

arXiv:1905.04181v1 [cs.LG] 10 N

IEEE Copyright Notice

© 2019 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Accepted to be Published in: The 13th IEEE International Conference on Self-Adaptive and Self-Organizing Systems (SASO 2019), June 16 – 20, 2019, Umeå, Sweden

Autonomous Management of Energy-Harvesting IoT Nodes Using Deep Reinforcement Learning

Abdulmajid Murad^{1*}, Frank Alexander Kraemer¹, Kerstin Bach², Gavin Taylor³

¹Department of Information Security and Communication Technology

²Department of Computer Science

Norwegian University of Science and Technology, Trondheim, Norway

³Department of Computer Science, United States Naval Academy, Annapolis, USA

*Corresponding author: abdulmajid.a.murad@ntnu.no

Abstract—Reinforcement learning (RL) is capable of managing wireless, energy-harvesting IoT nodes by solving the problem of autonomous management in non-stationary, resource-constrained settings. We show that the state-of-the-art policy-gradient approaches to RL are appropriate for the IoT domain and that they outperform previous approaches. Due to the ability to model continuous observation and action spaces, as well as improved function approximation capability, the new approaches are able to solve harder problems, permitting reward functions that are better aligned with the actual application goals. We show such a reward function and use policy-gradient approaches to learn capable policies, leading to behavior more appropriate for IoT nodes with less manual design effort, increasing the level of autonomy in IoT.

Index Terms—reinforcement learning, IoT, power management

I. INTRODUCTION

Digitizing domains like smart cities, precision agriculture, manufacturing, and healthcare requires the instrumentation of the physical world with wireless IoT nodes to sense environmental variables. The actual logic of an IoT device can be rather simple, as they usually make measurements, prepare the data and send it into a central fusion center for aggregation. However, challenges arise from the limited energy resources of the wireless and miniaturized nodes. Therefore, the choice of parameters like duty cycle (percentage of time a node is active) and sampling frequency is crucial since they have a dominant effect on the energy consumption and performance of the node. Currently, such parameters are often assigned manually to an entire family of IoT devices, and stay constant during deployment, or are determined by manually fine-tuned algorithms. However, IoT in its full scale presents a setting where nodes are placed in dynamic and non-stationary environments, so that a one-fits-all approach and manual design efforts are impossible. Future solutions should instead be based on autonomous nodes that optimize their behavior by learning within their specific environment.

Our vision for a more autonomous and hence scalable IoT is inspired by recent developments in deep reinforcement learning (RL), and its success in building capable autonomous

agents with little manual design effort in complex domains. Instead of manually adapting IoT nodes, we want engineers to only formulate the goal in the form of a reward function for the IoT nodes as close to the actual application goal as possible.

Reinforcement learning has indeed been used to control IoT nodes by various works, especially for power management [1]–[3]. However, these approaches use older RL techniques like SARSA and are constrained by coarsely discretized action and observation spaces, and tackle only simpler problems. Many approaches [4]–[6] thus formulate reward functions in terms of indirect indicators like energy budget. Those are often easier to learn, but are imperfectly correlated with actual goals. This not only adds manual design effort, but also doesn't live up to the potential of RL. Instead, RL should be used to sort out the complicated and node-specific technicalities on its own through learning.

In this paper, we explore deep RL as an end-to-end approach to enable autonomy within IoT nodes. For the training, we use a policy gradient method, namely Proximal Policy Optimization (PPO, [7]). We use neural networks as function approximators to learn state abstractions and effective policies directly from continuous input data. We also develop a reward function that closely reflects application goals. We automated the search for the best reward signal by defining it as hyperparameters from the application perspective. Additionally, we conduct a series of simulations to explore the influence of reward function design on a node's behavior.

To manage the computational complexity, we propose training the agent on a server as a part of device management [8] and update the node regularly based on a static or dynamic update interval [1]. To this end, we built a sensor simulator, called Sensor Gym (based on OpenAI Gym [9]), as a toolkit for training and comparing various RL algorithms in an IoT context.

The results in this paper show that modern and more successful deep RL techniques outperform the older ones, and that they make it possible to use more sensible reward functions that are based on IoT application goals. This leads to tractable learning, less manual design efforts, and hence scalable autonomy for IoT systems.

The rest of this paper is organized as follows: Sect. II provides a brief overview of deep reinforcement learning, and Sect. III presents related work that adopted RL to optimize energy-harvesting sensor nodes. In Sect. IV, we describe the proposed system setup, and in Sect. V we present a baseline for performance and suitability of PPO in an IoT by comparing it to an older RL technique and study the performance with a reward function based on energy neutrality. We then introduce a reward function closer to the application goals in Sect. VI, and present the results in Sect. VII.

II. DEEP REINFORCEMENT LEARNING

In reinforcement learning an agent learns to better perform a task by learning from its experiences in interacting with that task. Mathematically, the task is expressed with a Markov Decision Process (MDP). An MDP is a tuple $M = (\mathcal{S}, \mathcal{A}, \mathcal{P}, R, \gamma)$, where \mathcal{S} is the set of all states the problem might be in; \mathcal{A} is the set of all actions the agent might take in response; \mathcal{P} is a transition kernel describing the probabilities of transitioning between two states when performing an action ($p(s'|s, a) \forall s, s' \in \mathcal{S}, a \in \mathcal{A}$); $R: \mathcal{S} \rightarrow \mathbb{R}$ is a reward function (for example, positive numbers in states where the agent has achieved its goal, and negative numbers in states where the agent has harmed itself); and $\gamma \in (0, 1)$ is a discount factor, which is used to describe how much the agent prefers rewards in the short term to rewards in the long term. The behavior of the agent is expressed as a policy $\pi: \mathcal{S} \rightarrow \mathcal{A}$.

The quality of a policy can be measured using the value function of a state $V^\pi(s)$, or, the Q-value of a state-action pair $Q^\pi(s, a)$:

$$V^\pi(s) = R(s) + \int_{\mathcal{S}} p(s'|s, \pi(s)) V^\pi(s') ds' \quad (1)$$

$$Q^\pi(s, a) = R(s) + \int_{\mathcal{S}} p(s'|s, a) Q^\pi(s', \pi(s')) ds' \quad (2)$$

The goal of the agent is to learn a policy which produces high value for all states by collecting trajectory samples of the form (s_t, a_t, r_t, s_{t+1}) , where t denotes the timestep, and $s_t, s_{t+1} \in \mathcal{S}, a_t \in \mathcal{A}, r_t = R(s_t)$.

A variety of learning algorithms for RL exist. One traditional approach is to begin by learning Q-functions as accurately as possible using function approximation techniques; some are used in previous work applying RL to IoT, and will be introduced here.

Samples are themselves collected using a policy π_{smp} , where $a_t = \pi_{smp}(s_t)$. The Q-functions of this sampling policy can be learned using an algorithm called SARSA; this is known as *on-policy* learning. Alternatively, the Q-functions of an improved policy may be learned using an algorithm known as Q-learning; this is known as *off-policy* learning. In either case, given Q-function approximations, a new and likely-improved policy π^* is implied by $\pi^*(s) = \arg \max_{a \in \mathcal{A}} Q^\pi(s, a)$. With sufficient sampling, alterations to allow for exploration of the full state space, and (in the on-policy case) alternation between sampling, Q-function approximation and policy improvement steps, these approaches are

proven to converge to excellent policies. This is a convenient guarantee, but requires the state and action spaces to be finite and small enough to be easily tabulated. More complex approaches to value learning exist which allow state and action spaces to be continuous, but this results in a loss of these guarantees, and potentially much less capable policies. In recent years, these approaches have been further improved upon by approximating functions with deep neural networks [10], but they remain fundamentally difficult to apply successfully to continuous and complex domains.

In contrast with value-based methods, policy-search methods directly search for a parameterized policy π_θ . The result is a stochastic policy that learns a direct mapping from states to a probability distribution over actions, where preferred actions have a higher probability of being sampled. This output of a distribution makes it suitable for continuous domain tasks.

Policy gradient methods are the most prominent of policy-search methods. They optimize a policy by maximizing the value of sampled states $V^{\pi_\theta}(s_t)$ by performing gradient ascent over the parameters θ . Various approximations to this basic approach exist, with advantages in training time and sample efficiency. In this paper, we use Proximal Policy Optimization (PPO), one of the most successful of these approaches [7]. PPO is a trust region method, which seeks to iteratively maximize performance of π_θ compared to its previous iteration, without changing it too much. This allows for continuous, stable improvement, even in continuous and complex domains such as robotic running and difficult Atari games. This success in demanding domains make RL ready for application to real-world problems such as IoT.

III. RELATED WORK

Several works have already adopted RL within various contexts of IoT. An example is a dynamic power manager for energy-harvesting wireless sensor networks [11]. Here Hsu et al. used the Q-learning algorithm to train an agent to choose an action from four levels of duty cycles. Their environment state space is based on the distance from energy neutrality (i.e., the difference between harvested and consumed energy), the harvested energy, and the current battery level, while the reward function is based on the distance from energy neutrality and the current energy storage level. They extended this work in [12] to include quality of service (QoS), both in the state space and the reward function. Additionally, in [3], they included meeting a requested throughput along with the energy level in the reward function and added penalizing terms for overcharging, deep-discharge, and depletion of the energy storage. Furthermore, in [4] and [5], they used fuzzy decision processes to model an energy harvesting node as a fuzzy environment and used a modified Q-learning with fuzzy reward to train an RL agent.

Since the reward function formalizes the goal of an RL setup, Rioual et al. [13] discussed the choice of the reward function in the management of energy-harvesting IoT nodes but covered a limited range of design choices. To avoid intractable learning in old RL approaches when using an

application-level reward, most of the previous work used reward shaping, which is an alternative method to guide the learning process by rewarding the agent for achieving subgoals or developing an approximation to the desired behaviors. Unless the reward shaping function is based on a state-dependent potential, it may lead to learning suboptimal or locally-optimal policies [14]. Almost all previous work of using RL in IoT uses manually-designed shaped functions that produce arguably acceptable results without justifying how well the reward frames the application goal.

RL has also been used for power allocation in energy-harvesting communication systems. Ortiz et al. [15] used the SARSA algorithm with linear function approximation to learn a power allocation policy in two-hop communication. The objective of their policy is to maximize throughput of a communication system, but they designed a reward function based on the total power assigned to transmissions, which they claimed was correlated. They simulated a communication environment and approximated the state space with discrete features that indicate battery level and its constraints, harvested energy over an hour, characteristics of the communication link, data arrival process, and the data buffers at the communicating nodes. Similarly, Aoudia et al. [6], used an actor-critic method with linear function approximation to learn approximation for both policy and value function. They used a Gaussian policy to generate continuous values of bounded packet rates and summarized the state space by continuous values of the current residual energy. The objective is to maximize transmitted packet rate while sustaining perpetual operation; hence they designed a reward function to be a multiplication of normalized residual energy and the packet rate.

Shresthamali et al. [2] used a SARSA(λ) RL algorithm to develop adaptive power management for a solar-energy harvesting sensor node. To simulate a sensor node, they used a scaled-up version of a real sensor powered by a battery and a solar panel, and used solar radiation data to calculate hourly harvested energy. They designed a reward function based the distance from energy neutrality, defined as the difference between the current level of energy and the optimum battery level, which they calculated to be 60% of the battery's maximum level. They trained an agent in episodes of 24 hours and rewarded it at the end of a training episode. The state space consisted of discretized information about the battery level, the distance from energy neutrality, the harvested energy, and the weather forecast, which they approximate by calculating the total amount of energy harvested in a particular day.

Fraternali et al. [1] focused on the configuration of the sampling rate of indoor energy harvesting sensors. The objective of their agent is to maximize the number of samples while avoiding power failure, by designing a reward function that depends on the sampling rate and a penalty for power failure. They used Q-learning with a state space that included light intensity, the voltage level of the energy storage, and the current sampling rate. Another example of using RL for adaptive sampling can be found in [16]. Here Dias et al. proposed using Q-learning for adaptive sampling rate adaption to avoid

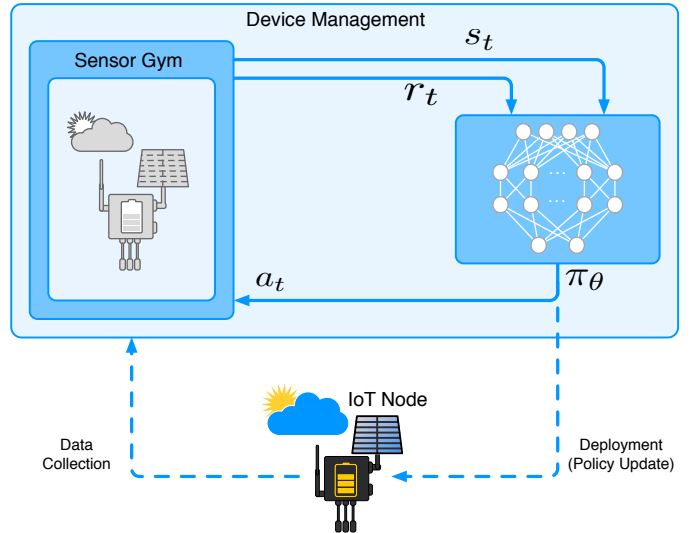


Fig. 1. System setup of agent-environment interaction

oversampling, while not missing environmental changes. They introduced a variable that describes the quality of measurement as the difference between two consecutive measurements and, depending on a specific application, this difference should be less than a threshold value. The action space is a range of possible sampling intervals, while the reward function is based on the transmission avoided and the quality of measurements.

IV. SYSTEM SETUP

In this paper, we propose and apply deep RL solutions for training RL agents to control IoT nodes autonomously. We use PPO as a policy gradient method for optimizing the agent's policy, using neural networks as function approximators. Executing the agent's policy corresponds to inference in a neural network which has become feasible even for computationally- and energy-constrained IoT nodes [17]. We propose to integrate the much more computation-intensive training phase on a server, as a part of the IoT device-management [8], illustrated in Fig. 1. The deployed agents can be updated regularly based on a static or dynamic update interval [1].

To train the agent, we built a simulator of a general sensor node with an energy-harvesting solar panel, an ideal energy buffer, and a load with variable duty-cycles. We base the simulator on the OpenAI Gym, which is a toolkit for developing and comparing RL algorithms [9]. We base our sensor specification on a realistic energy-harvesting IoT node. The simulated sensor has an energy buffer with a maximum capacity of $B_{max} = 40 \text{ W}$ and has a load that consumes energy during a time step according to its duty cycle. The maximum consumed energy in an hour is $Ec_{max} = 0.5 \text{ Wh}$, which corresponds to a duty cycle of $D_{max} = 100\%$ during that hour. We simulate the harvested energy by calculating the energy generated by a 6W solar panel using solar radiation data. At time step t , we represent the level of the energy buffer as B_t , the harvested energy as Eh_t , the consumed energy as

Ec_t , and the duty cycle, which is determined by the node's policy, as D_t . The consumed energy is calculated according to:

$$Ec_t = \begin{cases} 5D_t & \text{if } B_t > 5D_t \\ 0 & \text{if } B_t \leq 5D_t \end{cases} \quad (3)$$

The next level of the energy buffer is calculated according to:

$$B_{t+1} = \min(B_t + Eh_t - Ec_t, B_{max}). \quad (4)$$

To handle a continuous action space $\mathcal{A} = [0, D_{max}]$, the policy resulting from PPO is defined as a Gaussian distribution with mean $\mu(s_t)$ and standard deviation $\sigma(s_t)$, both of which are approximated by a neural network with parameters θ :

$$\pi_\theta(a_t|s_t) = \frac{1}{\sigma(s_t)\sqrt{2\pi}} \exp\left(-\frac{(a_t - \mu(s_t))^2}{2\sigma(s_t)^2}\right) \quad (5)$$

Thus, the actions are real numbers, chosen from the normal distribution in 5. Our neural network architecture has an output layer with two neurons which are the approximation of the mean and standard deviation of the Gaussian policy.

In Section V, we demonstrate that PPO outperforms older RL approaches on previously-proposed reward functions. Then, in Section VI, we demonstrate that PPO can develop good policies for reward functions that make the domain more difficult, but better capture desired behavior. We also illustrate that older RL approaches cannot produce similarly competent policies.

V. REWARD FUNCTION BASED ON ENERGY NEUTRALITY

To get a baseline for the performance and suitability, we first apply PPO in the same setting as Shresthamali et al. [2], who used the SARSA(λ) algorithm and designed a reward function based on the distance from energy neutrality. The concept of energy-neutrality was introduced by Kansal et al. [18], which states that a node is in energy-neutral operation if the consumed energy is less than or equal to the harvested energy. Accordingly, the distance from energy neutrality is redefined as the difference between the current level of the energy buffer B_t and the optimum buffer level B_0 to account for the variance in the harvested energy over a period:

$$Edist_t = |B_t - B_0|. \quad (6)$$

In [2], the reward function is formulated in terms of this distance to energy neutrality:

$$R_E(s_t) = \begin{cases} 500 & \text{if } Edist_t = 0 \text{ W h} \\ 500 - \frac{Edist_t}{10} & \text{if } 0 \text{ W h} < Edist_t \leq 1 \text{ W h} \\ 10 - \frac{Edist_t}{100} & \text{if } 1 \text{ W h} < Edist_t \leq 5 \text{ W h} \\ -500 & \text{if } 5 \text{ W h} < Edist_t \end{cases} \quad (7)$$

The action space is defined as discrete values of five duty cycles, which are

$$a_t \in \mathcal{A} = \{20\%, 40\%, 60\%, 80\%, 100\%\} \quad (8)$$

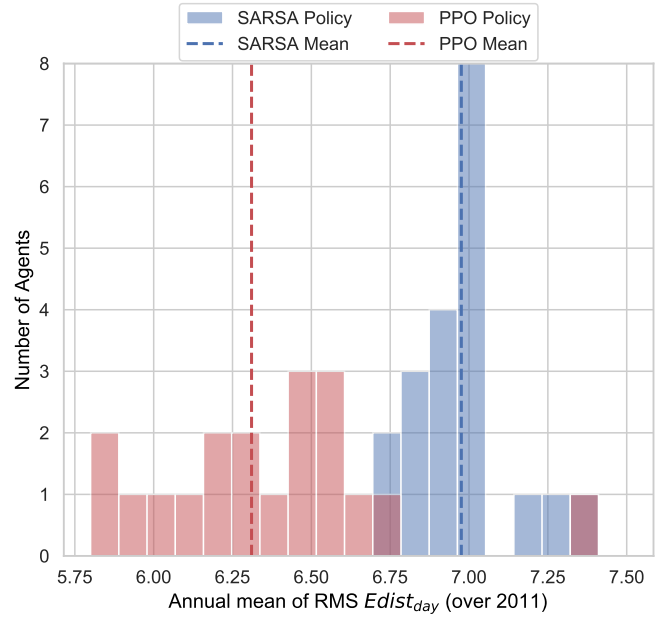


Fig. 2. Performance comparison of 20 agents trained with PPO and SARSA(λ)

The observation space contains vectors with four observations:

$$s_t = [B_t, Edist_t, Eh_t, Wf_{day}], \quad (9)$$

where Wf_{day} represents the weather forecast of the day. In [2], actual harvested energy data in a day are calculated before the training and used to give information about the expected weather. The agents can leverage this information to plan their energy expenditure accordingly.

While the state space requires discretization for SARSA with a hand-designed mapping, policy approximation has the advantage of handling continuous state spaces by generalization, hence eliminating another step of manual mapping an IoT problem to RL.

To compare the capabilities of SARSA and PPO, we trained 20 agents with PPO and 20 with SARSA with data of Tokyo 2010 and evaluated their policies on data of Tokyo 2011, using the same settings and data as [2]. Figure 2 represents a histogram for performance results of both policies. Here, the root mean square (RMS) of $Edist_{day}$ values are used to measure the deviation from energy-neutrality after a one-day window and the mean of daily deviation over the whole year to compare policies. We observe that the average PPO agent is considerably better than the best SARSA policy and that almost all PPO agents are better than the average SARSA agent. These excellent results support our motivation to use PPO for IoT agents.

We also take a more detailed look at the behavior of the agents by comparing their results over a one-week window. Figure 3 shows the performance results of PPO, SARSA policies, and linear programming when running them over one week of Tokyo weather data in February of 2011. The linear



Fig. 3. Comparison of PPO policy to SARSA(λ) and offline policy for a week of Tokyo 2011

programming policies are presented in [18], and use linear programming and acausal data (i.e., the exact future energy intake, hence not practical) to determine the optimal duty cycle given the energy neutrality constraint; this represents the upper limit of performance. The three policies started the week off with an energy level of $B_t = 60\%$ of B_{max} and try to maintain them at this level. $RMS\ Edist_{day}$ values are used to measure the deviation from energy-neutrality over the whole week, and we observe that the SARSA-trained policies have the highest deviation of (5.59%), PPO-trained policies have (5.43%) and the optimal policy based on linear programming has (4.17%).

Figure 3 illustrates a problem with the reward function based on distance from energy neutrality (7): The duty cycles of both RL policies are subject to high variance, often oscillating between highest and lowest duty cycle, instead of choosing a smoother course with more mid-range values. From an RL point-of-view, this maximizes the reward, but it is not a behavior appropriate for IoT nodes. Typically, we want to cover a phenomenon as continuously as possible, and spread out measurements over time, which corresponds to a smoother duty cycle.

Because PPO is so capable on this easier problem, in the next section we design a new, more difficult reward function,

which no longer rewards the intermediate goal of energy neutrality, but more clearly expresses the goals for the agent.

VI. REWARD FUNCTION BASED ON APPLICATION GOALS

Our objective in this section is to create a reward function that is closer to the actual application goals of an IoT system in order to encourage more desirable behavior. In particular, we want to depart from the inclusion of energy terms in the reward function, since energy is only a resource to be managed, but not a goal to be optimized. Instead, we want the agent to maximize the sum of the duty cycle over time Γ :

$$G_D = \sum_{t=0}^{\Gamma} D_t. \quad (10)$$

At the same time, IoT nodes must not empty their energy buffer completely. Not only would this prevent the node from taking any measurements until enough energy is harvested again, but it could also lead to the loss of data or leave the node in an undefined state. Since the sensor gym simulates

a failure when the idealized buffer is emptied, we want to minimize the occurrence of failures over time Γ

$$G_F = \sum_{t=0}^{\Gamma} \begin{cases} 1 & \text{if } B_t = 0 \\ 0 & \text{if } B_t > 0 \end{cases} \quad (11)$$

To ensure a smooth course of the duty cycle, corresponding to a continuous stream of measurements, we want to minimize the variance of the selected duty cycles, and so minimize

$$G_{Var} = \sum_{t=0}^{\Gamma} Var_t, \quad (12)$$

where the variance in the duty cycle is defined in an epoch t as the absolute difference to the previous duty cycle

$$Var_t = |D_t - D_{t-1}|. \quad (13)$$

We combine the conflicting objectives (10), (11), (12) into one reward function R_A :

$$R_A(s_t) = \begin{cases} D_t - \zeta [Var_t]^2 & \text{if } B_t > 0 \\ -F & \text{if } B_t = 0 \end{cases} \quad (14)$$

The agent receives its reward by maximizing the duty cycles D_t . To punish variance, we reduce the reward with the squared variance, scaled by a damping factor ζ . In the case of a failure, the reward is negative, using the punishment term F . The actual values of F and ζ are hyperparameters of the reward function, and we will have a detailed look at them in the next section.

We generalized the reward function in (14) to accommodate other RL techniques of sparse reward assignments. In these techniques, the agent gets a reward at the end of a training episode E . An episode consists of a set of epochs (time steps) that range from $t = 1$ to $t = T$, thus $E = \{1, 2, \dots, t, \dots, T\}$

$$R_A(s_T) = \sum_{t=1}^T \begin{cases} D_t - \zeta [Var_t]^2 & \text{if } B_t > 0 \\ -F & \text{if } B_t = 0 \end{cases} \quad (15)$$

If the reward is assigned at every epoch, $T = 1$ and (15) reduces to (14).

We considered the relevant attributes of the environment in the definition of state space. These attributes are observations of the current level in the energy buffer, current harvested energy, and the weather forecast of the whole episode. To simulate the harvested energy, we use solar radiation data to calculate the harvested energy in every hour. The weather forecast information is included to enhance performance by giving an estimate of the expected solar energy for that particular episode. This weather information can be acquired from external sources or real prediction algorithms with sufficient accuracy as the one presented in [19]. For our case study, we simulate weather information by calculating the total harvested energy in a particular day and introduce a 20% error to mimic inaccuracies in the weather forecast. We also include the previous duty cycle D_{t-1} in the state space, so the agent makes an informed decision to avoid high variability

in the duty cycle. The agent takes the observation as input summarized in a vector of four continuous values

$$s_t = [B_t, Eh_t, Wf_t, D_{t-1}]. \quad (16)$$

PPO allows us to use a continuous action space corresponding to setting the operational duty cycle of a node. This enables more accurate control of the consumed energy, and therefore the utility of a node. Therefore, the action space is

$$a_t \in \mathcal{A} = [D_{min}, D_{max}]. \quad (17)$$

Depending on IoT application requirements, the minimum value of the duty cycle can be set accordingly. For our case, we set it to $D_{min} = 0$.

VII. RESULTS AND DISCUSSION

We conducted a series of simulations to systematically explore the influence of the designed reward function R_A and training hyperparameters on the agent behavior. In total, we trained more than 300 agents on data of the year 2010 and tested their performance on data of the year 2011, using different values for the damping factor (ζ) for the reward function. For the PPO algorithm, we also trained with different values for the hyperparameters, namely the learning rate (α), batch size (number of state transitions used to calculate policy gradient),

discount factor (γ) and trace decay parameter (λ) of the advantage function [7]. In the following, we discuss the learning rate α and damping factor ζ in more detail.

The learning rate (α) has a significant impact on the learning process as shown in Fig. 4. The x-axis corresponds to the learning rate, which we chose to tune in the range of 10^{-5} to 10^{-1} . The y-axis in the upper plot corresponds to the yearly utilized energy (normalized by the maximum possible utilization), the y-axis in the lower plot corresponds to the mean-variance in duty cycle over the whole year. Each dot in the figure represents a single agent, and its color indicates the number of times the agent has emptied the energy buffer, i.e., failed. Training with low learning rates ($\alpha < 10^{-3}$) results in agents that learn poorly or not at all. These agents pick constant and low duty cycles, which results in low variance and few power failures, thus fulfilling application goals (11) and (12), but failing to achieve the goal of maximizing utilized energy in (10).

The choice of damping factor ζ in the reward function also has a significant impact on the resulting behavior. Figure 5 shows the results of all trained agents running over the whole year of 2011 for different damping factors ζ on the x-axis. We observe that agents trained with values of ζ in the approximate range ($10^{-2} < \zeta < 10^{-1}$) display the desired behavior of less variability and better energy utilization, which translates to an appropriate, continuous measurement coverage of the IoT node. We also observe that training with ζ in the approximate range ($\zeta \geq 10^{-1}$) leads to overdamped agents, which have less variance and fewer power failures, but also less energy utilization. Contrarily, most agents trained with values of ζ in the approximate range ($\zeta \geq 0.1$) are underdamped, with few

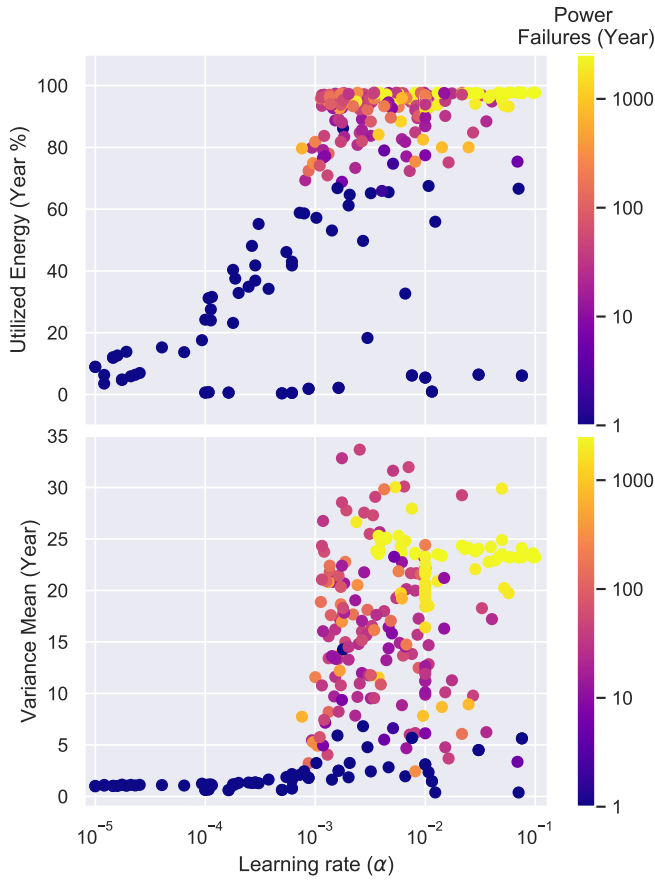


Fig. 4. Influence of the learning rate (α) on agent's learning process in terms of variance in duty cycle and utilized energy.

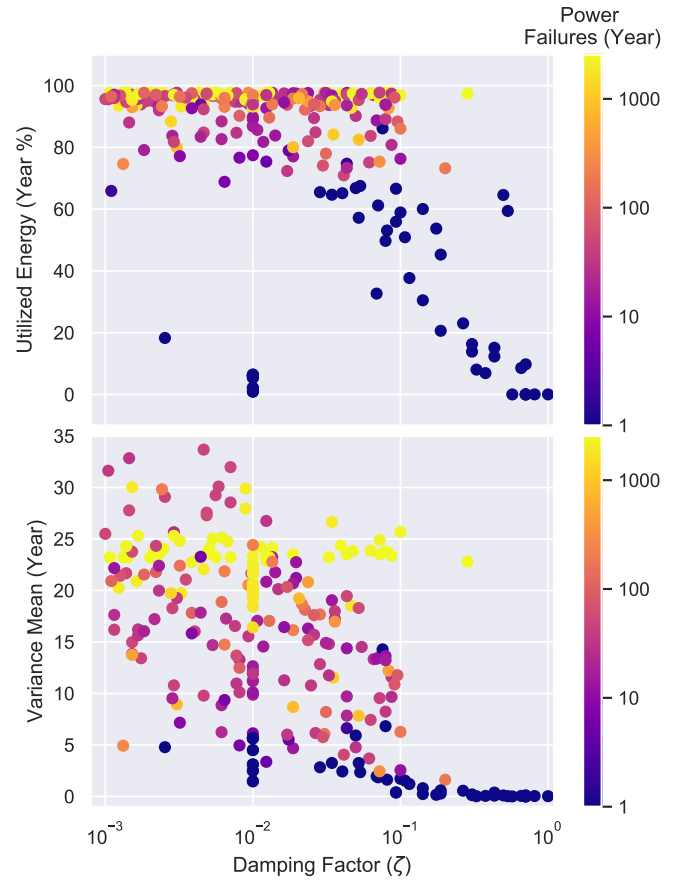


Fig. 5. Influence of the damping factor (ζ) in $R_A(t)$ on an IoT node's behavior in terms of utilized energy and variance in duty cycle.

exceptions as outliers. These underdamped agents have high utilized energy, but also higher variability in the duty cycle and more frequent power failures.

Depending on the choice of damping factor ζ and failure penalty F in the reward function R_A , agents try to maximize rewards by prioritizing among the conflicting objective (10), (11), (12). Accordingly, agents learn various policies that achieve different performance on each objective, as shown in Fig. 6. This enables choosing between different trade-off agents based on the context of individual IoT applications.

In comparison, using the same reward function R_A for SARSA-trained agents, achieves mediocre results at best, as shown in Fig. 6. The SARSA agents (shown as crosses) score poorly in terms of energy utilization and have a high variance. We attribute their relatively low failure rate to low utilization. This is as expected since the discretization necessary for SARSA does not offer the granularity needed for effective learning.

We also take a more detailed look at the behavior over a one-week window. Figure 7 shows the performance of four selected agents trained with different damping factors ζ . The first day yields much solar power, so all agents utilize this energy by setting their duty cycles to high values. Since the

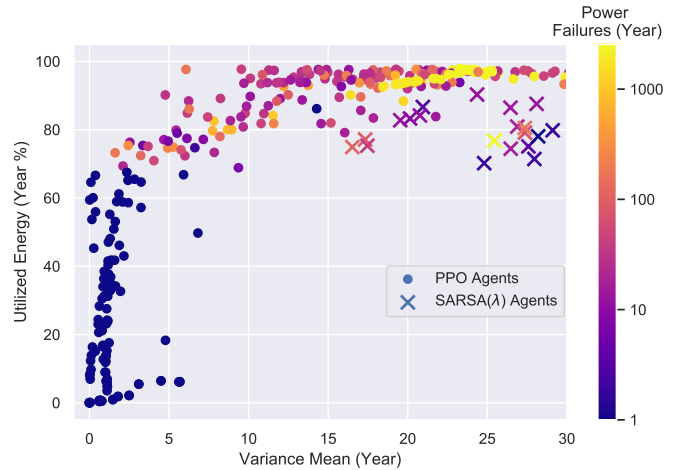


Fig. 6. Performance comparison of SARSA(λ) agents with PPO agents using the reward function R_A

weather forecast for the next day indicates less energy, all agents reduce their duty cycles. However, the pattern of the decrease differs according to the damping factor. Efficiently damped agents show the desired gradual decrease, while underdamped agents increase sharply and oscillate between

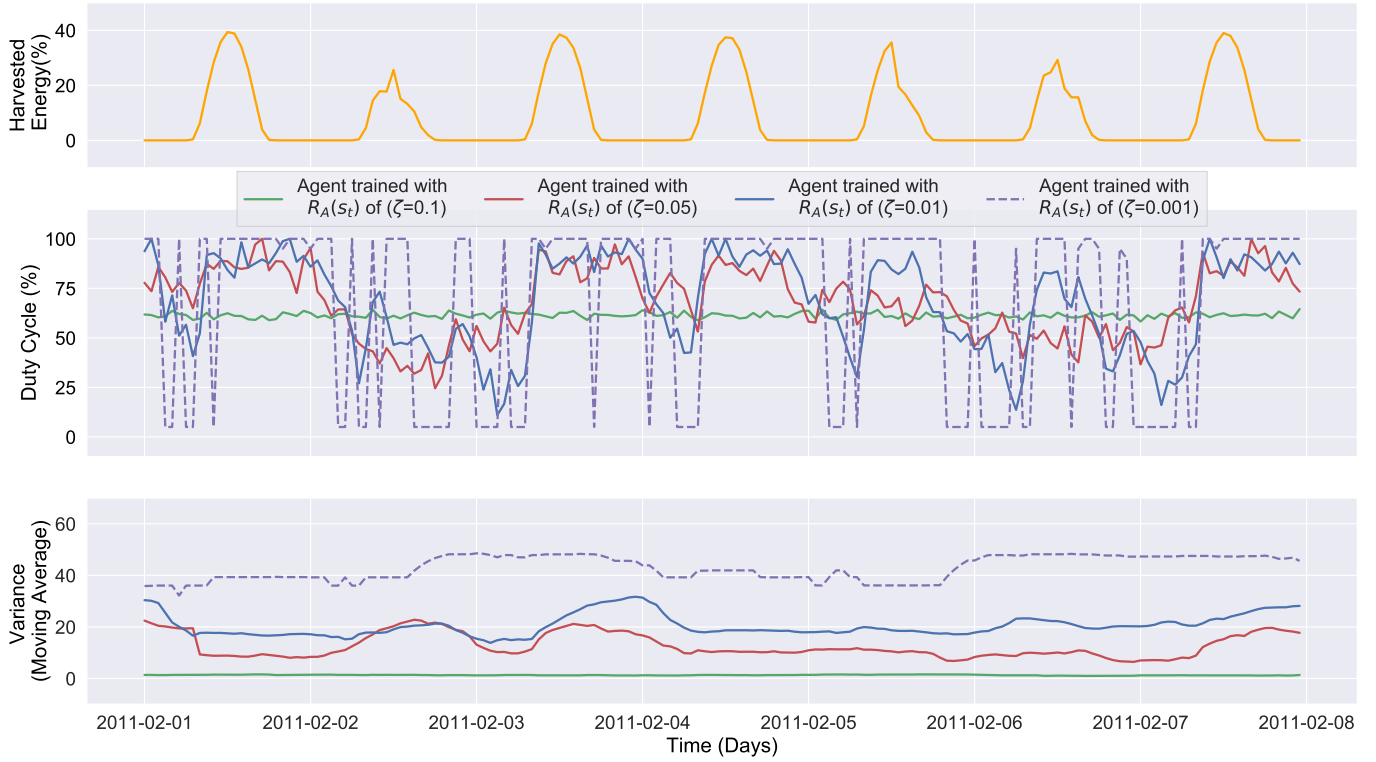


Fig. 7. Performance results over one week (Tokyo 2011) of four selected agents trained with different damping factors (ζ) in their reward function $R_A(t)$

TABLE I
SUMMARY OF PERFORMANCE RESULTS OVER A YEAR OF FOUR SELECTED AGENTS TRAINED WITH DIFFERENT DAMPING FACTORS (ζ)

ζ	Utilized Energy	Variance Mean	Power Failures
0.1	77%	2.5	11
0.05	94%	9.7	14
0.01	95%	12	17
0.001	98%	25.6	23

maximum and minimum duty cycles. When more energy is available again, we see analogous behavior when the duty cycles increase. The overdamped agents appear unaffected by the variance in energy supply, but they exhibit an overall low energy utilization. For example, the agent trained with $\zeta = 0.1$ avoids a high variance penalty by setting its duty cycle to a constant value, which leaves energy unutilized in times when much solar energy is available. Table I summarizes the overall performance of the agents in Fig. 7 over the whole year.

To investigate the effect of the neural network architecture of the policy approximation on performance, we trained with a different number of hidden layers and units. Our results show that shallower networks perform as well as deeper networks. We ended up using an architecture layout with two hidden layers of 64 units and a tanh activation function. The output layer has two units computing the mean and standard deviation of the Gaussian policy (5). Therefore, we substantiate that an

agent's policy can be approximated with a neural network that requires low computational effort and memory footprint, which makes it feasible for deployment also in resource-constrained sensor nodes.

VIII. CONCLUSION

Using reinforcement learning (RL) to manage constrained IoT nodes provides a path to learn optimal behavior without careful human attention. To the best of our knowledge, this is the first application of a policy-gradient method to this problem. We have shown that state-of-the-art policy-gradient RL methods such as PPO which use neural networks as function approximators are suitable for the use in IoT, that they outperform older RL approaches, and that they can solve more difficult problems which better describe and encourage desired behavior. Additionally, their suitability for continuous problems removes another manual design step of discretization. We also investigated the influence of hyperparameters on the resulting policies. While we have focused on the maximization of duty cycle with minimal variance, we believe there is potential to solve much more complex problems. The work presented hence leads to more autonomy in IoT systems, in which RL takes care of technicalities so that engineers can focus on the real application goals. This will allow a broader application of IoT in complex domains.

ACKNOWLEDGMENTS

This research was supported by the Office of Naval Research (N0001418WX01582) and the Department of Defense High Performance Computing Modernization Program.

REFERENCES

- [1] F. Fraternali, B. Balaji, and R. Gupta, "Scaling configuration of energy harvesting sensors with reinforcement learning," in *Proceedings of the 6th International Workshop on Energy Harvesting & Energy-Neutral Sensing Systems*. ACM, 2018, pp. 7–13.
- [2] S. Shresthamali, M. Kondo, and H. Nakamura, "Adaptive Power Management in Solar Energy Harvesting Sensor Node Using Reinforcement Learning," vol. 16, no. 5s, pp. 1–21, 2017.
- [3] R. C. Hsu, C. T. Liu, and H. L. Wang, "A reinforcement learning-based ToD provisioning dynamic power management for sustainable operation of energy harvesting wireless sensor node," *IEEE Transactions on Emerging Topics in Computing*, vol. 2, no. 2, pp. 181–191, 2014.
- [4] C. T. Liu and R. C. Hsu, "Dynamic power management utilizing reinforcement learning with fuzzy reward for energy harvesting wireless sensor nodes," *IECON Proceedings (Industrial Electronics Conference)*, pp. 2365–2369, 2011.
- [5] R. C. Hsu, T. H. Lin, S. M. Chen, and C. T. Liu, "Dynamic energy management of energy harvesting wireless sensor nodes using fuzzy inference system with reinforcement learning," *Proceeding - 2015 IEEE International Conference on Industrial Informatics, INDIN 2015*, pp. 116–120, 2015.
- [6] F. A. Aoudia, M. Gautier, and O. Berder, "RLMan: an Energy Manager Based on Reinforcement Learning for Energy Harvesting Wireless Sensor Networks," *IEEE Transactions on Green Communications and Networking*, vol. 2, no. 2, pp. 1–1, 2018.
- [7] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," *arXiv preprint arXiv:1707.06347*, pp. 1–12, 2017.
- [8] A. E. Braten and F. A. Kraemer, "Towards Cognitive IoT: Autonomous Prediction Model Selection for Solar-Powered Nodes," in *2018 IEEE International Congress on Internet of Things (ICIOT)*. Seattle, USA: IEEE, 2018, pp. 118–125.
- [9] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "OpenAI Gym," *arXiv preprint arXiv:1606.01540*, pp. 1–4, 2016.
- [10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [11] R. Chaoming Hsu, C. T. Liu, and W. M. Lee, "Reinforcement learning-based dynamic power management for energy harvesting wireless sensor network," in *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*, vol. 5579 LNAI. Springer, 2009.
- [12] R. C. Hsu, C.-T. Liu, K.-C. Wang, and W.-M. Lee, "Qos-aware power management for energy harvesting wireless sensor network utilizing reinforcement learning," in *2009 International Conference on Computational Science and Engineering*, vol. 2. IEEE, 2009, pp. 537–542.
- [13] Y. Rioual, Y. L. Moullec, J. Laurent, M. I. Khan, and J.-p. Diguët, "Reward Function Evaluation in a Reinforcement Learning Approach for Energy Management," in *2018 16th Biennial Baltic Electronics Conference (BEC)*. IEEE, 2018, pp. 1–4.
- [14] A. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," in *ICML*, 1999, vol. 99, pp. 278 —287.
- [15] A. Ortiz, H. Al-Shatri, X. Li, T. Weber, and A. Klein, "Reinforcement Learning for Energy Harvesting Decode-and-Forward Two-Hop Communications," *IEEE Transactions on Green Communications and Networking*, vol. 1, no. 3, pp. 309–319, 2017.
- [16] G. M. Dias, M. Nurchis, and B. Bellalta, "Adapting sampling interval of sensor networks using on-line reinforcement learning," in *2016 IEEE 3rd World Forum on Internet of Things (WF-IoT)*. IEEE, 2016, pp. 460–465.
- [17] "uTensor project," <https://github.com/utensor>, accessed: 2019-02-28.
- [18] A. Kansal, J. Hsu, S. Zahedi, and M. B. Srivastava, "Power management in energy harvesting sensor networks," *ACM Transactions on Embedded Computing Systems*, vol. 6, no. 4, pp. 32–es, 2007.
- [19] F. A. Kraemer, D. Ammar, A. E. Braten, N. Tamkittikhun, and D. Palma, "Solar energy prediction for constrained IoT nodes based on public weather forecasts," in *the Seventh International Conference*. New York, New York, USA: ACM Press, 2017, pp. 1–8.