

# Anomaly Detection via Controlled Sensing and Deep Active Inference

Geethu Joseph, Chen Zhong, M. Cenk Gursoy, Senem Velipasalar, and Pramod K. Varshney

*Department of Electrical Engineering and Computer Science*

*Syracuse University*

New York 13244, USA

Emails: {gjoseph, czhong03, mcgursoy, svelipas, varshney}@syr.edu.

**Abstract**—In this paper, we address the anomaly detection problem where the objective is to find the anomalous processes among a given set of processes. To this end, the decision-making agent probes a subset of processes at every time instant and obtains a potentially erroneous estimate of the binary variable which indicates whether or not the corresponding process is anomalous. The agent continues to probe the processes until it obtains a sufficient number of measurements to reliably identify the anomalous processes. In this context, we develop a sequential selection algorithm that decides which processes to be probed at every instant to detect the anomalies with an accuracy exceeding a desired value while minimizing the delay in making the decision and the total number of measurements taken. Our algorithm is based on *active inference* which is a general framework to make sequential decisions in order to maximize the notion of *free energy*. We define the free energy using the objectives of the selection policy and implement the active inference framework using a deep neural network approximation. Using numerical experiments, we compare our algorithm with the state-of-the-art method based on deep actor-critic reinforcement learning and demonstrate the superior performance of our algorithm.

**Index Terms**—Active hypothesis testing, anomaly detection, active inference, quickest state estimation, sequential decision-making, sequential sensing.

## I. INTRODUCTION

In many practical applications such as remote health monitoring using sensors, the goal is to identify the anomalies among a given set of functionalities of a system [1], [2]. Here, the system is equipped with multiple sensors and each sensor monitors a different, but not necessarily independent functionality (which we henceforth refer to as a process) of the system. The sensor sends its observations to the decision-making agent over a communication link, and the received observation may be distorted due to the unreliability in the sensor hardware and/or the noisy link (e.g., a wireless channel) between the sensor and the agent. Hence, the decision agent needs to probe each process multiple times before it declares one or more of the processes to be anomalous with the desired confidence. Repeatedly probing all the processes allows the agent to quickly find any potential system malfunction, but

this incurs a large cost (e.g., higher energy consumption that reduces the life span of the sensor network). Therefore, the agent uses the *controlled sensing* technique with which it probes a small subset of processes at every time instant. In this context, we address the question of how the agent sequentially chooses a subset of processes so that it accurately detects the anomalies with a minimum delay and a minimum number of sensor measurements.

A classical approach to solve the sequential sensor selection problem is based on the active hypothesis testing framework [3], [4] where the decision-making agent constructs a hypothesis corresponding to each of the possible states of the processes and determine which one of these hypotheses is true. Active hypothesis testing is a well-studied problem and several solution strategies have been proposed in the literature [5]–[9]. However, these approaches provide model-based algorithms which are designed under simplified modeling assumptions. This has motivated the researchers to design data-driven deep learning algorithms [3], [4], [10]. These algorithms are not only more flexible than traditional algorithms, but they also possess reduced computational complexity. The existing literature along these lines relies on the most fundamental reinforcement learning (RL) algorithms such as Q-learning [10] and actor-critic [3], [4]. However, recently a new framework called *active inference* has been shown to be a promising complement to the traditional RL approaches for several sequential decision-making problems [11]–[13]. Therefore, in this paper, we develop and implement a novel policy to select processes to obtain measurements at each step, inspired by the active inference approach.

The contributions of the paper are as follows: we first define the notion of *free-energy* based on the entropy associated with the estimate of the states of the processes and the cost of sensing. This allows us to reformulate the anomaly detection problem as an active inference problem in which the goal is to minimize the free energy. We then implement our algorithm using deep neural networks which are relatively less explored in the context of active inference. Our algorithm balances the model-based and the data-driven approaches of active inference. Specifically, we use the model-based posterior updates to tackle the uncertainties in the observations, and the data-driven neural network to handle the underlying statistical dependence between the processes.

The information, data, or work presented herein was funded in part by National Science Foundation (NSF) under Grant 1618615, Grant 1739748, Grant 1816732 and by the Advanced Research Projects Agency-Energy (ARPA-E), U.S. Department of Energy, under Award Number DE-AR0000940. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

The active inference approach has many similarities to the reinforcement-based algorithms, such as learning probabilistic models, exploration and exploitation of various actions, and efficient planning. So we compare our algorithm with the existing RL-based approach presented in [4] using numerical simulations. We observe that the delay in estimation is smaller for our method while the corresponding accuracy and cost of sensing are competitive to the performance of the RL-based method given in [4]. This advantage makes our active inference-based approach a better alternative to the existing RL-based method.

## II. ANOMALY DETECTION PROBLEM

We consider  $N$  random processes that are potentially statistically dependent. Each process is in one of the two states: normal (denoted by 0) or anomalous (denoted by 1). The states of these processes are denoted by a random vector  $\mathbf{s} \in \{0, 1\}^N$ . The goal of the work is to detect the anomalous processes out of the  $N$  processes, which is equivalent to estimating the random vector  $\mathbf{s}$ . The dependence pattern and the number of anomalous processes are unknown to the decision-making agent.

To estimate  $\mathbf{s}$ , the decision-making agent probes one or more processes at every time instant and obtains potentially erroneous observations of the corresponding entries of  $\mathbf{s}$ . Let the set of processes probed at time  $k$  be  $\mathcal{A}_k \in \mathcal{P}$  and the corresponding observation vector be  $\mathbf{y}_{\mathcal{A}_k}(k) \in \{0, 1\}^{|\mathcal{A}_k|}$ . Here,  $\mathcal{P}$  denotes the power set of  $\{1, 2, \dots, N\}$  without the null set ( $|\mathcal{P}| = 2^N - 1$ ). The observation corresponding to the  $i^{\text{th}}$  process at time  $k$ , denoted by  $\mathbf{y}_i(k) \in \{0, 1\}$ , obeys the following probabilistic model:

$$\mathbf{y}_i(k) = \begin{cases} \mathbf{s}_i & \text{with probability } 1 - p \\ 1 - \mathbf{s}_i & \text{with probability } p, \end{cases} \quad (1)$$

where  $p \in [0, 1]$  denotes the probability that the observation differs from the actual state of the process. We assume that given  $\mathbf{s}$ , the observations obtained across different time instants are jointly (conditionally) independent. Also, probing each process incurs a cost of sensing of  $\lambda \geq 0$ , i.e., the cost of sensing at time  $k$  is  $|\mathcal{A}_k| \lambda$ .

At each time  $k$ , the agent determines which processes to observe ( $\mathcal{A}_k$ ) until it declares the estimate of  $\mathbf{s}$  with the desired confidence. The selection policy is designed such that the stopping time  $K$  and the total cost of sensing  $\lambda \sum_{k=1}^K |\mathcal{A}_k|$  are minimized.

## III. ANOMALY DETECTION USING DEEP ACTIVE INFERENCE

The active inference framework relies on a normative theory of brain function based on its perception of the environment. At a high level, the active inference agent maintains a generative model that represents its perception. The generative model  $Q$  comprises a joint probability distribution on the state of the environment, the actions, and the corresponding observations. The generative model assigns higher probabilities to the states and actions that are favorable to the agent,

and therefore, it is biased towards the agent's preferences. Given a generative model, the agent inverts the model using the method of approximate Bayesian inference. To this end, it defines a variational distribution  $q$  that the agent controls. The distribution  $q$  is optimized by minimizing the Kullback-Leibler (KL) divergence between the distributions  $q$  and  $Q$ . Therefore, if we choose actions from the distribution  $q$ , they fulfill the agent's preferences. The KL divergence between the variational distribution and the generative model is called the variational free energy. In short, the goal of the active inference agent is to minimize its expected free energy (EFE) into the future up to the stopping time  $K$ . Next, we provide the details of the active inference framework in the context of anomaly detection.

### A. Environment

The environment of the active inference framework refers to the set of states, actions, and observations. In the context of our anomaly detection problem, we define the state of the active inference framework at time  $k$  as the posterior belief  $\pi(k)$  on the random vector  $\mathbf{s} \in \{0, 1\}^N$ . Since there are  $m = 2^N$  possible values for  $\mathbf{s}$ , the posterior belief is an  $m$ -dimensional vector  $\pi \in [0, 1]^m$ . Further, the actions refer to the selection of which processes to observe  $\mathcal{A}_k \in \mathcal{P}$ , and  $\mathbf{y}_{\mathcal{A}_k}$  denotes the observations.

We first note that at time  $k$ , the information available to the agent is the set of processes observed till time  $k$  and the corresponding observation vectors:  $\{\mathcal{A}_j, \mathbf{y}_{\mathcal{A}_j}\}_{j=1}^k$ . Using this information, the posterior belief vector  $\pi(k) \in [0, 1]^m$  can be computed in closed form as follows [4]:

$$\pi_i(k) = \frac{\pi_i(k-1) \prod_{a \in \mathcal{A}_k} [(1-p)\mathbb{1}_{\mathcal{E}_{a,k,i}} + p\mathbb{1}_{\mathcal{E}_{a,k,i}^c}]}{\sum_{i=1}^m \pi_i(k-1) \prod_{a \in \mathcal{A}_k} [(1-p)\mathbb{1}_{\mathcal{E}_{a,k,i}} + p\mathbb{1}_{\mathcal{E}_{a,k,i}^c}]}, \quad (2)$$

where  $\mathbb{1}$  is the indicator function and the event  $\mathcal{E}_{a,k,i} \triangleq \{\mathbf{y}_a(k) = \mathbf{s}_a | \mathcal{H} = i\}$  denotes the event that the observation obtained and the corresponding state are the same, when the index corresponding to the true value of  $\mathbf{s}$  is  $\mathcal{H} = i$ . Also, the event  $\mathcal{E}_{a,k,i}^c \triangleq \{\mathbf{y}_a(k) \neq \mathbf{s}_a | \mathcal{H} = i\}$  denotes the complement of  $\mathcal{E}_{a,k,i}$ . As a result, given the previous state  $\pi(k-1)$ , the action  $\mathcal{A}_k$  and the observation  $\mathbf{y}_{\mathcal{A}_k}$ , we can exactly compute the updated posterior belief  $\pi(k)$  using (2). Therefore, the generative model that learns the environment is a distribution on the actions and the observations:  $Q(\mathcal{A}_k, \mathbf{y}_{\mathcal{A}_k} | \pi(k-1))$ .

### B. Preferences

In this subsection, we consider the preferences of the agent that defines the generative model. Recall that our goal is to estimate the vector  $\mathbf{s}$  with confidence exceeding a specific level while minimizing the stopping time  $K$  and the cost of sensing  $\lambda \sum_{k=1}^K |\mathcal{A}_k|$ . Clearly, the best estimate of  $\mathbf{s}$  based on the posterior belief corresponds to  $i^*(k) \triangleq \arg \max_{i=1,2,\dots,m} \pi_i(k)$ ,

and the confidence associated with the estimation is  $\pi_{i^*(k)}(k)$ . Therefore, the agent terminates the detection algorithm when

$$\arg \max_{i=1,2,\dots,m} \pi_i(k) > \pi_{\text{upper}}, \quad (3)$$

where  $\pi_{\text{upper}}$  is the desired level of confidence. In short, the decision making relies only on the posterior belief  $\pi(k)$ . Also, as  $k$  increases, we get more observations and the posterior belief becomes more accurate. Therefore, the selection policy  $\mu$  is a function of the latest value of the posterior belief:  $\mu(\pi(k-1)) = \mathcal{A}_k$ .

Further exploring the objective of the policy design, we note that minimizing the stopping time is identical to driving the largest entry of  $\pi(k)$  to  $\pi_{\text{upper}}$  as soon as possible. We achieve this by minimizing the entropy  $H(\pi(K))$  of  $\pi(K)$  because the entropy is minimized when the largest entry of  $\pi(K)$  is 1 and the remaining entries are zeros. Here, the entropy is given by

$$H(\pi) = - \sum_{i=1}^m \pi_i \log(\pi_i). \quad (4)$$

We note that this approach is different from the Bayesian log likelihood ratio based-approach in [3], [4], [10]. Therefore, we define the instantaneous objective function that the agent aims to minimize at time  $k$  as follows:

$$r(k) = H(\pi(k)) - H(\pi(k-1)) + \lambda |\mathcal{A}_k|. \quad (5)$$

This definition ensures that the overall objective function is given by

$$\sum_{k=1}^K r(k) = H(\pi(K)) - H(\pi(0)) + \sum_{k=1}^K \lambda |\mathcal{A}_k|, \quad (6)$$

where minimizing  $H(\pi(K)) - H(\pi(0))$  minimizes the entropy in the posterior belief as  $H(\pi(0))$  is a constant, and minimizing  $\sum_{k=1}^K \lambda |\mathcal{A}_k|$  minimizes the total cost of sensing. The instantaneous objective function  $r(k)$  represents the preferences of the agent at time  $k$  and it is encoded into the generative model as the prior probability on the belief vector:

$$\begin{aligned} Q(\mathbf{y}_{\mathcal{A}_k} | \mathcal{A}_k, \pi(k-1)) \\ = \sigma(-H(\pi(k)) + H(\pi(k-1)) - \lambda |\mathcal{A}_k|), \end{aligned} \quad (7)$$

where  $\sigma(\cdot)$  is the softmax function. Also,  $\pi(k)$  is a function of  $\pi(k-1)$ ,  $\mathcal{A}_k$  and  $\mathbf{y}_{\mathcal{A}_k}$  due to (2). We also note that

$$Q(\mathbf{y}_{\mathcal{A}_k}, \mathcal{A}_k | \pi(k-1)) = Q(\mathbf{y}_{\mathcal{A}_k} | \mathcal{A}_k, \pi(k-1)) Q(\mathcal{A}_k | \pi(k-1)). \quad (8)$$

Therefore, the generative model is completely defined if we specify the distribution  $Q(\mathcal{A}_k | \pi(k-1))$ . This distribution is defined based on the EFE of the future as we discuss in the following subsection.

### C. Total expected free energy

The variational free energy  $F$  is the KL divergence between the variational distribution  $q(\mathcal{A} | \pi(k-1))$  and the generative model  $Q(\mathcal{A} | \pi(k-1))$ . Thus,

$$F(k) = \sum_{\mathcal{A} \in \mathcal{P}} q(\mathcal{A} | \pi(k-1)) \log \frac{q(\mathcal{A} | \pi(k-1))}{Q(\mathcal{A} | \pi(k-1))}. \quad (9)$$

The goal of the agent is to minimize the total free-energy of the expected trajectories into the future:

$$G(\mathcal{A}, \pi) = \sum_{j=k}^K \mathbb{E} \{ F(j) | \mathcal{A}_k = \mathcal{A}, \pi(k-1) = \pi \}. \quad (10)$$

In other words, the agent computes the expected free-energy of all paths into the future and probabilistically chooses an action that minimizes the expected free-energy. Therefore, a popular choice for the distribution over the actions assigned by the generative model is a Boltzmann distribution over the expected free energies [11], [14], [15]:

$$Q(\mathcal{A} | \pi(k-1)) = \sigma(-G(\mathcal{A}, \pi(k-1))), \quad (11)$$

where  $\sigma(\cdot)$  is again the softmax function, and  $G$  is given by (10).

So far, we have presented the conceptual aspects of our algorithm. We next discuss how to compute the expressions in (9) and (10).

### D. Deep-learning based implementation

We implement our algorithm using deep neural networks. We start with the computation of the free energy in (9):

$$\begin{aligned} F = -H(q(\mathcal{A} | \pi(k-1))) \\ - \sum_{\mathcal{A} \in \mathcal{P}} q(\mathcal{A} | \pi(k-1)) \log Q(\mathcal{A} | \pi(k-1)), \end{aligned} \quad (12)$$

where the entropy term  $H(q(\mathcal{A} | \pi(k-1)))$  is a function of the variational distribution  $q$  which is controlled by the agent. We implement this distribution using a neural network which we refer to as the *policy network*. The policy neural network takes the posterior belief  $\pi(k-1)$  as the input and outputs stochastic selection policy  $q_\theta \in [0, 1]^{m-1}$  which is a probability distribution on  $\mathcal{P}$  and parameterized by  $\theta$ . Therefore, the entropy term is computed using the entropy of the distribution outputted by the neural network. This neural network also gives the policy implemented by the agent, which is sampled from the distribution  $q$  learned at time  $k$ :

$$\mathcal{A}_k = \mu(\pi(k-1)) \sim q_\theta(\pi(k-1)). \quad (13)$$

Further, the second term in (12) can be determined using (11) and (10). From (10), the EFE for a single time-step can be approximated as follows [15]:

$$\begin{aligned} G(\mathcal{A}_k, \pi(k-1)) \approx -\log Q(\mathbf{y}_{\mathcal{A}_k} | \mathcal{A}_k, \pi(k-1)) \\ + \mathbb{E}_{\mathcal{A} \sim Q(\cdot | \pi(k))} \{ G(\mathcal{A}, \pi(k)) \}. \end{aligned} \quad (14)$$

Here, the first term is determined using (7). However, the second term in (14) involves explicit computation into the future values. Therefore, we learn a bootstrap estimate of this quantity using a neural network which we refer to as the *bootstrapped EFE-network*. Let  $G_\phi(\mathcal{A})$  denote this neural network where  $\phi$  is the parameter of the network. In other

words, the estimate of the neural network is the predicted value of the free-energy of the system. Thus, (14) reduces to

$$G(\mathcal{A}_k, \pi(k-1)) = H(\pi(k)) - H(\pi(k-1)) + \lambda |\mathcal{A}_k| + \mathbb{E}_{\mathcal{A} \sim Q(\cdot|\pi(k))} \{G_\phi(\mathcal{A}_{k+1}, \pi(k))\}. \quad (15)$$

Substituting (15) and (11) into (12) completes the derivation of the algorithm.

To summarize, our solution involves two neural networks  $q_\theta$  and  $G_\phi$  which represent the policy and the expected free-energy, respectively. At every time instant, we sample an action from the output distribution of the policy network  $q_\theta$  and obtain the corresponding observation  $\mathbf{y}_{\mathcal{A}_k}$ . Next, we compute the bootstrapped EFE estimate and the variational free energy using the neural networks and (12) and (15). Finally, the parameter  $\theta$  of the policy network is modified by minimizing the variational free energy  $F(k)$ . Similarly, the parameter  $\phi$  of the bootstrapped EFE-network is optimized by comparing EFE-network output with the value of the expected value  $G(\mathcal{A})$  calculated at time  $k$ . We use the  $\ell_2$ -norm of the difference between the two estimates:

$$L = \|G_\phi(\mathcal{A}) - G(\mathcal{A})\|^2. \quad (16)$$

The pseudo-code of the algorithm is summarized in Algorithm 1 below.

---

**Algorithm 1** Active inference for anomaly detection

---

**Initialization:** • Policy network  $q_\theta(a|\pi)$  with parameters  $\theta$   
• Bootstrapped EFE-network  $G_\phi(\pi; a)$  with parameters  $\phi$

- 1: **repeat**
- 2: Initialize the prior state  $\pi_0 \in [0, 1]^m$  (can be learned from the training data)
- 3: Time index  $k = 0$
- 4: **while**  $k < T$  and  $\max_i \pi_i > \pi_{\text{upper}}$  and  $k < T_{\text{max}}$  **do**
- 5: Choose action  $\mathcal{A}_k \sim q_\theta(\pi(k-1))$
- 6: Generate observations  $\mathbf{y}_{\mathcal{A}_k, k}$
- 7: Compute  $\pi(k+1)$  using (2)
- 8: Compute the bootstrapped EFE estimate  $G$  using (15)
- 9: Compute the variational free energy  $F$  using (11) and (12)
- 10: Update the policy network network by minimizing the variational free energy  $F$  with respect to  $\theta$
- 11: Update the bootstrapped EFE-network by minimizing the bootstrapping loss in (16) with respect to  $\phi$
- 12: Increase time index  $k = k + 1$
- 13: **end while**
- 14: **until**
- 15: Declare the estimate corresponding to  $\arg \max_i \pi_i$

---

#### IV. NUMERICAL RESULTS

In this section, we present numerical results comparing our algorithm with the actor-critic method in [4]. The simulation setup is similar to that in [4]. We choose the number of

processes as  $N = 3$  and thus,  $m = 2^N = 8$ . The probability of a process being normal is taken as  $q = 0.8$ . Here, the first and second processes are assumed to be statistically dependent, and the third process is independent of the other two. The correlation between the dependent processes is captured by the parameter  $\rho \in [0, 1]$ :

$$\mathbb{P}\{\mathbf{s}_1 = 0, \mathbf{s}_2 = 0\} = q^2 + \rho q(1 - q) \quad (17)$$

$$\mathbb{P}\{\mathbf{s}_1 = 0, \mathbf{s}_2 = 1\} = q(1 - q)(1 - \rho) \quad (18)$$

$$\mathbb{P}\{\mathbf{s}_1 = 1, \mathbf{s}_2 = 0\} = q(1 - q)(1 - \rho) \quad (19)$$

$$\mathbb{P}\{\mathbf{s}_1 = 1, \mathbf{s}_2 = 1\} = (1 - q)^2 + \rho q(1 - q). \quad (20)$$

Also, we assume that the crossover probability of the observations is  $p = 0.8$ , and the maximum number of time slots for each episode (trial or run) is  $T_{\text{max}} = 300$ .

For the active inference algorithm, we implement the policy neural network and the bootstrapped EFE-network with three layers and the ReLU activation function between consecutive layers. To update the parameters of the neural networks, we apply the Adam Optimizer, and we set the learning rates of the policy network and the bootstrapped EFE-network as  $10^{-6}$  and  $5 \times 10^{-6}$ , respectively. The implementation of the actor-critic method is the same as that in [4] except that we use the entropy based-reward function as defined in (5). Also, we choose the learning rates of the actor and critic networks as  $5 \times 10^{-4}$  and  $5 \times 10^{-3}$ , respectively.

The simulation results are presented in Figs. 1 to 3. Our observations from the numerical results are as follows:

- *Success rate:* In Fig. 1, we plot the success rates of the two algorithms as a function of the upper bound on the posterior  $\pi_{\text{upper}}$ . The success rate is defined as the ratio between the number of times the algorithm correctly identifies all the anomalous processes to the total number of trials. We observe that the success rates achieved by both algorithms are comparable in all the settings. Also, the success rate depends primarily on  $\pi_{\text{upper}}$  and it is almost insensitive to  $\lambda$  and  $\rho$ . This is intuitive because  $\pi_{\text{upper}}$  sets the confidence level with which the algorithms identify the anomalies, and therefore, for the same confidence level, the success rates achieved by the algorithms are almost the same.
- *Stopping time:* In Fig. 2, we show the variation of the stopping time  $K$  with  $\pi_{\text{upper}}$ . We see that the stopping time increases with  $\pi_{\text{upper}}$  in all cases, as a higher value of  $\pi_{\text{upper}}$  requires the algorithms to collect more observations before they make the decision regarding the anomalous processes. Also, we observe that the stopping time decreases with an increase in  $\rho$  for all values of  $\lambda$  and  $\pi_{\text{upper}}$ . This decrease is expected due to the fact that as the correlation increases, an observation corresponding to one of the dependent processes gives more information about the other. Consequently, the algorithms require fewer observations, and thus, a smaller stopping time, to achieve the same confidence level. Finally, we notice that the stopping time for the active inference algorithm is less than that of the actor-critic algorithm.

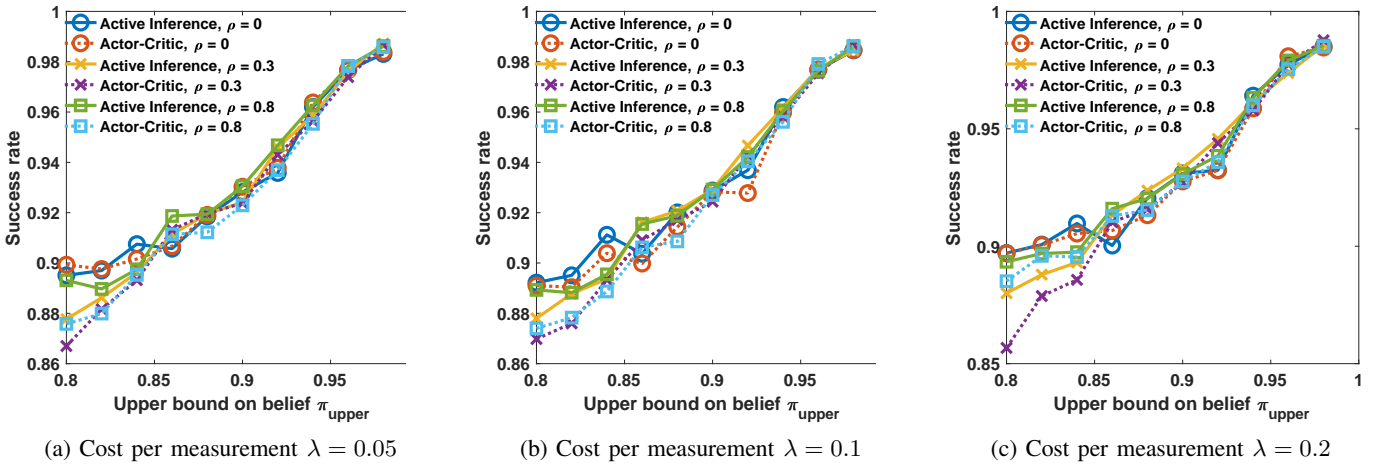


Fig. 1: Variation of the success rate of the active inference and the actor-critic algorithms when  $\pi_{\text{upper}}$ ,  $\lambda$  and  $\rho$  are varied.

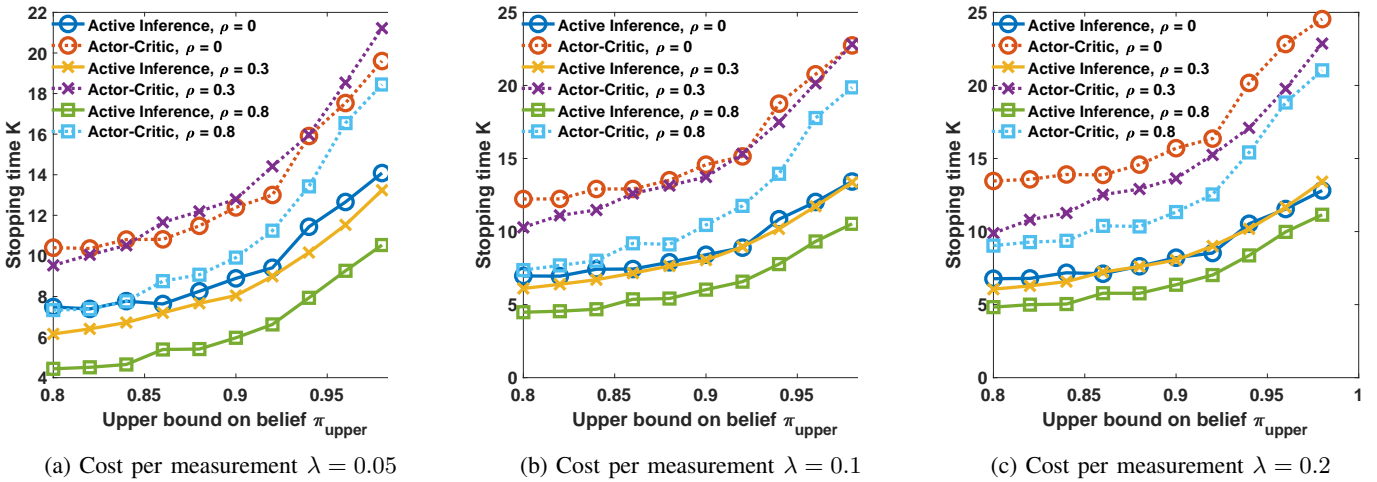


Fig. 2: Variation of the stopping time  $K$  of the active inference and the actor-critic algorithms when  $\pi_{\text{upper}}$ ,  $\lambda$  and  $\rho$  are varied.

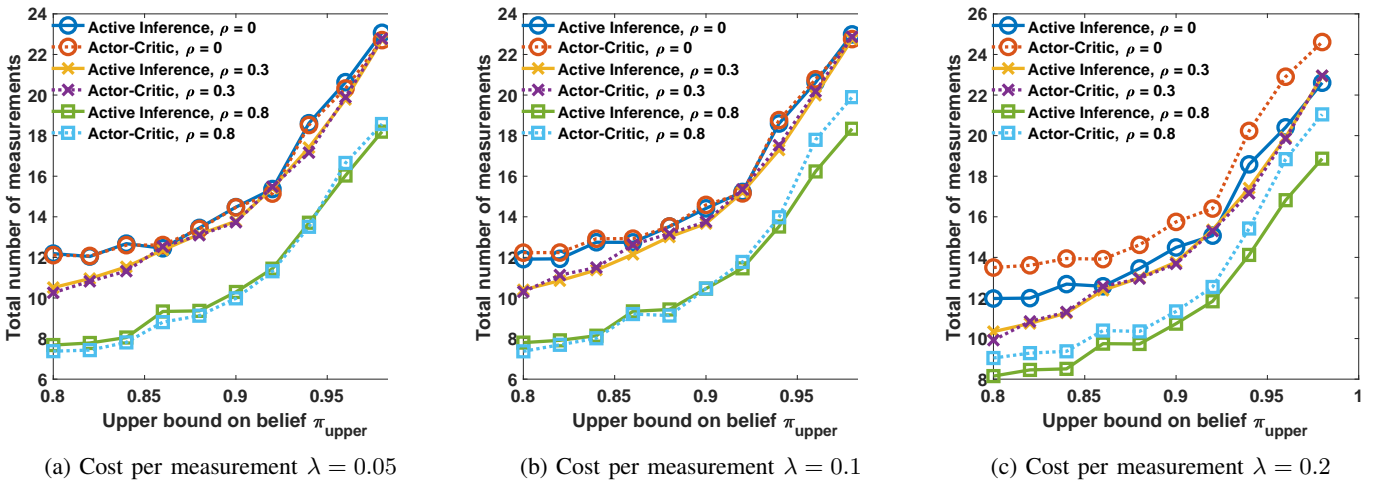


Fig. 3: Variation of the total number of measurements  $\sum_{k=1}^K |\mathcal{A}_k|$  of the active inference and the actor-critic algorithms when  $\pi_{\text{upper}}$ ,  $\lambda$  and  $\rho$  are varied.

- *Total number of measurements:* Fig. 3 compares the total number of measurements  $\sum_{k=1}^K |\mathcal{A}_k|$  obtained by the two algorithms in different settings. Clearly, the total number of measurements decreases with  $\rho$ , which is expected as mentioned above. Also, we infer that the total number of measurements obtained by both algorithms are similar in all the settings with the active inference algorithm collecting slightly fewer measurements compared to the actor-critic algorithm.

Thus, we conclude that the two algorithms achieve comparable success rates and incur a similar total cost of sensing, but the active inference algorithm has better stopping time compared to the actor-critic algorithm. This indicates that our algorithm identifies the anomalies faster than the actor-critic algorithm. Moreover, the stopping time of our algorithm does not vary much with  $\lambda$  while the stopping time of the actor-critic algorithm increases with  $\lambda$ . This implies that the actor-critic algorithm is more sensitive to the instantaneous cost of sensing  $\lambda |\mathcal{A}_k|$  than the total cost of sensing  $\sum_{k=1}^K \lambda |\mathcal{A}_k|$ . To elaborate, we note that both algorithms continue to acquire measurements until the desired level confidence level  $\pi_{\text{upper}}$  is achieved. However, since the actor-critic algorithm optimizes the average cost of sensing  $\frac{1}{K} \sum_{k=1}^K \lambda |\mathcal{A}_k|$ , as  $\lambda$  increases, it picks a fewer number of processes per time instant and this results in an increased stopping time. On the contrary, the average number of processes selected by our algorithm does not vary much with  $\lambda$ . Therefore, we achieve better performance by carefully designing the objective function using a novel entropy based-function and the total cost of sensing whereas the actor-critic algorithm optimizes the average change in entropy and the average cost of sensing.

## V. CONCLUSION

In this paper, we presented an anomaly detection algorithm using an active inference-based approach. We modeled the problem of anomaly detection as an active inference problem aiming at the detection accuracy exceeding a desired value while minimizing the delay and total cost of sensing. We designed a new objective function based on entropy and implemented the active inference algorithm using a deep learning-based approach. Through simulation results, we compared our algorithm with an algorithm based on the deep actor-critic method in terms of the success rate, stopping time, and total cost of sensing. The results demonstrated that our algorithm can detect the anomalies quicker (as indicated by the smaller stopping times) and achieves a competitive success rate with a similar cost of sensing as the actor-critic algorithm. However, we detect all the anomalous processes at a given time, assuming that the (normal or anomalous) behaviors of the processes remain unchanged until the agent makes a decision. Extending our algorithm to track any changes in the behavior of the processes over a longer time period is an interesting direction for future work.

## REFERENCES

- [1] W.-Y. Chung and S.-J. Oh, "Remote monitoring system with wireless sensors module for room environment," *Sensors Actuators B: Chemical*, vol. 113, no. 1, pp. 64–70, Jan. 2006.
- [2] A. Bujnowski, J. Ruminski, A. Palinski, and J. Wtrorek, "Enhanced remote control providing medical functionalities," in *Proc. Inter. Conf. Pervasive Comput. Tech Healthc. Workshops*, May 2013, pp. 290–293.
- [3] C. Zhong, M. C. Gurnoy, and S. Velipasalar, "Deep actor-critic reinforcement learning for anomaly detection," in *Proc. Globecom*, Dec. 2019.
- [4] G. Joseph, M. C. Gurnoy, and P. K. Varshney, "Anomaly detection under controlled sensing using actor-critic reinforcement learning," in *Proc. IEEE Inter. Workshop SPAWC*, May 2020.
- [5] H. Chernoff, "Sequential design of experiments," *Ann. Math. Stat.*, vol. 30, no. 3, pp. 755–770, Sep. 1959.
- [6] S. A. Bessler, "Theory and applications of the sequential design of experiments, k-actions and infinitely many experiments: Part I - theory," Stanford Univ CA Applied Mathematics and Statistics Labs, Tech. Rep., 1960.
- [7] S. Nitinawarat, G. K. Atia, and V. V. Veeravalli, "Controlled sensing for multihypothesis testing," *IEEE Trans. Autom. Control*, vol. 58, no. 10, pp. 2451–2464, May 2013.
- [8] M. Naghshvar, T. Javidi *et al.*, "Active sequential hypothesis testing," *Ann. Stat.*, vol. 41, no. 6, pp. 2703–2738, 2013.
- [9] B. Huang, K. Cohen, and Q. Zhao, "Active anomaly detection in heterogeneous processes," *IEEE Trans. Inf. Theory*, vol. 65, no. 4, pp. 2284–2301, Aug. 2018.
- [10] D. Kartik, E. Sabir, U. Mitra, and P. Natarajan, "Policy design for active sequential hypothesis testing using deep learning," in *Proc. Allerton*, Oct. 2018, pp. 741–748.
- [11] K. Friston, F. Rigoli, D. Ognibene, C. Mathys, T. Fitzgerald, and G. Pezzulo, "Active inference and epistemic value," *J. Cogn. Neurosci.*, vol. 6, no. 4, pp. 187–214, Oct. 2015.
- [12] K. Friston, T. FitzGerald, F. Rigoli, P. Schwartenbeck, and G. Pezzulo, "Active inference: A process theory," *Neural Comput.*, vol. 29, no. 1, pp. 1–49, Jan. 2017.
- [13] K. J. Friston, M. Lin, C. D. Frith, G. Pezzulo, J. A. Hobson, and S. Ondobaka, "Active inference, curiosity and insight," *Neural Comput.*, vol. 29, no. 10, pp. 2633–2683, Oct. 2017.
- [14] P. Schwartenbeck, J. Passecker, T. U. Hauser, T. H. FitzGerald, M. Kronbichler, and K. J. Friston, "Computational mechanisms of curiosity and goal-directed exploration," *Elife*, vol. 8, p. e41703, 2019.
- [15] B. Millidge, "Deep active inference as variational policy gradients," *J. Math. Psychol.*, vol. 96, p. 102348, Jan. 2020.