensor, and the sensors' ranges do not overlap. A Markov chain with an  $(n+1) \times (n+1)$  probability transition matrix  $P$  describes the motion of the target through the field of interest. The extra state is for an absorbing termination state of the Markov chain which is reached when the

object leaves the network. We let  $b_k$  denote the true target location at time  $k$ . It is further assumed that all information about the object trajectory is stored at some central unit and is used to determine the scheduling actions for the different sensors. It is worthwhile mentioning that alternatively one could also consider a case where the object might return to the network, causing the problem to never terminate. This would then necessitate a discounted cost (or average cost) formulation where cost incurred in the present has more weight than cost incurred at future stages.

We let  $u_{k,\ell}$  denote the action for sensor  $\ell$  at time  $k$ ;  $u_{k,\ell} = 1$  if sensor  $\ell$  is activated at time  $k + 1$  and 0 if the decision is to turn it off. The action vector at time  $k$ , denoted  $\mathbf{u}_k$ , is a binary vector of size  $n \times 1$ , one decision per sensor. In this simplistic model, we assume that the target is perfectly observable within the cell of an awake sensor or if it reaches the terminal state  $\tau$ , otherwise it is unobservable. Thus, the observation  $s_k$  at time  $k$  is defined according to

$$s_k = \begin{cases} b_k, & \text{if } b_k \neq \tau \text{ and } u_{k-1,b_k} = 1 \\ \varepsilon, & \text{if } b_k \neq \tau \text{ and } u_{k-1,b_k} = 0 \\ \tau, & \text{if } b_k = \tau \end{cases} \quad (1)$$

where  $\varepsilon$  stands for erasure. The observation model in (1) induces a well-defined probabilistic observation model  $p(s_k | b_k, \mathbf{u}_{k-1})$  such that the current observation depends on that actual target location and the scheduling action for the  $n$  sensors.

At each time step, the incurred cost is the sum of the energy and the tracking costs. An energy cost of  $c \in (0, 1]$  per unit time is incurred for every active sensor and a tracking cost of 1 for each time unit that the object is not observed. Once state  $\tau$  is reached the problem terminates and no further cost is incurred. In other words,  $\tau$  is an absorbing cost-free state; all  $n$  states are transient so that  $\tau$  is the only recurrence class of the Markov chain. Hence,

$$g(b_k, \mathbf{u}_{k-1}) = \mathbb{1}\{b_k \neq \tau\} \left( \mathbb{1}\{u_{k-1,b_k} = 0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_{k-1,\ell} = 1\} \right). \quad (2)$$

The parameter  $c$  is thus used to tradeoff energy consumption and tracking errors.

A drawback of our definition for the tracking cost in this simplified first level model is that cost may be incurred even if the object location is known through a process of elimination. Our definition of tracking cost is used to easily separate the problem into a set of simpler subproblems. However, note that not much is lost through this simplification since we require only one additional sensor awake per unit time to maintain zero tracking errors. Furthermore, in the models we introduce next in Sections II-B and II-C, this assumption is relaxed as we use a Hamming distance cost function where an error is incurred only if the estimated location is different from the actual target location.

### B. Overlapping Sensors With Discrete Observations Models

In this model, we continue to use a discrete model for the target transition but we redefine a new sensing model and cost

structure to account for the fact that sensors could have overlapping visibility regions. Within that model we further consider simple and probabilistic sensing.

1) *Overlapping Sensors With Simple Sensing*: In this case, the target is *perfectly* observed within the visibility region of *any* active sensor. Denote by  $R_\ell$  the set of locations in the visibility region of sensor  $\ell$  and by  $\mathcal{B}_i$  the set of sensors that observe location  $i$ . The observation at time  $k$  is given by

$$s_k = \begin{cases} b_k, & \text{if } b_k \neq \tau \text{ and } \exists j \in \mathcal{B}_{b_k} : u_{k-1,j} = 1 \\ \varepsilon, & \text{if } b_k \neq \tau \text{ and } u_{k-1,j} = 0, \quad \forall j \in \mathcal{B}_{b_k} \\ \tau, & \text{if } b_k = \tau. \end{cases} \quad (3)$$

Therefore, a tracking error is incurred if none of the sensors observing the current target location is active. Redefining the cost structure for this model:

$$g(b_k, \mathbf{u}_{k-1}) = \mathbb{1}\{b_k \neq \tau\} \left( \mathbb{1}\{u_{k-1,j} = 0, \forall j \in \mathcal{B}_{b_k}\} + \sum_{\ell=1}^n c \mathbb{1}\{u_{k-1,\ell} = 1\} \right) \quad (4)$$

2) *Overlapping Sensors With Probabilistic Sensing*: By probabilistic sensing, we account for observation uncertainty even if the target is within the visibility region of one or more active sensors. We assume

$$p(s_k | b_k, \exists j \in \mathcal{B}_{b_k} : u_{k-1,j} = 1) = \begin{cases} q, & s_k = b_k \\ \frac{1-q}{|\mathcal{R}_k|-1}, & s_k = i, \forall i \in \mathcal{R}_k \end{cases} \quad (5)$$

where

$$\mathcal{R}_k = \bigcap_{\substack{j \in \mathcal{B}_{b_k} \\ u_{k-1,j} = 1}} R_j \setminus \bigcup_{\substack{i \notin \mathcal{B}_{b_k} \\ u_{k-1,i} = 1}} R_i.$$

That is, the observation is uniformly distributed over the remaining locations  $\mathcal{R}_k$  (other than the true target location) that belong to the visibility regions of the set of awake sensors monitoring the true location  $b_k$ . The number of these locations is  $|\mathcal{R}_k|$  and is function of the control  $\mathbf{u}_{k-1}$  and the object state  $b_k$ . If the true target location does not belong to the visibility region of an awake sensor, we naturally exclude the visibility region of that sensor since no measurement is received from such a sensor. When  $\mathcal{R}_k$  is a singleton  $\{b_k\}$ , we set  $q = 1$ . A tracking error is incurred if the target is not directly observed and the uncertainty in the target location cannot be resolved.

### C. Continuous Observation, Continuous State and Arbitrary Cost Models

In this class of models, we allow for an arbitrary distribution of the observations given the current object location. Tracking cost is modeled through an arbitrary distance measure between the actual and the estimated object location. If we denote the set of possible object locations  $\mathcal{B}$ , we have  $\mathcal{B} = m + 1$ . Note that, in contrast to the simplistic model in Section II-A,  $m$  is different from  $n$  since object locations are arbitrary and we no longer assume one location corresponds to the sensing range of one particular sensor. The  $(m + 1)$ th state again corresponds to a termination state. Furthermore, the target can be moving on a continuous state space in which case  $m$  is  $\infty$ .

If the state space is discrete, then conditioned on the object state  $b_k$  at time  $k$ ,  $b_{k+1}$  has a probability mass function that is given by the  $b_k$ th row of the transition matrix  $P$ . If the state space is continuous,  $P$  is a kernel such that  $P(x, \mathcal{Y})$  is the probability that the next object location is in the set  $\mathcal{Y} \subset \mathcal{B}$  given the current object location is  $x$ . For simplicity of exposition, we focus on discrete state spaces. Also, we omit indexing time whenever the time evolution is well-understood to avoid cumbersome notation. We consider the following observation model for illustration; however, our approach is fairly general:

$$p(\mathbf{s}|\mathbf{b}, \mathbf{u}) = \prod_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{1}{2} \left( s_i - \frac{V}{(b - p_i)^2 + 1} \right)^2 \right) \times \mathbb{1}\{u_i = 1\} + \delta(s_i - \varepsilon) \mathbb{1}\{u_i = 0\} \right\} \quad (6)$$

where  $\mathbf{s}$  is an  $n \times 1$  continuous observation vector with the  $i$ th entry,  $s_i$ , representing the observation of sensor  $i$ ,  $p_i$ ,  $i = 1, \dots, n$ , is the position of the  $i$ th sensor,  $b$  is the target state,  $V$  is some positive constant,  $\varepsilon$  stands for erasure, and  $\delta(\cdot)$  is the Dirac Delta function. In (6), the observation of an active sensor is Gaussian with a mean received signal strength inversely proportional to the square of the distance between the sensor and the actual target location. The observation of an inactive sensor is just an erasure.

The estimated target location (given the entire history) is denoted by  $\hat{b}$ . We define the tracking error through an arbitrary bounded distance function  $d(b, \hat{b})$  between the actual and the estimated object locations, which can be the Hamming distance  $d(b, \hat{b}) = \mathbb{1}\{b \neq \hat{b}\}$  or the Euclidean distance for discrete and continuous state spaces, respectively. The control at each time step is the pair  $(\hat{b}_k, \mathbf{u}_k)$ . Since  $\hat{b}$  does not affect the state evolution, the optimal value for  $\hat{b}_k$  is the value that minimizes the tracking cost over a single time step given history up to time  $k$ , i.e.,

$$\hat{b}_k = \arg \min_{\hat{b}} E \left[ d(b_k, \hat{b}_k) | I_k \right] \quad (7)$$

where  $I_k$  denotes the information state, i.e., the total information available to the central controller at time  $k$ , which is given by

$$I_k = \{\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_k, \mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{k-1}\}.$$

In the case of Hamming cost, it follows that  $\hat{b}$  is simply the MAP decision, i.e.,  $\hat{b} = \arg \max_b p_k(b)$ , where  $\mathbf{p}_k$  is the posterior probability distribution of the target state.

#### D. Optimal Scheduling Policy

The design of an *optimal scheduling policy* depends on the history up to time  $k$ , i.e., the information state  $I_k$ . However, the posterior probability distribution,  $\mathbf{p}_k = \Pr[b_k | I_k]$ , of the target's state given  $I_k$  is a sufficient statistic for this class of partially observable processes. The distribution  $\mathbf{p}_k$ , also known as belief, summarizes all the information needed for optimal control. The sufficient statistic itself forms a Markov process

whose evolution can be obtained through Bayes' rule updates.<sup>1</sup> For example, the belief update equation for the simplistic model in Section II-A can be written as

$$\mathbf{p}_{k+1} = \begin{cases} \mathbf{e}_\tau, & \text{if } s_{k+1} = \tau \\ \mathbf{e}_{b_{k+1}}, & \text{if } u_{k,b_{k+1}} = 1 \\ [\mathbf{p}_k P]_{\{j: u_{k,j}=0\}}, & \text{if } u_{k,b_{k+1}} = 0 \end{cases} \quad (8)$$

where  $\mathbf{e}_i$  is a row vector with a 1 at the  $i$ th entry and 0 elsewhere. For an index set  $\mathcal{S}$ , the entries  $v'_i$  of a vector  $\mathbf{v}' = [\mathbf{v}]_{\mathcal{S}}$  are obtained from the vector  $\mathbf{v}$  as follows:

$$v'_i = \begin{cases} 0, & \text{if } i \notin \mathcal{S} \\ \frac{v_i}{\sum_{j \in \mathcal{S}} v_j}, & \text{if } i \in \mathcal{S} \end{cases}$$

Hence, the vector  $[\mathbf{p}_k P]_{\mathcal{S}}$  is the probability vector formed by setting the  $i$ th entry  $[\mathbf{p}_k P]_i$  of the vector  $\mathbf{p}_k P$  to 0,  $\forall i \notin \mathcal{S}$ , then normalizing the vector into a probability distribution. The set  $\{j : u_{k,j} = 0\}$  signifies the set of deactivated sensors. In other words, the updated belief for the model in Section II-A, is a point mass distribution concentrated at  $\tau$  if the object exits the network, and concentrated at  $b_{k+1}$  if the object is observed. When the object is unobservable, we eliminate the probability mass at all sensors that are awake, since the object cannot be at these locations, and normalize. The multi-valued function in (8), and equivalent Bayes' updates for the other models, define a transformation

$$\mathbf{p}_{k+1} = \phi(\mathbf{p}_k, s_{k+1}, \mathbf{u}_k) \quad (9)$$

mapping the current belief  $\mathbf{p}_k$ , the current control vector  $\mathbf{u}_k$ , and the future observation  $s_{k+1}$ , to a future belief.

The policy  $\mathbf{u}_k = \mu_k(I_k)$  is defined as a mapping from information states  $I_k$  to control actions  $\mathbf{u}_k$ . The goal is to design a policy that minimizes the expected sum of costs  $J$ , where

$$J(I_0, \mu_0, \mu_1, \dots) = E \left[ \sum_{k=1}^{\infty} g(b_k) | I_0 \right]. \quad (10)$$

$J$  is well-defined since  $g$  is upper bounded by  $cn + 1$  (regardless of the model) and the expected time till the object exits the network is finite. Note that the termination is inevitable, thus the objective is to reach the termination state with minimal expected cost. Hence, the scheduling policy is the solution of the minimization problem

$$J^* = \min_{\mu_0, \mu_1, \dots} J(I_0, \mu_0, \mu_1, \dots). \quad (11)$$

This POMDP problem falls within the class of infinite horizon stochastic shortest path problems. Noting that the termination state is observable, cost-free and absorbing, and that every policy is proper,<sup>2</sup> a stationary policy  $\mu^*(\cdot)$ , i.e., one which

<sup>1</sup>Equivalently, for a continuous state space, a sufficient statistic would be  $p_k(\mathcal{X}) = \Pr[b_k \in \mathcal{X} | I_k]$ . The updated belief  $p_{k+1}$  can be computed using standard Bayesian non-linear filtering as the posterior measure resulting from prior measure  $pP$  and observation  $s_{k+1}$ .

<sup>2</sup>A proper policy is a policy that leads to the termination state with probability one regardless of the initial state. In our problem, the scheduling policy does not affect the target motion and all policies are proper in the sense that there is a positive probability that the target will reach the termination state after a finite number of stages.

does not depend on  $k$ , is optimal in the class of all history-dependent policies and  $\mathbf{p}_k$  is a sufficient statistic for control [19], i.e.,  $u_k^* = \mu^*(\mathbf{p}_k)$ , is defined through a time-invariant mapping from the belief space to the action space.  $J$  can be written in terms of the sufficient statistic and the optimal policy can be obtained from the solution of the Bellman equation:

$$J(\mathbf{p}) = \min_{\mathbf{u} \in \{0,1\}^n} E[g(b, \mathbf{u})|\mathbf{p}, \mathbf{u}] + \sum_s p(s|\mathbf{p}, \mathbf{u}) J(\phi(\mathbf{p}, s, \mathbf{u})) \quad (12)$$

such that  $J(\mathbf{e}_\tau) = 0$ , where  $J(\cdot)$  is the value function for the POMDP and  $\phi$  is defined in (9). Note that we removed the time dependence due to the aforementioned time invariance property. For continuous observations, summation over  $s$  is replaced by an integration.

### III. APPROXIMATE SOLUTIONS AND LOWER BOUNDS

There are a number of algorithms for solving POMDPs exactly [20]–[22]. These algorithms rely on the powerful result of Sondik that the optimal value function for any POMDP can be approximated arbitrarily closely using a set of hyper-planes ( $\alpha$ -vectors) defined over the belief simplex [20]. This fact is the basis for exact value iteration based algorithms, such as the Witness algorithm [23] for computing the value function. The result is a value function parameterized by a number of hyperplanes (or vectors) whereby the belief space is partitioned into a finite number of regions. Each vector minimizes the value function over a certain region of the belief space and has a control action associated with it, which is the optimal control for the beliefs in its region.

To clarify, in value iteration we generally start with some initial estimate for  $J^*$  and repeatedly apply the transformation defined by the right-hand side of Bellman equation (12) until the sequence of cost functions converges. Let  $\{\alpha_i^{(k)}\}_{i=1}^{|J^{(k)}|}$  denote the set of vectors parameterizing the value function  $J^{(k)}$  after  $k$  iterations, where  $|J^{(k)}|$  is the total number of hyperplanes, and  $\alpha_i^{(k)}(b)$ , which is a hyperplane in the belief space, represents the value of executing the  $k$ -step policy associated with the  $i$ th vector starting from a state  $b$ . Hence, the value of executing the  $i$ th hyperplane policy starting from a belief state  $\mathbf{p}$  is simply the dot product of  $\alpha_i$  and  $\mathbf{p}$ :

$$J_i^{(k)}(\mathbf{p}) = \sum_b \mathbf{p}(b) \alpha_i^{(k)}(b) = \mathbf{p} \cdot \alpha_i^{(k)}.$$

Therefore, the value of the optimal  $k$ -step policy starting at  $\mathbf{p}$  is simply the minimum dot product over all hyperplanes, i.e.,

$$J^{*(k)}(\mathbf{p}) = \min_{\{\alpha_i^{(k)}\}} \mathbf{p} \cdot \alpha_i^{(k)}.$$

Hence,  $J^{*(k)}(\mathbf{p})$  is piecewise linear and concave. Some of the vectors (also known as policy trees) may be dominated by others in the sense that they are not optimal at any region in the belief simplex. Thus, many exact algorithms devise pruning mechanisms whereby a parsimonious representation with a minimal set of non-dominated hyperplanes is maintained [13].

Even though the aforementioned linearity/concavity property makes the policy search a great deal simpler, the exact computation is generally intractable except for relatively small problems. The two major difficulties for exact computation arise from the exponential growth of the vectors with the planning horizon and with the number of observations, and the inefficiencies related to identification of such vectors and subsequently pruning them. Namely, the number of hyperplanes grows double exponentially such that after  $k$  steps the number of hyperplanes is  $O(|\mathcal{U}|^{|\mathcal{S}|^k})$ , where  $|\mathcal{U}|$  and  $|\mathcal{S}|$  denote the cardinality of the control and observation spaces, respectively. Equivalently, the number of hyperplanes per iteration grows as

$$|J^{(k+1)}| = O(|\mathcal{U}| |J^{(k)}|^{|\mathcal{S}|}).$$

This has led to a number of approximations and suboptimal solutions techniques that trade off solution quality for speed.

*Remark III.1:* The intractability of the optimal solution for our problem is primarily due to the following reasons:

- i) The cost function is minimized over the simplex of probability distributions, i.e., the  $(m - 1)$ -dimensional belief simplex for  $m$ -state discrete state-space models, and the space of probability density functions for continuous state-space models.
- ii) The exponential explosion of the action space with the number of sensors ( $2^n$  actions).
- iii) The exponential growth of the  $\alpha$ -vectors with the planning horizon and with the number of observations, especially for continuous observation models.

#### A. Approximate Solutions

In this section, we outline our approximate solution methodologies for the different models introduced in Section II. First, we consider approximations where it is assumed that more information becomes available to the controller at future time steps. Policies based on the assumption that uncertainty in the current belief state will be gone after the next action were first introduced within the artificial intelligence community and known as  $Q_{\text{MDP}}$  policies [16], [23]. We show that under an observable-after-control assumption, our sensor scheduling problem decomposes into  $n$  simpler subproblems, one subproblem per sensor, for the *simplistic model* of Section II-A. These subproblems can then be solved exactly using policy iteration [19]. Furthermore, in this case, the  $Q_{\text{MDP}}$  solution gives us a lower bound on the optimal tracking-energy tradeoff. Unfortunately, this natural decomposition does not extend to the other class of models due to the inherent coupling of their tracking errors. However, based on the intuition gained from the simplistic model, we artificially decouple the scheduling problem for the more complicated models, and individually learn the tracking costs corresponding to each subproblem under the aforementioned  $Q_{\text{MDP}}$  assumption. This approach combines  $Q_{\text{MDP}}$  with reinforcement learning [24].

Second, we develop sensor scheduling strategies based on point-based approximations. Despite the fact that the generated  $Q_{\text{MDP}}$  based policies perform reasonably well, generally the

resulting policies would not take actions to gain information (an effect of the observable-after-control assumption), leading to situations wherein the belief state does not get updated appropriately. Furthermore, while decoupling the scheduling problem provides close-to optimal performance for uncoupled or lightly-coupled sensing and tracking models (see Section IV), it might come at the expense of reduction in solution quality for more realistic or heavily-coupled models. To that end, we develop point-based approximate scheduling policies. While our previous approach reduced complexity via decoupling and learning, the key idea here is to optimize the value function only for a small set of reachable beliefs  $\mathcal{P}$  and not over the entire belief simplex. Point-based methods have shown great potential for solving large scale POMDPs mostly for robotic applications [14], [15], [17], [25]. Pineau *et al.* [15] proposed point-based value iteration (PBVI) which performs point-based backups only at a discrete set of reachable belief points, that can be actually encountered by interacting with the environment. Developing a class of point-based algorithms, which mostly differ in the way the subset of belief points is chosen and the execution order of the backup operations over the selected belief points, has been the focus of recent algorithm-development research targeting large scale POMDPs. Perseus [17] is one such randomized point-based algorithm that maintains a fixed set of belief points. There, backup speedups can be obtained by exploiting the key observation that a single backup may improve the value of many belief points simultaneously. These algorithms were designed to deal with large state spaces, yet, two extra difficulties in the scheduling problem arise from the size of the action space (which is  $2^n$  for all models) and the observation space (for the models in Sections II-C). Regarding the dimensionality of the action space, we devise a strategy to sample actions based on the support of the beliefs and the sparse structure of the transition models. Intuitively speaking, an object can only move from one side of the network to the other side within time constraints rendering exponentially many scheduling actions irrational at certain times. Hence, instead of performing full updates including  $2^n$  actions, we perform the minimization over a reduced control space  $\mathcal{U}(\mathbf{p})$  for every  $\mathbf{p} \in \mathcal{P}$  (see Section III-C-1). When dealing with continuous or large observations, we combine that with a methodology that aggregates observations and uses aggregate observations for value iteration updates (Section III-C-2). At the core of the algorithm, we use Perseus [17], a variant of PBVI [15], whereby value iteration updates are not carried out for every sampled belief. Instead, the values for many belief points are improved simultaneously in one update. Fig. 1 depicts the structure of our point-based approximation, combining control space reduction and observation aggregation with point-based updates.

### B. $Q_{\text{MDP}}$ Based Scheduling Policies

Next, we consider our first class of policies based on the  $Q_{\text{MDP}}$  reduced future uncertainty assumption. First, we consider the simplistic model in Section II-A, then we use the intuition we developed from this model to devise similar policies for the other models. Since the POMDP is a stochastic shortest path

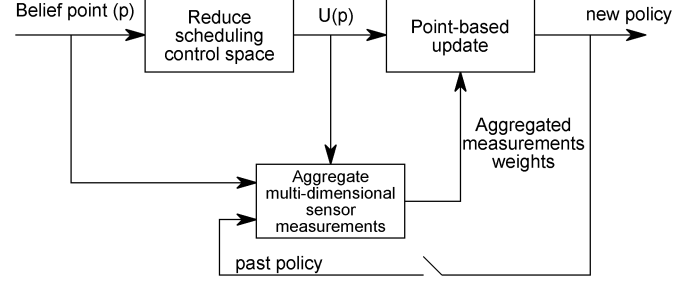


Fig. 1. Structure of the point-based scheduling approximation.

problem with an absorbing cost-free termination state, and the expected termination time is finite, the cost-to-go function for a given belief can be written as the minimum of the dot product of the belief vector and a set of hyperplanes ( $\alpha$  vectors)

$$\begin{aligned}
 J(\mathbf{p}) &= \min_{\{\alpha_i\}} \sum_b \alpha_i(b) \mathbf{p}(b) \\
 &= \min_{\mathbf{u} \in \{0,1\}^n} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) \right. \\
 &\quad + \sum_{s \in \{1, \dots, n, \varepsilon\}} \min_{\alpha_i} \sum_{b'} p(s|u, b') \\
 &\quad \left. \times \sum_b p(b'|b) \mathbf{p}(b) \alpha_i(b') \right\} \quad (13)
 \end{aligned}$$

where  $\{\alpha_i\}$  is the set of hyperplanes constituting the value function  $J$  and  $b'$  the future target state. In essence, the complexity of the Bellman equation (13) stems from the evolution of the belief  $\mathbf{p}_k$  in (8). We can see why (13) is hard to analyze if we further divide the second term in the summation into two terms depending on whether there is observability or there is an erasure,

$$\begin{aligned}
 J(\mathbf{p}) &= \min_{\mathbf{u} \in \{0,1\}^n} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) \right. \\
 &\quad + \sum_{b'} \mathbb{1}\{u_{b'}=1\} [\mathbf{p}P]_{b'} \min_{\{\alpha_i\}} \alpha_i(b') \\
 &\quad \left. + \min_{\{\alpha_i\}} \sum_{b'} \mathbb{1}\{u_{b'}=0\} [\mathbf{p}P]_{b'} \alpha_i(b') \right\}. \quad (14)
 \end{aligned}$$

To further clarify, we observe that

$$\begin{aligned}
 \sum_s p(s|\mathbf{u}, \mathbf{p}) J(\mathbf{p}_1) &= \sum_{i=1}^n \mathbb{1}\{u_i=1\} [\mathbf{p}P]_i J(\mathbf{e}_i) \\
 &\quad + \sum_{i=1}^n \mathbb{1}\{u_i=0\} [\mathbf{p}P]_i J([\mathbf{p}P]_{\{j:u_j=0\}}) \quad (15)
 \end{aligned}$$

and the minimization problem is coupled across the sensors as the second term in (15), which is due to nonobservability, depends on the action vector  $\mathbf{u}$ . The action of one sensor affects the belief evolution, therefore coupling the problem across sensors. Now, if we make the assumption that perfect observations would be available to the controller after taking a scheduling action, we obtain an approximate surrogate function, which can

be used to generate a suboptimal scheduling policy. Namely, we substitute  $p(s|u, b') = \delta(s - b')$  in (13). We get

$$J(\mathbf{p}) = \min_{\mathbf{u} \in \{0,1\}^n} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) + \sum_{b'} [\mathbf{p}P]_{b'} \min \alpha_i \cdot \mathbf{e}_{b'} \right\} \quad (16)$$

$$= \min_{\mathbf{u} \in \{0,1\}^n} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) + \sum_{b'} [\mathbf{p}P]_{b'} J(\mathbf{e}_{b'}) \right\}. \quad (17)$$

The terms in the summation in (17) only depend on the control action for each sensor. Furthermore, the belief evolution is independent of the scheduling action, wherefore the approximate recursion in (17) decomposes into separable terms, one per sensor. Hence, the value function and the scheduling policy for sensor  $\ell$ , under the observable-after-control assumption, can be obtained from the solution of per-sensor Bellman equation:

$$J^{(\ell)}(\mathbf{p}) = \min_{u_\ell \in \{0,1\}} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) + \sum_{b'} [\mathbf{p}P]_{b'} J^{(\ell)}(\mathbf{e}_{b'}) \right\}. \quad (18)$$

The POMDP problem is now decomposed into  $n$  separate simpler subproblems such that the total cost function is the sum of the per-sensor cost functions while the overall scheduling policy consists of the per-sensor policies applied in parallel. Each subproblem can be easily solved using standard policy iteration [19] with a simple minimization over a binary control action. Fundamentally, for the simplistic model, we are able to decompose the problem into  $n$  simpler subproblems due to the separability of the tracking cost into per-sensor costs. Note that the problem is still coupled due to the belief evolution in (8) yet that coupling is resolved under the observable-after-control assumption.

While separability holds for the simplistic model, this is not the case for the other models. Hence, we devise a strategy where we artificially decouple the problem into  $n$  simpler subproblems. To this end, we perform Monte Carlo simulations to determine appropriate values for the per-sensor tracking cost corresponding to each subproblem. For example, consider the continuous observation model of Section II-C. For simplicity of exposition, assume a discrete state space model with  $m$  possible object locations. In this case, we define a surrogate value function for the  $\ell$ th subproblem as follows:

$$J^\ell(\mathbf{p}) = \min_{\mathbf{u}} \left\{ \mathbb{1}\{u=0\} \sum_{i=1}^m \mathbf{p}(i) T(i, \ell) + \mathbb{1}\{u=1\} \sum_{i=1}^m c [\mathbf{p}P]_i + \sum_{i=1}^m [\mathbf{p}P]_i J^\ell(\mathbf{e}_i) \right\}, \quad \ell = 1, \dots, n \quad (19)$$

where  $T(i, \ell)$  captures the contribution of the  $\ell$ th sensor to the total tracking error when the target's previous state is  $i$  and is obtained via Monte Carlo simulations. Namely, the expected tracking cost can be evaluated by repeatedly simulating our system from time  $k-1$  to time  $k$  while changing the state of the  $\ell$ th sensor. Similarly, (19) can be generalized for continuous state spaces.

Even though the  $Q_{\text{MDP}}$  assumption leads to a separable problem and provides a lower bound on the optimal energy-tracking tradeoff for the simplistic model as we elaborate in Section III-D, the resulting scheduling policies are myopic, unlike the sleeping policies in [1]. This follows from the fact that under an observable-after-control assumption, the future cost term is independent of the control vector  $\mathbf{u}$ . Therefore, we consider more efficient, albeit more difficult, point-based approximations in the next section.

### C. Point-Based Approximate Policies

In the previous section, we described  $Q_{\text{MDP}}$  based policies, whereby issues i) and iii) in Remark III.1 are resolved since we only needed to solve the underlying Markov Decision Process to describe the full approximate surrogate value function. Decoupling the problem into one-per-sensor subproblems (naturally or artificially) further enabled us to address issue ii). Yet, we just argued in Sections III-A and III-B that the resulting scheduling policies are myopic and generally do not take control actions to gain information.

To that end, we develop point-based approximate scheduling policies. Instead of reducing complexity via artificial decoupling and learning, the key idea here is to optimize the value function only for specific reachable sampled beliefs and not over the entire belief simplex (addressing issue i) in Remark III.1). Such techniques have shown great potential for solving large scale POMDPs while significantly reducing complexity. Due to the large size of the control space, we also devise strategies to sample actions exploiting the sparsity of the beliefs and the problem structure (to address issue ii)). Moreover, observation aggregation is used for continuous observation models. Furthermore, since Perseus updates are not carried out for every sampled belief and multiple belief points are improved simultaneously, the number of  $\alpha$  vectors grows modestly with the number of iterations. This addresses issue iii) in Remark III.1.

For completeness, we first briefly outline the steps of Perseus and refer the reader to [17] and [18] for further details. Later, we discuss specific variations to the algorithm to address the dimensionality of the action and the observation spaces.

#### One iteration of Perseus

- 1) Sample a set of belief points  $\mathcal{P}$ . These beliefs are obtained by simulating the target motion through the field taking random actions and generating observation according to the observation models in (1), (3), (5), and (6).
- 2) Sample a belief point  $\mathbf{p} \in \mathcal{P}$  at random and compute the backup using (20a) and (20b),

$$\alpha = \arg \min_{\{\alpha_u^{\mathbf{p}}\}_{u \in \mathcal{U}}} \mathbf{p} \cdot \alpha_u^{\mathbf{p}} \quad (20a)$$

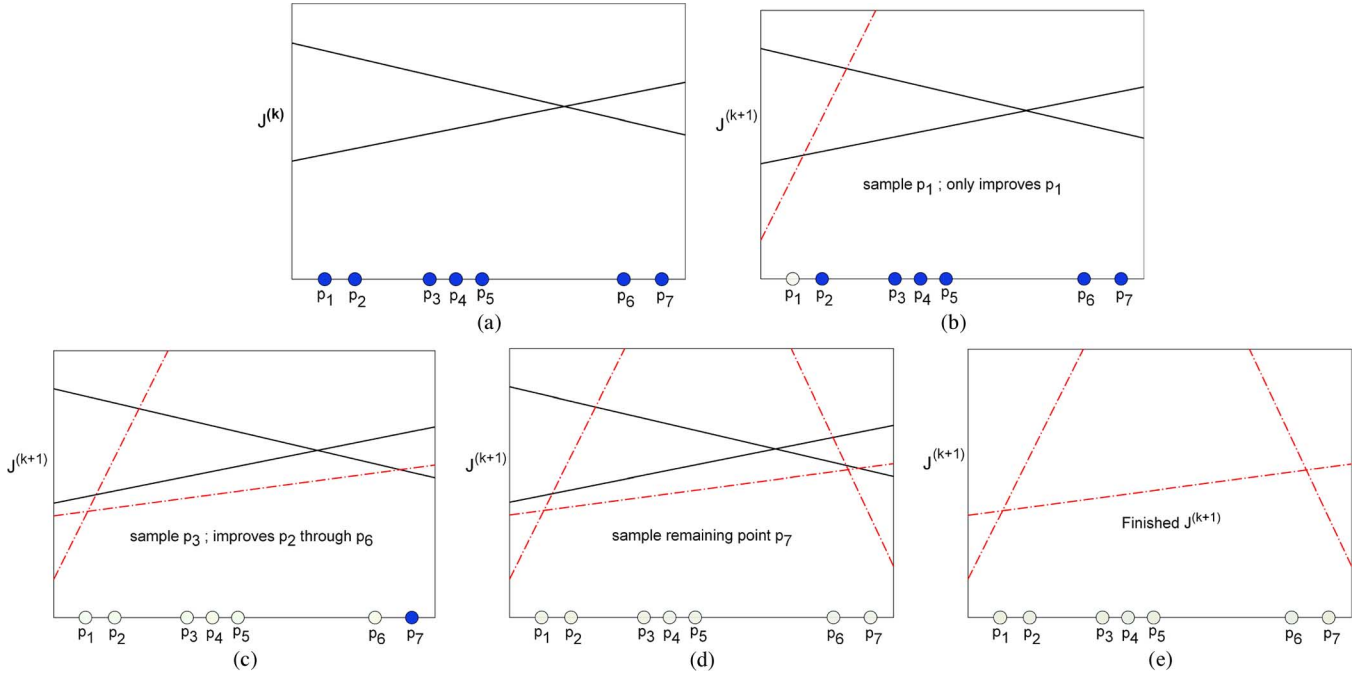


Fig. 2. One iteration of Perseus illustrating the progress of the algorithm. The  $x$  axis represents the belief space with circles representing the sampled belief set  $\mathcal{P} = \{p_1, \dots, p_7\}$ . The  $y$  axis is the value function at consecutive iterations, i.e.,  $J^{(k)}$  and  $J^{(k+1)}$ . Solid lines represent the hyperplanes in the  $k$ th iteration and dashed lines represent the newly added hyperplanes during the  $(k+1)$ th iteration. (a) The initial value function  $J^{(k)}$ ; (b)  $p_1$  is randomly selected and a new  $\alpha$  vector is added to  $J^{(k+1)}$ . This update step only happens to improve  $p_1$ . Dark circles represent belief points which did not yet improve; (c)  $p_3$  is sampled and a new hyperplane is added which improves the value for  $p_2$  through  $p_6$ ; (d) only  $p_7$  did not improve, thus  $p_7$  is sampled and a new hyperplane is added to  $J^{(k+1)}$ ; (e) all belief points improved,  $J^{(k+1)}$  is computed, the iteration ends.

where

$$\alpha_u^p = g(\mathbf{p}, \mathbf{u}) + \sum_s p(s|\mathbf{u}, \mathbf{p}) \min_{\alpha_i^{(k)}} \phi(\mathbf{p}, \mathbf{u}, s) \cdot \alpha_i^{(k)}. \quad (20b)$$

- 3) If  $\sum_b \mathbf{p}(b) \alpha(b) \leq J^{(k)}(\mathbf{p})$  then add new  $\alpha$  to  $J^{(k+1)}$  otherwise keep old hyperplane.
- 4) If  $\{\mathbf{p} \in \mathcal{P} : J^{(k+1)}(\mathbf{p}) > J^{(k)}(\mathbf{p})\} = \emptyset$ , i.e., the empty set, iteration is complete otherwise repeat from Step 1.

Fig. 2 illustrates the progress of one iteration of Perseus. The  $x$  axis represents the belief space with circles representing the sampled belief set  $\mathcal{P} = \{p_1, \dots, p_7\}$ . The  $y$  axis is the value function at consecutive iterations, i.e.,  $J^{(k)}$  (solid lines) and  $J^{(k+1)}$  (dashed lines). The figure displays the  $\alpha$  vectors and different steps illustrating the progress of the algorithm. The algorithm selects a belief point at random and updates the value function for that belief. Then a new update is carried out for a belief point randomly selected from the set of remaining beliefs, i.e., beliefs which did not improve in the previous step. The algorithm repeats till all belief points are updated. Solid lines represent the hyperplanes in the  $k$ th iteration and dashed lines represent the newly added hyperplanes during the  $(k+1)$ th iteration. In a way, the Perseus updates in POMDPs are the counterpart of asynchronous dynamic programming for MDPs [19] since the order of backup of the belief points is arbitrary and does not require full sweeps over the entire sampled belief set.

1) *Sampling Actions Based on the Support of the Belief:* Note that the update (20) involves a minimization over all control actions in  $|\mathcal{U}|$ . Even though one iteration of the algorithm is linear

in the cardinality  $|\mathcal{U}|$  of the control space,  $|\mathcal{U}|$  itself is exponential in the number of sensors, thus rendering the minimization infeasible for a relatively large sensor network.

The idea here is to exploit the structure of the scheduling/tracking problem. Since the target transition model is naturally sparse, we predict relatively small uncertainty regions for the target state at future time steps. More specifically, for every belief point in  $\mathcal{P}$ , we use prior information about the target transition model to project the future state of the target. This is particularly useful when the current belief vector is sparse leading to more restricted uncertainty regions. Subsequently, we restrict our attention to a *significant* subset of sensors, that is, the sensors of relevance to the particulars of the uncertainty region. Hence, we only consider scheduling actions involving scheduling different combinations of a reduced number of sensors which considerably reduces the control space for every belief in  $\mathcal{P}$ . If the number of significant sensors is still large, we randomly sample actions from the reduced control space. Note that the same intuition extends to more complex motion models wherein information about target speed, maneuver, and acceleration can be factored in to define the future uncertainty regions. Hence, instead of performing full updates including  $2^n$  actions, we perform the minimization over a reduced control space for every  $\mathbf{p} \in \mathcal{P}$ . Specifically, we redefine the point update equation as

$$\alpha = \arg \min_{\{\alpha_u^p\}_{u \in \mathcal{U}(\mathbf{p})}} \mathbf{p} \cdot \alpha_u^p \quad (21)$$

where  $\mathcal{U}(\mathbf{p})$  designates the reduced control space for the belief vector  $\mathbf{p}$ . Note that future iterations of the algorithm involving a particular belief point, ensure sufficient sampling of the relevant control actions in the reduced control space. This approach is



well suited to Perseus wherein the value for every belief point is guaranteed to improve over consecutive stages of the algorithm. It is worth mentioning that the observation and the cost models need to be computed on the fly for each sampled control action during the algorithm implementation.

2) *Observation Aggregation*: The point update (20) involves back-projecting all hyperplanes in the current iteration one step from the future and returning the vector that minimizes the value of the belief. Since this involves computing a cross sum by enumerating all possible combinations of alpha vectors for the different observations, a number of vectors which is exponential in the number of the observations is generated at each stage. The recursion has to be redefined to address continuous observation models. Looking carefully at (20), it is not hard to see that if different observations map to the same minimizing hyperplane, then they can be aggregated [26]. Hence, if we can partition the observation space into regions that map to the same hyperplane (possibly non contiguous), the continuous model is reduced to a corresponding discrete model. Integration is replaced by a summation over these partitions and the weighing probabilities are obtained by integrating the conditional density over these partitions. This is clarified in the following:

$$\begin{aligned} & \int_s \min_{\alpha_i} \sum_{b'} p(s|u, b') \sum_b p(b'|b) \mathbf{p}(b) \alpha_i(b') ds \\ &= \sum_j \int_{\mathcal{S}_j} \sum_{b'} p(s|u, b') \sum_b p(b'|b) \mathbf{p}(b) \alpha_j(b') ds \\ &= \sum_j \sum_{b'} [\mathbf{p}P]_{b'} \alpha_j(b') \int_{\mathcal{S}_j} p(s|u, b') ds \\ &= \sum_j \sum_{b'} [\mathbf{p}P]_{b'} \Pr[\mathcal{S}_j | \mathbf{u}, b'] \alpha_j(b'). \end{aligned} \quad (22)$$

To find the regions  $\mathcal{S}_j$  of aggregate observations, we need to solve for the boundaries, i.e., for each pair  $(i, j)$  of  $\alpha$  vectors we need to solve for  $\mathbf{s}$ :

$$\alpha_i \cdot \phi(\mathbf{p}, \mathbf{u}, \mathbf{s}) = \alpha_j \cdot \phi(\mathbf{p}, \mathbf{u}, \mathbf{s}) \quad (23)$$

where  $\phi(\mathbf{p}, \mathbf{u}, \mathbf{s}) = \mathbf{p}_1(b') \propto \sum_b \mathbf{p}(b) p(\mathbf{s}|b', \mathbf{u}) p(b'|b)$ . Hence, we need to solve:

$$\begin{aligned} & \sum_{b'} (\alpha_i(b') - \alpha_j(b')) [\mathbf{p}P]_{b'} \times \\ & \exp \left\{ -\frac{1}{2} \sum_{i: u_i=1} \left( s_i - \frac{V}{(b' - p_i)^2 + 1} \right)^2 \right\} = 0. \end{aligned} \quad (24)$$

After solving for the boundaries, we can readily define the regions:

$$\mathcal{S}_{j^*} = \left\{ \mathbf{s} | j^* = \arg \max_j \alpha_j \cdot \phi(\mathbf{p}, \mathbf{u}, \mathbf{s}) \right\}. \quad (25)$$

Now, the update step is simply

$$J(\mathbf{p}) = g(\mathbf{p}, \mathbf{u}^*) + \sum_j \sum_{b'} [\mathbf{p}P]_{b'} \Pr[\mathcal{S}_j | \mathbf{u}^*, b'] \alpha_j(b') \quad (26)$$

where  $\Pr[\mathcal{S}_j | \mathbf{u}^*, b'] = \int_{\mathcal{S}_j} p(\mathbf{s} | \mathbf{u}^*, b') d\mathbf{s}$ . Finding a closed form analytical solution for (24) is not feasible. Instead, we use Monte Carlo simulations to solve for the boundaries and get estimates of the weighing probabilities by sampling observations

from  $p(\mathbf{s} | \mathbf{u}, b')$  for different combinations of actions and target states.

#### D. Lower Bounds

We are able to derive lower bounds on the energy-tracking tradeoff for the simple as well as the continuous Gaussian observation models. For the simple model, the  $Q_{\text{MDP}}$  value function is itself a lower bound on the expected total cost since more information is available to the controller at future time steps given the reduced uncertainty assumption. To further clarify, observe that if we interchange the order of minimization and summation in the last term of (14), we obtain a lower bound on the optimal cost to go function. Hence, a lower bound can be obtained from the solution of the following equation:

$$\begin{aligned} J(\mathbf{p}) &= \min_{\mathbf{u} \in \{0,1\}^n} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) \right. \\ &\quad \left. + \sum_{b'} \mathbb{1}\{u_{b'}=1\} [\mathbf{p}P]_{b'} \min_{\{\alpha_i\}} \alpha_i(b') \right. \\ &\quad \left. + \sum_{b'} \mathbb{1}\{u_{b'}=0\} [\mathbf{p}P]_{b'} \min_{\{\alpha_i\}} \alpha_i(b') \right\} \\ &= \min_{\mathbf{u} \in \{0,1\}^n} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) \right. \\ &\quad \left. + \sum_{b'} [\mathbf{p}P]_{b'} \min_{\{\alpha_i\}} \alpha_i \cdot \mathbf{e}_{b'} \right\} \\ &= \min_{\mathbf{u} \in \{0,1\}^n} \left\{ \sum_{i=1}^n [\mathbf{p}P]_i \left( \mathbb{1}\{u_i=0\} + \sum_{\ell=1}^n c \mathbb{1}\{u_\ell=1\} \right) \right. \\ &\quad \left. + \sum_{b'} [\mathbf{p}P]_{b'} J(\mathbf{e}_{b'}) \right\}. \end{aligned} \quad (27)$$

Interchanging the order of the summation and minimization corresponds to a fully observable state after the next scheduling action, i.e., that the future belief is  $\mathbf{e}_{b'}$ . Hence, the  $Q_{\text{MDP}}$  value function is a lower bound on the cost function of the original problem.

Unfortunately, this is only true for the simplistic model and does not extend to the coupled models since the factored tracking cost in (19) need not be a lower bound on the true tracking cost. To obtain a lower bound on the optimal energy-tracking tradeoff for such models, we combine the observable-after-control assumption with a decomposable lower bound on the tracking cost which we derive next. Consider the continuous observation model with discrete state space. Given the current belief  $\mathbf{p}_k$  and a control vector  $\mathbf{u}_k$  the expected tracking cost can be written as

$$\begin{aligned} & E \left[ d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, \mathbf{u}_k \right] \\ &= \sum_{j=1}^m \Pr[\hat{b}_{k+1} \neq j | \mathbf{p}_k, \mathbf{u}_k, b_{k+1} = j] \Pr[b_{k+1} = j | \mathbf{p}_k, \mathbf{u}_k] \\ &= \sum_{i=1}^m \mathbf{p}_k(i) \sum_{j=1}^m p(b_{k+1} = j | b_k = i) \\ &\quad \times \Pr[\hat{b}_{k+1} \neq j | \mathbf{p}_k, \mathbf{u}_k, b_{k+1} = j]. \end{aligned} \quad (28)$$

Defining

$$P(E|H_j) \triangleq \Pr[\hat{b}_{k+1} \neq j | \mathbf{p}_k, \mathbf{u}_k, b_{k+1} = j]$$

which is a conditional error probability for a multiple hypothesis testing problem with  $m$  hypotheses, each corresponding to a different mean vector contaminated with white Gaussian noise. Conditioned on  $H_j$  the observation model is

$$H_j : \mathbf{s}(\ell) = (\mathbf{m}_j(\ell) + \mathbf{w}(\ell)) \mathbb{1}\{u_{k,\ell} = 1\} + \varepsilon \mathbb{1}\{u_{k,\ell} = 0\} \quad (29)$$

where  $\mathbf{s}(\ell)$  is the  $\ell$ th entry of an  $n \times 1$  vector  $\mathbf{s}$  denoting the received signal strength at the  $n$  sensors,  $\mathbf{m}_j$  is the mean received signal strength when the target is at state  $j$  ( $j$ th hypothesis), and  $\mathbf{w}$  is a zero mean white Gaussian Noise, i.e.,  $\mathbf{w} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ . According to (29), sensor  $\ell$  gets a Gaussian observation, which depends on the future target location, if activated at the next time step, and an erasure, otherwise. Since the current belief is  $\mathbf{p}_k$ , the prior for the  $j$ th hypothesis is  $\pi_j = [\mathbf{p}_k P]_j$ . The error event  $E$  can be written as the union of pairwise error regions as

$$p(E|H_j) = \Pr[\cup_{k \neq j} \zeta_{kj}] \quad (30)$$

where

$$\zeta_{kj} = \left\{ \mathbf{s} : L_{kj}(\mathbf{s}) > \frac{\pi_j}{\pi_k} \right\}$$

is the region of observations for which the  $k$ th hypothesis  $H_k$  is more likely than the  $j$ th hypothesis  $H_j$  and where

$$L_{kj} \triangleq \frac{f(\mathbf{s}|H_k)}{f(\mathbf{s}|H_j)}$$

denotes the likelihood ratio for  $H_k$  and  $H_j$ . Using standard analysis for likelihood ratio tests [27], [28], it is not difficult to show that

$$p(\zeta_{kj}|H_j) = Q\left(\frac{d_{kj}}{2} + \frac{\ln \frac{\pi_j}{\pi_k}}{d_{kj}}\right) \quad (31)$$

where,  $d_{kj}^2 = \frac{\Delta \mathbf{m}_{kj}^T \Delta \mathbf{m}_{kj}}{\sigma^2}$ ,  $\Delta \mathbf{m}_{kj} = \mathbf{m}_k - \mathbf{m}_j$ , and  $Q(\cdot)$  is the normal distribution  $Q$ -function. The quantity  $d_{kj}$  plays the role of distance between the two hypothesis and hence depends on the difference of their corresponding mean vectors and the noise variance  $\sigma^2$ . Note that, for different values of  $k$  and  $j$ ,  $\zeta_{kj}$  are not generally disjoint but allow us to lower bound the error probability in terms of pairwise error probabilities, namely, a lower bound can be written as

$$p(E|H_j) \geq \max_{k \neq j} p(\zeta_{kj}|H_j). \quad (32)$$

And we can readily lower bound the expected tracking error:

$$\begin{aligned} & E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, \mathbf{u}_k] \\ & \geq \sum_{i=1}^m \mathbf{p}_k(i) \sum_{j=1}^m p(b_{k+1} = j | b_k = i) \max_{k \neq j} p(\zeta_{kj}|H_j) \\ & = \sum_{i=1}^m \mathbf{p}_k(i) \sum_{j=1}^m p(b_{k+1} = j | b_k = i) \\ & \quad \times \max_{k \neq j} Q\left(\frac{d_{kj}}{2} + \frac{\ln \frac{\pi_j}{\pi_k}}{d_{kj}}\right). \end{aligned} \quad (33)$$

Next we separate out the effect of each sensor on the tracking error:

$$\begin{aligned} & E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, \mathbf{u}_k] \\ & \geq \mathbb{1}\{u_{k,\ell} = 1\} E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, \mathbf{u}_k = \mathbf{1}] \\ & \quad + \mathbb{1}\{u_{k,\ell} = 0\} E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, u_{k,i} = 1 \quad \forall i \neq \ell] \\ & \quad \text{for every } \ell \end{aligned} \quad (34)$$

where  $\mathbf{1}$  is the vector of all ones designating that all sensors will be active at the next time slot. The inequality in (34) follows from the fact that if we separate out the effect of the  $\ell$ th sensor, we get a better tracking performance when all the remaining sensors are awake. Since this holds for every  $\ell$ , a lower bound on the expected tracking error can be written as a convex combination of all sensors' contributions:

$$\begin{aligned} & E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, \mathbf{u}_k] \\ & \geq \sum_{\ell=1}^n \lambda_\ell(\mathbf{p}_k) \left\{ \mathbb{1}\{u_{k,\ell} = 1\} E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, \mathbf{u}_k = \mathbf{1}] \right. \\ & \quad \left. + \mathbb{1}\{u_{k,\ell} = 0\} E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, u_{k,i} = 1 \quad \forall i \neq \ell] \right\} \end{aligned} \quad (35)$$

where  $\sum_{\ell} \lambda_\ell(\mathbf{p}_k) = 1$ .

Let  $\mathbf{1}_{-\ell}$  denote a vector of length  $n$  with all entries equal to one except for the  $\ell$ th entry being zero. Then replacing from (33),

$$\begin{aligned} & E[d(\hat{b}_{k+1}, b_{k+1}) | \mathbf{p}_k, \mathbf{u}_k] \\ & \geq \sum_{\ell=1}^n \lambda_\ell(\mathbf{p}_k) \left\{ \mathbb{1}\{u_{k,\ell} = 1\} \sum_{i=1}^m \mathbf{p}_k(i) \sum_{j=1}^m p(b_{k+1} = j | b_k = i) \right. \\ & \quad \times \max_{k \neq j} Q\left(\frac{d_{kj}(\mathbf{1})}{2} + \frac{\ln \frac{\pi_j}{\pi_k}}{d_{kj}(\mathbf{1})}\right) + \mathbb{1}\{u_{k,\ell} = 0\} \\ & \quad \times \sum_{i=1}^m \mathbf{p}_k(i) \sum_{j=1}^m p(b_{k+1} = j | b_k = i) \\ & \quad \times \max_{k \neq j} Q\left(\frac{d_{kj}(\mathbf{1}_{-\ell})}{2} + \frac{\ln \frac{\pi_j}{\pi_k}}{d_{kj}(\mathbf{1}_{-\ell})}\right) \left. \right\}. \end{aligned} \quad (36)$$

To simplify notation, we define the following two quantities:

$$\begin{aligned} T_1(\mathbf{p}; i, \ell) & \triangleq \sum_{j=1}^m p(b_{k+1} = j | b_k = i) \\ & \quad \times \max_{k \neq j} Q\left(\frac{d_{kj}(\mathbf{1})}{2} + \frac{\ln \frac{[\mathbf{p}P]_j}{[\mathbf{p}P]_k}}{d_{kj}(\mathbf{1})}\right) \\ T(\mathbf{p}; i, \ell) & \triangleq \sum_{j=1}^m p(b_{k+1} = j | b_k = i) \\ & \quad \times \max_{k \neq j} Q\left(\frac{d_{kj}(\mathbf{1}_{-\ell})}{2} + \frac{\ln \frac{[\mathbf{p}P]_j}{[\mathbf{p}P]_k}}{d_{kj}(\mathbf{1}_{-\ell})}\right). \end{aligned}$$

Intuitively,  $T_1(\mathbf{p}; i, \ell)$  represents the contribution of sensor  $\ell$  to the total expected tracking cost when the underlying state is  $i$ , the belief is  $\mathbf{p}$ , and when all sensors are awake. On the other hand,  $T(\mathbf{p}; i, \ell)$  is the  $\ell$ th sensor contribution when it is inactive and all the other sensors are awake.

Now if we assume that the target will be perfectly observable after taking the scheduling action, a lower bound on the total cost can be readily obtained from the solution of the following Bellman equation:

$$J(\mathbf{p}) = \sum_{\ell} J^{(\ell)}(\mathbf{p}) \quad (37)$$

where

$$\begin{aligned} J^{(\ell)}(\mathbf{p}) &= \min_{u_{\ell} \in \{0,1\}} \left( \mathbb{1}\{u_{\ell}=1\} \left( \sum_b \mathbf{p}(b) \lambda_{\ell} T_1(\mathbf{p}; b, \ell) + c \sum_{i=1}^m [\mathbf{p}P]_i \right) \right. \\ &\quad \left. + \mathbb{1}\{u_{\ell}=0\} \sum_b \mathbf{p}(b) \lambda_{\ell} T(\mathbf{p}; b, \ell) \right. \\ &\quad \left. + \sum_{i=1}^m [\mathbf{p}P]_i J^{(\ell)}(\mathbf{e}_i) \right). \end{aligned} \quad (38)$$

Note that if we can solve the equation above for  $\mathbf{p} = \mathbf{e}_i$  for all  $i \in \{1, \dots, m\}$ , then it is straightforward to find the solution for all other values of  $\mathbf{p}$ . We therefore focus on specifying the value function at those points. Since this is the case, we further simplify our notation and use  $T(i, \ell)$  and  $\lambda(i, \ell)$  as shorthand for  $T(\mathbf{e}_i; i, \ell)$  and  $\lambda_{\ell}(\mathbf{e}_i)$ , respectively. We can see that a lower bound on the value function of sensor  $\ell$  can be obtained as a solution to the following minimization problem over  $u$ :

$$\begin{aligned} J^{(\ell)}(\mathbf{e}_b) &= \min \left\{ \lambda(b, \ell) T(b, \ell); \lambda(b, \ell) T_1(b, \ell) + c \sum_{i=1}^m [\mathbf{e}_b P]_i \right\} \\ &\quad + \sum_{i=1}^m [\mathbf{e}_b P]_i J^{(\ell)}(\mathbf{e}_i). \end{aligned} \quad (39)$$

Equation (39) together with (37) define a lower bound on the total expected cost. To further tighten the bound, we can now optimize over a matrix  $\Lambda$  for every value of  $c$ , where  $\Lambda(c)$  is an  $m \times n$  matrix with the  $(i, \ell)$  entry equal to  $\lambda(i, \ell)$ , i.e.,  $\Lambda(c) = \{\lambda(i, \ell)\}$ . Hence

$$\begin{aligned} J(\mathbf{e}_b) &= \max_{\Lambda(c)} \sum_{\ell=1}^n \left( \min \left\{ \lambda(b, \ell) T(b, \ell); \lambda(b, \ell) T_1(b, \ell) \right. \right. \\ &\quad \left. \left. + c \sum_{i=1}^m [\mathbf{e}_b P]_i \right\} + \sum_{i=1}^m [\mathbf{e}_b P]_i J^{(\ell)}(\mathbf{e}_i) \right) \\ &\text{subject to } \Lambda \mathbf{1}_n = \mathbf{1}_m \end{aligned} \quad (40)$$

where  $\mathbf{1}_m$  is a column vector of all ones of length  $m$ . The inner recursion can be solved to obtain a closed form solution for  $J^{(\ell)}(\mathbf{e}_b)$  as

$$\begin{aligned} J^{(\ell)}(\mathbf{e}_b) &= \sum_{j=0}^{\infty} \sum_{i=1}^m \min \left\{ [\mathbf{e}_b P^j]_i \lambda(i, \ell) T_1(i, \ell) \right. \\ &\quad \left. + c \sum_{k=1}^m [\mathbf{e}_b P^{j+1}]_k; [\mathbf{e}_b P^j]_i \lambda(i, \ell) T(i, \ell) \right\}. \end{aligned} \quad (41)$$

Since the problem is only constrained across the different sensors, we obtain a lower bound from the solution of the following optimization problem:

$$\begin{aligned} J &= \sum_{i=1}^m \max_{\lambda(i, \ell)} \sum_{\ell=1}^n \sum_{j=0}^{\infty} [\mathbf{e}_b P^j]_i \\ &\quad \times \min \left( \lambda(i, \ell) T_1(i, \ell) + c \sum_{k=1}^m [\mathbf{e}_b P]_k; \lambda(i, \ell) T(i, \ell) \right) \\ &\text{subject to } \sum_{\ell=1}^n \lambda(i, \ell) = 1 \quad \forall i = 1, \dots, m. \end{aligned} \quad (42)$$

We observe that for every  $i$  we are maximizing a concave piecewise linear function in  $\lambda(i, \ell)$ . We pose an equivalent convex optimization problem by realizing that the minimum of a set of concave functions is also concave. Since affine functions are concave, we can apply the technique here. Since the problem is unconstrained across the  $i$  dimension we focus on solving the max-min problem for a fixed  $i$ . The final solution can then be obtained by summing the objective function for  $m$  subproblems. For each  $\ell = 1, \dots, n$ , add a variable  $t_{\ell}$  to the optimization problem. Also for every  $\ell$ , append two constraints to the optimization problem. The constraints state the minimization over  $u_{\ell}$  implicitly, by requiring that  $\lambda(i, \ell) T_1(i, \ell) + c \sum_{k=1}^m [\mathbf{e}_b P]_k \geq t_{\ell}$  and  $\lambda(i, \ell) T(i, \ell) \geq t_{\ell}$ . The modified problem is therefore:

$$\begin{aligned} &\text{maximize}_{\lambda(i, \ell), t_{\ell}; \ell=1, \dots, n} \sum_{\ell=1}^n t_{\ell}, \\ &\text{subject to } \sum_{\ell=1}^n \lambda(i, \ell) \leq 1, \\ &\quad \lambda(i, \ell) T_1(i, \ell) + c \sum_{k=1}^m [\mathbf{e}_b P]_k \geq t_{\ell}, \\ &\quad \lambda(i, \ell) T(i, \ell) \geq t_{\ell}, \ell = 1, \dots, n \end{aligned} \quad (43)$$

which can be readily solved using standard convex optimization techniques [29]. Note that the Gaussian assumption was merely to find the expected tracking error in closed form. Our general approach could very well extend to non-Gaussian cases if we are given some lower bound on the expected tracking cost. The main ingredient of our approach, which leads to separability based on an all-awake assumption, remains valid.

#### IV. RESULTS AND SIMULATIONS

In this section, we show experimental results illustrating the performance of the proposed scheduling policies for the different models considered in this paper. In each simulation run, the object was initially placed at the center of the network and the simulation run concluded when the object reached the absorbing state  $\tau$ . We perform Monte Carlo runs to compute the average tracking and energy costs for different values of the energy parameter  $c$ . For the planning phase in case of point-based policies, beliefs are sampled by simulating multiple object trajectories through the sensor network. Each trajectory starts from a random state sampled from the initial belief, picking actions at random, until the target leaves the network. In our

TABLE I  
OBJECT MOVEMENT FOR A NETWORK OF 41 SENSORS WITH  
SIMPLE COST AND SENSING MODELS

Change in Position	-3	-2	-1	0	1	2	3
Probability	0.23	0.10	0.01	0.33	0.06	0.05	0.22

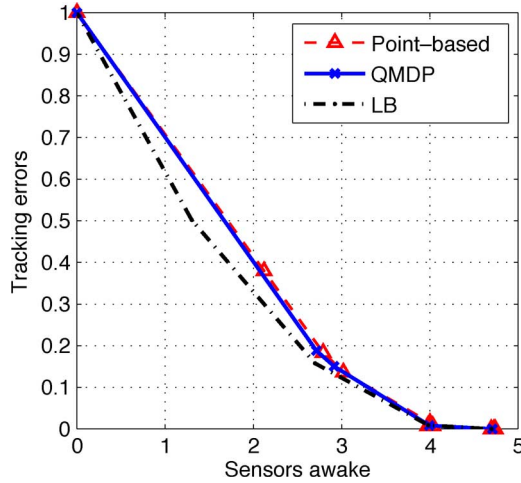


Fig. 3. Energy-tracking tradeoff for a one-dimensional network of 41 sensors with the simplistic sensing and cost model in Section II-A.

simulations backups are performed until the difference between successive value function estimates is below a small threshold. Other alternatives for convergence criteria could also be used such as tracking the number of policy changes between consecutive backups.

First, we consider the simple model in Section II-A with a linear network of 41 sensors. The object can move anywhere from three steps to the left to three steps to the right in each time step. The distribution for these movements is given in Table I. The change in position indicates movement by a corresponding number of steps to the right or to the left. Fig. 3 shows the tradeoff curve between the number of active sensors per unit time and the tracking error per unit time using the point-based and the  $Q_{MDP}$  policies. The figure also shows a lower bound on the optimal performance (see Section III-D). It is clear that both policies lead to tradeoffs that closely approach the lower bound. The  $Q_{MDP}$  policy gets even closer to the lower bound at small tracking errors since the observable-after-control assumption is more meaningful in this regime. In Fig. 4, we show convergence results for the point-based algorithm with reduced control space minimization. The top left subplot displays the convergence of the sum cost of all the belief points in  $\mathcal{P}$ ; the top right shows the expected cost averaged over many trajectories; the bottom left subplot shows the number of hyperplanes constituting the value function as a function of time; the bottom right subplot shows the number of policy changes versus time, i.e., the number of belief points for which the optimal action changed over two consecutive iterations of the algorithm.

Fig. 6 displays the average cost and the tradeoff curves for the network in Fig. 5 with a probabilistic observation model. The cost per unit time is the average ratio of the total energy plus tracking cost and the time the object spends in the network before reaching the termination state. The network is composed of 12 sensors and 20 object locations with the shown connectivity

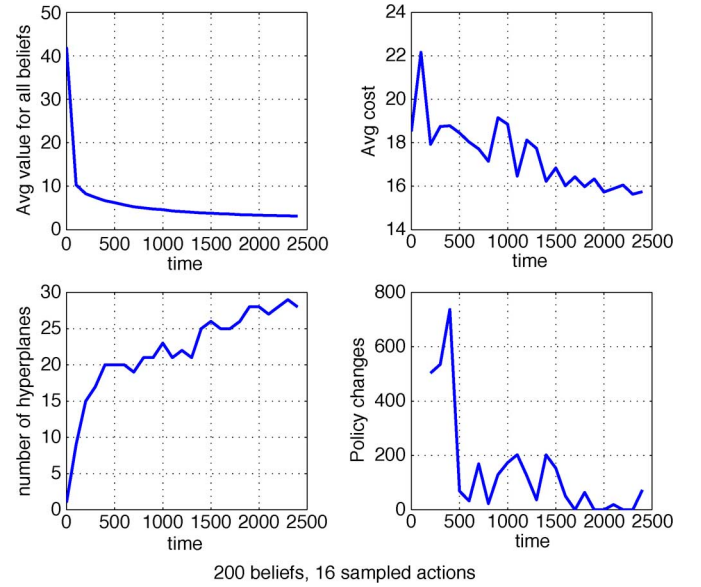


Fig. 4. Convergence results for the point-based algorithm for a one-dimensional network of 41 sensors with the simplistic sensing and cost model in Section II-A.

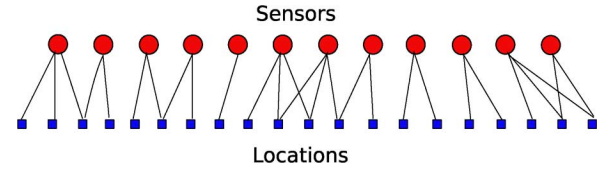


Fig. 5. A sensor network with overlapping sensing ranges (12 sensors and 20 object locations). An edge connects a sensor to a given location if this location falls within the sensing range of that sensor.

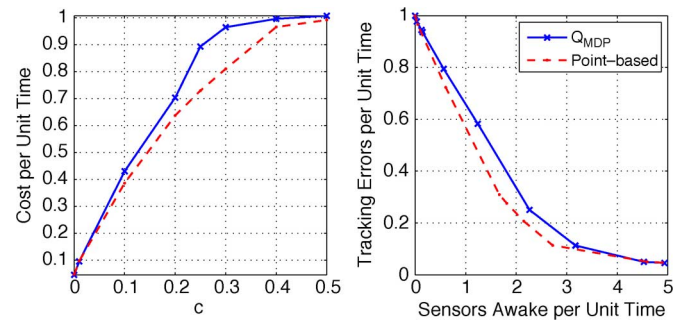


Fig. 6. Overlap model.

TABLE II  
OBJECT MOVEMENT FOR THE OVERLAP NETWORK AND CONTINUOUS

Change in Position	0	1	2	3
Probability	0.3125	0.2344	0.0938	0.0156

such that the observation range for the different sensors overlap. The object moves according to a random walk anywhere from three steps to the left to three steps to the right in each time step. The distribution of these movements is given in Table II. For the locations close to the boundaries, i.e., when less than three steps are available on the right or left, the remaining probability is absorbed in the transition to the termination state. Since the tracking error for this model is inherently coupled across

TABLE III  
 SENSOR LOCATIONS FOR NETWORK B

Sensor	1	2	3	4	5	6	7	8	9	10
Location	1.36	1.61	3.91	8.09	11.96	13.39	13.52	13.66	16.60	18.68

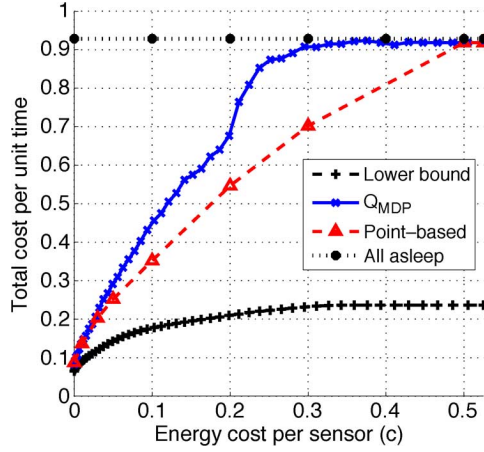


Fig. 7. Continuous observation model: total cost versus energy cost per sensor.

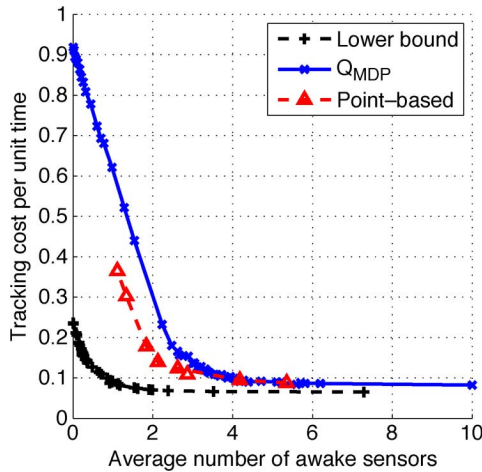


Fig. 8. Continuous observation model: energy-tracking tradeoff.

sensors, the global point-based policy clearly outperforms the learning-based  $Q_{MDP}$  policy.

Next, we consider a network of 10 sensors where object locations are located on integers from 1 to 21. The observation for each awake sensor is continuous and Gaussian as in (6) with  $V = 10$ . The locations of the sensors are given in Table III and the object moves according to the random walk defined in Table II. For every object state and every scheduling action in the reduced control space, we sample 50 observations to construct estimates of the weight probabilities and compute the aggregate observation boundaries. Up to 32 actions are sampled from the reduced control space. In this setup, the belief set consists of 500 sampled belief vectors and we assume a Hamming error cost. Figs. 7 and 8 show the performance of the different policies for the continuous observation model. It is shown that the point-based scheduling policy outperforms the  $Q_{MDP}$  policy. We further show a lower bound on the optimal performance

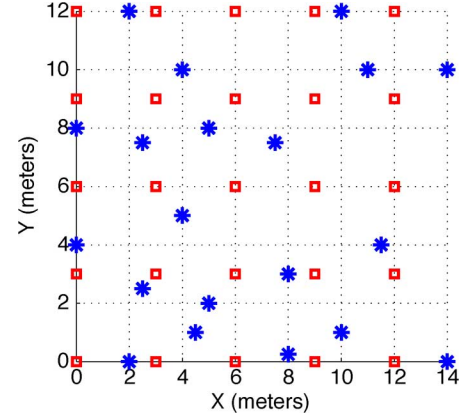
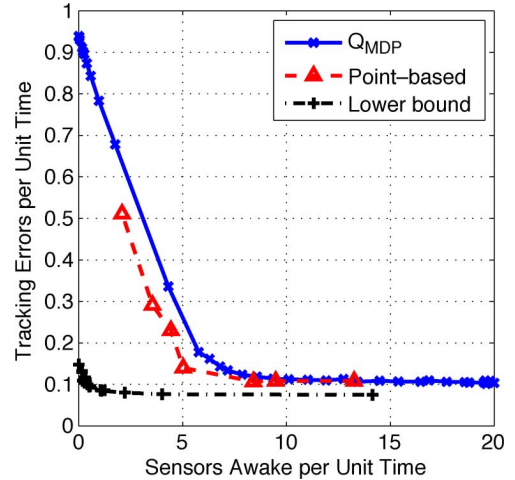


Fig. 9. 2-D network with 20 sensors (stars) and 25 possible object locations (squares).


 Fig. 10. Energy-tracking tradeoff of the  $Q_{MDP}$  and point-based scheduling policies for a 2-D network with continuous observations and Hamming cost.

tradeoff. The lower bound is loose especially in the high tracking error regime since the derived bound on per-sensor tracking errors assumes all other sensors are awake. However, we can exactly compute the saturation point for the optimal scheduling policy since every policy has to eventually meet the all-asleep performance curve, shown in Fig. 7, when the energy cost per sensor is high. At that point, all sensors are inactive and hence the target estimate can only be based on prior information. Our results are not restricted to 1-D networks but easily apply to 2-D networks. Namely, Fig. 10 shows the energy-tracking tradeoff of the  $Q_{MDP}$  and point-based policies in addition to a lower bound on optimal performance for the 2-D network of Fig. 9 with continuous observations and Hamming cost. The entries of the object transition matrix are generated randomly with the restriction that at any state the object can only move to its neighboring locations or remain at its current state. This simulation shows similar

trends to the previously observed results. The point-based policies outperform the  $Q_{MDP}$  approach at the expense of an increase in the offline computational complexity of the planning phase. Furthermore, the lower bound is reasonably tight in the low tracking error regime.

## V. CONCLUSIONS

In this paper we studied the problem of tracking an object moving randomly through a dense network of wireless sensors. We devised approximate strategies for scheduling the sensors to optimize the tradeoff between tracking performance and energy consumption for a wide range of models. First, we proposed policies that rely on an observable-after-control assumption ( $Q_{MDP}$  policies). Key to this solution is the decoupling of the optimization problem into per-sensor subproblems combined with simulation-based learning of individual tracking costs for each subproblem. Second, we developed point-based sensor scheduling strategies, which optimize the value function over a small set of reachable beliefs within the belief simplex. Based on the belief support and the sparsity of the transition models, we developed a methodology to sample actions from reduced control spaces. This was combined with observation aggregation to address the complexity of the observation space for continuous observations models. In some cases we derived lower bounds on the optimal tradeoff curves. While being suboptimal, the generated scheduling policies often provide close-to-optimal energy-tracking tradeoffs. Developing distributed scheduling strategies when no central controller is available is an area for future research. Another interesting challenge is when the statistics for object movement are unknown or partially known.

## REFERENCES

- [1] J. A. Fuemmeler and V. V. Veeravalli, "Smart sleeping policies for energy efficient tracking in sensor networks," *IEEE Trans. Signal Process.*, vol. 56, no. 5, pp. 2091–2101, May 2008.
- [2] D. A. Castanon, "Approximate dynamic programming for sensor management," in *Proc. 36th IEEE Conf. Decision Control*, 1997, vol. 2, pp. 1202–1207.
- [3] J. L. Williams, J. W. Fisher, and A. S. Willsky, "Approximate dynamic programming for communication-constrained sensor network management," *IEEE Trans. Signal Process.*, vol. 55, no. 8, pp. 4300–4311, Aug. 2007.
- [4] C. Kreucher, K. Kastella, and A. Hero, "Sensor management using an active sensing approach," *IEEE Trans. Signal Process.*, vol. 85, no. 3, pp. 607–624, Mar. 2005.
- [5] C. Kreucher, K. Kastella, and A. Hero, "A bayesian method for integrated multitarget tracking and sensor management," in *Proc. 6th Int. Conf. Inf. Fusion*, 2003, pp. 704–711.
- [6] A. Logothetis and A. Isaksson, "On sensor scheduling via information theoretic criteria," in *Proc. Amer. Control Conf.*, San Diego, CA, 1999, pp. 2402–2406.
- [7] J. Liu, Reich, and F. Zhao, "Collaborative in-network processing for target tracking," *EURASIP J. Appl. Signal Process.*, vol. 4, pp. 378–391, Mar. 2003.
- [8] Y. He and E. K. Chong, "Sensor scheduling for target tracking: A Monte Carlo sampling approach," *Digit. Signal Process. (Special Issue on DASP 2005)*, vol. 16, no. 5, pp. 533–545, Sep. 2006 [Online]. Available: <http://www.sciencedirect.com/science/article/B6WDJ-4FS8943-1/2/5d9aa547ea82e8874c16a8f29cb55936>
- [9] F. Zhao, J. Shin, and J. Reich, "Information-driven dynamic sensor collaboration," *IEEE Signal Process. Mag.*, vol. 19, no. 3, pp. 61–72, Mar. 2002.
- [10] W. Xiao, L. Xie, J. Lin, and J. Li, "Multi-sensor scheduling for reliable target tracking in wireless sensor networks," in *Proc. ITST Int. Conf. ITS Telecommun.*, 2007, pp. 996–1000.
- [11] M. Kalandros and L. Pao, "Covariance control for multisensor systems," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 38, no. 4, pp. 1138–1157, Oct. 2002.
- [12] A. S. Chhetri, D. Morrell, and A. Papandreou-Suppappola, "Non-myopic sensor scheduling and its efficient implementation for target tracking applications," *EURASIP J. Appl. Signal Process.*, vol. 2006, Jan. 2006 [Online]. Available: <http://dx.doi.org/10.1155/ASP/2006/31520>
- [13] G. Monahan, "A survey of partially observable Markov decision processes: Theory, models and algorithms," *Manag. Sci.*, vol. 28, pp. 1–16, Jan. 1982.
- [14] M. Hauskrecht, "Value-function approximations for partially observable Markov decision processes," *J. Artif. Intell. Res. (JAIR)*, vol. 13, pp. 33–94, 2000.
- [15] J. Pineau, G. Gordon, and S. Thrun, "Point-based value iteration: An anytime algorithm for POMDPs," in *Proc. Int. Joint Conf. Artif. Intell. (IJCAI)*, 2003, pp. 1025–1032.
- [16] M. L. Littman, A. R. Cassandra, and L. P. Kaelbling, "Learning policies for partially observable environments: Scaling up," in *Proc. 12th Int. Conf. Mach. Learn.*, 1995, pp. 362–370.
- [17] M. T. J. Spaan and N. Vlassis, "Perseus: Randomized point-based value iteration for POMDPs," *J. Artif. Intell. Res. (JAIR)*, vol. 24, pp. 195–220, 2005.
- [18] J. M. Porta, N. Vlassis, M. T. J. Spaan, and P. Poupart, "Point-based value iteration for continuous POMDPs," *J. Mach. Learning Research*, vol. 7, pp. 2329–2367, 2006.
- [19] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, 3rd ed. Belmont, MA: Athena Scientific, 2007.
- [20] E. J. Sondik, "The Optimal control of partially observable Markov processes," Ph.D. dissertation, Stanford Univ., Stanford, CA, 1971.
- [21] H. T. Cheng, "Algorithms for partially observable Markov decision processes," Ph.D. dissertation, Univ. of British Columbia, Vancouver, BC, Canada, 1988.
- [22] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra, "Planning and acting in partially observable stochastic domains," *Artif. Intell.*, vol. 101, pp. 99–134, May 1998.
- [23] A. Cassandra, M. L. Littman, and N. L. Zhang, "Incremental pruning: A simple, fast, exact algorithm for partially observable Markov decision processes," in *Proc. 13th Annu. Conf. Uncertainty Artif. Intell.*, 1997, pp. 54–61, Morgan Kaufmann.
- [24] L. P. Kaelbling, M. L. Littman, and A. Moore, "Reinforcement learning: A survey," *J. Artif. Intell. Res.*, vol. 4, pp. 237–285, 1996.
- [25] N. Roy and G. Gordon, "Exponential family PCA for belief compression in POMDPs," *Adv. Neural Inf. Process. Syst.*, vol. 15, 1995.
- [26] J. Hoey and P. Poupart, "Solving POMDPs with continuous or large discrete observation spaces," in *Proc. 19th Int. Joint Conf. Artif. Intell. (IJCAI)*, San Francisco, CA, 2005, pp. 1332–1338, Morgan Kaufman.
- [27] H. V. Poor, *An Introduction to Signal Detection and Estimation*, 2nd ed. New York: Springer-Verlag, Feb. 1994.
- [28] B. C. Levy, *Principles of Signal Detection and Parameter Estimation*. Boston, MA: Springer, 2008.
- [29] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, Mar. 2004.



**George K. Atia** (S'01–M'04) received the B.Sc. and M.Sc. degrees from Alexandria University, Egypt, in 2000 and 2003, respectively, and the Ph.D. degree from Boston University, MA, in 2009, all in electrical and computer engineering.

He joined the Department of Electrical and Computer Engineering at the University of Illinois at Urbana-Champaign in fall 2009, where he is currently a Postdoctoral Research Associate in the Coordinated Science Laboratory. His research interests include wireless communications, statistical signal processing and information theory.

Dr. Atia is the recipient of many awards, including the Outstanding Graduate Teaching Fellow of the Year Award in 2003–2004 from the Electrical and Computer Engineering Department at Boston University, the 2006 College of Engineering Dean's Award at the BU Science and Engineering Research Symposium, and the Best Paper Award at the International Conference on Distributed Computing in Sensor Systems (DCOSS) in 2008.



**Venugopal V. Veeravalli** (M'92–SM'98–F'06) received the B.Tech. degree (Silver Medal Hons.) from the Indian Institute of Technology, Bombay, in 1985, the M.S. degree from Carnegie Mellon University, Pittsburgh, PA, in 1987, and the Ph.D. degree from the University of Illinois at Urbana-Champaign, in 1992, all in electrical engineering.

He joined the University of Illinois at Urbana-Champaign in 2000, where he is currently a Professor in the Department of Electrical and Computer Engineering, and a Research Professor in the Coordinated Science Laboratory. He served as a Program Director for communications research at the U.S. National Science Foundation in Arlington, VA from 2003 to 2005. He has previously held academic positions at Harvard University, Rice University, and Cornell University, and has been on sabbatical at MIT, IISc Bangalore, and Qualcomm, Inc. His research interests include distributed sensor systems and networks, wireless communications, detection and estimation theory, and information theory.

Prof. Veeravalli is a Distinguished Lecturer for the IEEE Signal Processing Society for 2010–2011. He has been on the Board of Governors of the IEEE Information Theory Society. He has been an Associate Editor for Detection and Estimation for the IEEE TRANSACTIONS ON INFORMATION THEORY and for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS. Among the awards he has received for research and teaching are the IEEE Browder J. Thompson Best Paper Award, the National Science Foundation CAREER Award, and the Presidential Early Career Award for Scientists and Engineers (PECASE).



**Jason A. Fuemmeler** (S'97–M'00) received the B.E.E. degree in electrical engineering from the University of Dayton, Dayton, OH, in 2000 and the M.S. and Ph.D. degrees in electrical engineering from the University of Illinois at Urbana-Champaign in 2004 and 2008, respectively.

He has been awarded an NSF graduate research fellowship and a Vodafone fellowship. He is currently employed in the Advanced Technology Center, Rockwell Collins, performing research in electronic warfare and wireless communications.