
Decentralized Multi-Agent Active Search for Sparse Signals: Supplementary Material

Ramina Ghods¹

Arundhati Banerjee¹

Jeff Schneider¹

¹School of Computer Science, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA

ABSTRACT

This appendix includes derivations for the proposed Thompson Sampling algorithms Laplace-TS, SPATS and LATSI, discussion and example on failure mode of Laplace-TS, theoretical proofs for Theorem 1 and Theorem 2 and additional numerical results.

Notation Lowercase and uppercase boldface letters represent column vectors and matrices, respectively. For a matrix \mathbf{A} , its transpose \mathbf{A}^T , and the k th row and ℓ th column entry is $A_{k,\ell}$. For a vector \mathbf{a} , the k th entry is a_k , and the sub-vector containing the i th to j th entries (excluding j) is $\mathbf{a}[i : j]$ (Python notation). The ℓ_1 and ℓ_2 -norm of \mathbf{a} are denoted by $\|\mathbf{a}\|_1$ and $\|\mathbf{a}\|_2$, and $|\mathbf{a}|$ represents its absolute value applied element-wise. The Kronecker product is \otimes , and the trace operator is $\text{tr}(\cdot)$. The $N \times N$ identity matrix is denoted by \mathbf{I}_N . $\text{diag}(\mathbf{a})$ is a square matrix with \mathbf{a} on the main diagonal. \emptyset denotes an empty set and for a set \mathcal{S} , $|\mathcal{S}|$ shows the number of elements in that set. $\mathcal{A} \wedge \mathcal{B}$ depicts the logical AND between two sets \mathcal{A} and \mathcal{B} . We use symbols \triangleq and $\mathbb{1}(\cdot)$ as the symbols for mathematical definition and indicator function, respectively.

1 LAPLACE-TS: MULTI-AGENT THOMPSON SAMPLING WITH SPARSE PRIOR

We will now derive the posterior sampling and design stages of Laplace-TS algorithm introduced in Section 3.2. As detailed in this section, we use a Laplace sparse prior $p(\boldsymbol{\beta}) = \frac{1}{(2b)^n} \exp(-\frac{\|\boldsymbol{\beta}\|_1}{b})$ and likelihood distribution $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}_{t-1})$ for this problem.

Posterior Sampling In this stage (also referred to as inference stage in active learning algorithms), each agent j regularly updates its posterior distribution of $\boldsymbol{\beta}$ given the data set \mathbf{D}_t^j available to it. Using Bayes rule we have:

$$p(\boldsymbol{\beta}|\mathbf{D}_t^j) = p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{y}) = \frac{1}{Z} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) p(\boldsymbol{\beta}).$$

Unfortunately, computing the normalization factor Z above is intractable for a Bayesian likelihood distribution with a Laplace prior. To compute this posterior, we borrow ideas from Figueiredo [2003]. In particular, we substitute the Laplace prior distribution with a zero-mean Gaussian prior per entry $p(\beta_i|\tau_i) = \mathcal{N}(0, \tau_i)$, where its variance τ_i has an exponential hyper prior $p(\tau_i) = \frac{\eta}{2} \exp(-\frac{\eta}{2}\tau_i)$ with $\tau_i \geq 0$.

Integrating out τ_i to compute $p(\beta_i)$, we see that this new prior is equivalent to the original Laplace prior per entry β_i :

$$p(\beta_i) = \int p(\beta_i|\tau_i)p(\tau_i) d\tau_i = \frac{\sqrt{\eta}}{2} \exp(-\sqrt{\eta}|\beta_i|).$$

Here, $\sqrt{\eta} = 1/b$ is the scaling hyperparameter in Laplace distribution and for best performance needs to be tuned given the sparsity rate of the original signal $\boldsymbol{\beta}$. If we were able to observe variance τ_i , the posterior distribution of $\boldsymbol{\beta}$ given data set \mathbf{D}_t^j would become:

$$p(\boldsymbol{\beta}|\mathbf{D}_t^j, \boldsymbol{\tau}) = \frac{1}{Z} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) p(\boldsymbol{\beta}|\boldsymbol{\tau})$$

$$\begin{aligned}
&= \frac{1}{Z} \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_{t-1}) \times \mathcal{N}(0, \text{diag}(\boldsymbol{\tau})) \\
&= \mathcal{N}(\boldsymbol{\mu}_\beta(\boldsymbol{\tau}), \boldsymbol{\Sigma}_\beta(\boldsymbol{\tau})),
\end{aligned} \tag{8}$$

with,

$$\begin{aligned}
\boldsymbol{\Sigma}_\beta(\boldsymbol{\tau}) &= \left((\text{diag}(\boldsymbol{\tau}))^{-1} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} \right)^{-1}, \\
\boldsymbol{\mu}_\beta(\boldsymbol{\tau}) &= \frac{1}{\sigma^2} \boldsymbol{\Sigma}_\beta(\boldsymbol{\tau}) \mathbf{X}^T \mathbf{y}.
\end{aligned} \tag{9}$$

Here, $\boldsymbol{\tau}$ is the vector containing all elements of τ_i for $i = 1, \dots, n$. Since we cannot observe $\boldsymbol{\tau}$, we use Gibbs sampling [Gelfand and Smith, 1990] to iteratively sample $\boldsymbol{\tau}$ and $\boldsymbol{\beta}$ from their conditional posterior distributions. For the given prior, Bae and Mallick [2004] has computed the conditional distributions to iteratively sample $\boldsymbol{\tau}$ and $\boldsymbol{\beta}$ as follows:

$$\begin{aligned}
\boldsymbol{\beta}^* &\sim \mathcal{N}(\boldsymbol{\mu}_\beta(\tilde{\boldsymbol{\tau}}), \boldsymbol{\Sigma}_\beta(\tilde{\boldsymbol{\tau}})) \\
\tilde{\tau}_i^{-1} &\sim \text{InvGauss}\left(\frac{\sqrt{\eta}}{\beta_i^*}, \eta\right), \quad i = 1, \dots, n.
\end{aligned} \tag{10}$$

Design In this stage, agent j computes the expected reward $\lambda^+(\boldsymbol{\beta}^*, \mathbf{D}_t^j, \mathbf{x})$ and optimize it for best sensing action \mathbf{x}_t . Let us use the reward function $\lambda(\boldsymbol{\beta}^*, \mathbf{D}_t^j \cup (\mathbf{x}, y)) = -\|\boldsymbol{\beta}^* - \hat{\boldsymbol{\beta}}(\mathbf{D}_t^j \cup (\mathbf{x}, y))\|_2^2$. In order to compute the expected reward in (3), we need to design an estimator $\hat{\boldsymbol{\beta}}(\mathbf{D}_t^j \cup (\mathbf{x}, y))$ whose expectation we can compute. While there has been many sparse recovery algorithms proposed in the literature [Marques et al., 2018], for many of these well-known thresholding or iterative algorithms (e.g. [Pati et al., 1993, Needell and Tropp, 2010, Blumensath and Davies, 2009, Daubechies et al., 2004, Maleki, 2011]) computing the expectation in (3) is not tractable. Alternatively, we propose using maximum a posteriori (MAP) estimate of the posterior $p(\boldsymbol{\beta} | \mathbf{D}_t^j, \boldsymbol{\tau})$ in (8) where we estimate $\boldsymbol{\tau}$ using Expectation-Maximization (EM) proposed by Figueiredo [2003] with the E-step and M-step for $p = 1, \dots, P$ iterations as follows:

$$\begin{aligned}
\text{E-step} : \boldsymbol{\tau}^{(p)} &= |\hat{\boldsymbol{\beta}}^{(p-1)}|/\eta \\
\text{M-step} : \hat{\boldsymbol{\beta}}^{(p)} &= \boldsymbol{\mu}_\beta(\boldsymbol{\tau}^{(p)}).
\end{aligned} \tag{11}$$

Since many elements of $\hat{\boldsymbol{\beta}}^{(p)}$ and consequently $\boldsymbol{\tau}^{(p)}$ are expected to be zero, to avoid inverting $\text{diag}(\boldsymbol{\tau})$ we can rewrite the M-step as follows: $\hat{\boldsymbol{\beta}}^{(p)} = \frac{1}{\sigma^2} \boldsymbol{\Sigma}^{(p)} \mathbf{X}^T \mathbf{y}$, with $\boldsymbol{\Sigma}^{(p)} = \mathbf{U} (\mathbf{I}_n + \frac{1}{\sigma^2} \mathbf{U} \mathbf{X}^T \mathbf{X} \mathbf{U})^{-1} \mathbf{U}$ and $\mathbf{U} = (\text{diag}(\boldsymbol{\tau}^{(p)}))^{1/2}$.

Since the posterior is Gaussian, its MAP estimate is the mean of the distribution, i.e. $\hat{\boldsymbol{\beta}}^{(P)}$ in (11). Adding (\mathbf{x}, y) to the data set \mathbf{D}_t^j , our MAP estimate becomes:

$$\hat{\boldsymbol{\beta}}(\mathbf{D}_t^j \cup (\mathbf{x}, y)) = \underbrace{\mathbf{U} \left(\sigma^2 \mathbf{I}_n + \mathbf{U} \begin{bmatrix} \mathbf{X}^T & \mathbf{x} \end{bmatrix} \begin{bmatrix} \mathbf{X} \\ \mathbf{x}^T \end{bmatrix} \mathbf{U} \right)^{-1} \mathbf{U} \begin{bmatrix} \mathbf{X}^T & \mathbf{x} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ y \end{bmatrix}}_{\triangleq \mathbf{Q}}.$$

Now, we can derive the expected reward in (3) as follows:

$$\begin{aligned}
\lambda^+(\boldsymbol{\beta}^*, \mathbf{D}_t^j, \mathbf{x}) &= \mathbb{E}_{y|\mathbf{x}, \boldsymbol{\beta}^*} \left[-\|\boldsymbol{\beta}^* - \hat{\boldsymbol{\beta}}(\mathbf{D}_t^j \cup (\mathbf{x}, y))\|_2^2 \right] = -\mathbb{E}_{y|\mathbf{x}, \boldsymbol{\beta}^*} \left[\|\mathbf{Q} \mathbf{X}^T \mathbf{y} + \mathbf{Q} \mathbf{x} y - \boldsymbol{\beta}^*\|_2^2 \right] \\
&= -\|\mathbf{Q} \mathbf{X}^T \mathbf{y} - \boldsymbol{\beta}^*\|_2^2 - \|\mathbf{Q} \mathbf{x}\|_2^2 \mathbb{E}_{y|\mathbf{x}, \boldsymbol{\beta}^*} [y^2] - 2 (\mathbf{Q} \mathbf{X}^T \mathbf{y} - \boldsymbol{\beta}^*)^T \mathbf{Q} \mathbf{x} \mathbb{E}_{y|\mathbf{x}, \boldsymbol{\beta}^*} [y] \\
&= -\|\mathbf{Q} \mathbf{X}^T \mathbf{y} - \boldsymbol{\beta}^*\|_2^2 - \|\mathbf{Q} \mathbf{x}\|_2^2 (\sigma^2 + (\mathbf{x}^T \boldsymbol{\beta}^*)^2) - 2 (\mathbf{Q} \mathbf{X}^T \mathbf{y} - \boldsymbol{\beta}^*)^T \mathbf{Q} \mathbf{x} \boldsymbol{\beta}^*.
\end{aligned} \tag{12}$$

With the posterior sampling and design stage developed, we can now run Algorithm 1 for the active search problem in Section 2. Let us call this algorithm Laplace-TS. In the next section, we will evaluate the performance of this algorithm.

2 FAILING PERFORMANCE OF LAPLACE-TS IN SINGLE AGENT SETTING

Figure 7 illustrates simulation results of the full recovery performance of Laplace-TS compared to a simple point algorithm that exhaustively searches the entire action space one location at a time with one agent. Here, we are plotting full recovery rate of vector $\boldsymbol{\beta}$ with size $n = 128$ as a function of number of measurements T for two sparsity rates of $k = 1$ and $k = 5$. The curves show mean and standard error of full recovery over 50 random trials. This figure demonstrates that, unfortunately,

Laplace-TS with one agent leads to poor performance that is on par with a point-wise algorithm that exhaustively searches all locations one at a time.

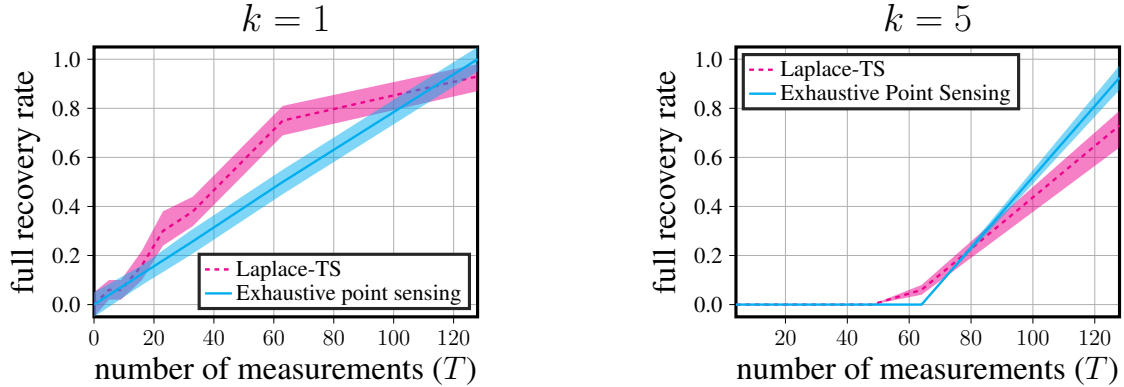


Figure 7: Full recovery rate of Laplace-TS with single agent for sparsity rate $k = 1$

We can associate this poor performance with one of the failure modes of TS discussed in Sec. 8.2 of the tutorial by Russo et al. [2018]. According to this tutorial, TS faces a dilemma when solving certain kinds of active learning problems. One such scenario are problems that require a careful assessment of information gain. In general, by optimizing the expected reward, TS always restricts its actions to those that have a chance in being optimal. However, in certain active learning problems, suboptimal actions can carry additional information regarding the parameter of interest.

We provide the following example to better understand the situation in our problem formulation. Let us assume there is no noise in the sensing model ($\epsilon_t = 0$) and that there is only one non-zero element ($k = 1$). Under such conditions, the problem amounts to finding the location of the non-zero element in β which we call \tilde{i} . Thus, for every action \mathbf{x}_t , the observation $y_t = \mathbf{x}_t^T \beta$ is non-zero if \tilde{i} is in the support of \mathbf{x}_t and is zero, otherwise. Clearly, a binary search can locate \tilde{i} in $\log(n)$ steps. TS, however, in each time step estimates the sample β^* with $\hat{\beta}$ and picks the sensing action $\mathbf{x}_t = \hat{\beta}$. Or, in case of perfect estimation, we will have $\mathbf{x}_t = \hat{\beta} = \beta^*$. So, unless β^* is the true β , observation \mathbf{y}_t would be zero and TS will only manage to eliminate support of β^* from the list of possibly correct supports. Therefore, TS ends up eliminating wrong supports one location at a time which explains its on par performance to an exhaustive point sensor. Similarly, for $k > 1$ TS is always limited to picking sensing actions \mathbf{x}_t that are in the support of the sparse estimates of the samples. Adding the region sensing constraint will only aggravate this problem by further shrinking this support set. In the next section, we will modify Laplace-TS and propose two algorithms that can bypass this failure mode.

3 SPATS: THOMPSON SAMPLING WITH BLOCK SPARSE PRIOR

In what follows, we will derive the posterior sampling and design stages for SPATS algorithm proposed in Section 4.1. SPATS is an asynchronous multi-agent TS algorithm using a block sparse prior for the problem formulation in Section 2.

As discussed in Section 4.1, we use the following prior and likelihood distributions. The likelihood function is $p(\mathbf{y}|\mathbf{X}, \beta) = \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}_{t-1})$ from (2). For prior distribution, as discussed in Section 4.1 we use a block sparse prior $p(\beta) = \mathcal{N}(\mathbf{0}_{n \times 1}, \Sigma_0)$, where:

$$\Sigma_0 = \text{diag}(\gamma) \otimes \mathbf{B}, \quad (13)$$

with $\gamma \in \mathbb{R}^M$ and $\mathbf{B} \in \mathbb{R}^{L \times L}$ ($M = n/L$) as hyperparameters.

Posterior Sampling Similar to (8) and (9) in computing the posterior distribution for agent j , the posterior with a Gaussian likelihood and Gaussian prior becomes:

$$p(\beta|\mathbf{D}_t^j, \gamma, \mathbf{B}) = \mathcal{N}(\mu_\beta(\gamma, \mathbf{B}), \Sigma_\beta(\gamma, \mathbf{B})), \quad (14)$$

with,

$$\Sigma_\beta(\gamma, \mathbf{B}) = ((\Sigma_0)^{-1} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X})^{-1},$$

$$\boldsymbol{\mu}_\beta(\boldsymbol{\gamma}, \mathbf{B}) = \frac{1}{\sigma^2} \boldsymbol{\Sigma}_\beta(\boldsymbol{\gamma}, \mathbf{B}) \mathbf{X}^\top \mathbf{y}.$$

Using Expectation-Maximization proposed by Zhang and Rao [2011], we can estimate the hyperparameters $\boldsymbol{\gamma}$ and \mathbf{B} for $p = 1, \dots, P$ iterations as follows.

$$\begin{aligned} \mathbf{E}\text{-step} : \quad & \boldsymbol{\mu}_\beta^{(p)} = \boldsymbol{\mu}_\beta(\boldsymbol{\gamma}^{(p-1)}, \mathbf{B}^{(p-1)}), \\ & \boldsymbol{\Sigma}_\beta^{(p)} = \boldsymbol{\Sigma}_\beta(\boldsymbol{\gamma}^{(p-1)}, \mathbf{B}^{(p-1)}) \\ \mathbf{M}\text{-step} : \quad & \gamma_m^{(p)} = \frac{\text{tr}(\mathbf{B}^{-1}(\boldsymbol{\Sigma}_\beta^m + \boldsymbol{\mu}_\beta^m (\boldsymbol{\mu}_\beta^m)^\top))}{L}, \quad (m = 1, \dots, M) \\ & \mathbf{B}^{(p)} = \frac{1}{M} \sum_{m=1}^M \frac{\boldsymbol{\Sigma}_\beta^m + \boldsymbol{\mu}_\beta^m (\boldsymbol{\mu}_\beta^m)^\top}{\gamma_m} \end{aligned} \quad (15)$$

with $\boldsymbol{\mu}_\beta^m = \boldsymbol{\mu}_\beta^{(p)}[(m-1) \times L : m \times L]$, and $\boldsymbol{\Sigma}_\beta^m = \boldsymbol{\Sigma}_\beta^{(p)}[(m-1) \times L : m \times L, (m-1) \times L : m \times L]$.

Design Next, for the estimator of agent j , $\hat{\boldsymbol{\beta}}(\mathbf{D}_t^j)$, we use the MAP estimator given by $\boldsymbol{\mu}_\beta^{(P)}$. Particularly, we have:

$$\hat{\boldsymbol{\beta}}(\mathbf{D}_t^j \cup (\mathbf{x}, y)) = \underbrace{\left(\sigma^2 \boldsymbol{\Sigma}_0^{-1} + [\mathbf{X}^\top \ \mathbf{x}] \begin{bmatrix} \mathbf{X} \\ \mathbf{x}^\top \end{bmatrix} \right)^{-1}}_{\mathbf{Q}} [\mathbf{X}^\top \ \mathbf{x}] \begin{bmatrix} \mathbf{y} \\ y \end{bmatrix} \quad (16)$$

with $\boldsymbol{\Sigma}_0 = \text{diag}(\boldsymbol{\gamma}^{(P)}) \otimes \mathbf{B}^{(P)}$. Hence, the reward function for this estimator is given by (12) with \mathbf{Q} from (16).

For best results, we need to initialize parameter $\boldsymbol{\Sigma}_0$ relative to signal power.

4 PROOF FOR THEOREM 1

Since this theorem focuses on single agent settings, both proposed TS algorithms in the theorem have the theoretical guarantees as well as the limits of TS bounds as summarized in [Russo et al., 2018, Ch. 8]. As discussed in this chapter, there are different theoretical approaches proposed in literature for regret bounds on TS. For this proof, we use regret bounds computed via information theory proposed in Russo and Van Roy [2016]. According to this paper, the resulting bounds better describe the benefits of prior distribution which helps point to failure modes of TS and potential solutions. Defining π_t as the distribution of actions for an online optimization policy at time t , Russo and Van Roy [2016] shows that for this policy, the expected regret defined in Theorem 1 is bounded by:

$$\mathbb{E}[\text{Reg}(T)] \leq \sqrt{T \bar{\Psi}_t(\pi) \mathbf{H}(\mathbf{x}^*)}, \quad (17)$$

where, $\mathbf{H}(\mathbf{x}^*)$ is the entropy for optimal action \mathbf{x}^* and $\bar{\Psi}_t(\pi)$ is the average expected information ratio defined as follows.

$$\bar{\Psi}_t(\pi) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\pi_t} [\Psi_t(\pi_t)] \quad \& \quad \Psi_t(\pi_t) \triangleq \frac{\Delta_t(\pi_t)^2}{\mathcal{I}_t(\pi_t)}, \quad (18)$$

where, $\Delta_t(\pi_t)$ and $\mathcal{I}_t(\pi_t)$ are the single-period expected regret and information gain over policy π_t at time step t defined below. For the single-period expected regret we have:

$$\Delta_t(\pi_t) = \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \Delta_t(\mathbf{x}) \quad \& \quad \Delta_t(\mathbf{x}) = \mathbb{E}[\mathcal{R}_t(\mathbf{x}^*, \boldsymbol{\beta}) - \mathcal{R}_t(\mathbf{x}, \boldsymbol{\beta})], \quad (19)$$

where, $\mathcal{R}_t(\mathbf{x}, \boldsymbol{\beta})$ is the reward function for action \mathbf{x} at time t . And, for information gain, we have:

$$\mathcal{I}_t(\pi_t) = \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x}), \quad (20)$$

where, $\mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x})$ is the mutual information between optimal action \mathbf{x}^* and observation $y = (\mathbf{x}^T \boldsymbol{\beta})^2$ at time t .

To compute the regret bound in (17), we need to first focus on computing single-period expected regret and information gain in (19) and (20). Computing these terms for a general problem formulation can be difficult. However, for an online optimization problem with a finite prior distribution and finite action set, [Russo and Van Roy, 2017, Algorithm 1] has provided the necessary formulation to compute these terms. Since our problem formulation falls under the category of this framework, we can take advantage of their formulations as follows. For single-period expected regret, we have:

$$\Delta_t(\mathbf{x}) = \mathcal{R}(\mathbf{x}^*) - \sum_{\beta} p(\beta) \mathcal{R}(\mathbf{x}, \beta), \quad (21)$$

and, for the mutual information we have:

$$\mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x}) = \sum_{\mathbf{x}^*, y} p(\mathbf{x}^*, y|\mathbf{x}) \log \left(\frac{p(\mathbf{x}^*, y|\mathbf{x})}{p(\mathbf{x}^*)p(y|\mathbf{x})} \right), \quad (22)$$

where,

$$\begin{aligned} p(\mathbf{x}^*) &= \sum_{\beta \in \mathcal{B}(\mathbf{x}^*)} p(\beta | \mathbf{D}_t^1), \\ p(y|\mathbf{x}) &= \sum_{\beta} p(\beta | \mathbf{D}_t^1) p(y|\mathbf{x}, \beta), \\ p(\mathbf{x}^*, y|\mathbf{x}) &= p(\mathbf{x}^*) \sum_{\beta \in \mathcal{B}(\mathbf{x}^*)} p(y|\mathbf{x}, \beta), \end{aligned} \quad (23)$$

with $\mathcal{B}(\mathbf{x}) = \{\beta : \mathbf{x} = \arg \max_{\tilde{\mathbf{x}}} \mathbb{E}_{y|\tilde{\mathbf{x}}, \beta} [\mathcal{R}(\tilde{\mathbf{x}}, \beta)]\}$ as the set containing all vectors β for which action \mathbf{x} is optimal.

Next, we will compute this regret bound for our TS algorithms with 1-sparse and block sparse prior introduced in Theorem 1. Our emphasis here is on understanding the benefits of a block sparse prior distribution. Hence, to simplify the discussion we assume both algorithms are aware of the true 1-sparse prior for β .

1-SPARSE PRIOR

Let us here establish two parameters i^* and \tilde{i} that we will use throughout the proof. 1) Since the true parameter β is 1-sparse, we assume the location of the non-zero element of true β is i^* . Given the uniform 1-sparse prior on true β , i^* can be either of indices 1 through n with probability $1/n$. 2) Since our TS algorithm here assumes a 1-sparse prior distribution, the only feasible sensing actions in \mathcal{X} are 1-sparse with one nonzero element at location \tilde{i} .

With two parameters i^* and \tilde{i} , we can now compute the single-period expected regret in (21) for $t = 1$ as:

$$\Delta_1(\mathbf{x}) = \Delta_1(\tilde{i}) = 1 - \sum_{i^*=1}^n 1/n \times \mathbb{1}(\tilde{i} = i^*),$$

Hence, the single-period expected regret $\Delta_1(\pi_1)$ over our policy becomes:

$$\Delta_1(\pi_1) = \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \Delta_1(\mathbf{x}) \stackrel{(a)}{=} \sum_{\tilde{i}=1}^n 1/n \Delta_1(\tilde{i}) = \frac{n-1}{n}$$

Here, (a) follows from the fact that our TS algorithm at $t = 1$ will pick all sensing actions with the nonzero element $j = 1$ through $j = n$ equally likely. Next, we need to compute single-period expected regret in (21) for $t = 2$. For $t = 2$, there are two policies for our TS algorithm depending on the action \mathbf{x}_1 chosen in previous time step. If $\mathbf{x}_1 = \mathbf{x}^*$, then our TS algorithm will be exploiting the same optimal action at time $t = 2$ and set $\mathbf{x}_2 = \mathbf{x}^*$ for which $\Delta_2(\mathbf{x}) = \Delta_2(\pi_2) = 0$. The probability of finding \mathbf{x}^* at $t = 1$ is $1/n$. However, if $\mathbf{x}_1 \neq \mathbf{x}^*$, then our TS algorithm will be exploring the $n - 1$ feasible actions that are not observed. Let us assume \tilde{i}_1 shows the location of nonzero element in \mathbf{x}_1 , then:

$$\Delta_2(\mathbf{x}) = \Delta_2(\tilde{i}) = 1 - \sum_{i^*=1, i^* \neq \tilde{i}_1}^n \frac{1}{n-1} \times \mathbb{1}(\tilde{i} = i^*),$$

Hence, for both policies the single-period expected regret $\Delta_2(\pi_2)$ becomes:

$$\Delta_2(\pi_2) = \begin{cases} 0, & \text{prob} = \frac{1}{n} \\ \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \Delta_2(\mathbf{x}) = \sum_{\tilde{i}=1, \tilde{i} \neq \tilde{i}_1}^n \frac{1}{n-1} \Delta_2(\tilde{i}) = \frac{n-2}{n-1}, & \text{prob} = \frac{n-1}{n} \end{cases}$$

Similarly, we can conclude that for any $t > 2$, we will have:

$$\Delta_t(\pi_t) = \begin{cases} 0, & \text{prob} = \frac{t-1}{n} \\ \frac{n-t}{n-t+1}, & \text{prob} = \frac{n-t+1}{n} \end{cases} \quad (24)$$

We will next compute the information gain in (22) for time step t . Recall that our TS algorithm have two policies. One policy π_t is to keep exploiting the optimal action if it has been discovered in the previous $t-1$ measurements \mathbf{D}_t^1 . Since $\Delta_t(\pi_t) = 0$ for this policy, we have $\Psi_t(\pi_t) = \frac{\Delta_t(\pi_t)^2}{\mathcal{I}_t(\pi_t)} = 0$ and hence there is no need to compute the information term $\mathcal{I}_t(\pi_t)$ for this scenario. The second policy of our TS method is used if optimal action \mathbf{x}^* has not been found in the previous $t-1$ measurements. For this scenario, at time t the TS algorithm picks the sensing action \mathbf{x}_t out of $n-t+1$ measurements that have not been observed yet. For this policy, we can compute the probabilities in (23) as follows. First, for $p(\mathbf{x}^*)$, we have:

$$p(\mathbf{x}^*) = \sum_{\beta \in \mathcal{B}(\mathbf{x}^*)} p(\beta = \mathbf{x}^* | \mathbf{D}_t^1) = \begin{cases} \frac{1}{n-t+1}, & \text{if } \mathbf{x}^* \notin \mathbf{D}_t^1 \wedge \mathbf{x}^* \in \mathcal{X} \\ 0, & \text{otherwise} \end{cases}$$

Second, for $p(y|\mathbf{x})$ we have:

$$\begin{aligned} p(y|\mathbf{x}) &= p(y|\tilde{i}) = \begin{cases} \sum_{i^*=1, i^* \neq \{\tilde{i}_1, \dots, \tilde{i}_{t-1}\}}^n \frac{1}{n-t+1} p(i^* \neq \tilde{i}), & \text{if } y = 0, \\ \sum_{i^*=1, i^* \neq \{\tilde{i}_1, \dots, \tilde{i}_{t-1}\}}^n \frac{1}{n-t+1} p(i^* = \tilde{i}), & \text{if } y = 1, \end{cases} \\ &= \begin{cases} \frac{n-t}{n-t+1} & \text{if } y = 0, \\ \frac{1}{n-t+1} & \text{if } y = 1, \\ 0 & \text{otherwise,} \end{cases} \end{aligned}$$

where, $\{\tilde{i}_1, \dots, \tilde{i}_{t-1}\}$ denotes the location of the nonzero elements in the previous $t-1$ measurements \mathbf{D}_t^1 . Third, for $p(\mathbf{x}^*, y|\mathbf{x})$ we get:

$$p(\mathbf{x}^*, y|\mathbf{x}) = \frac{1}{n-t+1} \sum_{\beta=\mathbf{x}^*} p(y|\mathbf{x}, \beta) = \begin{cases} \frac{1}{n-t+1} & \text{if } y = 0 \wedge \mathbf{x}^* \neq \mathbf{x}, \\ \frac{1}{n-t+1} & \text{if } y = 1 \wedge \mathbf{x}^* = \mathbf{x}, \\ 0 & \text{otherwise.} \end{cases}$$

Using these three probability densities, the information gain in (22) becomes:

$$\begin{aligned} \mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x}) &= \sum_{\substack{\mathbf{x}^* \in \mathcal{X}, \mathbf{x}^* \notin \mathbf{D}_t^1 \\ \mathbf{x}^* \neq \mathbf{x}, y=0}} \frac{1}{n-t+1} \log \left(\frac{\frac{1}{n-t+1}}{\frac{1}{n-t+1} \times \frac{n-t}{n-t+1}} \right) \\ &+ \sum_{\substack{\mathbf{x}^* \in \mathcal{X}, \mathbf{x}^* \notin \mathbf{D}_t^1 \\ \mathbf{x}^* = \mathbf{x}, y=1}} \frac{1}{n-t+1} \log \left(\frac{\frac{1}{n-t+1}}{\frac{1}{n-t+1} \times \frac{1}{n-t+1}} \right) \\ &= \frac{n-t}{n-t+1} \log \left(\frac{n-t+1}{n-t} \right) + \frac{1}{n-t+1} \log(n-t+1), \end{aligned}$$

resulting in:

$$\mathcal{I}_t(\pi_t) = \sum_{\mathbf{x} \in \mathcal{X}, \mathbf{x}^* \notin \mathbf{D}_t^1} \frac{1}{n-t+1} \mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x}) = \mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x}). \quad (25)$$

Now, using (24) and (25), we can compute the average expected information ratio in (18) as:

$$\bar{\Psi}_t(\pi) = \frac{1}{T} \sum_{t=1}^T \left[\frac{t-1}{n} \times 0 + \frac{n-t+1}{n} \times \frac{\left(\frac{n-t}{n-t+1} \right)^2}{\frac{n-t}{n-t+1} \log \left(\frac{n-t+1}{n-t} \right) + \frac{1}{n-t+1} \log(n-t+1)} \right],$$

which together with an upperbound $H(\mathbf{x}^*) < \log(|\mathcal{X}|)$ and (17) gives the following regret bound:

$$\mathbb{E}[\text{Reg}(T)] \leq \left(\log(|\mathcal{X}|) \sum_{t=1}^T \frac{\frac{(n-t)^2}{n(n-t+1)}}{\frac{n-t}{n-t+1} \log\left(\frac{n-t+1}{n-t}\right) + \frac{1}{n-t+1} \log(n-t+1)} \right)^{1/2}.$$

Lastly, the regret bound we computed above is only applicable when $T \leq n$. Since there are n 1-sparse actions available for a vector β with length n , after $T = n$ actions the algorithm has surely found the optimal action. As a result, $\Delta_t(\pi_t) = 0$ for $t \geq n$ which gives us the regret bound:

$$\mathbb{E}[\text{Reg}(T)] \leq \left(\log(|\mathcal{X}|) \sum_{t=1}^{\min\{T, n-1\}} \frac{\left(1 - \frac{t}{n}\right) \left(1 - \frac{1}{n-t+1}\right)}{\left(\frac{n-t-1}{n-t} \log\left(\frac{n-t}{n-t-1}\right) + \frac{1}{n-t} \log(n-t)\right)} \right)^{1/2} \quad (26)$$

BLOCK SPARSE PRIOR

We will again start with establishing two parameters i^* and m_t that we will use throughout the proof. 1) We assume the location of the non-zero element of true β is i^* . Given the uniform 1-sparse prior on true β , i^* can be either of indices 1 through n with probability $1/n$. 2) Since our TS algorithm here assumes a 1-block sparse prior distribution, the feasible sensing actions in \mathcal{X} at time step t are $\mathbf{x} = [\gamma_1 \mathbf{b}, \dots, \gamma_{M_t} \mathbf{b}]^T$ where only one of parameters γ_1 through γ_{M_t} is nonzero and $\mathbf{b} = [\frac{1}{\sqrt{L_t}}, \dots, \frac{1}{\sqrt{L_t}}]$ are the blocks with length L_t . We call the location of the nonzero block m_t which can be either of indices 1 through $M_t = n/L_t$ with probability $1/M_t$ at time step t .

With two parameters i^* and m_t , we can now compute the single-period expected regret $\Delta_1(\pi_1)$ in (19) for $t = 1$ as:

$$\begin{aligned} \Delta_1(\pi_1) &= \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \Delta_1(\mathbf{x}) \stackrel{(b)}{=} \sum_{m_1=1}^{M_1} \frac{1}{M_1} \Delta_1(m_1) \\ &\stackrel{(c)}{=} \sum_{m_1=1}^{M_1} \frac{1}{M_1} \left(1 - \sum_{i^*=1}^n \frac{1}{n} \times \mathbb{1}(i^* \in \{m_1 L_1 + 1, \dots, m_1 L_1 + L_1\}) \times \frac{1}{L_1} \right) = \frac{n-1}{n}, \end{aligned}$$

here, (b) follows from the fact that our TS algorithm at $t = 1$ will pick all sensing actions with the nonzero blocks $m_1 = 1$ through $m_1 = M_1$ equally likely, and (c) follows from (21). For $t > 2$, there are two policies depending on whether i^* has been discovered to belong to any of the elements of block m_1 or not. If i^* does not belong to block m_1 , then:

$$\begin{aligned} \Delta_2(\pi_2) &= \sum_{\substack{\mathbf{x} \in \mathcal{X} \\ \mathbf{x} \notin \mathbf{D}_t^1}} \pi(\mathbf{x}) \left(1 - \sum_{\substack{i^*=1 \\ i^* \notin \{m_1 L_1 + 1, \dots, m_1 L_1 + L_1\}}}^n \frac{1}{n - L_1} \times \mathbb{1}(i^* \in \{m_2 L_2 + 1, \dots, m_2 L_2 + L_2\}) \times \frac{1}{L_2} \right) \\ &= 1 - \frac{1}{n - L_1}. \end{aligned}$$

Or if i^* does belong to block m_1 , similarly we $\Delta_2(\pi_2) = 1 - \frac{1}{L_1}$. For block length $L_1 = n/2$, we will have the regret $\Delta_2(\pi_2) = 1 - \frac{1}{n/2}$ for both policies which is the smallest regret value. Similarly, for any $t > 2$, we have:

$$\Delta_t(\pi_t) = 1 - \frac{1}{n - \sum_{t'=1}^{t-1} \frac{n}{2^{t'}}}, \quad (27)$$

as long as $L_t = \frac{n}{2^t}$.

We will next compute the information gain in (22) for time step $t = 1$. For that, we will first compute the three probability densities in (23). For $p(\mathbf{x}^*)$, we have:

$$p(\mathbf{x}^*) = \sum_{\beta \in \mathcal{B}(\mathbf{x}^*)} p(\beta = \mathbf{x}^*) = \begin{cases} \frac{1}{n}, & \text{if } \mathbf{x}^* \text{ is 1-sparse} \\ 0, & \text{otherwise} \end{cases}$$

which follows the fact that the optimal action is 1-sparse. Second, for $p(y|\mathbf{x})$ we have:

$$p(y|\mathbf{x}) = p(y|m_1) = \begin{cases} \sum_{i^*=1}^n \frac{1}{n} p(i^* \notin \{m_1 L_1 + 1, \dots, m_1 L_1 + L_1\}), & \text{if } y = 0, \\ \sum_{i^*=1}^n \frac{1}{n} p(i^* \in \{m_1 L_1 + 1, \dots, m_1 L_1 + L_1\}), & \text{if } y = 1/L_1, \end{cases}$$

$$= \begin{cases} \frac{n-L_1}{n} & \text{if } y = 0, \\ \frac{L_1}{n} & \text{if } y = 1/L_1, \\ 0 & \text{otherwise,} \end{cases}$$

Third, for $p(\mathbf{x}^*, y|\mathbf{x})$ we get:

$$p(\mathbf{x}^*, y|\mathbf{x}) = \frac{1}{n} \sum_{\beta=\mathbf{x}^*} p(y|\mathbf{x}, \beta) = \begin{cases} \frac{1}{n} & \text{if } y = 0 \wedge \mathbf{x}^T \mathbf{x}^* = 0, \\ \frac{1}{n} & \text{if } y = 1/L_1 \wedge \mathbf{x}^T \mathbf{x}^* \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Using these three probability densities, the information gain in (22) becomes:

$$\begin{aligned} \mathcal{I}_1(\mathbf{x}^*; y|\mathbf{x}) &= \sum_{\substack{\mathbf{x}^* \in \mathcal{X} \\ \mathbf{x}^T \mathbf{x}^* \neq 0, y=1/L_1}} \frac{1}{n} \log \left(\frac{\frac{1}{n}}{\frac{1}{n} \times \frac{L_1}{n}} \right) + \sum_{\substack{\mathbf{x}^* \in \mathcal{X} \\ \mathbf{x}^T \mathbf{x}^* = 0, y=0}} \frac{1}{n} \log \left(\frac{\frac{1}{n}}{\frac{1}{n} \times \frac{n-L_1}{n}} \right) \\ &= \frac{L_1}{n} \log \left(\frac{n}{L_1} \right) + \frac{n-L_1}{n} \log \left(\frac{n}{n-L_1} \right). \end{aligned}$$

The information term above is maximized when $L_1 = n/2$ which gives $\mathcal{I}_1(\mathbf{x}^*; y|\mathbf{x}) = \log(2)$.

Next, let us discuss information gain for $t = 2$. Recall that for $t = 2$, there are two policies depending on whether i^* has been discovered to belong to elements of block m_1 or not. Similar computations to that of $t = 1$ shows that for the two policies we have:

$$\mathcal{I}_2(\mathbf{x}^*; y|\mathbf{x}) = \begin{cases} \frac{L_2}{L_1} \log \left(\frac{L_1}{L_2} \right) + \frac{L_1-L_2}{L_1} \log \left(\frac{L_1}{L_1-L_2} \right), & \text{if } i^* \in \{m_1 L_1 + 1, \dots, m_1 L_1 + L_1\} \\ \frac{L_2}{n-L_1} \log \left(\frac{n-L_1}{L_2} \right) + \frac{n-L_1-L_2}{n-L_1} \log \left(\frac{n-L_1}{n-L_1-L_2} \right), & \text{if } i^* \notin \{m_1 L_1 + 1, \dots, m_1 L_1 + L_1\} \end{cases}$$

By setting $L_1 = n/2$ and $L_2 = n/4$, both policies will be maximized to give $\mathcal{I}_2(\mathbf{x}^*; y|\mathbf{x}) = \log(2)$. Similarly, for any $t > 2$, the information gain $\mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x})$ is maximized with:

$$\mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x}) = \log(2), \quad (28)$$

as long as $L_t = \frac{n}{2^t}$. The fact that the varying block length $L_t = \frac{n}{2^t}$ maximizes information gain $\mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x})$ and minimizes regret term $\Delta_t(\pi_t)$ shows that this varying block length is the optimal approach for our TS algorithm finding a 1-sparse signal of interest. This result in essence describes our strategy for varying block length $L_t = \frac{L_{t-1}}{2}$ in SPATS algorithm.

Finally, using (27) and (28) with the upperbound $H(\mathbf{x}^*) < \log(|\mathcal{X}|)$, we can compute the expected regret in (17) as:

$$\mathbb{E}[\text{Reg}(T)] \leq \left(\log(|\mathcal{X}|) \sum_{t=1}^T \left[\frac{\left(1 - \frac{1}{n - \sum_{t'=1}^{t-1} \frac{n}{2^{t'}}}\right)^2}{\log(2)} \right] \right)^{1/2},$$

Lastly, the regret bound we computed above is only applicable when $T \leq \log_2(n)$. Since the algorithm has a varying block length of $L_t = \frac{n}{2^t}$, after $T = \log_2(n)$ the block length is reduced to 1 which certainly includes nonzero element of β . As a result, after $T = \log_2(n)$ actions the algorithm has surely found the optimal action. This result means that $\Delta_t(\pi_t) = 0$ for $t > \log_2(n)$ which gives us the regret bound:

$$\mathbb{E}[\text{Reg}(T)] \leq \left(\log(|\mathcal{X}|) \sum_{t=1}^{\min\{T, \log_2(n)\}} \left(1 - \frac{1}{n - \left(\sum_{t'=1}^{t-1} \frac{n}{2^{t'}}\right)}\right)^2 / \log(2) \right)^{1/2}. \quad (29)$$

5 PROOF FOR THEOREM 2

Unlike Theorem 1, we use a direct approach to compute expected regret in this proof. Similar to Appendix 4, to simplify the discussion, we assume both algorithms are aware of the prior distribution on true parameter β . We will first compute the expected regret for the single agent algorithm.

SINGLE AGENT

Since the true parameter β is 1-sparse, we assume the location of the non-zero element is i^* . Recall the equation for expected regret:

$$\mathbb{E}[\text{Reg}(T)] = \mathbb{E}_{i^*, \mathbf{x}_1, \dots, \mathbf{x}_T} \left[\sum_{t=1}^T [\mathcal{R}(\mathbf{x}^*, \beta) - \mathcal{R}(\mathbf{x}_t, \beta)] \right] \stackrel{(a)}{=} \sum_{t=1}^T \mathcal{R}(\mathbf{x}^*, \beta) - \sum_{t=1}^T \mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta)], \quad (30)$$

where, (a) follows from assuming i^* is a fixed unknown location without loss of generality. To compute the first term $\mathcal{R}(\mathbf{x}^*, \beta)$, note that the optimal action \mathbf{x}^* is the 1-sparse vector with i^* as the nonzero element, i.e. $\mathcal{R}(\mathbf{x}^*, \beta) = 1$.

To compute the second term $\mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta)]$, remember that \mathbf{x}_t for TS is selected by maximizing the reward function given the available measurements \mathbf{D}_t^1 . For a single agent setting, as depicted in Figure 2a, we have $\mathbf{D}_t^1 = \{(\mathbf{x}_{t'}, y_{t'}) | t' = 1, \dots, t-1\}$. Since our algorithm has a 1-sparse prior, only feasible sensing actions are 1-sparse. If $\mathbf{x}^* \in \mathbf{D}_t^1$, then TS will be exploiting the same optimal action at time t and set $\mathbf{x}_t = \mathbf{x}^*$ for which $\mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta) | \mathbf{x}^* \in \mathbf{D}_t^1] = \mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta) | \mathbf{x}_t = \mathbf{x}^*] = 1$. On the other hand, if $\mathbf{x}^* \notin \mathbf{D}_t^1$, then TS will select action \mathbf{x}_t out of $n - t + 1$ actions that are not observed as part of \mathbf{D}_t^1 . Hence, the reward $\mathcal{R}(\mathbf{x}_t, \beta)$ is 1 with probability $\frac{1}{n-t+1}$ and 0 otherwise, i.e. $\mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta) | \mathbf{x}^* \notin \mathbf{D}_t^1] = \frac{1}{n-t+1}$.

Putting the information for the first and second term in (30), we have:

$$\begin{aligned} \mathbb{E}[\text{Reg}(T)] &= \sum_{t=1}^T \mathcal{R}(\mathbf{x}^*, \beta) - \sum_{t=1}^T \mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta)] \\ &= T - \sum_{t=1}^T \left[p(\mathbf{x}^* \in \mathbf{D}_t^1) \times \mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta) | \mathbf{x}^* \in \mathbf{D}_t^1] \right. \\ &\quad \left. + p(\mathbf{x}^* \notin \mathbf{D}_t^1) \times \mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta) | \mathbf{x}^* \notin \mathbf{D}_t^1] \right] \\ &\stackrel{(b)}{=} T - \sum_{t=1}^T \left[\frac{t-1}{n} \times 1 + \frac{n-t+1}{n} \times \frac{1}{n-t+1} \right] \\ &= T - \frac{T(T+1)}{2n}, \end{aligned}$$

where, (b) follows from $p(\mathbf{x}^* \in \mathbf{D}_t^1) = 1 - p(\mathbf{x}^* \notin \mathbf{D}_t^1) = \frac{t-1}{n}$ given that there are $t-1$ distinct measurements out of n distinct possible actions until \mathbf{x}^* is found.

Lastly, the regret term we computed above is only applicable when $T \leq n$. Since there are n 1-sparse actions available for a vector β with length n , after $T = n$ actions the algorithm has surely found the optimal action. As a result, $\mathcal{R}(\mathbf{x}^*, \beta) - \mathcal{R}(\mathbf{x}_t, \beta) = 1 - 1 = 0$ for $t > n$ which gives us:

$$\mathbb{E}[\text{Reg}(T)] = T_n - \frac{T_n(T_n + 1)}{2n}, \quad T_n = \min\{T, n\}.$$

MULTI-AGENT

Now, we will compute the expected regret for TS with asynchronous multi-agent setting. Similar to single agent, the equation for expected regret follows from (30) with i^* as the location of the nonzero element of β . Again, for the first term we have $\mathcal{R}(\mathbf{x}^*, \beta) = 1$ where the optimal action \mathbf{x}^* is the 1-sparse vector with i^* as the nonzero element. For the second term $\mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta)]$, we know that if $\mathbf{x}^* \in \mathbf{D}_t^j$, then TS will be exploiting the same optimal action at time t and set $\mathbf{x}_t = \mathbf{x}^*$ for which $\mathbb{E}_{\mathbf{x}_t} [\mathcal{R}(\mathbf{x}_t, \beta) | \mathbf{x}^* \in \mathbf{D}_t^j] = 1$. On the other hand, if $\mathbf{x}^* \notin \mathbf{D}_t^j$, then TS will select action \mathbf{x}_t out of actions that are not

observed as part of \mathbf{D}_t^j . Precisely:

$$\begin{aligned} \mathbb{E}[\text{Reg}(\mathbf{T})] &= T - \sum_{t=1}^T \left[p(\mathbf{x}^* \in \mathbf{D}_t^j) \times \mathbb{E}_{\mathbf{x}_t} \left[\mathcal{R}(\mathbf{x}_t, \boldsymbol{\beta}) | \mathbf{x}^* \in \mathbf{D}_t^j \right] \right. \\ &\quad \left. + p(\mathbf{x}^* \notin \mathbf{D}_t^j) \times \mathbb{E}_{\mathbf{x}_t} \left[\mathcal{R}(\mathbf{x}_t, \boldsymbol{\beta}) | \mathbf{x}^* \notin \mathbf{D}_t^j \right] \right]. \end{aligned} \quad (31)$$

However, computing $p(\mathbf{x}^* \in \mathbf{D}_t^j) = 1 - p(\mathbf{x}^* \notin \mathbf{D}_t^j)$ and $\mathbb{E}_{\mathbf{x}_t} \left[\mathcal{R}(\mathbf{x}_t, \boldsymbol{\beta}) | \mathbf{x}^* \notin \mathbf{D}_t^j \right]$ is not as straight forward as the single agent setting. In a single agent setting, $|\mathbf{D}_t^j| = t - 1$ and all $t - 1$ actions are distinct until \mathbf{x}^* is found. In an asynchronous multi-agent setting, $|\mathbf{D}_t^j| < t - 1$ and some of the actions in \mathbf{D}_t^j might be the same. Specifically, as illustrated in Figure 2b an action t starts before all previous $t - 1$ actions are completed and as a result there is a chance that action t will equal the previous actions that were not included in \mathbf{D}_t^j . While we cannot precisely compute $p(\mathbf{x}^* \in \mathbf{D}_t^j)$ in this asynchronous multi-agent setting, in the following lemma we compute a lower bound for it with the proof provided in Appendix 5.1.

Lemma 1. *Consider the active search problem in Theorem 2. For an asynchronous multi-agent TS algorithm, $p(\mathbf{x}^* \in \mathbf{D}_t^j)$ is bounded by:*

$$\begin{aligned} p(\mathbf{x}^* \in \mathbf{D}_t^j) &\geq 1 - \frac{(n-1)^t}{n^t}, \quad \text{if } t < J \\ p(\mathbf{x}^* \in \mathbf{D}_t^j) &\geq 1 - \frac{(n-1)^{J-1}(n-t+2J-1)}{n^J}, \quad \text{if } t \geq J \end{aligned}$$

Using the bound in Lemma 1 along with a naïve bound $\mathbb{E}_{\mathbf{x}_t} \left[\mathcal{R}(\mathbf{x}_t, \boldsymbol{\beta}) | \mathbf{x}^* \notin \mathbf{D}_t^j \right] \geq 0$, we can upper bound the expected regret in (31) as follows:

$$\begin{aligned} \mathbb{E}[\text{Reg}(\mathbf{T})] &\leq T - \sum_{t=1}^T \left[p(\mathbf{x}^* \in \mathbf{D}_t^j) \times \mathbb{E}_{\mathbf{x}_t} \left[\mathcal{R}(\mathbf{x}_t, \boldsymbol{\beta}) | \mathbf{x}^* \in \mathbf{D}_t^j \right] \right] \\ &= T - \sum_{t=1}^J \left[1 - \frac{(n-1)^t}{n^t} \times 1 \right] - \sum_{t=J+1}^T \left[1 - \frac{(n-1)^{J-1}(n-t+2J-1)}{n^J} \times 1 \right] \\ &\leq T - J + \sum_{t=1}^J \left[\frac{(n)^t}{n^t} \right] - (T - J) + \sum_{t=J+1}^T \left[\frac{(n)^{J-1}(n-t+2J-1)}{n^J} \right] \\ &= T - \frac{T(T+1)}{2n} + \frac{T(2J-1)}{n} - \frac{3J^2 - 3J}{2n} \\ &\leq T - \frac{T(T+1)}{2n} + \frac{T(2J-1)}{n}. \end{aligned}$$

Lastly, the regret term we computed above is only applicable when $T \leq n + J$. If $T > n + J$, then $|\mathbf{D}_t^j| = n$ which includes all n 1-sparse actions available for a vector $\boldsymbol{\beta}$ with length n . Hence, after $T = n + J$ actions, the algorithm has surely found the optimal action which results in $\mathcal{R}(\mathbf{x}^*, \boldsymbol{\beta}) - \mathcal{R}(\mathbf{x}_t, \boldsymbol{\beta}) = 1 - 1 = 0$, and consequently:

$$\mathbb{E}[\text{Reg}(T)] \leq T_n - \frac{T_n(T_n+1)}{2n} + \frac{T_n(2J-1)}{n}, \quad T_n = \min\{T, n+J\}$$

5.1 PROOF FOR LEMMA 1

Proof. Since we need to upper bound the expected regret in (30), we will lower bound $p(\mathbf{x}^* \in \mathbf{D}_t^j)$. To compute this bound, we need to better understand \mathbf{D}_t^j . For that, we will look at the example illustrated in Figure 2b. For this example, \mathbf{D}_1^1 through \mathbf{D}_{12}^1 have the following measurement sets:

$$\begin{aligned} \mathbf{D}_1^1 &= \emptyset & \mathbf{D}_7^3 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3, 2, 5\}\} \\ \mathbf{D}_2^2 &= \emptyset & \mathbf{D}_8^2 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3, 2, 5, 6\}\} \\ \mathbf{D}_3^3 &= \emptyset & \mathbf{D}_9^1 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3, 2, 5, 6, 4\}\} \end{aligned}$$

$$\begin{aligned}
\mathbf{D}_4^1 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1\}\} & \mathbf{D}_{10}^2 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3, 2, 5, 6, 4, 8\}\} \\
\mathbf{D}_5^3 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3\}\} & \mathbf{D}_{11}^3 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3, 2, 5, 6, 4, 8, 7\}\} \\
\mathbf{D}_6^2 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3, 2\}\} & \mathbf{D}_{12}^1 &= \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3, 2, 5, 6, 4, 8, 7, 9\}\}
\end{aligned}$$

Our first observation is that $|\mathbf{D}_t^j| = t - J$ which applies to any asynchronous example [Kandasamy et al., 2018]. Next, let us compute $p(\mathbf{x}^* \in \mathbf{D}_t^j)$ as an example for \mathbf{D}_7^3 :

$$\begin{aligned}
p(\mathbf{x}^* \in \mathbf{D}_7^3) &= 1 - p(\mathbf{x}^* \notin \mathbf{D}_7^3) = 1 - p\left((\mathbf{x}^* \neq \mathbf{x}_1) \wedge (\mathbf{x}^* \neq \mathbf{x}_3) \wedge (\mathbf{x}^* \neq \mathbf{x}_2) \wedge (\mathbf{x}^* \neq \mathbf{x}_5)\right) \\
&= 1 - \left(p(\mathbf{x}^* \neq \mathbf{x}_1) \times p(\mathbf{x}^* \neq \mathbf{x}_3 | \mathbf{x}^* \neq \mathbf{x}_1) \times p(\mathbf{x}^* \neq \mathbf{x}_2 | \mathbf{x}^* \neq \mathbf{x}_1, \mathbf{x}_3) \right. \\
&\quad \left. \times p(\mathbf{x}^* \neq \mathbf{x}_5 | \mathbf{x}^* \neq \mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_2)\right) \\
&\stackrel{(c)}{=} 1 - \left(p(\mathbf{x}^* \neq \mathbf{x}_1) \times p(\mathbf{x}^* \neq \mathbf{x}_3 | \emptyset) \times p(\mathbf{x}^* \neq \mathbf{x}_2 | \emptyset) \times p(\mathbf{x}^* \neq \mathbf{x}_5 | \mathbf{x}^* \neq \mathbf{x}_1, \mathbf{x}_3)\right) \\
&\stackrel{(d)}{=} 1 - \frac{n-1}{n} \times \frac{n-1}{n} \times \frac{n-1}{n} \times \frac{n-3}{n-2}. \tag{32}
\end{aligned}$$

Here, (c) follows from the fact that on picking \mathbf{x}_3 and \mathbf{x}_2 , they depend on empty sets of measurements \mathbf{D}_3^3 and \mathbf{D}_2^2 . And, on choosing sensing action \mathbf{x}_5 , the algorithm has access to measurements $\mathbf{D}_5^3 = \{(\mathbf{x}_{t'}, y_{t'}) | t' = \{1, 3\}\}$. Consequently, in (d) we see that $\mathbf{x}_1, \mathbf{x}_3$ and \mathbf{x}_2 have n possible actions to pick, while \mathbf{x}_5 will pick out of $n - 2$ actions that had excluded the ones in \mathbf{D}_5^3 .

As evident in (32), computing $p(\mathbf{x}^* \in \mathbf{D}_t^j)$ depends on the specifics of every example and is different for all \mathbf{D}_t^j . However, this probability can be lower-bounded for any asynchronous setting as follows. We know that \mathbf{D}_t^j has a cardinality of $t - J$. As a result, each asynchronous example translates to a set \mathbf{D}_t^j with permutation of $t - J$ measurements picked out of all $t - 1$ possible measurements $\{(\mathbf{x}_{t'}, y_{t'}) | \{t' = 1, \dots, t - 1\}\}$. Out of all possible asynchronous examples, $\mathbf{D}_t^W = \{(\mathbf{x}_{t'}, y_{t'}) | \{t' = 1, \dots, t - J\}\}$ is the worst case scenario in terms of finding the optimal action. In other words, for this measurement set, the algorithm has used the least amount of information to decide on each of the actions \mathbf{x}_1 through \mathbf{x}_t . Hence, the asynchronous example with \mathbf{D}_t^W has the lowest probability to pick the optimal sensing action \mathbf{x}^* and can be used to lower-bound $p(\mathbf{x}^* \in \mathbf{D}_t^j)$:

$$p(\mathbf{x}^* \in \mathbf{D}_t^j) \geq p(\mathbf{x}^* \in \mathbf{D}_t^W)$$

Now, all we have to do is compute $p(\mathbf{x}^* \in \mathbf{D}_t^W)$. For $t < J$, we have:

$$\begin{aligned}
p(\mathbf{x}^* \in \mathbf{D}_t^W) &= 1 - p(\mathbf{x}^* \notin \mathbf{D}_t^W) = 1 - \left(p(\mathbf{x}^* \neq \mathbf{x}_1) \times p(\mathbf{x}^* \neq \mathbf{x}_2 | \emptyset) \times \dots \times p(\mathbf{x}^* \neq \mathbf{x}_t | \emptyset)\right) \\
&= 1 - \underbrace{\frac{n-1}{n} \times \frac{n-1}{n} \times \dots \times \frac{n-1}{n}}_{t \text{ terms}} = 1 - \frac{(n-1)^t}{n^t}
\end{aligned}$$

For $t \geq J$, we have:

$$\begin{aligned}
p(\mathbf{x}^* \in \mathbf{D}_t^W) &= 1 - p(\mathbf{x}^* \notin \mathbf{D}_t^W) \\
&= 1 - \left(p(\mathbf{x}^* \neq \mathbf{x}_1) \times p(\mathbf{x}^* \neq \mathbf{x}_2 | \emptyset) \times \dots \times p(\mathbf{x}^* \neq \mathbf{x}_J | \emptyset) \right. \\
&\quad \times p(\mathbf{x}^* \neq \mathbf{x}_{J+1} | \mathbf{x}^* \neq \mathbf{x}_1) p(\mathbf{x}^* \neq \mathbf{x}_{J+2} | \mathbf{x}^* \neq \mathbf{x}_1, \mathbf{x}_2) \\
&\quad \left. \times \dots \times p(\mathbf{x}^* \neq \mathbf{x}_{t-J} | \mathbf{x}^* \neq \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t-2J})\right) \\
&= 1 - \underbrace{\frac{n-1}{n} \times \frac{n-1}{n} \times \dots \times \frac{n-1}{n}}_{J \text{ terms}} \times \frac{n-2}{n-1} \times \frac{n-3}{n-2} \times \dots \times \frac{n-(t-2J)-1}{n-(t-2J)} \\
&= 1 - \frac{(n-1)^{J-1} (n-t+2J-1)}{n^J}
\end{aligned}$$

■

6 ADDITIONAL DETAILS ON LATSI

In this section, we provide additional details on deriving LATSI algorithm we proposed in Section 4.2. First, let us do a short review on IDS and RSI algorithms:

REVIEW OF IDS (INFORMATION DIRECTED SAMPLING)

There are certain online optimization problems such as our active search problem where traditional multi-armed bandit algorithms such as TS and Upper Confidence Bound (UCB) fail due to a careless assessment of the information gain (see Appendix 2). IDS is an online optimization algorithm proposed by Russo and Van Roy [2017] that addresses this problem by introducing a novel reward function that balances between expected single-period regret and a measure of information gain. Specifically, defining π_t^{IDS} as the action sampling distribution of IDS at time t , at each time step t , IDS is computed by:

$$\pi_t^{\text{IDS}} = \arg \min_{\pi \in \mathcal{D}(\mathcal{X})} \left\{ \Psi_t(\pi) \triangleq \frac{\Delta_t(\pi)}{\mathcal{I}_t(\pi)} \right\},$$

where they call $\Psi_t(\pi)$ the information ratio of the action sampling distribution π , and $\Delta_t(\pi)$ and $\mathcal{I}_t(\pi)$ are the single-period expected regret and information gain at time t defined below. For the single-period expected regret, we have:

$$\Delta_t(\pi) = \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \mathbb{E}[\mathcal{R}_t(\mathbf{x}^*) - \mathcal{R}_t(\mathbf{x})],$$

where, $\mathcal{R}_t(\mathbf{x})$ is the reward function for action \mathbf{x} at time t . And, for information gain, we have:

$$\mathcal{I}_t(\pi) = \sum_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) \mathbf{I}_t(\mathbf{x}^*; y|\mathbf{x}),$$

where, $\mathcal{I}_t(\mathbf{x}^*; y|\mathbf{x})$ is the mutual information between optimal action \mathbf{x}^* and observations y at time t if sensing action \mathbf{x} is chosen.

Once IDS computes the action sampling distribution at time t , it will choose an action by randomly sampling this distribution.

REVIEW OF RSI (REGION SENSING INDEX)

RSI is a single agent active search algorithm designed to locate sparse signals by actively making data-collection decisions. Similar to our problem formulation in Section 2, RSI makes a practical assumption that at each time step the agent senses a contiguous region of the space. To decide on their next action, RSI at each time step chooses the sensing action \mathbf{x}_t that maximizes the mutual information between the next observation y_t and the signal of interest β , i.e.

$$\mathbf{x}_t = \arg \max_{\mathbf{x}} I(\beta; y|\mathbf{x}, \mathbf{D}_t^1),$$

where, the mutual information is computed using posterior distribution $p(\beta|\mathbf{D}_t^1) = p_0(\beta) \prod_{\tau=1}^{t-1} p(y_\tau|\mathbf{x}, \beta)$ with a k -sparse uniform prior $p_0(\beta)$ and same likelihood distribution as in (1).

Unfortunately, computing the mutual information $I(\beta; y|\mathbf{x}, \mathbf{D}_t^1)$ has high complexity for sparsity rates of $k > 1$. In order to reduce this complexity for $k > 1$ RSI recovers the support of β by repeatedly applying RSI assuming $k = 1$.

DERIVING LATSI

Recall how IDS algorithm solves for the failure mode of TS by using a novel reward function that is the ratio of the expected single-period regret and information gain. Inspired by this algorithm, we propose revising the expected reward $\lambda^+(\beta^*, \mathbf{D}_t^j, \mathbf{x})$ for agent j in (12) for Laplace-TS by adding a measure of information gain to it. For the information gain, we use the mutual information $I(\beta; y|\mathbf{x}, \mathbf{D}_t^j)$ computed in RSI algorithm. Thus, we can say that LATSI is the Laplace-TS algorithm where the design stage has been replaced by:

$$\mathbf{x}_t = \arg \max_{\mathbf{x}} \mathbf{R}^+(\beta^*, \mathbf{D}_t^j, \{\mathbf{x}, y\}),$$

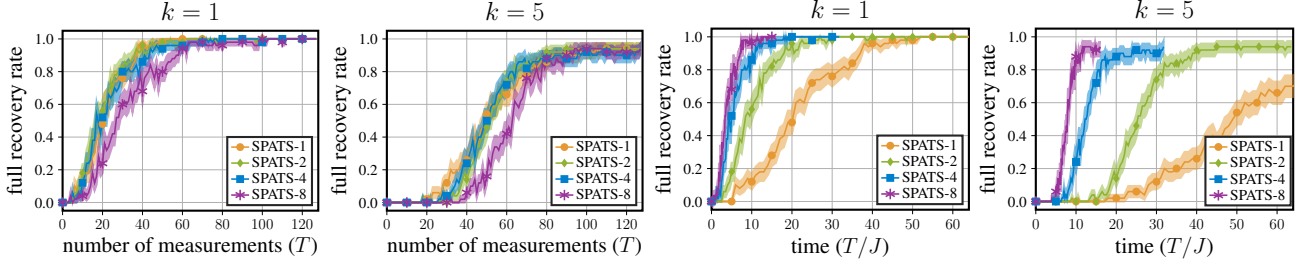


Figure 8: Full recovery rate of SPATS with 1, 2, 4 and 8 agents for $n = 128$, $k = 1, 5$

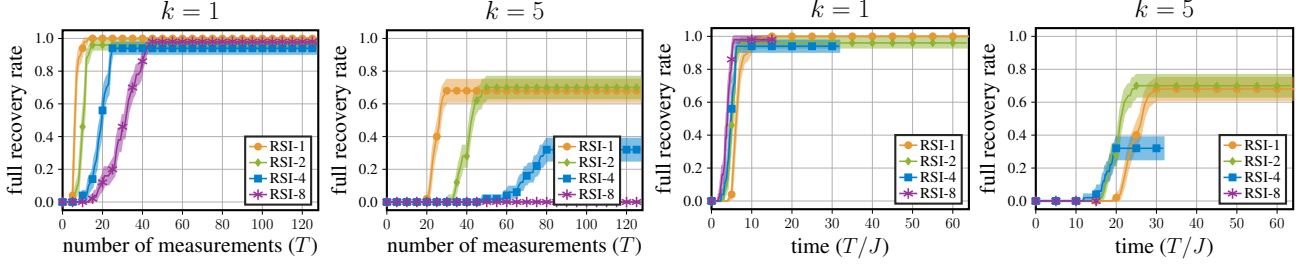


Figure 9: Full recovery rate of RSI with 1, 2, 4 and 8 agents for $n = 128$, $k = 1, 5$

where,

$$\mathbf{R}^+ = \frac{I(\beta^*; y|\mathbf{x}, \mathbf{D}_t^j)}{\text{average } I \text{ over all } \mathbf{x}} + \alpha \times \frac{\lambda^+(\beta^*, \mathbf{D}_t^j, \mathbf{x})}{\text{average } \lambda^+ \text{ over all } \mathbf{x}}, \quad (33)$$

using I from Ma et al. [2017] and λ^+ in (12). Here, we normalize and add a tuning parameter α to best control the importance of each term on the overall reward \mathbf{R}^+ .

7 ADDITIONAL NUMERICAL RESULTS

We now provide additional simulation results to support our analysis in Section 5.

1-DIMENSIONAL SEARCH SPACE

In this section we provide additional results for a 1-dimensional search space ($d = 1$), where we estimate a k -sparse signal β of length $n = 128$ with two sparsity rates of $k = 1, 5$. Here, we use the same set of parameters as those outlined in Section 5. Figure 8, 9 and 10 show the performance of SPATS, RSI and LATSI respectively. Our observations overall are similar to those in Section 5. Concretely, we observe that both SPATS and LATSI show an improvement in performance with time as more agents become available. The efficiency of this asynchronous search by SPATS and LATSI becomes more pronounced with higher sparsity rate k where there are more targets in the search space. On the other hand, multi-agent RSI's performance worsens as we increase the number of agents (especially with larger sparsity rate k). This observation is similar to that in Section 5 and is explained by the lack of randomness in the information theoretic reward function of RSI.

2-DIMENSIONAL SEARCH SPACE

In this section we provide additional results for a 2-dimensional search space ($d = 2$), where we estimate a k -sparse signal β of length $n = 16 \times 16$ and two sparsity rates of $k = 1, 5$. The same set of parameters outlined in Section 5 are followed. Figure 11, 12 and 13 show the performance of SPATS, RSI and LATSI respectively, in the multi-agent setting. Our observations overall are similar to those in Section 5. The probabilistic decision-making nature of SPATS ensures that its performance multiplies by its number of agents. RSI relies on information theoretic decision making (no randomness).

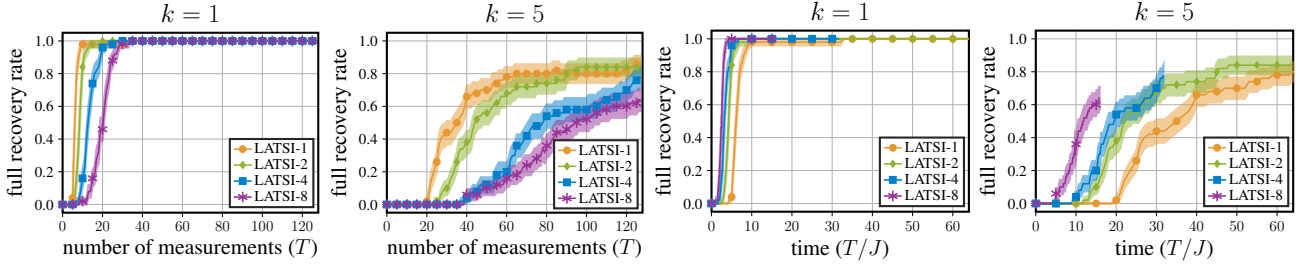


Figure 10: Full recovery rate of LATSI with 1, 2, 4 and 8 agents for $n = 128$, $k = 1, 5$

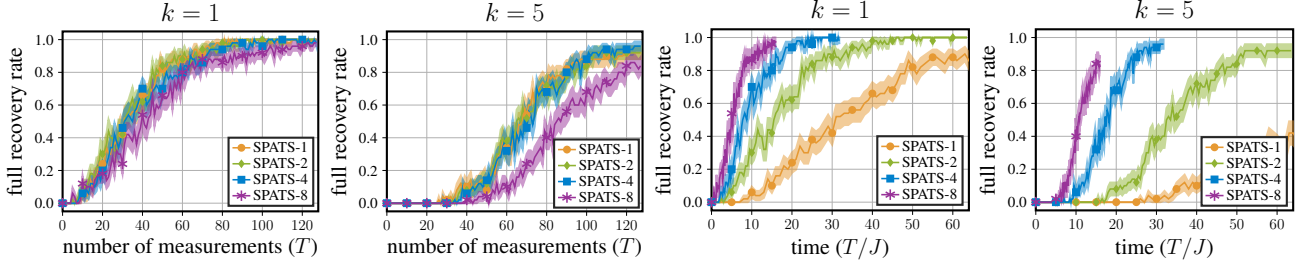


Figure 11: Full recovery rate of SPATS with 1, 2, 4 and 8 agents for $n = 16 \times 16$, $k = 1, 5$

Hence, this lack of randomness in its reward function together with the poor approximation of mutual information for $k > 1$, results in a poor performance generally worsening with more agents. The multi-agent performance of LATSI shows similar trends as SPATS given its partly probabilistic reward function. However, LATSI is not as efficient as SPATS given its poor approximation of mutual information due to the larger size of the environment. These results further support the analysis in Section 5.

1-DIMENSIONAL VS 2-DIMENSIONAL SEARCH SPACE

In Figure 14, we compare the performance of SPATS and LATSI for 1-dimensional (1d) and 2-dimensional (2d) search spaces with the same length $n = 128$. We compare a 1d search space of length $n = 128$ with a 2d search space of length $n = 8 \times 16$ (i.e. flattened length $n = 128$) for two sparsity rates $k = 1, 5$. The same set of parameters outlined in Section 5 are followed. As evident in this figure, while 1-dimensional SPATS and LATSI have a similar performance for $k = 1$, with larger k they perform better in a 2d space rather than in 1d. The reason for this behavior is our region sensing assumption. Specifically, region sensing in the 2-dimensional grid expands the available action space and allows for a larger feasible action set compared to the 1-dimensional grid. Mainly, the larger action set helps our algorithms to choose better sensing actions in the 2-dimensional space, particularly when multiple targets are present.

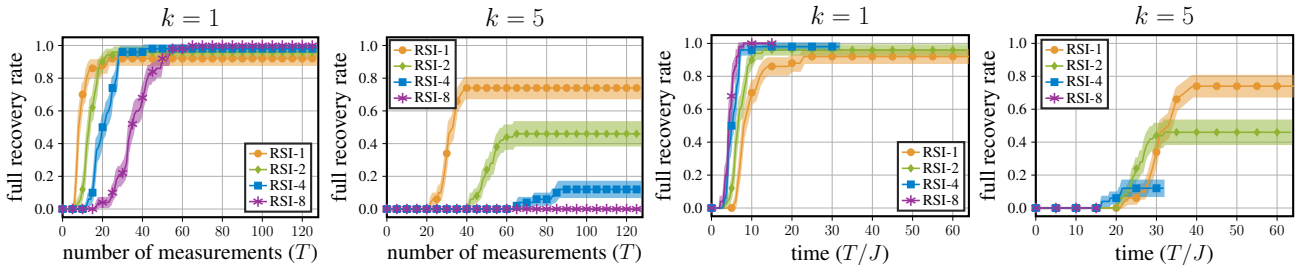


Figure 12: Full recovery rate of RSI with 1, 2, 4 and 8 agents for $n = 16 \times 16$, $k = 1, 5$

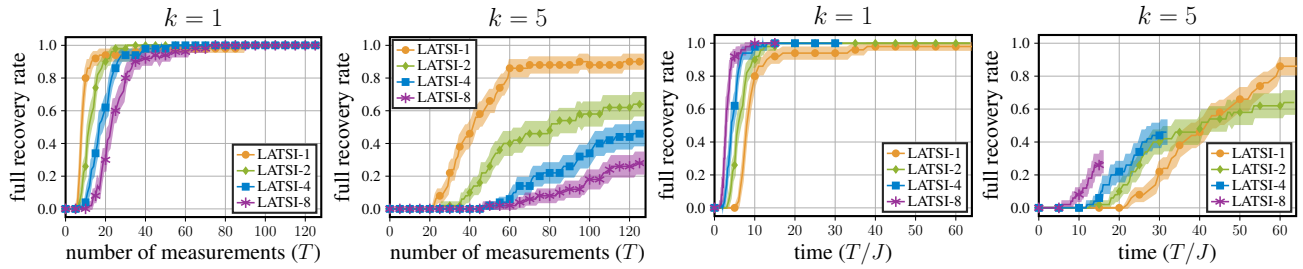


Figure 13: Full recovery rate of LATSI with 1, 2, 4 and 8 agents for $n = 16 \times 16$, $k = 1, 5$

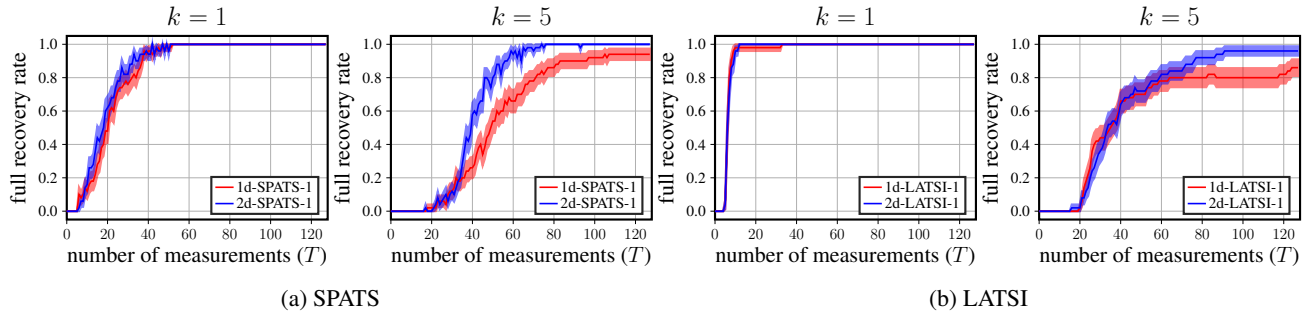


Figure 14: Full recovery rate in 1-dimensional ($n = 128$) vs 2-dimensional ($n = 8 \times 16$) search space with 1 agent for sparsity rates $k = 1, 5$

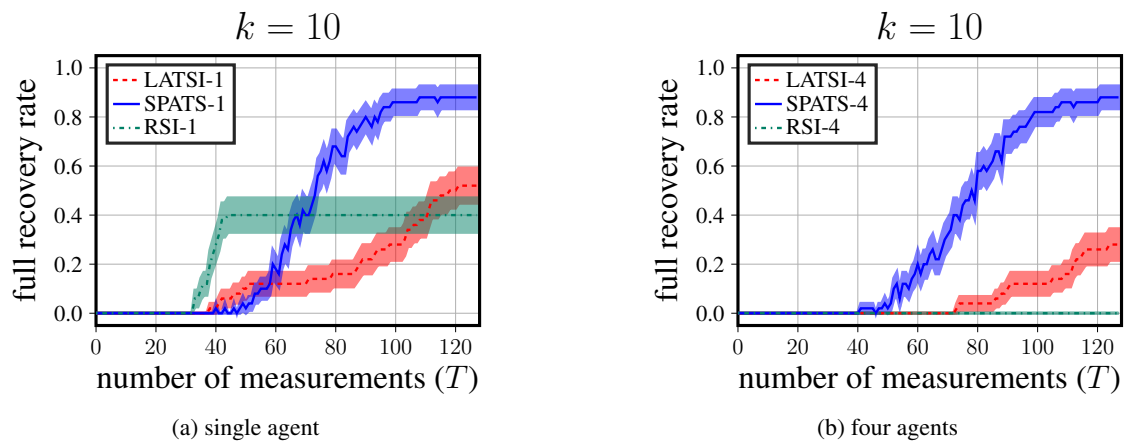


Figure 15: Full recovery rate of SPATS, LATSI and RSI for 1 and 4 agents for sparsity rate $k = 10$

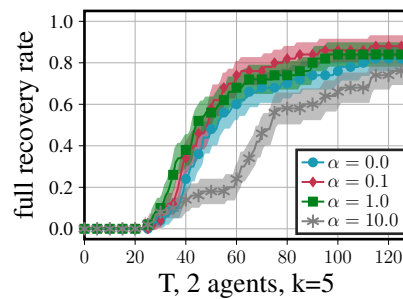


Figure 16: Effect of α in full recovery performance of LATSI with multiple agents and multiple targets

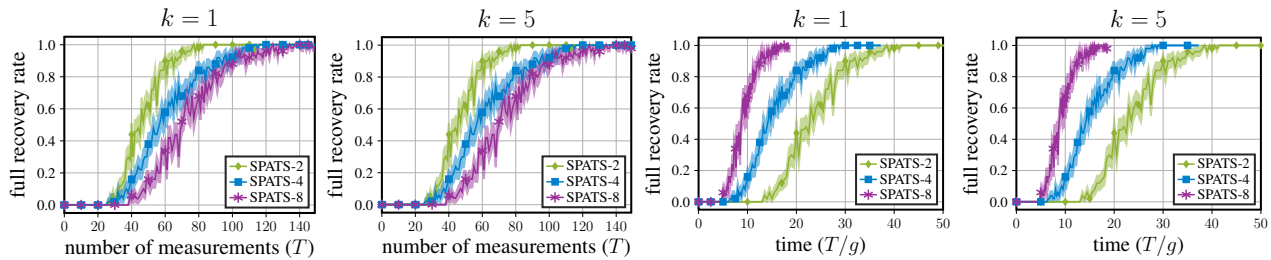


Figure 17: Full recovery rate of SPATS in the presence of communication failure with 2, 4 and 8 agents for sparsity rates $k = 1, 5$

SPARSITY RATE

In this section, we provide an additional set of results for sparsity rate of $k = 10$ in a 1-dimensional search space of length $n = 128$. Figure 15 illustrates the full recovery rate of SPATS, LATSI and RSI for one and four agents, respectively. Overall, we see that increasing the number of targets from $k = 5$ to $k = 10$ reduces the performance of all algorithms. Both figures confirm our previous observations that for sparsity rates of $k > 1$, SPATS algorithm outperforms LATSI and RSI. In the case of single agent, the poor performance of LATSI and RSI can be traced back to poor approximation of mutual information for $k > 1$ (see Section 5). In the four-agent scenario, RSI and LATSI are significantly worse than SPATS due to lack of randomness in information-greedy approaches (see Section 5 for details).

SENSITIVITY ANALYSIS FOR LATSI

Recall tuning parameter α as the scaling factor in (33) which is used to combine Laplace-TS and RSI's reward functions in LATSI's reward function \mathbf{R}^+ . As evident by (33), increasing the value of α in \mathbf{R}^+ increases the weight of the reward computed by Laplace-TS. We experiment with different values of $\alpha \in \{0, 0.1, 1, 10\}$ to determine its effect on the full recovery rate performance. Since it is our aim to exploit parallelization, we perform the simulations with 2 agents and $k = 5$ in a 1-dimensional search space of length $n = 128$. Plotting the results in Figure 16, we observe that the algorithm's performance is robust for a wide range of α . Based on this observation, we chose $\alpha = 1$ for all experiments with LATSI.

ROBUSTNESS TO UNRELIABLE INTER-AGENT COMMUNICATION

In this section, we provide experimental results showing the robustness of SPATS to unreliable communication among agents in asynchronous multi-agent active search. Using the same parameters and setting described in Section 5, we focus on a 2-dimensional search space where we estimate a k -sparse signal β with length $n = 8 \times 16$ and two sparsity rates of $k = 1, 5$. Assuming J number of agents are actively searching for the targets and asynchronously communicating their measurements to each other, we introduce unreliability in inter-agent communication by randomly dropping, for each agent, up to three of its last received new measurement messages from other agents. We measure the mean and standard error of the full recovery rate as a function of the number of measurements T as well as time over 50 random trials. We vary the number of agents J between 2, 4 and 8. Figure 17 shows that the agents are able to fully recover all targets even in this additionally introduced unreliable communication setup. The probabilistic nature of SPATS allows the agents to maintain their performance even when they cannot access all measurements from other agents.

References

- Kyoungwa Bae and Bani K Mallick. Gene selection using a two-level hierarchical Bayesian model. *Bioinformatics*, 20(18):3423–3430, 2004.
- Thomas Blumensath and Mike E Davies. Iterative hard thresholding for compressed sensing. *Applied and computational harmonic analysis*, 27(3):265–274, 2009.
- Ingrid Daubechies, Michel Defrise, and Christine De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11), 2004.

- Mário AT Figueiredo. Adaptive sparseness for supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 25(9):1150–1159, 2003.
- Alan E Gelfand and Adrian FM Smith. Sampling-based approaches to calculating marginal densities. *Journal of the American statistical association*, 85(410), 1990.
- Kirthevasan Kandasamy, Akshay Krishnamurthy, Jeff Schneider, and Barnabás Póczos. Parallelised bayesian optimisation via thompson sampling. In *International Conference on Artificial Intelligence and Statistics*, 2018.
- Yifei Ma, Roman Garnett, and Jeff Schneider. Active search for sparse signals with region sensing. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Arian Maleki. Approximate message passing algorithms for compressed sensing. *a degree of Doctor of Philosophy, Stanford University*, 2011.
- Elaine Crespo Marques, Nilson Maciel, Lirida Naviner, Hao Cai, and Jun Yang. A review of sparse recovery algorithms. *IEEE Access*, 7:1300–1322, 2018.
- Deanna Needell and Joel A Tropp. CoSaMP: iterative signal recovery from incomplete and inaccurate samples. *Communications of the ACM*, 53(12):93–100, 2010.
- Yagyensh Chandra Pati, Ramin Rezaifar, and Perinkulam Sambamurthy Krishnaprasad. Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Asilomar*. IEEE, 1993.
- Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling. *Journal of Machine Learning Research (JMLR)*, 17(1), 2016.
- Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling. *Operations Research*, 66(1): 230–252, 2017.
- Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial on Thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- Zhilin Zhang and Bhaskar D Rao. Sparse signal recovery with temporally correlated source vectors using sparse Bayesian learning. *IEEE Journal of Selected Topics in Signal Processing*, 5(5):912–926, 2011.