
Dynamic Observation Policies in Observation Cost-Sensitive Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Reinforcement learning (RL) has been shown to learn sophisticated control policies
2 for complex tasks including games, robotics, heating and cooling systems and
3 text generation. The action-perception cycle in RL, however, generally assumes
4 that a measurement of the state of the environment is available at each time step
5 without a cost. In applications such as materials design, deep-sea and planetary
6 robot exploration and medicine, however, there can be a high cost associated with
7 measuring, or even approximating, the state of the environment. In this paper, we
8 survey the recently growing literature that adopts the perspective that an RL agent
9 might not need, or even want, a costly measurement at each time step. Within
10 this context, we propose the Deep Dynamic Multi-Step Observationless Agent
11 (DMSOA), contrast it with the literature and empirically evaluate it on OpenAI
12 gym and Atari Pong environments. Our results, show that DMSOA learns a better
13 policy with fewer decision steps and measurements than the considered alternative
14 from the literature.

15 1 Introduction

16 In many applications of reinforcement learning (RL), such as materials design, computational
17 chemistry, deep-sea and planetary robot exploration and medicine (13; 15; 2), there is a high cost
18 associated with measuring, or even approximating, the state of the environment. Thus, the RL system
19 as a whole faces observation costs in the environment, along with processing and decision making
20 costs in the agent. On both sides, the costs result from a cacophony of factors including the use
21 of energy, systems and human resources. In this work, we propose the Deep Dynamic Multi-Step
22 Observationless Agent (DMSOA), the first RL agent in its class to reduce both measurement and
23 decision making costs.

24 Since standard RL agents require a large number of *state-action-reward-next state* interactions with
25 the environment during policy learning and application, the measurement and decision making costs
26 can be very high. Traditionally, these underlying costs are hidden from the agent. Indeed, little
27 consideration has been given to the idea that the agent might not need, or even want, a potentially
28 costly observation at each time step. In the real-world, however, agents (animal or artificial) are
29 limited by their resources. To save time and energy, decision making associated with common or
30 predictable tasks is believed to be conducted re-actively or based on fast, low-resource systems.
31 Only with deliberate cognitive intervention are the slower, resource-intensive planning systems used
32 (22; 8).

33 Recently, there have been a number of interesting conference papers, workshop discussions and
34 theses discussing how to address this problem in RL (4; 15; 5; 11; 6; 9; 20). The general approach is
35 to augment RL by *a*) assign an intrinsic cost to measure the state of the environment, and *b*) provide

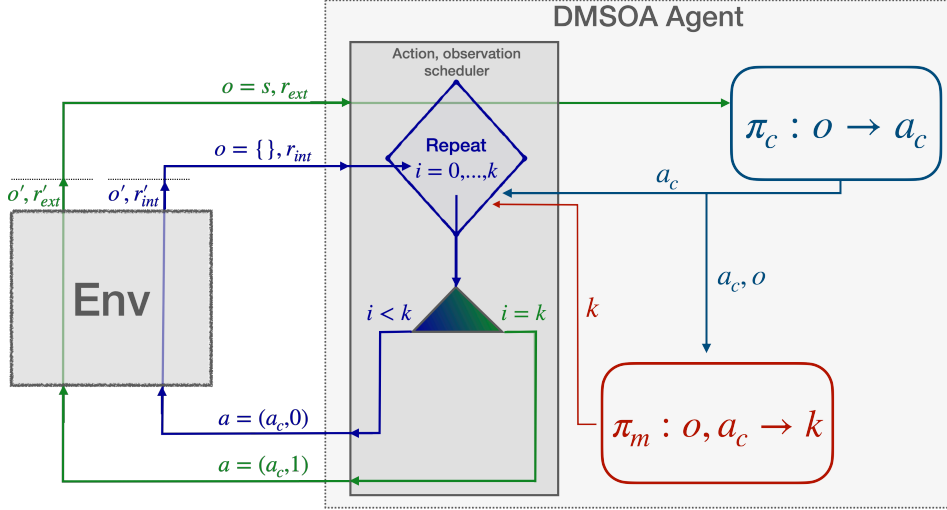


Figure 1: Illustration of the DMSOA Framework.

36 the agent with the flexibility to decide if the next state should be measured. Together, these provide a
 37 mechanism and a learning signal to encourage the agent to reduce its intrinsic measurement costs
 38 relative the to the explicit control rewards it receives.

39 When the agent opts not to measure the state of the environment, it must make its next control
 40 decision based on stale information or an estimate. Thus, at the highest level, this constitutes a
 41 partially observable Markov decision process (POMDP). Learning a POMDP, however, is much more
 42 difficult than learning a Markov decision process (MDP) with RL (17). When designed as shown in
 43 Fig. 1, the agent’s experience is composed of fully observable measurements and partially observable
 44 estimates of the state of the environment. Thus, the problem is related to mixed observable Markov
 45 decision processes (MOMDPs), which are an easier sub-class of POMDPs (16). The authors in
 46 (15) denoted this class of RL problem as action-contingent, noiselessly observable Markov decision
 47 processes (AC-NOMDPs). Distinct from an MOMDP, action-contingent in AC-NOMDP relates
 48 the fact that the agent explicitly chooses between measuring and not measuring the environment.
 49 “Noiselessly observable” relates to the fact that when the agent decides to measure at a cost, the
 50 state is fully observable. Although previous works have used different terms to refer to this class of
 51 problem, we believe that AC-NOMDP is the most descriptive of the underlying dynamics and use it
 52 throughout this paper.

53 Recently, (15) provided a theoretical analysis of the advantages of AC-NOMDP of over a general
 54 POMDP formulation and found a significant improvement in efficiency for RL with explicit obser-
 55 vation costs and actions. All other previous works have carried out limited empirical evaluations in
 56 which the proposed algorithm is compared to a baseline MDP (11; 5; 9; 20). Moreover, the previous
 57 analyses primarily relied on just a few, or even single, experiments on OpenAI gym classic control
 58 environments and grid-worlds (11; 5; 20).

59 In this work, we compare DMSOA to the one-step memory-based observationless agent (OSMBOA)
 60 recently proposed in (5). Our contribution serves to expand the state-of-the-art in AC-NOMDP
 61 algorithms and the understanding of how different classes of algorithms impact the observation
 62 behaviour. OSMBOA was selected for its demonstrated effectiveness and easy of use. At each time
 63 step, OSMBOA selects a control action and makes a decision about measuring the next state of the
 64 environment. If no measurement is made, the agent’s next control action is selected based on its
 65 fixed-size internal memory of the last measured state(s). In contrast, DMSOA selects a control action
 66 and the number of times to apply the action before measuring the next state. Therefore, DMSOA
 67 learns to reduce both observation costs and decision making costs by dynamically applying its control
 68 action multiple times.

69 To facilitate fair comparison, we implement both agents as double DQN (25) with prioritized
 70 experience replay (19). We evaluate the agents in terms of the accumulated extrinsic control reward
 71 and by the reduction in the number of observations and decision steps made on Atari Pong and

72 OpenAI gym (7). The results show that the proposed method learns a better control policy, requires
73 fewer measurements of the environment and decision steps. Moreover, DMSOA has less variance
74 across independent training runs.

75 The remainder of the paper is organized as follows. In the next section, Section 2, we outline the
76 related work. Section 3 formalizes the AC-NOMDP problem and Section 4 presents the proposed
77 algorithms. Section 5 provides the experimental setup. The results are shown in Section 6 and
78 discussed in Section 7. Our concluding remarks are in Section 8.

79 **2 Related Work**

80 This work fits into a small but growing sub-area of RL in which observations are optional at each time
81 step and have an explicit cost to the agent when they are made. In the subsection immediately below,
82 we provide an overview of methods recently applied to AC-NOMDP. Following that, we discuss the
83 literature directly related to the proposed DMSOA algorithm.

84 **2.1 AC-NOMDP Methods**

85 In the existing work, the authors in (4; 15; 20) proposed tabular Q-learning-based algorithms for
86 AC-NOMDPs and (11; 15; 5; 9) proposed deep RL based methods. In (11), the authors modify TRPO,
87 and in (15; 9), actor-critic frameworks with a recurrent neural networks are used. In (5), the authors
88 provide a wrapper class that modifies the underlying environment by expanding the observation and
89 action spaces to facilitate any off-the-shelf deep RL algorithm to work in the AC-NOMDP setting.

90 Through our analysis of the literature, we have identified 4 key questions addressed when developing
91 for AC-NOMDPs. These are: 1) the mechanism by which the agent expresses its desire not to
92 measure, 2) how the observation is supplemented when no measurement is made, 3) how the agent is
93 encouraged to reduce its reliance of costly measurements, and 4) how the agent is constructed.

94 The most common way to handle question 1) is by expanding the action space. In the case of
95 discrete actions, (4; 15; 5; 9; 20) expanded the action space to action tuples: $\langle \text{control actions} \rangle \times$
96 $\langle \text{measure, don't measure} \rangle$. Alternatively, in (11), the agent specifies the control action plus a sample
97 purity value $q \in \mathbb{R}^1$, where a larger q triggers a less accurate measurement with lower associated
98 costs.

99 With respect to question 2), when the agent does not request a fresh measurement in (11; 4; 15; 9),
100 the environment sends a Null state observation or an observation composed of zeros. In (4), the
101 agent uses an internal statistical model to estimate the next state and in (15; 9) the agents utilize
102 a deep recurrent networks for estimating belief states and encoded states, respectively. In (5), the
103 agent utilizes a fixed-size memory of recent measurements when no measurement is made. To reduce
104 partial observability, each observation is augmented with a flag indicating whether or not it is the
105 result of a fresh measurement of the environment. Since the agent in (11) adjusts the noise level
106 rather than turning on and off measurements, it makes its next action selection purely based on the
107 noisy measurement returned.

108 Question 3) relates to the rewards structure. This is generally divided into intrinsic rewards (or costs),
109 which are used to encourage the agent to reduce its reliance on costly measurements and extrinsic
110 rewards that push the agent to achieve the control objective. At each time step in (11; 4; 15; 9), an
111 intrinsic cost is subtracted from the extrinsic reward if the agent measures the state. Alternatively, in
112 (21) a positive intrinsic reward is added when the agent foregoes a measurement. A critical point that
113 remains unclear in the literature is how to acquire the extrinsic reward when no measurement is made.
114 In most cases, this is simply assumed to be available. We argue that if no measurement is made, the
115 extrinsic reward cannot be known. As a result, in this work only the intrinsic portion of the reward is
116 provided when no measurement is made.

117 The final question relates to the architecture of the agent. In (4; 15; 5), a single agent policy select
118 both the control action and the measurement behaviour. Alternatively, in (11; 9; 20), separate policies
119 are learned to determine the control actions and measurement actions. In addition, (4; 15; 9) learn
120 models for estimating the next state.

121 **2.2 Works Related to DMSOA**

122 The proposal of (20) is most algorithmically related to the DMSOA. In it, the author demonstrated
 123 the potential of dynamic action repetition for RL with observation costs using tabular q-learning. The
 124 agent learns to forego a sequence of one or more measurements in predictable regions of the state
 125 space by repeatedly applying the same action. Their proposed method is found to require fewer
 126 measurement steps to reach the goal than the MDP baseline. However, it is only suitable for discrete
 127 state and actions spaces, and was only evaluated on grid-world problems. In this work, we show
 128 how action repetition and measurement skipping can be implemented in deep RL for continuous and
 129 image-based observation spaces.

130 DMSOA is a method that aims to improve the efficiency of RL. To this end, it is weakly related to
 131 other techniques to improve the sampling efficiency (19; 10; 18). The classic sample efficiency work,
 132 however, aims to reduce the overall number of training steps needed to learn a suitable policy, rather
 133 than reducing the measurement or decision steps made by the agent.

134 DMSOA utilizes concepts from the RL literature on dynamic frame skipping to repeatedly apply
 135 the selected control action (12). In DMSOA, however, the agent’s measurement skipping policy is
 136 shaped by the intrinsic reward. Moreover, unlike frame skipping applications, which are concerned
 137 with processing speed not measurement costs, DMSOA does not have access to privileged extrinsic
 138 control rewards from intermediate steps. This makes the problem more challenging.

139 In addition, DMSOA has a connection to the options framework (24), and particularly dynamic
 140 options (1). Similar to the options framework, at each decision point the DMSOA agent chooses to
 141 apply a sequence of actions that will transition the agent through multiple states. In DMSOA, the
 142 agent’s policy selects a single control action and the number of times to apply the action in order to
 143 reduce its measurement costs while still achieving the control objective. Through the incorporation of
 144 measurement costs, the agent is able to learn how many times the action should be applied in order to
 145 arrive at the next meaningful state. For DMSOA, a meaningful state is one for which the information
 146 provided by it is greater than the cost to measure it.

147 **3 Problem Setup**

148 An AC-NOMDP is defined by $\langle S, A, \mathcal{O}, \mathcal{P}, \mathcal{R}_{ext}, \mathcal{R}_{int}, \mathcal{P}_{s_0}, \gamma \rangle$ where S is the state space, $A =$
 149 $\langle A_c \times A_m \rangle$ is the set of action tuples composed of control actions a_c and binary measurement actions
 150 $a_m \in \{0, 1\}$ that specify if an observation of the next state is requested. The observation space \mathcal{O}
 151 is related to S by the observation emission function $p(o|s', a)$ (more on this below). $\mathcal{P} : S \times A \times S \rightarrow \mathbb{R}$
 152 denotes the transition probabilities, $\mathcal{R}_{ext} : S \times A_c \rightarrow \mathbb{R}$ denotes the extrinsic reward function,
 153 $\mathcal{R}_{int} : S \times A_m \rightarrow \mathbb{R}$ denotes the intrinsic reward function that encourages the agent to reduce the
 154 number of measurements it makes. The r_{int} value is typically set to slightly outweigh the r_{ext} value
 155 to achieve the balance between the need for information to solve the control problem and the cost of
 156 information. If r_{int} is very large or very small relative to r_{ext} , the agent may never measure at the
 157 cost of solving the control objective or always measure and fail to reduce the observation costs.

158 The function $\mathcal{P}_{s_0} : S \rightarrow \mathbb{R}$ denotes the probability distribution over the initial state and $\gamma \in (0, 1]$
 159 is the discount factor. The observation emission function $p(o|s', a)$ specifies the probability of
 160 observing $o \in \mathcal{O}$ given the action a in state s . Unlike the more general POMDP, the observation
 161 space in an AC-NOMDP is limited to $\mathcal{O} = S \cup \{empty\}$, where *empty* is the missing measurement
 162 of the environment. In this setup, the potential probabilities of $p(o|s', a_m = 1) \in \{0, 1\}$, with
 163 $p(o|s', a_m = 1) = 1$ if and only if $o = s'$ and $p(o|s', a_m = 1) = 0$ for all $o \neq s'$. In contrast,
 164 $p(o|s', a_m = 0) = 1$ if and only if $o = stale$.

165 The agent learns a policy $\pi(o) : \mathcal{O} \rightarrow A$ that maps observations to action tuples. The initial
 166 observation $o_0 = s_0$ contains a fresh measurement of the environment. The control action selected
 167 by the agent is applied in the environment and the underlying state transitions according to \mathcal{P} . At
 168 each time step, the reward, r_t is the intrinsic reward, $r_t = r_{(int,t)}$, if $a_m = 0$, otherwise the extrinsic
 169 reward, $r_t = r_{(ext,t)}$, is given in response to the state s_t and control action a_c selected by the agent.
 170 When the measurement action, $a_m = 1$, is selected, the agent receives a fresh measurement of the
 171 underlying state $o_{t+1} = s_{t+1}$ of the environment. Alternatively, when $a_m = 0$, the agent does not
 172 obtain a fresh measurement and the next action must be selected based on the agent’s internal mechanism,
 173 such as an internal memory or model.

174 The OSMBOA agent selects one control and one measurement action, $\langle a_{c,t}, a_{m,t} \rangle$ per time step,
 175 whereas the DMSOA agent moves from decision point to decision point with a frequency less than
 176 or equal to the environment’s clock. At each decision point, the DMSOA agent selects a control
 177 action and the number of times to apply it, k . The state is only measured on the k^{th} application (e.g.
 178 $\langle a_c, a_m = 0 \rangle_t, \dots, \langle a_c, a_m = 0 \rangle_{t+(k-1)}, \langle a_c, a_m = 1 \rangle_{t+k}$).

179 The agent’s objective is to learn a policy π that maximizes the discounted expected costed return
 180 which incorporates both the intrinsic and extrinsic rewards:

$$J(\pi) = \mathbb{E}_{a_t \sim \pi, s_t \sim P} \left[\sum_t \gamma^t r(s_t, a_t) \right], \quad (1)$$

181 where $\gamma < 1$ is the discount factor. In this work, we focus on deep Q-learning based solutions (14)
 182 combined with standard improvements for better convergence and stability (23; 25; 19). Although
 183 this work examines problems with discrete action spaces, the proposed algorithms can be modified
 184 for continuous action spaces.

185 4 Deep Dynamic Multi-Step Observationless Agent

186 The Deep Dynamic Multi-Step Observationless Q-learning Agent (DMSOA) for noiselessly observ-
 187 able RL environments with explicit observation costs is presented in Figure 1. The framework has
 188 three key components: the control policy $\pi_c : o \rightarrow a_a$ that maps the observation to a control action,
 189 the measurement skipping policy $\pi_m : o, a_c \rightarrow k$ that maps the observation and selected control ac-
 190 tion to $k \in \{1, \dots, K\}$ the number of steps to apply a_c to the environment, and the action-observation
 191 scheduler. The action-observation scheduler applies the action pair $(a_c, 0)$ $k - 1$ times and collects
 192 the intrinsic rewards r_{int} from the environment. On the k^{th} iteration, it applies the action pair $(a_c, 1)$,
 193 records the extrinsic reward r_{ext} , and passes the new observation to π_c and π_m . The extrinsic reward
 194 is equal to the control policy reward for applying a_c and arriving in the measured state after step $i = k$.
 195 The intrinsic reward is $r_{int} \in \{0, c\}$, where c is a bonus (ie. “cost saving”) given to the agent when it
 196 chooses not to measure. To ensure the agent is motivated to omit measurements whenever possible,
 197 we set $c \geq r_{ext}^{max}$. The optimal setting of c will depend on the application and the requirements of the
 198 domain.

199 In this work, the policies are implemented as deep Q networks (DQN), however, other forms of policy
 200 learning could be utilized. The agent’s objective is to maximize the costed rewards $\sum_{t=0}^{\infty} \gamma^t r_t$. To
 201 achieve this we learn parameterized value functions $Q_c(o; \theta)$ and $Q_m(o, a; \zeta)$ as feed-forward deep
 202 neural networks. As described above, for an m -dimensional observation space and an n -dimensional
 203 action space, Q_c is a mapping from an m -dimensional observation to an n -dimensional vector of
 204 action values. The function Q_m is a mapping from an $m + 1$ -dimensional observation-action to a
 205 K -dimensional vector of measurement values. In the case of image data, each channel is augmented
 206 with the action details. The argmax of each output indicates the action to apply and the number of
 207 times to apply it.

208 During training, the experience tuples $(o_t, a_{(c,t)}, a_{(m,t)}, r_t, o_{t+1})$ are stored in a prioritized experience
 209 replay buffer. To improve stability, target networks θ^- and ζ^- for Q_c and Q_m are copied from θ and
 210 ζ every τ steps. The target for the control network is:

$$Y_i^{Q_c} \equiv r_t + \gamma Q_c(o_{t+1}, \text{argmax}_a Q_c(o_{t+1}, a; \theta_t); \theta_t^-). \quad (2)$$

211 For the same update step, the target for the measurement network is:

$$Y_i^{Q_m} \equiv r_t + \gamma Q_m((o_{t+1}, a_{(c,t+1)}), \text{argmax}_a Q_m((o_{t+1}, a_{(c,t+1)}), a; \zeta_t); \zeta_t^-). \quad (3)$$

212 The corresponding losses are:

$$\mathcal{L}_i^{Q_c}(\theta_i) = \mathbb{E}_{(o_t, a_{(c,t)}) \sim \mathcal{D}} [(Y_i^{Q_c} - Q_c(o_t, a_{(c,t)}; \theta_i))^2], \quad (4)$$

213 and

$$\mathcal{L}_i^{Q_m}(\zeta_i) = \mathbb{E}_{(o_t, a_{(c,t)}, a_{(m,t)}) \sim \mathcal{D}} [(Y_i^{Q_m} - Q_m((o_t, a_{(c,t)}), a_{(m,t)}; \zeta_i))^2] \quad (5)$$

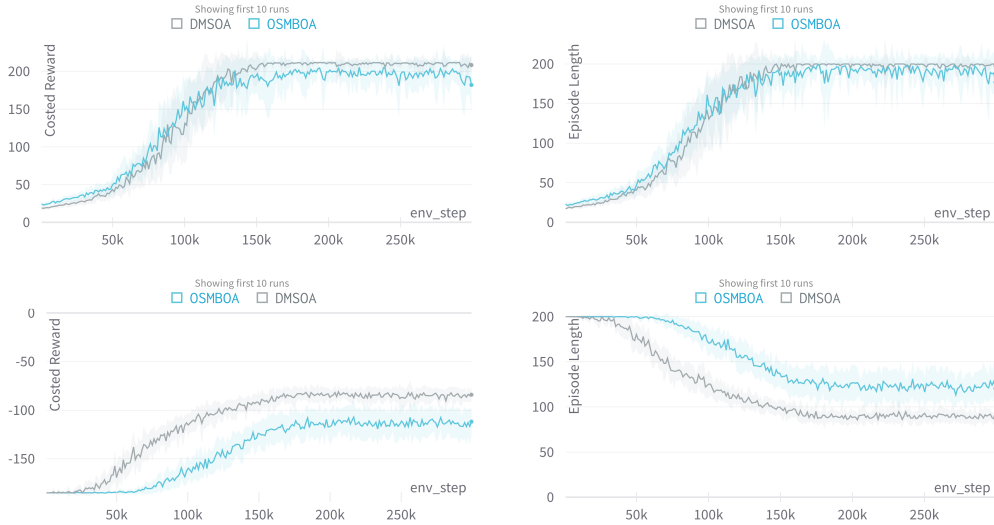


Figure 2: Mean and standard deviation of performance on Cartpole (upper) and Acrobot (lower). Left: costed reward and right: episode length. In both cases, DMSOA has a higher mean costed reward, and is superior in terms of the control objective (longer episodes on Cartpole and shorter episodes on Acrobot.)

214 5 Experimental Setup

215 In this section, we compare the performance of DMSOA to OSMBOA. In order to highlight the
 216 differences in the measurement behaviour of each method, we implement both with double DQN
 217 and a prioritized replay buffer. The hyper-parameters were selected via grid search with 3 random
 218 trials. For the evaluation, we report the mean and standard deviation of the reward during training
 219 and the observation behaviour of the best policy. Each agent is reinitialized with 20 difference seeds
 220 and trained on the OpenAI gym environments Cartpole, Acrobot, Lunar Lander and Atari Pong. The
 221 experiments were run on CentOS with Intel Xeon Gold 6130 CPU and 192 GB memory. In addition,
 222 a NVIDIA V100 GPU was used in the training of the Atari agent.

223 6 Results

224 Figure 2 shows the mean and standard deviation for each agent on the Cartpole and Acrobot environ-
 225 ments. The aim in the Cartpole environment is for the agent to operate a cart such that a vertical pole
 226 remains balanced for as long as possible. The extrinsic reward is set to 1 and the intrinsic reward is
 227 set to 1.1. We truncate each episode at a maximum of 200 time steps. In the Acrobot environment,
 228 the objective is to apply torque to flip an arm consisting of two actuated links connected linearly
 229 above a target height in as few steps as possible. The agent receives an extrinsic reward of -1 or an
 230 intrinsic reward of -0.85 at each time step. Each episode ends after 200 steps or when the arm is
 231 successfully flipped over the line.

232 DMSOA learns a policy for both environments that pro-
 233 duces a higher costed reward than OSMBOA. This indicates that DMSOA requires fewer measurements whilst
 234 carrying out the control policy. In addition, the standard deviation is lower indicating more stability across independent
 235 training runs. The episode length plots on the right
 236 show that DMSOA learned policies to keep the Cartpole
 237 upright longer and flip the Acrobot over the goal faster.
 238 Thus, DMSOA outperforms OSMBOA on all accounts.

Table 1: Ratio of steps with measurements to steps without measurements of the converged policy during training.

Env.	DMSOA	OSMBOA
Cartpole	1:1.27	1:0.37
Acrobot	1:0.45	1:1.03
Lunar Lander	1:1.56	1:0.33

241 The objective of the Lunar Lander environment is to fire
 242 the lander’s rockets such that it lands squarely in the target area. In general, the results show shows

243 that DMSOA has significantly more successful landings than OSMBOA. Similarly to Cartpole
244 and Acrobot, our analysis finds that DMSOA is more efficient in the number of decision steps
245 and observations. More details on the environment along with performance plots are available in
246 Appendix A Figure 5

247 Table 1 shows the ratio of the number of steps made without measuring for each measurement made.
248 On Cartpole and Lunar Lander, DMSOA makes more than one step without measuring for each
249 measurement step, whereas on Acrobot it makes an average of 0.5 non-measuring steps for each
250 measuring step. This suggests that the dynamics of Acrobot are less predictable, causing DMSOA to
251 measure more frequently on Acrobot. Interestingly, Acrobot is the only environment where OSMBOA
252 does better than a 1:1 ratio.

253 6.1 Examination of Measurement Policies

254 Figure 3 shows the measurement behaviour of the best OSMBOA (left) and DMSOA (right) policies
255 for Cartpole (top), Acrobot (middle) and Lunar Lander (lower) environments. Each row specifies a
256 1-episode roll-out of the best policy. Each column in the OSMBOA plots is the environment time
257 step during the episode. For OSMBOA, the number of decision steps is equivalent to the number of
258 steps in the environment. In contrast, each column in the DMSOA plots corresponds to a decision by
259 the agent, with one or more environment time steps associated with it. In addition to highlighting
260 the measurement efficiency, this also shows the decision efficiency. On Cartpole, DMSOA makes
261 approximately 70 action selections (decisions) per episode of 200 environment steps (the mean steps
262 per episode are shown in Figure 2).

263 For OSMBOA, an orange cell indicates that a fresh measurement of the environment and blue
264 specifies that no measurement was requested at corresponding time step. In the case of DMSOA, the
265 colour indicates the number of consecutive steps that were taken without a fresh measurement. Blue
266 indicates that a measurement is made after the control action is applied once, yellow indicates that a
267 measurement is made after the control action is applied twice and red indicates that a measurement is
268 made after the control action is applied three times.

269 The distinct pattern in each plot suggests the different capabilities of each class of AC-NOMDP
270 agent, along with the fact that each environment is unique in terms of its dynamics and complexity.
271 The consistent measurement patterns for OSMBOA and DMSOA on Cartpole suggest that the
272 environment has very regular dynamics. OSMBOA switches between selecting the next action from a
273 freshly measured observation and selecting it from a stale observation. Alternatively, DMSOA learns
274 to apply an action 3 times before measuring. This clearly demonstrates the potential of DMSOA to
275 take more environment steps without measuring than OSMBOA.

276 Acrobot and Lunar Lander show much more complex measurement behaviour. For both OSMBOA
277 and DMSOA on Acrobot during approximately the first 3/4s of the each episode they display a pattern
278 of frequently measuring followed by briefly not measuring. Our analysis finds that both OSMBOA
279 and DMSOA skip measurements while the arcobot is in the lower left region of the observation space.
280 This is roughly where the momentum of the Acrobot shifts from heading way from the goal, back
281 towards the goal. In this area, it is deemed safe to apply torque back towards the goal without
282 observing. In the last quarter of each episode, both agents take more steps without measuring. It
283 is noteworthy that in most episodes OSMBOA takes significantly more steps without observing
284 than DMSOA. However, this has a negative impact on the total number of action decisions made by
285 OSMBOA on route to achieving to goal. This particularly visible in episodes 1 and 4. DMSOA does
286 not to suffer from similar behaviour.

287 On the Lunar Lander environment both methods take few or no measurements near the end of the
288 episode when the agent is close to landing. In addition, DMSOA repeatedly takes 2-3 steps before
289 measuring at the beginning of each episode, whereas OSMBOA repeatedly measures early in each
290 episode. Each method measures frequently during the middle of the episode as agent attempts to
291 direct the lander safely towards the landing area. OSMBOA generally alternates between measuring
292 and not measuring at each time step, whereas DMSOA typically takes many measurement steps
293 followed by 1 to 2 steps without measuring before returning to measuring again. Similar to Acrobot,
294 once OSMBOA estimates that it is on target to reach the goal it commits to never measuring again.
295 When this estimate is erroneous, the leads to much longer episodes than necessary and the risk of
296 crashing the ship.

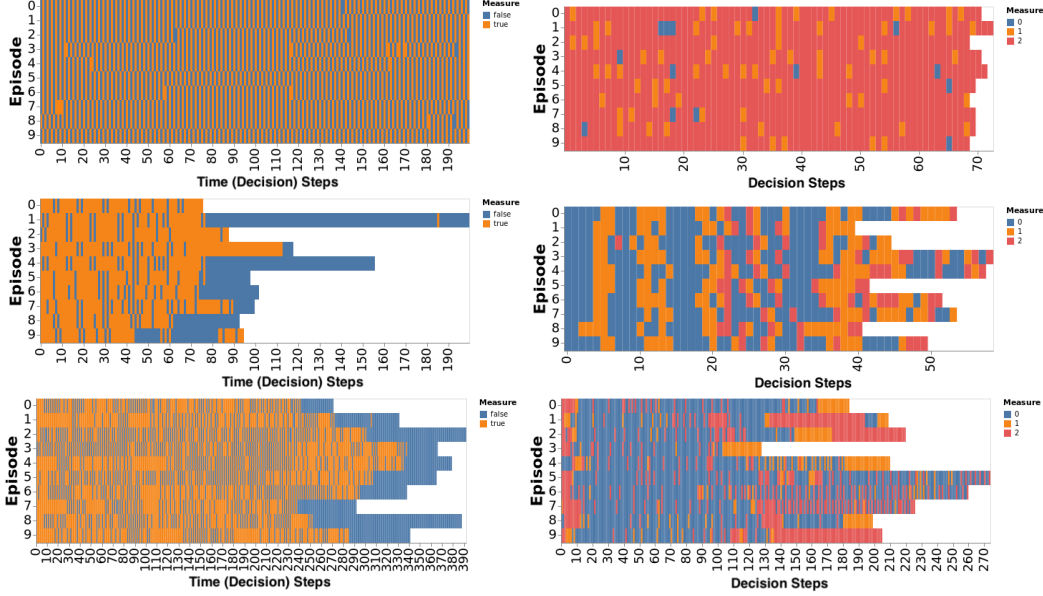


Figure 3: Measurement behaviour of the best OSMBOA (left) and DMSOA (right) policies. Top: Cartpole, centre: Acrobot, lower: Lunar Lander. For OSMBOA, blue indicates that no measurement was made and orange indicates that a measurement was made. In the case of DMSOA, blue indicates that a measurement is made after one step, orange indicates that a measurement is made after two steps, and red indicates that a measurement is made after three steps. This figure reveals the very distinct measurement behaviour between the two classes of AC-NOMDP agents.

297 6.2 Image-Based RL Results

298 The objective in the pong Atari game is to bounce the ball off of your paddle and past the opponents
 299 paddle into its goal (3). The action space is 6-dimensional including do nothing, fire, move right,
 300 move left, fire right and fire left. The observation space is a (210, 160, 3) image. In the case of
 301 OSMBOA, a 210 by 1 vector of ones or zeros is added to each channel to indicate if the observation is
 302 fresh or stale. The agent gets an extrinsic reward of 1 for winning a match and 0 for each intermediate
 303 step. Each episode is composed of 21 matches and the intrinsic reward is 0.001.

304 The results in Figure 4 show learns a policy
 305 to achieve a higher costed reward. In general,
 306 DMSOA learns to be a better Pong player and
 307 requires fewer measurements. Additional results
 308 showing DMSOA’s advantage in terms of wins
 309 and intrinsic rewards can be found in Appendix
 310 A Figure 6.

311 From our analysis for the measurement policies
 312 of each agent, we found that both learn to mea-
 313 sure less frequently when the ball is travelling
 314 away from their paddle. Alternatively, if the
 315 ball is near their paddle or the opponents pad-
 316 dle, each agent measures more frequently. Inline
 317 with the observations on Acrobot and Lunar Lan-
 318 der, when OSMBOA reaches a state from which
 319 it expects to win the match, it switch to not mea-
 320 suring for the remainder of the match. If the prediction is correct, it can achieve a greater reduction
 321 in measurements than DMSOA. If it is wrong, however, OSMBOA general loses the match. An
 322 erroneous prediction of this nature is particularly risky in a complex and dynamic environment.

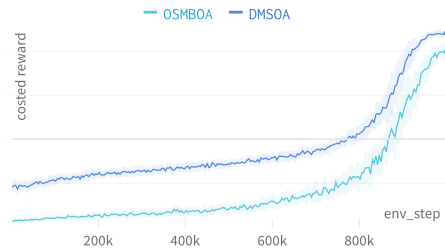


Figure 4: Mean and standard deviation of the costed reward on the Atari Pong environment. DMSOA learns a policy that for a better costed reward.

323 7 Discussion

324 The results indicate that DMSOA has a clear advantage over OSMBOA in terms of its convergence
325 rate and the reduction in measurements and decision steps. We believe that the control action
326 repetition capabilities of DMSOA improve its exploration of the environment and its understanding
327 of the implications (positive and negative) of taking multiple steps without measuring. This helps it
328 to quickly converge to good control and measurement policies. In addition, the fact that DMSOA’s
329 multi-step action sequences always ends with a measurement of the final state provides it with a good
330 grounding from which to select the next control action. On the other hand, because the extrinsic reward
331 for intermediate steps is not available, there is the potential for more noise in the reward signal for
332 longer DMSOA action repetition trajectories. Due to the fact that OSMBOA is limited to one-step
333 action, noise in the reward is less of a concern. Although, DMSOA appears to handle the noisy
334 reward signal, future work should examine this in more detail.

335 For unshaped (or uniform) reward environments, such as Cartpole, Acrobot and Pong, setting the
336 intrinsic reward is simple and the agent is insensitive to the value so long as it is slightly larger than
337 the extrinsic reward. Alternative, the intrinsic reward requires fine tuning on environments with
338 complex reward shaping such as Lunar Lander. As heuristic, we suggest starting the fine tuning from
339 the mean of the extrinsic reward collected over multiple random walks in the environment.

340 In multiple environments, we found that OSMBOA commits to not measuring towards the end of
341 each episode. This is surprising since if OSMBOA takes more than one step without measuring, it
342 enters a partially observable state. This is akin to playing the game with its eyes closed. The agent is,
343 thus, unaware if any unexpected event occurs. On Acrobot and pong, this resulted in it not achieving
344 the goal, or it taking much longer than otherwise necessary. An example of this is seen in episodes 1
345 and 4 of the Acrobot plot in Figure 3.

346 8 Conclusion

347 In this work, we consider the problem of RL for environments where agent’s decision making and
348 measuring of the state of the environment have explicit costs, namely AC-NOMDPs. We provide
349 the first survey of methods recently proposed for AC-NOMDPs. Building on the existing work, we
350 propose DMSOA, an RL algorithm learns a control and a measurement policy to reduce measurement
351 and decision steps. Our empirical results confirm the previously published results for OSMBOA on
352 Cartpole, Acrobot and Lunar Lander, and show that OSMBOA is also capable on the more complex,
353 image-based Atari Pong environment. However, we find that our proposed method DMSOA learns a
354 *better control policy* than OSMBOA, and *requires fewer costly measurement and decision steps*.

355 This demonstrates the great potential to reduce measurement and decision costs associated with RL
356 by allowing the agent to take control of its action and observation behaviour. We expect this to be a
357 necessary capability of RL agents applied in many real-world applications. The next steps that we
358 envision are developing more sophisticated loss functions for DMSOA, incorporating recurrency into
359 the network to deal with time, expanding the analysis to additional methods and more realistic setting
360 such as self-driving chemistry where observation can be costly and potentially destructive (2).

361 References

- 362 [1] Barto, A.G., Mahadevan, S.: Recent advances in hierarchical reinforcement learning. *Discrete*
363 *event dynamic systems* **13**(1-2), 41–77 (2003)
- 364 [2] Beeler, C., Subramanian, S.G., Sprague, K., Chatti, N., Bellinger, C., Shahan, M., Paquin, N.,
365 Baula, M., Dawit, A., Yang, Z., et al.: Chemgymrl: An interactive framework for reinforcement
366 learning for digital chemistry. *arXiv preprint arXiv:2305.14177* (2023)
- 367 [3] Bellemare, M.G., Naddaf, Y., Veness, J., Bowling, M.: The arcade learning environment: An
368 evaluation platform for general agents. *Journal of Artificial Intelligence Research* **47**, 253–279
369 (2013)
- 370 [4] Bellinger, C., Coles, R., Crowley, M., Tamblyn, I.: Active Measure Reinforcement Learning
371 for Observation Cost Minimization. In: *Proceedings of the Canadian Conference on*

- 372 Artificial Intelligence. Canadian Artificial Intelligence Association (CAIAC) (jun 8 2021),
373 <https://caiac.pubpub.org/pub/3hn8s5v9>
- 374 [5] Bellinger, C., Drozdyuk, A., Crowley, M., Tamblin, I.: Balancing Information with Observation
375 Costs in Deep Reinforcement Learning. In: Proceedings of the Canadian Conference on
376 Artificial Intelligence. Canadian Artificial Intelligence Association (CAIAC) (may 27 2022),
377 <https://caiac.pubpub.org/pub/0jmy7gpd>
- 378 [6] Bellinger, C., Drozdyuk, A., Crowley, M., Tamblin, I.: Scientific discovery and the cost of
379 measurement – balancing information and cost in reinforcement learning. In: ICML 2nd Annual
380 AAAI Workshop on AI to Accelerate Science and Engineering (AI2ASE) (Feb 13 2023)
- 381 [7] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., Zaremba, W.:
382 Openai gym. arXiv preprint arXiv:1606.01540 (2016)
- 383 [8] Daniel, K.: Thinking fast and slow. United States of America (2011)
- 384 [9] Fu, Y.: The Cost of OPS in Reinforcement Learning. Master’s thesis, University of California,
385 Berkeley (2021)
- 386 [10] Gal, Y., McAllister, R., Rasmussen, C.E.: Improving pilco with bayesian neural network
387 dynamics models. In: Data-Efficient Machine Learning workshop, ICML. vol. 4, p. 34 (2016)
- 388 [11] Koseoglu, M., Özcelikkale, A.: How to miss data?: Reinforcement learning for environments
389 with high observation cost. In: 2020 International Conference on Machine Learning (ICML)
390 Workshop, Wien, Österrike, 12-18 juli (2020)
- 391 [12] Lakshminarayanan, A., Sharma, S., Ravindran, B.: Dynamic action repetition for deep rein-
392 forcement learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 31
393 (2017)
- 394 [13] Mills, K., Ronagh, P., Tamblin, I.: Finding the ground state of spin hamiltonians with reinforcement learning. *Nature Machine Intelligence* **2**(9), 509–517 (Sep 2020). <https://doi.org/10.1038/s42256-020-0226-x>, <https://doi.org/10.1038/s42256-020-0226-x>
- 398 [14] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., et al.: Human-level control through deep reinforcement learning. *nature* **518**(7540), 529–533 (2015)
- 401 [15] Nam, H.A., Fleming, S., Brunskill, E.: Reinforcement learning with state observation costs in action-contingent noiselessly observable markov decision processes. *Advances in Neural Information Processing Systems* **34**, 15650–15666 (2021)
- 404 [16] Ong, S.C., Png, S.W., Hsu, D., Lee, W.S.: Planning under uncertainty for robotic tasks with mixed observability. *The International Journal of Robotics Research* **29**(8), 1053–1068 (2010)
- 406 [17] Papadimitriou, C.H., Tsitsiklis, J.N.: The complexity of markov decision processes. *Mathematics of operations research* **12**(3), 441–450 (1987)
- 408 [18] Pathak, D., Agrawal, P., Efron, A.A., Darrell, T.: Curiosity-driven exploration by self-supervised prediction. In: International conference on machine learning. pp. 2778–2787. PMLR (2017)
- 410 [19] Schaul, T., Quan, J., Antonoglou, I., Silver, D.: Prioritized experience replay. arXiv preprint arXiv:1511.05952 (2015)
- 412 [20] Shann, T.Y.A.: Reinforcement learning in the presence of sensing costs. Master’s thesis, University of British Columbia (2022). <https://doi.org/http://dx.doi.org/10.14288/1.0413129>, <https://open.library.ubc.ca/collections/ubctheses/24/items/1.0413129>
- 415 [21] Sharma, S., Srinivas, A., Ravindran, B.: Learning to repeat: Fine grained action repetition for deep reinforcement learning. arXiv preprint arXiv:1702.06054 (2017)
- 417 [22] Simon, H.A.: Bounded rationality. *Utility and probability* pp. 15–18 (1990)

- 418 [23] Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. MIT press (2018)
- 419 [24] Sutton, R.S., Precup, D., Singh, S.: Between mdps and semi-mdps: A framework for temporal
420 abstraction in reinforcement learning. *Artificial intelligence* **112**(1-2), 181–211 (1999)
- 421 [25] Van Hasselt, H., Guez, A., Silver, D.: Deep reinforcement learning with double q-learning. In:
422 Proceedings of the AAAI conference on artificial intelligence. vol. 30 (2016)

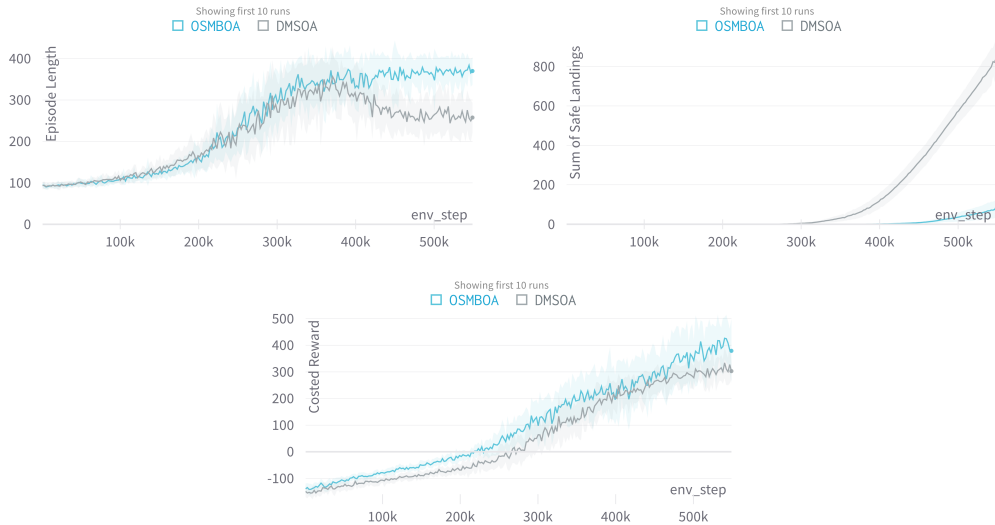


Figure 5: Mean and standard deviation of the performance on the Lunar Lander environment. Top left: episode length, top right: sum of the number of successful landings, lower: costed reward. DMSOA and OSMBOA achieve similar mean costed rewards. DMSOA, however, learns to successfully land the ship more frequently and in fewer steps.

423 A Appendix: Results

424 A.1 Lunar Lander

425 The results for the Lunar Lander environment are presented in Figure 5. The Lunar Lander environ-
 426 nment is a rocket trajectory optimization problem (7). The objective is to fire the lander’s rockets such
 427 that it lands squarely in the target area. The fuel supply is infinite, but the best policy uses it sparingly.
 428 The environment has four discrete actions available: do nothing, fire left orientation engine, fire main
 429 engine, fire right orientation engine. The intrinsic reward is 0.1. The extrinsic reward is -0.3 for
 430 firing the main engine and -0.03 for side engines, the reward is also scaled by the lander’s distance
 431 from the landing pad. Ten points are added to the extrinsic reward for each leg that is in contact with
 432 the ground, and an additional 100 points are added for landing, while 100 points are subtracted for
 433 crashing. The episode ends when the agent lands or crashes, or is truncated after a maximum of 400
 434 time steps.

435 The plot on the top left in Figure 5 shows that mean episode length is longer for OSMBOA than
 436 DMSOA, and the top right plot shows that DMSOA has significantly more successful landings. This
 437 indicates that DMSOA learns a policy that quickly navigates the ship to a safe landing. The lower
 438 plot shows that OSMBOA has a slightly higher costed reward. As suggested by the first two plots,
 439 this is due to the fact that it takes longer to land and not because the policies is superior.

440 A.2 Image-Based RL Results

441 The objective in the pong Atari game is to bounce the ball off of your paddle and past the opponents
 442 paddle into its goal (3). The action space is 6-dimensional including do nothing, fire, move right,
 443 move left, fire right and fire left. The observation space is a (210, 160, 3) image. In the case of
 444 OSMBOA, a 210 by 1 vector of ones or zeros is added to each channel to indicate if the observation is
 445 fresh or stale. The agent gets an extrinsic reward of 1 for winning a match and 0 for each intermediate
 446 step. Each episode is composed of 21 matches and the intrinsic reward is 0.001.

447 The results in Figure 6 show that DMSOA wins significantly more matches than OSMBOA (top left),
 448 achieves a higher costed reward (top right) and more intrinsic reward (lower). Thus, DMSOA learns
 449 to be a better Pong player and requires fewer measurements. Due the longer episodes and training
 450 times, a measurement behaviour plot similar to Figure 3 is not feasible within the confines of this
 451 paper.

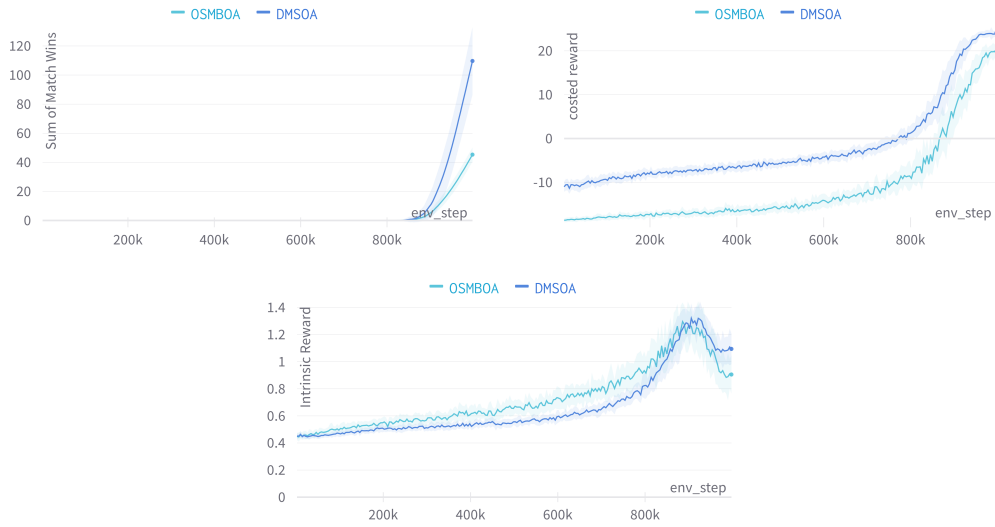


Figure 6: Mean and standard deviation of the performance on the Atari Pong environment. Top left: sum of the number of wins, top right: costed reward, lower: intrinsic reward. DMSOA learns a policy that wins more games with a better costed reward and fewer measurements.

452 From our analysis for the measurement policies of each agent, we found that both learn to measure
 453 less frequently when the ball is travelling away from their paddle. Alternatively, if the ball is near their
 454 paddle or the opponents paddle, each agent measures more frequently. Inline with the observations on
 455 Acrobot and Lunar Lander, when OSMBOA reaches a state from which it expects to win the match,
 456 it switch to not measuring for the remainder of the match. If the prediction is correct, it can achieve a
 457 greater reduction in measurements than DMSOA. If it is wrong, however, OSMBOA general loses
 458 the match. An erroneous prediction of this nature is particularly risky in a complex and dynamic
 459 environment.