# Adaptive Search for Sparse Targets with Informative Priors

Gregory Newstadt, Eran Bashan, Alfred O. Hero III

## I. ABSTRACT

This works considers the problem of efficient energy allocation of resources in a continuous fashion in determining the location of targets in a sparse environment. We extend the work of Bashan [1] to analyze the use of non-uniform prior knowledge for the location of targets. We show that in the best-case scenario (i.e., when the known prior knowledge is also the underlying prior), then we can get significant gains (several dB) by using a two-level piecewise uniform prior over using the uniform prior that is assumed in [1]. Moreover, even when we have uncertainty in our prior knowledge, we show that we can always do at least as well as the uniform alternative in terms of worst-case and expected gains. In future work, we plan to extend our analysis to general piecewise uniform priors in order to develop multistage (i.e., greater than 2) adaptive energy allocation policies.

## II. INTRODUCTION

In many situations, it might be desirable to allocate a limited amount of energy to a small region of interest (ROI) within a larger environment by using adaptive sampling techniques. For example, consider the problem of minimizing communication costs when tracking a target in a distributed sensor network. Clearly, when a node in our sensor network is far from our previous estimate of the target, we would like to reduce the communication from that node in order to preserve its battery life. On the other hand, we would like all sensors within a region near the previous target estimate to be used to estimate the target's position at the next time step. For this purpose, then, we can use adaptive sampling to provide both smart sensor management and improved tracking performance.

In another application, we may be interested in locating and estimating the content of tumors in early cancer detection. In this situation, the ROI consists of the tumor location, which in early cancer detection is assumed to be much smaller in area than the entire image. Moreover, the total amount of energy used

by CT scans and X-rays for this purpose is limited by safety constraints. Thus, an adaptive sampling scheme could be used to allocate energy efficiently only to the regions where a tumor may exist.

Lastly, we may be interested in the detection and estimation of airplanes in an airport landing field using active radar systems. We assume that airplanes are much more likely to approach from one direction than some others. In this case, it would be desirable to search for airplanes in an optimal manner that reduces the amount of energy spent searching the low probability regions and augments the amount of energy spent searching the high probability regions.

In this work, we consider the problem of estimating a sparse ROI where we may have prior knowledge for the locations of the targets. Bashan [1] showed that under a total energy constraint, an adaptive resource allocation policy (ARAP) can be used to form a two-stage energy allocation policy that optimally allocates energy according to a uniform prior with respect to a suitable cost function. ARAP was developed for a general prior, but all of the previous analysis focused on the uniform case (i.e., targets distributed uniformly across the entire signal), where analysis was straightforward.
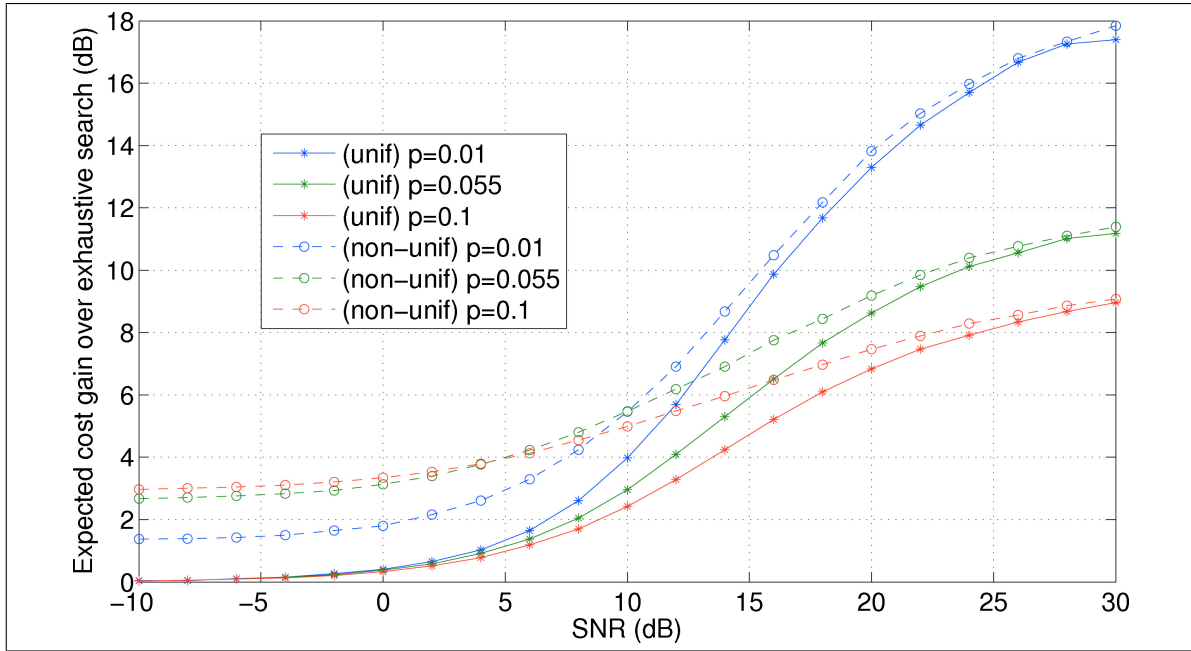


Fig. 1.   This plot compares the performance of energy policies using uniform prior information for the locations of targets (circles) versus the performance of energy policies using non-uniform prior information (stars). For low SNR, we see enhanced performance for all sparsity levels when using the non-uniform prior information. For high SNR, we see that the performance of either adaptive policy converges to $-10 \log p$, (where $p$ is the sparsity level), the predicted asymptotic gain.

This work, then, focuses on relaxing the assumption of a uniform prior in detecting/estimating the ROI. Moreover, we are interested in developing intuition for when to use our prior knowledge as opposed to

just assuming a uniform prior (i.e., how much margin for error do we have). Figure 1 shows the gain in expected cost over using an exhaustive search as a function of SNR. The circles represent the gains when using a uniform prior, while the stars represent the gains when using a two-level piecewise uniform prior. We see that for low SNR values, we can get significant gains when using non-uniform prior knowledge.

It should be noted that much of the previous work in resource allocation has been directed towards inhomogeneous signals [2] [3] [4] [5], though in this work we consider signals with a small ROI. Thus, our signals could be considered sparse in the sense that the ratio of the ROI to the entire signal is small. In compressed sensing [6] [7], the goal is also to reconstruct a sparse signal with reduced energy. Resource allocation and sensor management is considered for choosing between a discrete set of events (where to point the sensor, which sensing modality, etc.) in [8]. Resource allocation for multi-target tracking is also considered in [9] and [10]. In this work, though, we focus on the continuous allocation of energy. For an extended literature review, see section II-A.

The rest of the paper is organized as follows: Section II-A provides an extended literature review of resource allocation and adaptive sampling research. Section II-B specifically details the work of Bashan. Section III includes a formal statement of the problem. Section IV presents a performance analysis for when to use possibly incorrect prior knowledge rather than the uniform prior. Section V provides conclusions and future work. Lastly, appendices are provided to present the derivations of the analytical expressions and densities used for the analysis in this work.

### A. Extended literature review

Adaptive sampling (also referred to as active learning) was studied by Castro, Nowak, and Willet [3] [5] [4] in the context of estimating inhomogeneous functions from noisy measurements. They developed a method called backcasting that involved an initial "preview" step that distributed energy uniformly to estimate borders, followed by a "refinement" step to refine the boundary points. It was found that fast convergence rates in terms of the mean square error (MSE) could be obtained if the complex regions are relatively small compared to the rest of the signal. Castro also showed [2] that compressive sampling techniques can approach the performance of reconstruction of piecewise constant functions with adaptive sampling for high SNR. Castro et al. considered different applications that were characterized by inhomogeneous signals for which adaptive sampling would be efficient. However, in this work we consider general signals only restricted to have a small ROI.

Since we consider signals with a small ROI, we can refer to our signals as sparse. Compressive sampling to reconstruct sparse signals has been researched extensively over the past years. Tropp and

Gilbert [6] showed that m-sparse signals can be recovered with high probability while using many fewer measurements than the size of the signal. Moreover, efficient algorithms exist (orthogonal matching pursuit, etc.) that are both easy to implement and converge to the solution fast. However, compressive sampling does not lend itself easily to incorporating Bayesian knowledge into determining the location of sparse signals. Bayesian compressive sensing (BCS) [7] incorporates relevance vector machines that assume a prior on the signal in order to estimate both the signal and the confidence in the reconstructed signal. Moreover, BCS provides a means to determine if a sufficient number of measurements have been taken. However, since BCS relies on relevance vector machines [11], the prior information is governed by hyperparameters of a Bayesian linear model, and thus does not reflect the true prior information[1].

Resource allocation is considered in the context of sensor scheduling/management by Kastella [8] and Kreucher [9] [10]. Sensor management considers the problem of choosing between a discrete set of actions, such as choosing which cell to search at the next time step and in what mode. Kastella showed that under a myopic strategy, pointing the sensor to the cell that maximizes the discrimination gain (based on the Kullback-Leibler information) decreases the target misdetection probability. Kreucher et al. show that by combining sensor management with the joint multi-target probability density (JMPD) for target tracking, one can predict which measurement provides the most information gain. In our work, we consider distributing resources continuously.

*B. Review of ARAP*

Bashan [1] explored the problem of estimating a sparse ROI where prior knowledge is known in the form of a Bayesian prior on the existence of a target at a given cell. In particular, if we assume that targets may exist at a discrete set of $Q$ points, then the prior knowledge can be represented as the set of Bernoulli probabilities $\{p_1, p_2, \ldots, p_Q\}$, where $\Pr(I_i = 1) = p_i$ and $I_i$ is an indicator variable for the existence of a target at cell $i$.

Bashan considered random measurements at time $t$, $Y(t) = \{y_1(t), y_2(t), \ldots, y_Q(t)\}$, where $Y(t)$ depends on the energy allocation policy at that time[2], $\{\lambda(i,t)\}_{i=1}^{Q}$, the target locations, $\{I_i\}_{i=1}^{Q}$, and the random returns from each cell, $\{\theta_i\}_{i=1}^{Q}$. Moreover, Bashan considered $\lambda(i,t)$ to be a deterministic mapping from the past observations of $Y(1)$ through $Y(t-1)$. Lastly, optimal solutions were derived

---

[1]In essence, this model assumes that the posterior is a multivariate Gaussian. In our work, we do not wish to restrict our attention to just this case.

[2]$\lambda(i,t)$ is the amount of energy allocated to cell $i$ at time $t$

for both 1-stage and 2-stage energy allocation policies (henceforth referred to as ARAP), where the total energy was constrained to $\lambda_T$:

$$\sum_{t=1}^{T}\sum_{i=1}^{Q}\lambda(i,t) = \lambda_T \tag{1}$$

Optimality was defined through the cost function[3]:

$$J = \sum_{i=1}^{Q}\frac{\nu I_i + (1-\nu)(1-I_i)}{\sum_{t=1}^{T}\lambda(i,t)} \tag{2}$$

where $\nu$ controls the percentage of energy devoted to the ROI or to its complement. In particular, since $J$ is a random variable, the optimal energy allocation policies minimized $E[J]$. Note that for $T = 1$, the energy allocation policies are just deterministic, since they don't depend on random measurements. Bashan showed that the optimal allocation for $T = 1$ is given by:

$$\lambda(i,1) = \frac{\lambda_T \sqrt{p_i}}{\sum\limits_{j=1}^{Q}\sqrt{p_j}} \tag{3}$$

However, the problem is much more interesting when we consider $T > 1$. It was shown that given the energy allocation at $T = 1$, the optimal allocation of resources for $T = 2$ is a quantity that can be computed in $\mathcal{O}(Q)$ time. In particular, it was found that given $\{\lambda(i,1)\}_{i=1}^{Q}$ and $Y(1)$

$$\lambda(i,2) = \left( \frac{\lambda_T - \sum\limits_{j=1}^{Q}\lambda(j,1)}{\sum\limits_{j=k_0+1}^{Q}\sqrt{w_{(j)}}}\sqrt{w_{(i)}} - \lambda(i,1) \right) I(i > k_0) \tag{4}$$

where $w_i$ is a realization of the random variable $W_i = \Pr(I_i = 1|Y(1))$, $w_{(i)}$ is an ordered version of $w_i$, and $k_0$ defines a cutoff point based upon that ordering. Noting that $\lambda(i,2)$ depends on ordered random variables, Bashan also developed a suboptimal allocation policy whose performance paralleled closely that of the optimal policy, where

$$\lambda(i,2) = \left( \lambda_T - \sum_{j=1}^{Q}\lambda(j,1) \right)\frac{\sqrt{w_i}}{\sum\limits_{j=1}^{Q}\sqrt{w_j}} \tag{5}$$

---

[3]Bashan provided several reasons for the selection of this particular cost function. He showed that the cost function was lower bounded by a value that, and the bound could be attained by using an intuitive optimal resource allocation. Moreover, Bashan showed that by minimizing the cost function, one also minimized the CRB lower bound of estimating a deterministic quantity $\theta_i I_i$ in additive Gaussian noise, as well as uniformly minimizing the Chernoff bound on the misdetection probability over the ROI [1].

In both the optimal and suboptimal resource allocation policies, the second stage allocations, $\{\lambda(i,2)\}_{i=1}^{Q}$ determined by equation (4) or (5), is a function of the first stage allocations, $\{\lambda(i,1)\}_{i=1}^{Q}$. Thus, in general, the minimization problem involves $Q$ degrees of freedom in determining the first stage allocations (which consequently determine the second stage allocations). Since $Q$ is large by assumption, this minimization is infeasible for general prior information.

Bashan proposed two solutions to this problem. First, for uniform prior information (i.e., $p_i = \Pr(I_i = 1) = p$ for all $i$), $\lambda(i,1)$ is constant for all $i$. Therefore, the minimization problem reduces to determining the percentage of energy to allocate at the first and second stages, which can be done by grid search to provide the optimal allocation.

Second, in the general prior case, Bashan proposed to use a myopic approach, in the sense that one should distribute a total effort of $(\alpha\lambda_T, 1 - \alpha\lambda_T)$ to each stage for $\alpha \in (0,1)$, and then distribute effort optimally within each step according to equations (3-5). Lastly, one should grid search over $\alpha$ to determine the best allocation.

This research is primarily concerned with the case where the prior information is not uniform for multiple reasons. First, in many applications we know more than just a general sparsity constraint. For example, in active radar imaging, geopolitical constraints (mountains, oceans, country borders, etc.) may restrict the locations of targets. The next section describes one possible application where this applies directly. Second, for the case of $T > 2$, the measurements at times 1 through $T - 1$ will provide prior information that is non-uniform in just about all possible cases. Therefore, in this work, we aim to provide a theoretical basis for non-uniform prior information in order to potentially create adaptive resource allocation policies for multi-stage effort allocations. This work considers simple models for these non-uniform priors, and we show that under reasonable conditions, one can gain significant improvement over using uniform priors. Moreover, we discuss the extension to general priors and possible future work.

### C. Motivating application: Active radar in an airport landing field

Consider the problem of using an active radar system to detect and estimate the location of airplanes at an airport landing field. In this situation, we could expect that airplanes will approach the landing field from certain directions with much higher probability than others (due to geographic and safety constraints, for instance). In Figure 2 we show a possible configuration of probability regions for this application, where $p_0 \leq p_1 \leq p_2 \leq p_3$ and $p_0$ is much smaller than $p_1$.

In such a situation, we would expect that ARAP will provide significantly better results if we use this prior knowledge rather than just a uniform allocation, since we can allocate more resources to the higher
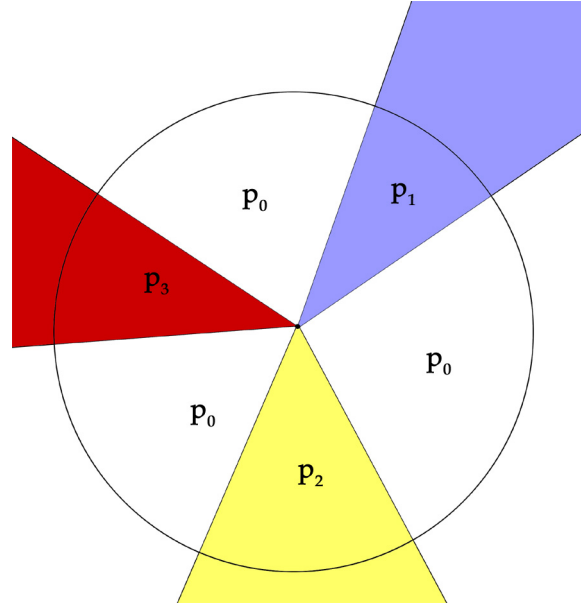
Fig. 2. Possible probability regions for the approach direction of airplanes at a landing field. We would expect that for many instances of this application, the prior information would lead to highly non-uniform applications. Thus, we would have $p_1, p_2, p_3 >> p_0$. As will be seen later in this report, significant gains can be had over a uniform alternative if we choose to use (possibly incorrect) non-uniform prior information for these applications.

probability regions and less resources to the lower probability regions.

Indeed, the more extensive our knowledge is about the locations of targets, the better our performance through ARAP. In the limiting case, we propose prior knowledge of the form:

$$\Pr(I_i = 1) = \begin{cases} 1, & i \in \Psi \\ 0, & i \notin \Psi \end{cases} \tag{6}$$

Clearly, in this case, we would allocate all of our resources to the ROI, and none outside of it. This energy allocation policy is denoted as $\Lambda_o$ in [1] and is the optimal resource allocation policy with respect to the cost function. However, since we don't know the ROI, this is clearly not a feasible solution.

From this discussion, we can easily see that the more knowledge we know about the locations of the targets, the better our performance gains will be in ARAP with respect to just using the uniform prior. On the other hand, as the amount of prior knowledge increases, so does the restrictiveness of our model. Thus, if our prior knowledge is inaccurate, we risk missing or badly estimating targets that lie in regions that we assumed to be improbable. For this reason, the following work also considers the problem of determining when it is advisable to use possibly incorrect prior knowledge versus the uniform alternative.

## III. PROBLEM STATEMENT

Let $\mathcal{X} = \{1, 2, \ldots, Q\}$ be a discrete space containing $Q$ cells equipped with a probability measure $P$. Let $\Psi \subseteq \mathcal{X}$ be a ROI that we are interested with an associated indicator function:

$$I_i = \begin{cases} 1, & i \in \Psi \\ 0, & \text{else} \end{cases} \tag{7}$$

and $\{p_i = \Pr(I_i = 1)\}_{i=1}^{Q}$ is an associated set of prior probabilities. Let $\lambda(i, t)$ be the amount of energy allocated to cell $i$ at time $t$. Let $T$ be the number of stages in the energy allocation policy. Then, the total amount of energy is constrained to $\lambda_T$ according to equation (1) with $0 \leq \lambda(i, t) \leq \lambda_T$. Let

$$\Lambda_i = \sum_{t=1}^{T} \lambda(i, t) \tag{8}$$

be the total energy allocated to cell $i$. Then let $\Lambda = \{\Lambda_1, \Lambda_2, \ldots, \Lambda_Q\}$ be the associated energy allocation policy, and define the cost function as:

$$J(\Lambda) = \sum_{i=1}^{Q} \frac{\nu I_i + (1 - \nu)(1 - I_i)}{\Lambda_i} \tag{9}$$

Let measurements at time $t$ be defined as

$$y_i(t) = \sqrt{\lambda(i, t)} \theta_i(t) I_i + \gamma_i \tag{10}$$

where $\theta_i(t)$ is the random return from cell $i$ and normally distributed with mean $\mu_\theta$ and variance $\sigma_\theta^2$, and $\nu_i(t)$ is normally distributed with zero mean and unit variance. We assume that $\nu_i(t)$ is independent for varying $i$ and $t$, and $\theta_i(t)$ is independent for varying $i$, but possibly dependent for different $t$. Note that we also are assuming that we are dealing with a static scenario, since the indicator function does not vary with time.

### A. Prior knowledge parameterization families

Let us assume that we are dealing with simple two-level piecewise uniform priors. Under this assumption, very few parameters actually describe all of our knowledge. In fact, four parameters, $(n_0, n_1, p_0, p_1)$ describe this knowledge, where

$$\Pr(I_i = 1) = \begin{cases} p_0, & i \in \{1, 2, \ldots, n_0\} \\ p_1, & i \in \{n_0 + 1, \ldots, n_0 + n_1\} \end{cases} \tag{11}$$

Let $(n_0, n_1, p_0, p_1)$ accurately describe the locations of targets, and let $(\hat{n}_0, \hat{n}_1, \hat{p}_0, \hat{p}_1)$ be the prior knowledge that we know. We will make a couple of assumptions in order to make our analysis easy. These include:

1) The number of cells is constant, so that $Q = n_0 + n_1 = \hat{n}_0 + \hat{n}_1$.

2) The number of expected targets over the entire region is known, where

$$E[|\Psi|] = Qp_{unif} = p_0 n_0 + p_1 n_1 = \hat{p}_0 \hat{n}_0 + \hat{p}_1 \hat{n}_1. \tag{12}$$

3) The low probability regions will assigned the same value $p_0 = \hat{p}_0$.

4) The high probability region (i.e., $i \in \{n_0 + 1, n_0 + 2, \ldots, Q\}$) may be underestimated or overestimated, but not missed altogether.

Note that the first two assumptions do not require any additional knowledge than in the uniform prior case. The third assumption results from the fact that in the applications that we are interested in, $p_0 << p_1$ and so in terms of our cost function, only the high probability region is really important to us. Thus, we simplify our parameterization by setting $\hat{p}_0 = p_0$. The last assumption is really a statement that we have some confidence in our prior knowledge. Clearly, if our prior knowledge is so bad that we don't include any of the high probability region, then we should use the uniform prior instead.

Let $\mathbf{g} = (\hat{n}_0, \hat{n}_1, p_0, \hat{p}_1)$ be a particular parameterization of our prior knowledge. Let $p_{unif}$ be the fixed sparsity level. Then, let us define a family of prior knowledge parameterizations, $\mathbf{G}$ to be

$$\mathbf{G} = \left\{ (\hat{n}_0, \hat{n}_1, p_0, \hat{p}_1) \mid \hat{n}_0 p_0 + \hat{n}_1 \hat{p}_1 = Qp_{unif}, \ \hat{n}_0 + \hat{n}_1 = Q \right\}, \tag{13}$$

where $\hat{n}_1$ is uniformly distributed over a discrete set of values:

$$\Pr(\hat{n}_1 = k) = \begin{cases} (n_1^{\max} - n_1^{\min})^{-1}, & k \in \left\{ n_1^{\min}, n_1^{\min} + 1, \ldots, n_1^{\max} \right\} \\ 0, & else \end{cases} \tag{14}$$

Note that in this family, $\hat{n}_1$ will define $\hat{n}_0$ and $\hat{p}_1$ as well due to the constraints on $\mathbf{G}$ and fixed $p_0$. In particular, we know

$$\hat{n}_1 = \hat{n}_1 \tag{15}$$

$$\hat{n}_0 = Q - \hat{n}_1 \tag{16}$$

$$\hat{p}_0 = p_0 \tag{17}$$

$$\hat{p}_1 = \frac{Qp_{unif} - \hat{n}_0 \hat{p}_0}{\hat{n}_1} = \frac{Q(p_{unif} - p_0) + \hat{n}_1 p_0}{\hat{n}_1} \tag{18}$$

Figure 3 presents a possible implementation of $\mathbf{G}$ for fixed $Q$ and $p_0$. The red curve represents the prior information with the smallest high probability region. Note that since the number of expected targets is equal for all elements of $\mathbf{G}$, this parameterization also has the highest probability in region 1. Thus, this is the *most* non-uniform prior, in the sense that $p_1$ and $p_0$ have the largest absolute difference from $p_{unif}$. The blue curve, on the other hand, represents the prior information with the largest high probability

region (and smallest $p_1$), thus being the *least* non-uniform prior in $\mathbf{G}$. The maximum difference between $n_1^{\max}$ and $n_1^{\min}$ will be determined by the application and will be representative of the confidence that the user has in their prior knowledge.
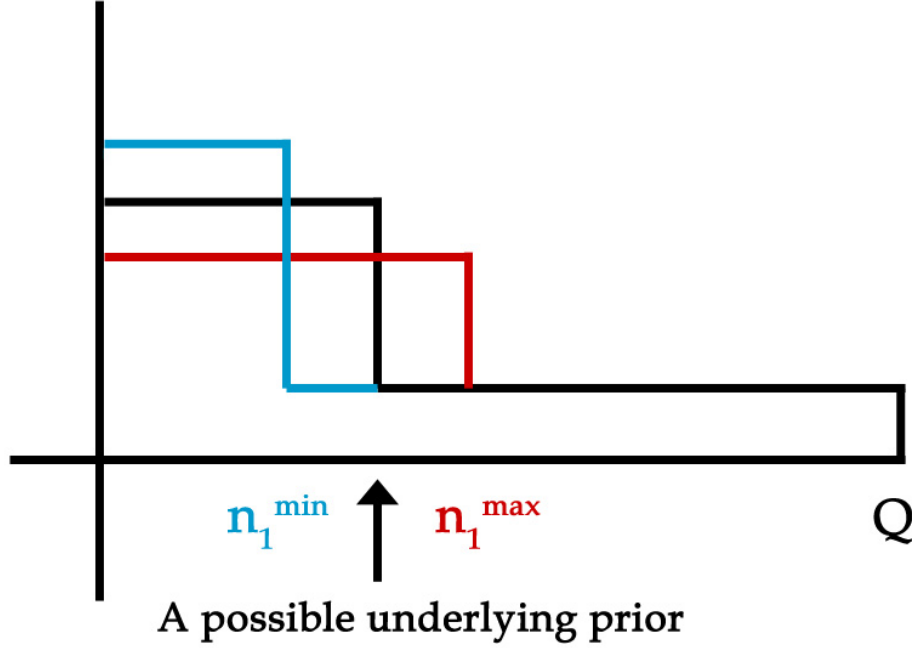


Fig. 3. A possible implementation of a two-level family of prior parameterizations for a given sparsity level. Denote $n_i$ to be the number of cells in region $i$, and $p_i$ to be the sparsity level within region $i$. In this analysis, we consider families of two-level uniform prior information, where $p_0$ is fixed to be a small value (i.e., a low probability region), the size of the high probability region $n_1$ is distributed uniformly over a discrete set of values in $\{n_1^{\min}, \ldots, n_1^{\max}\}$, and the number of expected targets is constant across all parameterizations. Under these constraints (discussed in Section III-A), the parameter $n_1$ determines all of the other descriptive parameters through equations (15-18).

It should be noted that our prior knowledge is **NOT** a Bayesian prior on the locations of targets, since it does not sum to one. Rather, our knowledge is a set of Bernoulli probabilities for each cell in $\mathcal{X}$. For this reason, we will henceforth refer to our knowledge only as parameterizations, even though they contain Bayesian information.

## IV. PERFORMANCE ANALYSIS

### A. Performance measures

For the purposes of comparison, we will draw a distinction between the underlying model denoted as $\boldsymbol{g}^* = \{p_i\}_{i=1}^{Q}$, and the assumed prior denoted as $\hat{\boldsymbol{g}} = \{\hat{p}_i\}_{i=1}^{Q}$. Note that in general, allocations formed by ARAP will depend on both the assumed and underlying priors, since the first stage allocation depends only on the assumed prior, while subsequent allocations depend on both models. Moreover, $\boldsymbol{u} := \{\hat{p}_i = p\}_{i=1}^{Q}$ will be referred to as a uniform prior and will be used as a base for comparison.

We will consider as a performance measure, the function that is minimized by ARAP, $E[J(\Lambda)|\boldsymbol{g}^*, \hat{\boldsymbol{g}}, \alpha]$, where $\alpha$ is assumed to be known. In practice, $\alpha$ can be computed by doing a line search over the minimum expected cost for each $\boldsymbol{g}^*$ and $\hat{\boldsymbol{g}}$.

Two quantities are explicitly computed in this analysis: (1) the expected cost, $C(\hat{\boldsymbol{g}}; \boldsymbol{g}^*)$, and (2) the worst-case gain when compared with using a uniform prior, $K(\hat{\boldsymbol{g}})$. Expressions for these quantities are given below.

$$C(\hat{\boldsymbol{g}}; \boldsymbol{g}^*) = E[J(\Lambda)|\boldsymbol{g}^*, \hat{\boldsymbol{g}}, \alpha] \tag{19}$$

$$K(\hat{\boldsymbol{g}}) = \max_{\boldsymbol{g}^* \in \boldsymbol{G}} -10\log \frac{C(\hat{\boldsymbol{g}}; \boldsymbol{g}^*)}{C(\boldsymbol{u}; \boldsymbol{g}^*)} \tag{20}$$

### B. Results and Discussion

*1) Derivation of density on $(\Lambda_i|\mathbf{g})$:* The details of this derivation are shown in Appendix A. However, it is important to note a couple of points. First, $\lambda(i, 1)$ is deterministic, but $\lambda(i, 2)$ depends on random measurements and therefore is random. Moreover, since $\Lambda_i = \lambda(i, 1) + \lambda(i, 2)$, it suffices to find the density on $\lambda(i, 2)$.

For optimal ARAP, $\lambda(i, 2)$ is a random variable that depends on ordered statistics, as well as the measurement and target models. This makes it unwieldy to find a useful analytical expression for this density. Thus, we consider suboptimal ARAP, where $\lambda_2(i)$ is defined as

$$\lambda(i, 2) = \frac{(1 - \alpha)\lambda_T \sqrt{W_{iy}}}{\sum\limits_{j=1}^{Q} \sqrt{W_{jy}}}, \tag{21}$$

where $W_{iy} = \Pr(I_i = 1|y(1))$. Using our knowledge of the measurement and target models, deriving a density on $W_{iy}$ is straightforward. Moreover, since $Q$ is very large, the denominator can be approximated by the central limit theorem (CLT) or strong law of large numbers. Combining these points, we are able

to derive an explicit expression for the density on $\lambda_2(i)$ with few approximations. Moreover, this density has been verified to be accurate by performing Kolmogorov-Smirnov tests on various instances of $\mathbf{g} \in \mathbf{G}$. However, the expression for the density involves integrals that are impossible (or at least, very difficult) to evaluate analytically. Thus, for the results presented next, numerical approximations to the integrals were used.

*2) Description of* $\mathbf{G}$ *used in the simulations:* We considered five parameterization families who differed only in their sparsity level, $p_{unif}$. In particular, we considered $p_{unif} = \{0.01, 0.0325, 0.055, 0.0775, 0.10\}$. Since the considered $\mathbf{G}$ differ only by sparsity level, we will abuse notation slightly by referring to the family only by its sparsity level.

Moreover, we fixed $p_0 = 0.002$, $n_1^{\min} = 1,500$ and $n_1^{\max} = 4,000$ for all parameterization families, and set the SNR to 8 dB. For a discussion on the choice of these parameters, see Section IV-B6. Lastly, we set the distribution parameters for the random cell returns to $\mu_\theta = 1$ and $\sigma_\theta^2 = 0.0625$.

It is important to note that we refer to the following results as simulations, even though they based on evaluating the analytical expressions derived in Appendices A and B. However, since we must rely on numerical approximations to some of the integrals, we denote our results as simulations to emphasize that we are approximating the theoretical results.

*3) Performance of uniform parameterization:* Recall that $C(\mathbf{u}, \mathbf{g}^*)$ is the expected cost of using the uniform parameterization in ARAP when the underlying parameterization is actually $g^*$. Figure 4(a) presents $C(\mathbf{u}, \mathbf{g}^*)$ versus $\mathbf{g}^*$ for all of the considered sparsity levels [4]. Clearly, we see that for each sparsity level, $C(\mathbf{u}, \mathbf{g}^*)$ remains nearly constant as a function of $g^*$. Figure 4(b) plots $C(\mathbf{u}, \mathbf{g}^*)$ for $p_{unif} = 0.055$, from which we can conclude that the expected cost is not constant everywhere. However, the variations are insignificant enough within a particular sparsity level to ignore their dependence on $\mathbf{g}^*$.

*4) Best-case performance: when the underlying parameterization is known:* The best we can do should intuitively occur when our energy allocation policy is derived by using the underlying parameterization. Figure 5 plots the gains in expected cost when the underlying parameterization is known with respect to the uniform parameterization alternative.

We see that the performance gain increases as $n_1$ decreases. Moreover, since $C(\mathbf{u}, \mathbf{g}^*)$ is approximately constant over $\mathbf{G}$, we can conclude that the expected cost decreases as $n_1$ decreases. This is equivalent

[4]Recall that for fixed $p_0$ and $Q$, we can equivalently represent $\mathbf{g}^* = (p_0, \hat{p}_1, \hat{n}_1, \hat{n}_0)$ with just $\hat{n}_1$, due to the constraints of $\mathbf{g}^* \in \mathbf{G}$. Therefore, in all remaining plots, we will refer to a particular $\mathbf{g}^*$ only by its associated $n_1$.
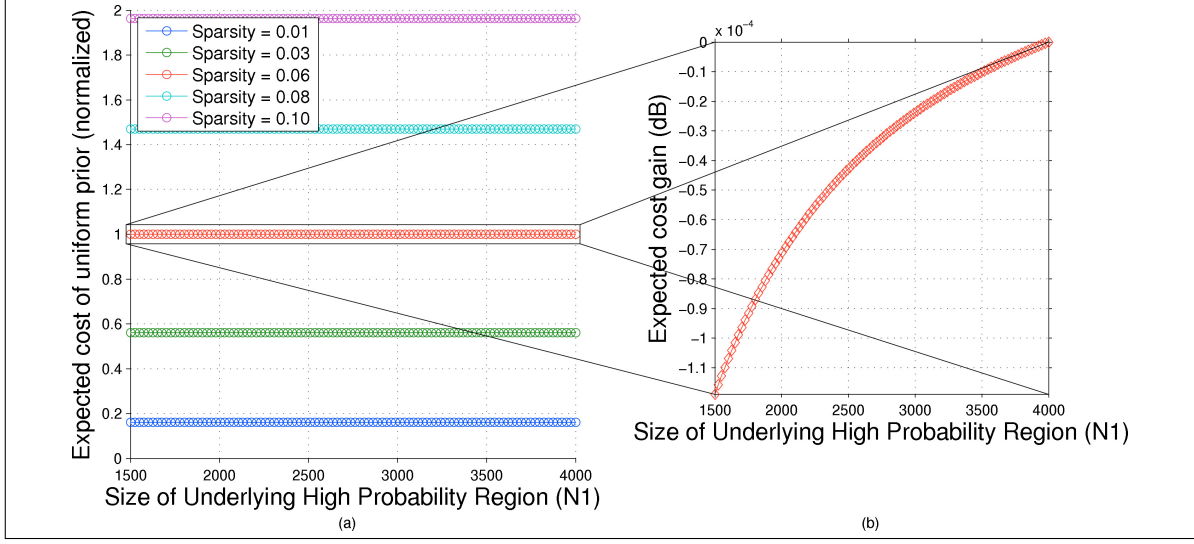
Fig. 4. We display the expected cost of using the uniform parameterization in ARAP as a function of the underlying parameterization in our family of prior knowledge, **G**. The results are plotted against $\hat{n}_1$, which is a sufficient parameter for describing our prior knowledge by equations (15-18). In (a), the expected cost is plotted for several sparsity levels. We have scaled the axes by the maximum expected cost for $p_{unif} = 0.055$ in order to easily interpret the results. We see that across all possible prior parameterizations, the expected cost when using a uniform prior is nearly constant. Moreover, the expected cost looks to be a largely linear function of sparsity level. In (b), we zoom in on the behavior for $p_{unif} = 0.055$. We see that although the expected cost is not uniform everywhere, the maximum variation is less than 1.2e-4 dB (i.e., less than 0.003% difference), which is practically insignificant in our analysis.

to saying that the expected cost is inversely proportional to the uniformity of the parameterization; i.e., as the parameterization becomes more non-uniform, its performance increases.

Note also that performance gains increase as the sparsity level increases. This point is interesting, in particular, because this is the reverse pattern from the results in [1]. However, the explanation for this discrepancy is simple: In [1], Bashan considered the alternative energy allocation policy to be an exhaustive search where higher percentages of energy is wasted as $p_{unif}$ decreases. On the other hand, in these simulations, we consider the alternative energy allocation policy to be ARAP with the uniform parameterization. Thus, as $p_{unif}$ decreases, there are fewer targets for which we can adaptively sample, and the margin for performance gains decrease. Section IV-B6 discusses the phenomenon in more detail.

Since our primary goal of this research is to determine when it makes sense to use a non-uniform parameterization over the uniform alternative, we are still mostly interested in the gains shown in Figure 5 (as opposed to comparing to an exhaustive search). However, it should be remarked that a 2 dB gain over the uniform alternative truly represents a much larger gain over an exhaustive search policy.

Table 1 shows the expected and worst-case gains with respect to the uniform alternative as well as the exhaustive search policy. We see that we get expected gains over an exhaustive search of more than 4 dB for all sparsity levels.
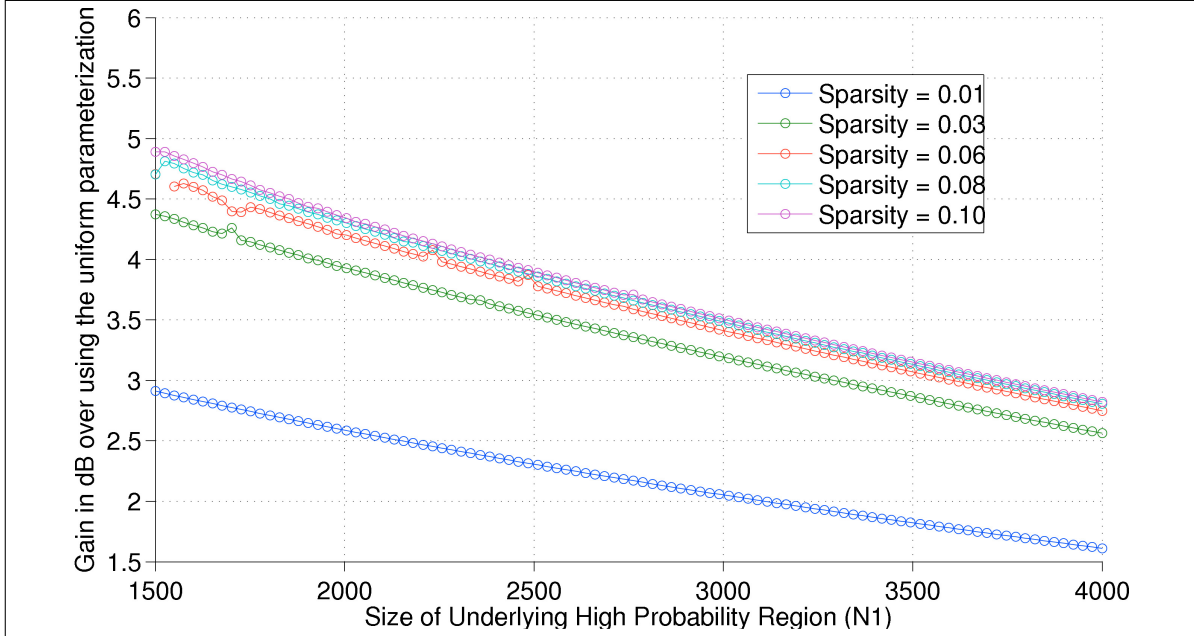


Fig. 5. We show the gain in expected cost over the uniform alternative in the best-case scenario, when the underlying prior parameterization is known. We see that as $\hat{n}_1$ decreases for all sparsity levels (and hence, the prior becomes more non-uniform), the gains in expected cost over the uniform alternative increase. It is important to observe that these gains are shown with respect to another adaptive resource allocation policy that already has significant gains over an exhaustive search. See Table 1 for additional discussion. Lastly, note that the gains over the uniform alternative increase as $p_{unif}$ increases, which is the reverse of the pattern found by Bashan [1]. In the limiting case, as $p_{unif} \to 0$, we expect there to be no gain by using non-uniform prior knowledge over a uniform alternative. Section IV-B6 discusses an explanation for this phenomenon.

*5) Practical case: when the underlying parameterization is approximated:* Now we consider the practical case, when we have to approximate the underlying parameterization with an element from $\mathbf{G}$. To get an intuitive sense for the best approximate element, we chose 5 possible approximate parameterizations for each sparsity level. Figure 6 shows the gain in dB over the uniform alternative as a function of the underlying parameterization for the case of $p_{unif} = 0.055$. We see several interesting and intuitive patterns here. First, the gain of an approximate parameterization never outperforms the true parameterization gain. Second, when an approximate parameterization overestimates the high probability region of the underlying parameterization, the gain is approximately equal to the best case gain of the approximate parameterization. This leads us to a lemma:

TABLE I

GAINS WHEN UNDERLYING PRIOR IS KNOWN

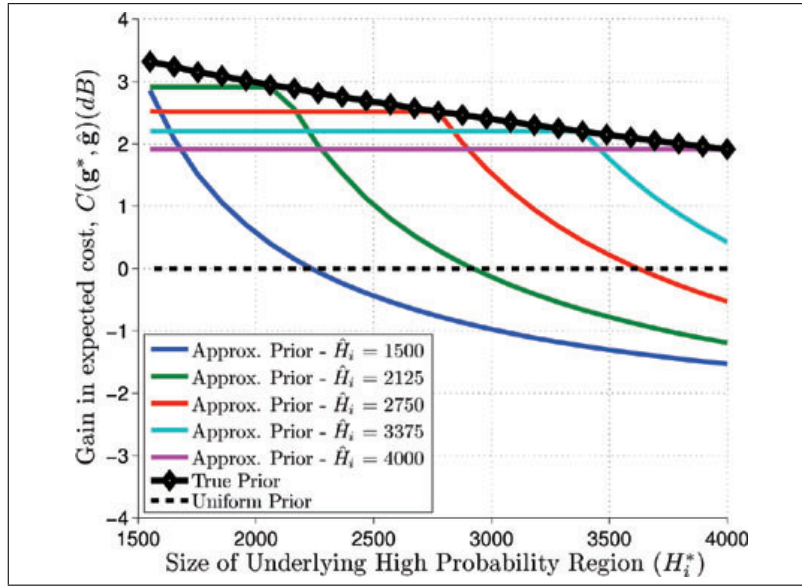| Sparsity Level | Uniform Alternative | | Exhaustive Search | |
|---|---|---|---|---|
| | Expected Gain (dB) | Min. Gain (dB) | Expected Gain (dB) | Min. Gain (dB) |
| 0.01 | 2.2047 | 1.6119 | 4.6261 | 4.2294 |
| 0.0325 | 3.3999 | 2.5622 | 5.4071 | 4.8460 |
| 0.055 | 3.6257 | 2.7486 | 5.3962 | 4.8028 |
| 0.0755 | 3.7165 | 2.8053 | 5.2823 | 4.6785 |
| 0.10 | 3.7503 | 2.8217 | 5.1538 | 4.5431 |



Fig. 6.    For a fixed sparsity level ($p_{unif} = 0.055$), we show the gain in dB over the uniform alternative when using several prior parameterizations to approximate all other possible underlying prior parameterizations. This figure presents one of the key results of this work: when $\hat{n}_1 > n_1$, then the gain in expected cost is nearly constant and equal to $C(\hat{\mathbf{g}}, \hat{\mathbf{g}})$, the gain in expected cost when $\hat{\mathbf{g}}$ is also the underlying prior. By Lemma 2, we can conclude that minimax solution to equation (20) will always be the prior parameterization with the largest $n_1 = n_1^{\max}$. This basic result tells us that we only need to construct our family of prior knowledge parameterizations, $\mathbf{G}$, so that $n_1^{\max}$ is greater than the uncertainty we have in our prior knowledge on the location of the probability regions. If this is done appropriately, then we can always do at least as well as the uniform alternative (which is just a special case with $n_1^{\max} = Q$).

*Lemma 1. Let $\hat{\mathbf{g}} = (p_0, \hat{p}_1, \hat{n}_1, \hat{n}_0) \in \mathbf{G}$ be an approximate parameterization. Let $\mathbf{g}^* = (p_0, p_1^*, n_1^*, n_0^*) \in \mathbf{G}$ be an underlying parameterization. If $n_1^* < \hat{n}_1$, then*

$$C(\hat{\mathbf{g}}, \mathbf{g}^*) \approx C(\hat{\mathbf{g}}, \hat{\mathbf{g}}) \tag{22}$$
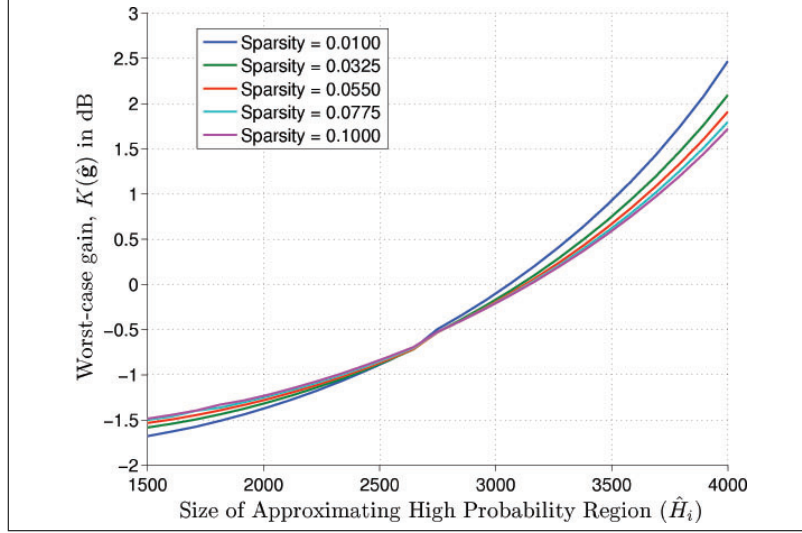
Fig. 7. We display here the worst-case gains over the uniform alternative for several potential approximating prior parameterizations. For each sparsity level, five approximating prior parameterizations were chosen that were evenly distributed across $\{n_1^{\min}, \ldots, n_1^{\max}\}$. We see that the worst-case gains are maximized when $n_1 = n_1^{\max}$ for all sparsity levels, which agrees with Lemma 2.

We are working on proving this formally, but currently we leave it as heuristic knowledge. Moreover, we see that when $\hat{n}_1 < n_1^*$ (i.e., the high probability region is underestimated), then the gains decay rapidly and drop below zero in some cases.

Figure 7 shows the worst-case gains as a function of the approximating parameterization (i.e., $K(\hat{g})$). We see that for all sparsities, the worst-case gain is maximized by the element with the highest $n_1$ value. In other words, the element that overestimates the high probability region for all other parameterizations in $\mathbf{G}$ should be used as the approximating prior if the worst-case criterion is used. This leads us to another lemma:

*Lemma 2. Let $\hat{\mathbf{g}} = (p_0, \hat{p}_1, \hat{n}_1, \hat{n}_0) \in \mathbf{G}$ be an approximate parameterization. Let $\mathbf{g}^* = (p_0, p_1^*, n_1^*, n_0^*) \in \mathbf{G}$ be an underlying parameterization. Then the $\hat{\mathbf{g}}$ that minimizes equation (20) is given by $\hat{\mathbf{g}}$ with $\hat{n}_1 < n_1^*$ for all $\mathbf{g}^* \in \mathbf{G}$.*

Once again we leave this as heuristic knowledge for now. However, we see that the worst-case gain shown in Figure 7(b) is a monotonically increasing function in $n_1$. The lemma follows from this observation.

*6) Note on simulations:* As mentioned earlier, we fixed particular values for $n_1^{\min}$ and $n_1^{\max}$ for the simulations described above. These values were chosen to satisfy two constraints:

$$\hat{p}_1 \geq \gamma p_{unif} \tag{23}$$

$$\hat{p}_1 \leq \beta < 1 \tag{24}$$

For $\gamma > 1$, the first constraint enforces that all possible parameterizations are non-uniform. In particular, for fixed $p_0$ and $p_{unif}$, we see that $\gamma$ is inversely proportional to $n_1^{\max}$. We see intuitively that as $n_1^{\max}$ approaches $Q$, $\gamma$ approaches 1. For the simulations above, we set $n_1^{\max} = 4,000$ so that $\gamma > 1.5$ for all sparsity levels (note that as the sparsity level increases, so does $\gamma$).

The second constraint enforces that our prior knowledge is valid for the applications we are looking at. Clearly, since $\hat{p}_1 \leq 1$, since it represents a Bernoulli probability. Moreover, for our simulations we set $\beta < 0.5$, so that we don't violate the sparsity assumption that motivates this research.

Another point is important with regard to selecting $n_1^{\max}$ in particular. We have seen in our previous discussion that the approximate element with maximum worst-case gain is the parameterization with $\hat{n}_1 = n_1^{\max}$ (in the case of uniformly distributed elements in $\mathbf{G}$, this generalizes to the element that maximizes expected gain as well). Moreover, that gain is approximately equal to the best-case gain of the parameterization with $n_1^* = n_1^{\max}$. Since we have shown that the best-case performance increases as $n_1^*$ decreases, we can conclude that as $n_1^{\max}$ decreases, our worst-case (and possibly expected) gain will increase.

Also, since the performance in Figure 5 is approximately linear, we can conclude that our uncertainty in $n_1^{\max}$ will be have a inversely proportional relationship to the worst-case performance. This relationship may depend on several parameters, such as SNR, $Q$, and $p_0$, and it may be worthwhile to investigate this further.

We also decided to fix the SNR to 8 dB for the presented simulations. Bashan [1] showed that for very low SNR, there was little room for performance gain over an exhaustive search alternative. This intuitively makes sense, since at very low SNR, all measurements will provide little information since they will be dominated by noise.

Now consider the situation when we have very high SNR. Bashan showed that as SNR and $Q$ approach infinity, the gain of using ARAP with a uniform energy allocation at the first step converges to $-10 \log p_{unif}$. Moreover, we can easily see that this result generalizes if assume any prior information (i.e., not just uniform). This leads us to a third Lemma.

*Lemma 3: The gain of using ARAP with any prior will asymptotically (in SNR) approach* $-10 \log p_{unif}$ *if the number of expected targets is given by* $Q p_{unif}$.

*Proof:* See Appendix C, which proves a two-fold result. First, ARAP is consistent in the sense that the posterior probabilities converge to $\Pr(I_i = 1)$ in probability with asymptotic SNR. Second, the gain asymptotically approaches the optimal gain.

Thus, as SNR approaches infinity, we would expect that ARAP will perform nearly identically regardless of the assumed prior information. This is corroborated by Figure 1, where we see that the gain in expected cost asymptotically approaches $-10 \log p_{unif}$.

Combining these thoughts, we see that there should be an optimal SNR operating point for which we can obtain the largest performance gains over using the uniform parameterization alternative. To illustrate this point, Monte-Carlo simulations were performed for several priors who differed only in their sparsity level. Figure 8 shows the gains in MSE[5], respectively, of using the true parameterization over using the uniform parameterization. We see that when minimizing MSE we have the largest margin for gain with SNR values between 8 and 16 dB.

We will now discuss the interesting behavior with respect to decreasing gains in expected cost over the uniform alternative as $p_{unif} \to 0$. Let us consider the limiting situation where $p_{unif}$ is very small. Note that in this case, our measurements will be dominated by noise, since our measurement model gives

$$y_i(1) = \sqrt{\lambda(i, 1)} \theta_i I_i + n_i \tag{25}$$

where $\theta_i$ and $n_i$ are normally distributed, and $E[\theta_i] \le 1$ (i.e., it doesn't amplify the signal on average). Thus, in order to get good performance, we must have high values of SNR. From our previous discussion, though, we know that when SNR is high, the performances of all ARAP policies with arbitrary priors will converge to the same value.

For an illustrative example, consider screening for terrorists in an airport. In this situation, we might expect there to be a terrorist in with probability $1e-7$. We could provide a non-uniform prior by assuming that foreigners were 3x more likely to be terrorists than citizens of that country. However, since the likelihood of either of these events is still remarkably small, the very large majority of $W_i = \Pr(I_i = 1)$ will be nearly zero, regardless of our prior knowledge. Thus, the second stage allocation will skewed to favor measurements from cells with targets (i.e., decreasing $p_{unif}$ is similar to increasing the SNR).

---

[5]MSE is calculated in an identical way to [1]; i.e., we use a naive Bayes estimator based on the measurements at both time steps to get $\hat{\theta}_i$, an estimate of $\theta_i$, and then calculate the MSE accordingly
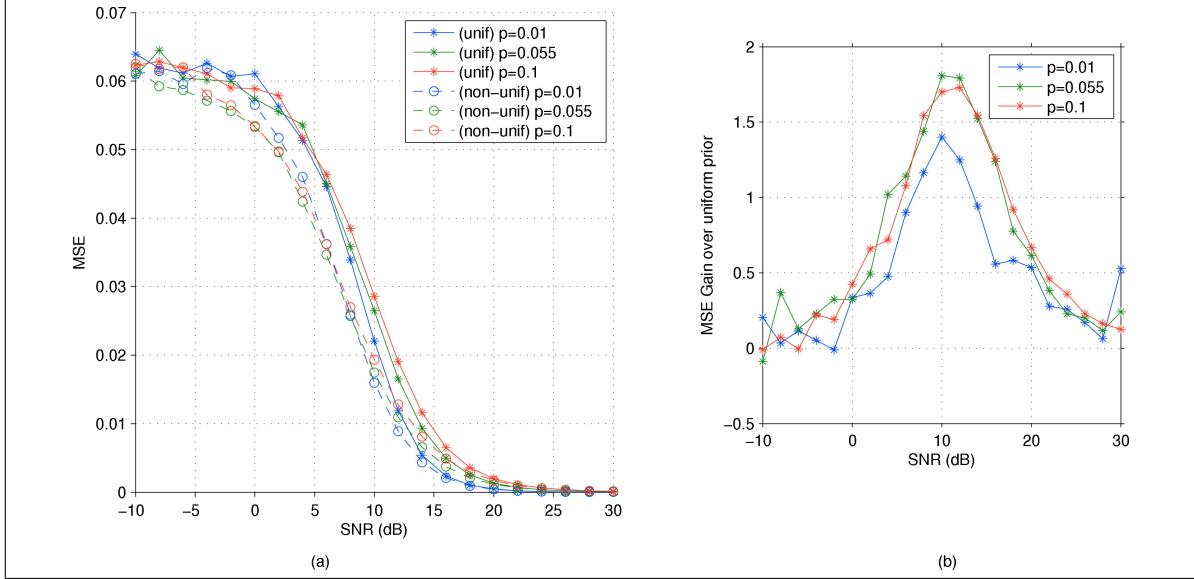
Fig. 8. We display here the results of Monte-Carlo simulations for MSE when comparing using a non-uniform prior versus the uniform alternative. In (a), we show the absolute MSE versus SNR for three sparsity levels ($p_{unif} = 0.01, 0.055, 0.1$). In (b), we plot the $dB$ gain of using the non-uniform prior over the uniform alternative. We see that for small SNR values or large SNR values, the performances are very similar. However, these simulations show that for SNR values in the range of 8-16 dB, we can get significant gains when using non-uniform prior information.

Therefore, all *adaptive* energy allocation policies will do very well when comparing against an exhaustive search as $p_{unif}$ decreases.

## V. CONCLUSIONS AND FUTURE WORK

In this research, we have extended the work of Bashan regarding optimal allocation of resources using adaptive sampling. In particular, we have looked at the case where our prior knowledge can be represented as two-level piecewise uniform that is drawn from a family of parameterizations describing that knowledge. Analytical expressions have been derived for the density on energy allocation policies obtained from suboptimal ARAP, as well as the expected cost.

We have shown that if the underlying parameterization is known, then significant gains over using uniform prior knowledge can be attained. These gains were seen to decrease in an approximately linear fashion with the size of the high probability region, and the slope of the decrease was not dependent on the sparsity level.

In any practical application, however, the underlying parameterization is unknown. In this case, we have shown heuristically that we can find a particular parameterization that maximizes the worst-case gain

and the expected gain under certain assumptions. This "minimax" solution depends solely on a single parameter $(n_1^{\max})$ of the overarching family of parameterizations, and thus, we have observed that the uncertainty in our prior knowledge is inversely proportional to the maximum worst-case performance; i.e., a linear increase in our uncertainty (represented through an increase in $n_1^{\max}$) will result in a proportional decrease in the maximum worst-case performance.

Our future work consists of several goals. In the near future, we would like to prove the lemmas described in this paper, as well as apply our analysis (at least in a Monte-Carlo sense) to more realistic cost functions, such as MSE and detection probability of error. Moreover, we would like to explore the performance difference between using suboptimal ARAP (for which analysis was tractable) and optimal ARAP. More importantly, we would like to show that we can extend our analysis to more general families of prior knowledge, so that we can broaden ARAP to a general multi-stage energy allocation policy. Monte-Carlo simulations have led us to believe that general priors can be decomposed into approximately piecewise uniform priors, for which the analysis in this paper would be useful.

## APPENDIX A. DERIVATION OF PRIOR DENSITY ON ENERGY ALLOCATION POLICIES

*A. Assumptions:*

For this derivation, we will assume that the following are known a priori:

1) $\mathbf{g}^* = (n_0^*, n_1^*, p_0^*, p_1^*)$, the underlying prior distribution parameters for which the targets follow the model:

$$\Pr(I_i = 1) = \begin{cases} p_0^*, & i \in \{1, 2, \ldots, n_0^*\} \\ p_1^*, & i \in \{n_0^* + 1, \ldots, n_0^* + n_1^*\} \end{cases} \tag{26}$$

2) $\mathbf{g} = (n_0, n_1, p_0, p_1)$, the assumed prior distribution parameters.

3) $i$, the cell whose energy allocation prior we are calculating. Note that this defines $p_i^*$, the underlying Bernoulli probability for the existence of a target in this cell, and $p_i$, the assumed Bernoulli probability for the existence of a target in this cell.

4) $\lambda_T$, the total amount of energy allocated across all time steps

5) $\alpha$, the percentage of energy allocated at the first time step (i.e., $\sum_{j=1}^{Q} \lambda(i, 1) = \alpha \lambda_T$).

6) $Q = n_0^* + n_1^* = n_0 + n_1$ is very large (in the thousands).

7) $n_0, n_1, n_0^*, n_1^*$ are quite large (at least in the hundreds)

8) The measurements follow the model:

$$y_i(t) = \sqrt{\lambda(i, 1)} \theta_i I_i + \nu_i(t) \tag{27}$$

where $\theta_i \sim \mathcal{N}(\mu_\theta, \sigma_\theta^2)$, $\nu_i(t) \sim \mathcal{N}(0, 1)$, $\theta_i$ and $\nu_i(t)$ are independent of each other, and $\nu_i(t)$ is independent across time.

9) $\{\lambda(i, 1)\}$, the deterministic set of energy allocations at the first time step, based on the assumed prior information, where

$$\lambda(i, 1) = \frac{\alpha \lambda_T \sqrt{p_i}}{\sum\limits_{j=1}^{Q} \sqrt{p_j}} \tag{28}$$

We also assume that we are allocating energy at the second step based on the suboptimal ARAP algorithm, in order to ignore ordered statistics. Note that [1] found that suboptimal ARAP to have nearly the same performance as optimal ARAP. Under this assumption, we have

$$\lambda(i, 2) = \frac{(1 - \alpha) \lambda_T \sqrt{W_i}}{\sum\limits_{j=1}^{Q} \sqrt{W_j}} = \frac{(1 - \alpha) \lambda_T Z_i}{\sum\limits_{j=1}^{Q} Z_j} \tag{29}$$

where $Z_i = \sqrt{W_i} = \sqrt{\Pr(I_i = 1 | \{y_j(1)\}_{j=1}^{Q})}^6$.

## B. Monotonic transformation

First, we note that $Z_i$ is a function of the measurements $\{y_j(1)\}_{j=1}^{Q}$ only through $y_i(1)$. Then, from Bayes rule, we can rewrite

$$W_i = \Pr(I_i = 1 | y_i(1)) = \frac{\Pr(y_i(1) | I_i = 1) \Pr(I_i = 1)}{\Pr(y_i(1) | I_i = 1) \Pr(I_i = 1) + \Pr(y_i(1) | I_i = 0) \Pr(I_i = 0)} \tag{30}$$

$$= \frac{1}{1 + \frac{\Pr(y_i(1) | I_i = 0) \Pr(I_i = 0)}{\Pr(y_i(1) | I_i = 1) \Pr(I_i = 1)}} \tag{31}$$

From our assumptions above, we have $\Pr(I_i = 1) = p_i = 1 - \Pr(I_i = 0)$ where $p_i$ is equal to $p_1$ or $p_0$ depending on $i$. Moreover, from our measurement model, we know

$$y_i(1) | I_i \sim \mathcal{N}(\sqrt{\lambda(i, 1)} \mu_\theta I_i, 1 + \lambda(i, 1) \sigma_\theta^2 I_i) \tag{32}$$

Therefore, we have

$$Z_i = \left( 1 + \frac{1 - p_i}{p_i} \sqrt{1 + \lambda(i, 1) \sigma_\theta^2} \exp\left\{ -\frac{1}{2(1 + \lambda(i, 1) \sigma_\theta^2)} \left( y_i^2(1) \lambda(i, 1) \sigma_\theta^2 + 2\mu_\theta \sqrt{\lambda(i, 1)} y_i(1) - \lambda(i, 1) \mu_\theta^2 \right) \right\} \right)^{-0.5} \tag{33}$$

---

[6]In [1], $W_i$ is given in a more general form with a tunable parameter for determining how much energy should be devoted to the ROI as well as its complement. We are looking at the specific case where all energy is devoted to estimating the ROI only, which gives this specific form for $W_i$

To make this derivation easier, let us define another variable

$$X_i = -\frac{1}{2(1 + \lambda(i,1)\sigma_\theta^2)}\left(y_i^2(1)\lambda(i,1)\sigma_\theta^2 + 2\mu_\theta\sqrt{\lambda(i,1)}y_i(1) - \lambda(i,1)\mu_\theta^2\right) \tag{34}$$

Then, we see that we can write $Z_i$ as a monotonic transformation of $X_i$, so that $Z_i = h(X_i)$, where

$$h(x) = \left(1 + e^x\frac{1 - p_i}{p_i}\sqrt{1 + \lambda(i,1)\sigma_\theta^2}\right)^{-0.5} \tag{35}$$

Since $h$ is monotonic, we can write our density on $Z_{(i)}$ as

$$f_{Z_i}(z) = f_{X_i}(h^{-1}(z)) \cdot |(h^{-1})'(z)| \tag{36}$$

$$= f_{X_i}\left(\ln(z^2 - 1) + \ln\left(\frac{p_i}{1 - p_i} \cdot \frac{1}{\sqrt{1 + \lambda(i,1)\sigma_\theta^2}}\right)\right) \cdot \frac{2}{z^3 - z} \tag{37}$$

*C. Density on $X_i$*

From equation (34), we see that $X_i$ is a quadratic function of $y_i(1)$. In particular, we can write

$$X_i = ay_i(1)^2 + by_i(1) + c \tag{38}$$

for

$$a = -\frac{\lambda(i,1)\sigma_\theta^2}{2(1 + \lambda(i,1)\sigma_\theta^2)} \tag{39}$$

$$b = -\frac{2\mu_\theta\sqrt{\lambda(i,1)}}{2(1 + \lambda(i,1)\sigma_\theta^2)} \tag{40}$$

$$c = \frac{\lambda(i,1)\mu_\theta^2}{2(1 + \lambda(i,1)\sigma_\theta^2)} \tag{41}$$

Noting that $a$ is negative, we see that $X_i$ is a concave function, and has a maximum at

$$\frac{dX_i}{dy_i(1)} = 0 \leftrightarrow y_i(1) = -b/2a = -\frac{2\mu_\theta\sqrt{\lambda(i,1)}}{\lambda(i,1)\sigma_\theta^2} \tag{42}$$

Thus, $F_{X_i}(x) = \Pr(X_i \leq x) = 0$ for $x \geq -\frac{2\mu_\theta\sqrt{\lambda(i,1)}}{\lambda(i,1)\sigma_\theta^2}$. Equivalently, $f_{X_i}(x) = 0$ for $x \geq -\frac{2\mu_\theta\sqrt{\lambda(i,1)}}{\lambda(i,1)\sigma_\theta^2}$.
For any other $x$, we have

$$F_{X_i}(x) = \Pr(X_i \leq x) \tag{43}$$

$$= \Pr(ay_i(1)^2 + by_i(1) + c - x \leq 0) \tag{44}$$

$$= \Pr(r_1(x) \leq y_i(1) \leq r_2(x)) \tag{45}$$

where $\{r_1(x), r_2(x)\} = \frac{-b \mp \sqrt{b^2 + 4a(x-c)}}{2a}$ are the roots of $ay_i(1)^2 + by_i(1) + c - x = 0$, and $r_1(x) \leq r_2(x)$. Noting that $y_i(1)$ is conditionally Gaussian given $I_i$, we have the following expression for the density on $y_i(1)$:

$$
\begin{align}
f_{y_i(1)}(y) &= f_{y_i(1)|I_i}(y|1)p_i^* + f_{y_i(1)|I_i}(y|0)(1 - p_i^*) \tag{46} \\
&= \frac{p_i^*}{\sqrt{2\pi(1 + \lambda(i,1)\sigma_\theta^2)}} \exp\left\{-\frac{(y - \sqrt{\lambda(i,1)}\mu_\theta)^2}{2(1 + \lambda(i,1)\sigma_\theta^2)}\right\} + \frac{1 - p_i^*}{\sqrt{2\pi}} \exp\left\{-\frac{y^2}{2}\right\} \tag{47}
\end{align}
$$

where $p_i^*$ is the underlying Bernoulli probability of the location of a target (i.e., the ideal probability, not the assumed one). With this density, we can define

$$
f_{X_i}(x) = \frac{d}{dx} F_{X_i}(x) = \frac{d}{dx} \int_{r_1(x)}^{r_2(x)} f_{y_i(1)}(y) dy \tag{48}
$$

Lastly, by Leibniz's rule, we see that for $x < -\frac{2\mu_\theta \sqrt{\lambda(i,1)}}{\lambda(i,1)\sigma_\theta^2}$ we have

$$
f_{X_i}(x) = \frac{f_{y_i(1)}(r_1(x)) + f_{y_i(1)}(r_2(x))}{\sqrt{b^2 + 4a(x - c)}} \tag{49}
$$

yielding the final density on $X_i$ as

$$
f_{X_i}(x) = \begin{cases} \frac{f_{y_i(1)}(r_1(x)) + f_{y_i(1)}(r_2(x))}{\sqrt{b^2 + 4a(x-c)}}, & x < -\frac{2\mu_\theta \sqrt{\lambda(i,1)}}{\lambda(i,1)\sigma_\theta^2} \\ \\ 0, & \text{else} \end{cases} \tag{50}
$$

Putting equations (37) and (50) together provide the density on $Z_i$.

### D. Beta density approximation

Looking at the shape of $f_{Z_i}(z)$, it was suggested that we could possibly use the beta density defined by

$$
f(k; \alpha, \beta) = \frac{k^{\alpha-1}(1 - k)^{\beta-1}}{B(\alpha, \beta)} \tag{51}
$$

for random variable K. where $B(\alpha, \beta)$ is the beta function. For this distribution, we know that

$$
\begin{align}
E[K] &= \frac{\alpha}{\alpha + \beta} \tag{52} \\
E[K^2] &= \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)} \tag{53}
\end{align}
$$

Since our earlier derivations give us analytical solutions to find $E[Z]$ and $E[Z^2]$, we can use the previous two equations to solve a system of 2 equations with 2 unknowns, yielding the beta density parameters, $\alpha$ and $\beta$. However, when this is done, we have found that Kolmorogov-Smirnov tests rejected these densities

as being good approximations to $f_{Z_i}$ (and hence $f_{\lambda_2(i)}$ or $f_{\Lambda_i}$). Thus, we conclude that we cannot use the beta density for our analysis.

On the other hand, an eye test shows that the shape of the beta density still resembles $f_{Z_i}$, which we may find useful for theoretical proofs further in our analysis.

*E. Density on $\lambda(i, 2)$*

In general, the denominator of equation (29) involves the summation of thousands of independent random variables. An explicit density could theoretically be calculated as the convolution of the densities of each of the r.v.'s, but this is not possible in practice. However, since $Q$ is large, we can exploit laws of large numbers in order to create an approximation to the density. First, we note that

$$\sum_{j=1}^{Q} Z_j = \sum_{j \in R_{00}} Z_j + \sum_{j \in R_{01}} Z_j + \sum_{j \in R_{10}} Z_j + \sum_{j \in R_{11}} Z_j, \tag{54}$$

where we define the regions $R_{mn}$ to be the region where the underlying model has the Bernoulli probability $p_m$, while the assumed model has the Bernoulli probability $p_n$. Within each region, the $Z_j$ are identically distributed and independent. Moreover, we can reasonably assume the following:

1) Either a region will have a very large number of i.i.d. random variables, in which case a law of large numbers will apply, or

2) The region will have few i.i.d. random variables, but the overall sum will be dominated by the other partial sums. Thus, even a bad approximation to this partial sum will result in little overall error.

*1) CLT approximation:* Let us assume that $E[Z_j | j \in R_{mn}] = \mu_{mn}$, $Var[Z_j | j \in R_{mn}] = \sigma_{mn}^2$, and $n_{mn} = |R_{mn}|$. Then, by the CLT we have

$$\frac{1}{\sqrt{n_{mn}}} \sum_{j \in R_{mn}} \frac{Z_j - \mu_{mn}}{\sigma_{mn}} \sim \mathcal{N}(0, 1) \tag{55}$$

or equivalently

$$\sum_{j \in R_{mn}} Z_j \sim (n_{mn}\mu_{mn}, n_{mn}\sigma_{mn}^2) \tag{56}$$

Thus, our denominator is normally distributed since it is the linear combination of normal variables:

$$\sum_{j=1}^{Q} Z_j \sim \mathcal{N}(\mu_{den}, \sigma_{den}^2) \tag{57}$$

where $\mu_{den} = \sum_{\substack{m \in \{0,1\} \\ n \in \{0,1\}}} \mu_{mn}$ and $\sigma_{den}^2 = \sum_{\substack{m \in \{0,1\} \\ n \in \{0,1\}}} \sigma_{mn}^2$ Since both the numerator and denominator is random in equation 29, we must compute the density of a ratio of two dependent variables [7]. Let us consider the ratio

$$A = \frac{\beta Z}{T} \tag{58}$$

Then we know that

$$f_A(a) = \frac{d}{da} \Pr(A \leq a) = \frac{d}{da} \left\{ \int_{-\infty}^{\infty} \Pr\left(\frac{\beta Z}{T} \leq a | T = t\right) f_T(t) dt \right\} \tag{59}$$

$$= \frac{d}{da} \left\{ \int_{-\infty}^{\infty} \Pr\left(Z \leq \frac{at}{\beta} | T = t\right) f_T(t) dt \right\} \tag{60}$$

$$= \int_{-\infty}^{\infty} \frac{d}{da} \left\{ F_{Z|T}\left(\frac{at}{\beta} | t\right) \right\} f_T(t) dt \tag{61}$$

$$= \int_{-\infty}^{\infty} \frac{t}{\beta} f_{Z|T}\left(\frac{at}{\beta} | t\right) f_T(t) dt \tag{62}$$

$$= \int_{-\infty}^{\infty} \frac{t}{\beta} f_{T|Z}\left(t | \frac{at}{\beta}\right) f_Z\left(\frac{at}{\beta}\right) dt \tag{63}$$

$$\tag{64}$$

From this, we see see that for $T = \sum_{j=1}^{Q} Z_j$ and $\beta = (1 - \alpha)\lambda_T$, we have

$$f_{\lambda(i,2)}(\lambda) = \int_{-\infty}^{\infty} \frac{t}{(1-\alpha)\lambda_T} f_{T|Z}\left(t \Big| \frac{\lambda t}{(1-\alpha)\lambda_T}\right) f_Z\left(\frac{\lambda t}{(1-\alpha)\lambda_T}\right) dt \tag{65}$$

Note that given $Z_i = z$, $T = \sum_{j=1}^{Q} Z_j \sim \mathcal{N}(\mu_{den} + (z - \mu_i), \sigma_{den}^2 - \sigma_i^2)$ for $E[Z_i] = \mu_i$ and $Var[Z_i] = \sigma_i^2$, so all of the quantities in the equation above are known.

*2) Strong law of large numbers approximation:* The CLT gives a good approximation to the distribution of the denominator. However, using this method also requires that we calculate an additional integral that is not easily computable (or numerically approximated). Since we have thousands of observations, however, we can use the SLLN to approximate the denominator as just the mean of all of the elements of the sum. Then,

$$f_{\lambda(i,2)}(\lambda) = \frac{\mu_{den}}{(1-\alpha)\lambda_T} f_{Z_i}\left(\frac{\mu_{den}\lambda}{(1-\alpha)\lambda_T}\right) \tag{66}$$

*3) Chebyshev's law of large numbers:* This is actually identical in form to the CLT above, except that we use Chebyshev's justification for using the mean, since from his notes *On mean values*, he shows that as the number of independent variables goes to infinity, the sample mean approaches the average of the means of the variables almost surely.

---

[7]Note that we can make the approximation that the denominator is independent of the numerator. This does not reduce the complexity greatly though, so we include the full analysis here

APPENDIX B. DERIVATION OF CLOSED-FORM EXPRESSION FOR EXPECTED COST

*F. A note on $W_i = \Pr(I_i = 1|Y(1))$*

Recall that we defined $W_i = Z_i^2 = \Pr(I_i = 1|Y(1))$. In the absence of a true model, we had to approximate the true quantity $\Pr(I_i = 1) = p_i^*$ with $\Pr(I_i = 1) \approx p_i$. Let $M_i^2 = \Pr(I_i = 1|Y(1))$, where we note that $Z_i^2 \approx M_i^2$ if $p_i \approx p_i^*$. Then, we can write

$$M_i^2 = \left(1 + \frac{\Pr(y_i(1)|I_i = 0)\Pr(I_i = 0)}{\Pr(y_i(1)|I_i = 1)\Pr(I_i = 1)}\right)^{-1} \tag{67}$$

$$= \left(1 + \frac{1 - p_i^*}{p_i^*} \cdot \frac{\Pr(y_i(1)|I_i = 0)}{\Pr(y_i(1)|I_i = 1)}\right)^{-1} \tag{68}$$

and

$$Z_i^2 = \left(1 + \frac{1 - p_i}{p_i} \cdot \frac{\Pr(y_i(1)|I_i = 0)}{\Pr(y_i(1)|I_i = 1)}\right)^{-1} \tag{69}$$

Clearly, we see that if we let $\gamma = \frac{\Pr(y_i(1)|I_i=0)}{\Pr(y_i(1)|I_i=1)}$, then

$$Z_i^2 = \left(1 + \frac{1 - p_i}{p_i}\gamma\right)^{-1} \tag{70}$$

or equivalently

$$\gamma = \frac{p_i}{1 - p_i}(Z_i^{-2} - 1) \tag{71}$$

So that we can write $M_i^2$ explicitly in terms of $Z_i$:

$$M_i^2 = \left(1 + \frac{1 - p_i^*}{p_i^*} \cdot \frac{p_i}{1 - p_i}(Z_i^{-2} - 1)\right)^{-1} = \left(1 + \beta(Z_i^{-2} - 1)\right)^{-1} = \frac{Z_i^2}{\beta + (1 - \beta)Z_i^2} \tag{72}$$

for $\beta = \frac{1 - p_i^*}{p_i^*} \cdot \frac{p_i}{1 - p_i}$. Let us define $g(x)$ as

$$g(x) = \frac{x^2}{\beta + (1 - \beta)x^2} \tag{73}$$

Then, clearly we have $M_i^2 = g(Z_i)$, which will be useful in the next section.

*G. An expression for the expected cost*

Now let's explore the quantity

$$E[J(\Lambda)] = E\left[\sum_{j=1}^{Q} \frac{I_j}{\Lambda_j}\right] = \sum_{j=1}^{Q} E\left[\frac{I_j}{\Lambda_j}\right] \tag{74}$$

We assumed that $I_j$ are independent for all $j$. Thus, we can write

$$E[J(\Lambda)] = \sum_{j=1}^{Q} E\left[\frac{I_j}{\Lambda_j}\right] = \sum_{\substack{n \in \{0,1\} \\ m \in \{0,1\}}} \left(\sum_{j \in R_{mn}} E\left[\frac{I_j}{\Lambda_j}\right]\right) \tag{75}$$

where all $j \in R_{mn}$ are identically distributed. Once the expectation has been taken, the quantity is no longer random. Thus, we can find a final expression for the expectation as:

$$E[J(\Lambda)] = \sum_{j=1}^{Q} E\left[\frac{I_j}{\Lambda_j}\right] = \sum_{\substack{n \in \{0,1\} \\ m \in \{0,1\}}} n_{mn} E\left[\frac{I_j}{\Lambda_j}\bigg| j \in R_{mn}\right] \tag{76}$$

Now consider the quantity inside the sum. We can find an equivalent representation as

$$E\left[\frac{I_j}{\Lambda_j}\bigg| j \in R_{mn}\right] = E\left[E\left[\frac{I_j}{\Lambda_j}\bigg| Y(1), j \in R_{mn}\right]\right] \tag{77}$$

$$= E\left[\frac{E[I_j|Y(1)]}{\Lambda_j}\bigg| j \in R_{mn}\right] \tag{78}$$

$$= E\left[\frac{\Pr(I_j = 1|Y(1))}{\Lambda_j}\bigg| j \in R_{mn}\right] \tag{79}$$

$$= E\left[\frac{R_j^2}{\Lambda_j}\bigg| j \in R_{mn}\right] \tag{80}$$

Dropping the conditioning on region for clarification in notation, we see that we can write

$$E\left[\frac{I_j}{\Lambda_j}\right] = E\left[\frac{R_j^2}{\Lambda_j}\right] = E\left[\frac{g(Z_j)}{\Lambda_j}\right] \tag{81}$$

$$= \int_{-\infty}^{\infty} E\left[\frac{g(Z_j)}{\Lambda_j}\bigg| \Lambda_j = k\right] f_{\Lambda_j}(k)dk \tag{82}$$

$$= \int_{-\infty}^{\infty} \frac{1}{k} E\left[g(Z_j)\bigg| \Lambda_j = k\right] f_{\Lambda_j}(k)dk \tag{83}$$

Now we note that $\Lambda_j = \lambda(j,1) + \lambda_j(2)$ so that $f_{\Lambda_j}(k) = f_{\lambda_j(2)}(k - \lambda(j,1))$. Making the change of variables $k' = k - \lambda(j,1)$, we can rewrite the integral as:

$$E\left[\frac{I_j}{\Lambda_j}\right] = \int_{-\infty}^{\infty} \frac{1}{k' + \lambda_j(1)} E\left[g(Z_j)\bigg| \lambda_j(2) = k'\right] f_{\lambda_j(2)}(k')dk' \tag{84}$$

Now we just need to deal with the conditional expectation. In particular, since we are dealing with expectation of a function of $Z_j$ given $\lambda_j(2)$, we just need the distribution on $Z_j|\lambda_j(2)$. Clearly, by equation 29, we can write

$$(1-\alpha)\lambda_T Z_i = \lambda(i,2)\sum_{j=1}^{Q} Z_j \tag{85}$$

$$Z_i\left((1-\alpha)\lambda_T - \lambda(i,2)\right) = \lambda(i,2)\sum_{\substack{j=1 \\ j \neq i}}^{Q} Z_j \tag{86}$$

$$Z_i = \frac{\lambda(i,2)}{(1-\alpha)\lambda_T - \lambda_j(2)}\sum_{\substack{j=1 \\ j \neq i}}^{Q} Z_j \tag{87}$$

Thus, if we know $\lambda(i, 2)$, then $Z_i$ is random only through the summation term. This quantity can be easily seen to be the same quantity approximated in sections 1.5.1 through 1.5.3. Thus, either $Z_i|\lambda(i, 2)$ will be normally distributed with mean $\beta(\lambda(i, 2)) \cdot (\mu_{den} - \mu_i)$ and variance $\beta^2(\lambda(i, 2)) \cdot (\sigma_{den}^2 - \sigma_i^2)$ (CLT approximation) for

$$\beta(\lambda) = \frac{\lambda}{(1 - \alpha)\lambda_T - \lambda}, \tag{88}$$

or it will be a deterministic quantity with value $\beta(\lambda(i, 2)) \cdot (\mu_{den} - \mu_i)$ (Strong Law of Large Numbers)[8]. In the former case, we have

$$E\left[g(Z_j)|\lambda_j(2) = k'\right] = \int_{-\infty}^{\infty} \frac{g(z)}{\sqrt{2\pi\beta^2(k')(\sigma_{den}^2 - \sigma_i^2)}} \exp\left\{-\frac{(z - \beta(k')(\mu_{den} - \mu_i))^2}{2\beta^2(k')(\sigma_{den}^2 - \sigma_i^2)}\right\} dz \tag{89}$$

In the latter case

$$E\left[g(Z_j)|\lambda_j(2) = k'\right] = g\left(\beta(k') \cdot (\mu_{den} - \mu_i)\right) \tag{90}$$

Plugging either of these into equation (84) yields the quantity $E[I_j/\Lambda_j]$ for any $j$. Plugging these quantities into equation (76) yields the expected cost.

## APPENDIX C. PROOF OF ASYMPTOTIC SNR RESULT

We will show that as SNR goes to infinity, the gain of using optimal ARAP approaches $-10\log p_{unif}$ where $p_{unif} = E\left[\frac{|\Psi|}{Q}\right]$. Note that we do not make any assumptions on $p_i$. If $p_i = q$ for all $i$, then we have a uniform prior knowledge.

*a) Consistency:* Define $\mathcal{H}_0$ to be the event that a target does not exist at cell $i$ (i.e., $I_i = 0$) and $\mathcal{H}_1$ to be its complement event (i.e., $I_i = 1$). Then we see that the posterior probabilities $\{p_{I_i|\mathbf{Y}(1)}\}$ can be rewritten as

$$p_{I_i|\mathbf{y}(1)} = \frac{p_i f_1(\mathbf{y}(1))}{p_i f_1(\mathbf{y}(1)) + (1 - p_i) f_0(\mathbf{y}(1))} \tag{91}$$

Under $\mathcal{H}_0$, we know that $\Pr(I_i = 1) = p_i = 0$, which by the above equation yields $p_{I_i|\mathbf{y}(1)} = 0 = p_i$. Otherwise, we have

$$p_{I_i|\mathbf{y}(1)} = \frac{1}{1 + \frac{1 - p_i}{p_i} \frac{f_0(\mathbf{y}(1))}{f_1(\mathbf{y}(1))}}, \tag{92}$$

where $f_j(\mathbf{y}(1))$ denotes the probability density function (pdf) of a specific observation conditioned on $\mathcal{H}_j$ for $j = 0, 1$, and noting $p_i > 0$. Moreover, according to equation (32) we see that $\mathbf{y}(1)$ is Gaussian when conditioned on either $\mathcal{H}_0$ or $\mathcal{H}_1$, and $\frac{1 - p_i}{p_i} \frac{f_0(\mathbf{y}(1))}{f_1(\mathbf{y}(1))} \propto \frac{1}{\text{LRT}}$, where the likelihood ratio test (LRT) is

---

[8]As before, $\mu_i = E[Z_i]$ and $\sigma_i^2 = Var[Z_i]$.

between two Gaussian distributions. Therefore, the posterior probability $p_{I_i|\boldsymbol{y}(1)}|\mathcal{H}_1 \to 1$ since it follows the performance of an LRT for a simple binary hypothesis testing problem with a Gaussian distribution. Recall that an LRT is the uniformly most powerful test for distinguishing between two Gaussians with $\mu_1 > \mu_0$, where $\mu_r$ is the mean conditioned on $\mathcal{H}_r$ and $r = 0, 1$. Therefore, under either $\mathcal{H}_0$ or $\mathcal{H}_1$, we have shown $p_{I_i|\boldsymbol{y}(1)} \to p_i$, which concludes the proof of consistency.

*b) Asymptotic gain:* Now that we know $p_{I_i|\boldsymbol{y}(1)} \to p_i$, we can show that the gain using any general prior approaches the optimal gain as SNR increases asymptotically. Following the proof of a similar result in [1], we know that we can order the posterior probabilities so that $w_i$ (i.e., a realization of $W_i$) is equal to 0 for all $\tau(i) \le \tilde{k}$ and 1 for all $\tau(i) > \tilde{k}$, where $\tilde{(\cdot)}$ represents an ordering transformation. In this case, $\tilde{k} = Q - \sum_{i=1}^{Q} I(w_i = 1)$. Thus, we see that $|\Psi| = Q - \tilde{k}$.

Assume that we know the 1st stage allocations, $\{\lambda_1(i)\}_{i=1}^{Q}$. Then, from [1], equation (136), we know that

$$\lambda_2(i) = \left( \frac{\lambda_T - \sum_{j=1}^{\tilde{k}} \lambda_1(\tau(j))}{\sum_{j=\tilde{k}+1}^{Q} \sqrt{w_{\tau(j)}}} \sqrt{w_{\tau(i)}} - \lambda_1(\tau(i)) \right) I(\tau(i) > \tilde{k}) \tag{93}$$

But we know that $w_{\tau(j)} = 0$ for $j \le \tilde{k}$ and $w_{\tau(j)} = 1$ for $j > \tilde{k}$, giving:

$$\lambda_2(i) = \left( \frac{\lambda_T - \sum_{j=1}^{\tilde{k}} \lambda_1(\tau(j))}{Q - \tilde{k}} - \lambda_1(\tau(i)) \right) I(\tau(i) > \tilde{k}) \tag{94}$$

Since, $\Lambda_i = \lambda_1(i) + \lambda_2(i)$, we have

$$\Lambda_i = \begin{cases} \lambda_1(i), & \tau(i) \le \tilde{k} \\ \frac{1}{Q-\tilde{k}} \left( \lambda_T - \sum_{j=1}^{\tilde{k}} \lambda_1(\tau(j)) \right), & else \end{cases} \tag{95}$$

So the expected cost is

$$E[J(\Lambda)] = E\left[ \sum_{i=1}^{Q} \frac{I_i}{\Lambda_i} \right] = E\left[ \sum_{i=1}^{\tilde{k}} \frac{I_{\tau(i)}}{\Lambda_{\tau(i)}} \right] + E\left[ \sum_{i=\tilde{k}+1}^{Q} \frac{I_{\tau(i)}}{\Lambda_{\tau(i)}} \right] \tag{96}$$

Note in the first sum, we know that $w_{\tau(i)} = 0$ and $W_i \to I_i$. Thus, for high SNR, the first sum is equal to 0, yielding:

$$\lim_{SNR \to \infty} E[J(\Lambda)] = E\left[ \sum_{i=\tilde{k}+1}^{Q} \frac{1}{\Lambda_{\tau(i)}} \right] = E\left[ \sum_{i=\tilde{k}+1}^{Q} \frac{Q - \tilde{k}}{\lambda_T - \sum_{j=1}^{\tilde{k}} \lambda_1(\tau(j))} \right] \tag{97}$$

Since the term inside of the expectation operator is monotonically increasing in $\{\lambda_1(\tau(j))\}_{j=1}^{\tilde{k}}$, the expected cost will be minimized when all of $\lambda_1(i) = 0$ (note that the expected cost does not depend on $\tau(i) > \tilde{k}$). Thus, we have

$$\lim_{SNR \to \infty} E[J(\Lambda)] = E\left[\sum_{i=\tilde{k}+1}^{Q} \frac{Q - \tilde{k}}{\lambda_T}\right] = E\left[\frac{(Q - \tilde{k})^2}{\lambda_T}\right] = E\left[\frac{|\Psi|^2}{\lambda_T}\right] \tag{98}$$

Noting that $|\Psi|$ is a Binomially distributed variable with mean $Qp_{unif}$, we arrive at our final expression for the asymptotic cost:

$$\lim_{SNR \to \infty} E[J(\Lambda)] = \frac{(Qp_{unif})^2 + Qp_{unif}(1 - p_{unif})}{\lambda_T} \tag{99}$$

which we now recognize as the minimal cost from [1]. Thus, we have $G(\Lambda) = G(\Lambda_0) = -10\log p_{unif}$, where $\Lambda_0$ is the optimal energy allocation policy that distributes energy equally across the ROI, and $G(\cdot)$ is the gain in expected cost over the exhaustive search.

## REFERENCES

[1] E. Bashan, R. Raich, and A. O. H. III, "Optimal two-stage search for sparse targets using convex criteria," *IEEE Transaction Signal Processing*, vol. 56, pp. 5389–5402, November 2008.

[2] R. Castro, J. Haupt, and R. Nowak, "Compressed sensing vs. active learning," in *2006 International Conference on Acoustics, Speech and Signal Processing*, pp. 820–823.

[3] R. Castro, R. Willett, and R. Nowak, "Faster rates in regression via active learning," *Advances in Neural Information Processing Systems*, vol. 18, p. 179, 2006.

[4] R. Willett, A. Martin, and R. Nowak, "Backcasting: adaptive sampling for sensor networks," in *Proceedings of the 3rd international symposium on Information processing in sensor networks*, pp. 124–133, ACM New York, NY, USA, 2004.

[5] R. Nowak, U. Mitra, and R. Willett, "Estimating inhomogeneous fields using wireless sensor networks," *IEEE Journal on Selected Areas in Communications*, vol. 22, no. 6, pp. 999–1006, 2004.

[6] J. Tropp and A. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Transactions on Information Theory*, vol. 53, no. 12, p. 4655, 2007.

[7] S. Ji, Y. Xue, and L. Carin, "Bayesian compressive sensing," *IEEE Transactions on] Signal Processing, IEEE Transactions on [see also Acoustics, Speech, and Signal Processing*, vol. 56, no. 6, pp. 2346–2356, 2008.

[8] K. Kastella, L. Syst, and M. St Paul, "Discrimination gain to optimize detection and classification," *IEEE Transactions on Systems, Man and Cybernetics, Part A*, vol. 27, no. 1, pp. 112–116, 1997.

[9] C. Kreucher, K. Kastella, and A. Hero Iii, "Sensor management using an active sensing approach," *Signal Processing*, vol. 85, no. 3, pp. 607–624, 2005.

[10] C. Kreucher, K. Kastella, and A. Hero III, "Multitarget tracking using the joint multitarget probability density," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 41, no. 4, pp. 1396–1414, 2005.

[11] M. Tipping, "Sparse Bayesian learning and the relevance vector machine," *Journal of Machine Learning Research*, vol. 1, no. 3, pp. 211–244, 2001.