

# Learning to Search for Targets

- A Deep Reinforcement Learning Approach for Visual Search in Unknown Environments

---

*Inlärd sökning efter mål*

**Oskar Lundin**

Supervisor : Sourabh Balgi

Examiner : Jose M. Peña

External supervisor : Fredrik Bissmarck

## Upphovsrätt

Detta dokument hålls tillgängligt på Internet - eller dess framtida ersättare - under 25 år från publiceringsdatum under förutsättning att inga extraordinära omständigheter uppstår.

Tillgång till dokumentet innebär tillstånd för var och en att läsa, ladda ner, skriva ut enstaka kopior för enskilt bruk och att använda det oförändrat för ickekommersiell forskning och för undervisning. Överföring av upphovsrätten vid en senare tidpunkt kan inte upphäva detta tillstånd. All annan användning av dokumentet kräver upphovsmannens medgivande. För att garantera äktheten, säkerheten och tillgängligheten finns lösningar av teknisk och administrativ art.

Upphovsmannens ideella rätt innefattar rätt att bli nämnd som upphovsman i den omfattning som god sed kräver vid användning av dokumentet på ovan beskrivna sätt samt skydd mot att dokumentet ändras eller presenteras i sådan form eller i sådant sammanhang som är kränkande för upphovsmannens litterära eller konstnärliga anseende eller egenart.

För ytterligare information om Linköping University Electronic Press se förlagets hemsida <http://www.ep.liu.se/>.

## Copyright

The publishers will keep this document online on the Internet - or its possible replacement - for a period of 25 years starting from the date of publication barring exceptional circumstances.

The online availability of the document implies permanent permission for anyone to read, to download, or to print out single copies for his/hers own use and to use it unchanged for non-commercial research and educational purpose. Subsequent transfers of copyright cannot revoke this permission. All other uses of the document are conditional upon the consent of the copyright owner. The publisher has taken technical and administrative measures to assure authenticity, security and accessibility.

According to intellectual property law the author has the right to be mentioned when his/her work is accessed as described above and to be protected against infringement.

For additional information about the Linköping University Electronic Press and its procedures for publication and for assurance of document integrity, please refer to its www home page: <http://www.ep.liu.se/>.

## Abstract

The abstract resides in file `Abstract.tex`. Here you should write a short summary of your work.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Pellentesque in massa suscipit, congue massa in, pharetra lacus. Donec nec felis tempor, suscipit metus molestie, consectetur orci. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Curabitur fermentum, augue non ullamcorper tempus, ex urna suscipit lorem, eu consectetur ligula orci quis ex. Phasellus imperdiet dolor at luctus tempor. Curabitur nisi enim, porta ut gravida nec, feugiat fermentum purus. Donec hendrerit justo metus. In ultrices malesuada erat id scelerisque. Sed sapien nisi, feugiat in ligula vitae, condimentum accumsan nisi. Nunc sit amet est leo. Quisque hendrerit, libero ut viverra aliquet, neque mi vestibulum mauris, a tincidunt nulla lacus vitae nunc. Cras eros ex, tincidunt ac porta et, vulputate ut lectus. Curabitur ultricies faucibus turpis, ac placerat sem sollicitudin at. Ut libero odio, eleifend in urna non, varius imperdiet diam. Aenean lacinia dapibus mauris. Sed posuere imperdiet ipsum a fermentum.

Nulla lobortis enim ac magna rhoncus, nec condimentum erat aliquam. Nullam laoreet interdum lacus, ac rutrum eros dictum vel. Cras lobortis egestas lectus, id varius turpis rhoncus et. Nam vitae auctor ligula, et fermentum turpis. Morbi neque tellus, dignissim a cursus sed, tempus eu sapien. Morbi volutpat convallis mauris, a euismod dui egestas sit amet. Nullam a volutpat mauris. Fusce sed ipsum lectus. In feugiat, velit eu fermentum efficitur, mi ex eleifend ante, eget scelerisque sem turpis nec augue.

Vestibulum posuere nibh ut iaculis semper. Ut diam justo, interdum quis felis ac, posuere fermentum ex. Fusce tincidunt vel nunc non semper. Sed ultrices suscipit dui, vel lacinia lorem euismod quis. Etiam pellentesque vitae sem eu bibendum. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Pellentesque scelerisque congue ullamcorper. Sed vehicula sodales velit a scelerisque. Pellentesque dignissim lectus ipsum, quis consectetur tellus rhoncus a.

Nunc placerat ut lectus vel ornare. Sed nec dictum enim. Donec imperdiet, ipsum ut facilisis blandit, lacus nisi maximus ex, sed semper nisl metus eget leo. Nunc efficitur risus ac risus placerat, vel ullamcorper felis interdum. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Duis vitae felis vel nibh sodales fringilla. Donec semper eleifend sem quis ornare. Proin et leo ut dolor consectetur vehicula. Lorem ipsum dolor sit amet, consectetur adipiscing elit.

Nunc dignissim interdum orci, sit amet pretium nibh consectetur sagittis. Aenean a eros id risus aliquam placerat nec ut lectus. Curabitur at quam in nisi sodales imperdiet in at erat. Praesent euismod pulvinar imperdiet. Nam auctor mattis nisi in efficitur. Quisque non cursus ipsum, consequat vehicula justo. Fusce varius metus et nulla rutrum scelerisque. Praesent molestie elementum nulla a consequat. In at facilisis nisi, convallis molestie sapien. Cras id ullamcorper purus. Sed at lectus sit amet dolor finibus suscipit vel et purus. Sed odio ipsum, dictum vel justo sit amet, interdum dictum justo. Quisque euismod quam magna, at dignissim eros varius in. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas.

# Acknowledgments

Acknowledgments.tex

# Contents

<b>Abstract</b>	<b>iii</b>
<b>Acknowledgments</b>	<b>iv</b>
<b>Contents</b>	<b>v</b>
<b>List of Figures</b>	<b>vii</b>
<b>List of Tables</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Aim . . . . .	2
1.3 Research questions . . . . .	2
1.4 Delimitations . . . . .	3
<b>2 Theory</b>	<b>4</b>
2.1 Background . . . . .	4
2.1.1 Visual Search and Attention . . . . .	4
2.1.2 Active Vision and Object Search . . . . .	5
2.1.3 Deep Learning . . . . .	5
2.1.3.1 Feedforward Neural Network . . . . .	6
2.1.3.2 Convolutional Neural Network . . . . .	6
2.1.3.3 Recurrent Neural Network . . . . .	6
2.1.4 Reinforcement Learning . . . . .	7
2.1.4.1 Partially Observable Markov Decision Processes . . . . .	7
2.1.4.2 Policies and Value Functions . . . . .	8
2.1.4.3 Design Challenges . . . . .	9
2.1.4.4 Deep Reinforcement Learning . . . . .	9
2.1.4.5 Proximal Policy Optimization . . . . .	10
2.2 Related Work . . . . .	10
2.2.1 Search with Reinforcement Learning . . . . .	10
2.2.1.1 Sequential Visual Attention . . . . .	11
2.2.1.2 Active Object Detection . . . . .	11
2.2.1.3 Visual Navigation . . . . .	12
2.2.2 Memory Architectures . . . . .	13
2.2.3 Generalization and Inductive Bias . . . . .	14
2.2.4 Evaluation of Agents . . . . .	16
<b>3 Method</b>	<b>17</b>
3.1 Problem Statement . . . . .	17
3.2 Environments . . . . .	17
3.2.1 Observations, Actions and Reward . . . . .	18

3.2.2	Gaussian Environment . . . . .	19
3.2.3	Terrain Environment . . . . .	19
3.2.4	Camera Environment . . . . .	20
3.3	Approach . . . . .	21
3.3.1	Architecture . . . . .	21
3.3.2	Training . . . . .	23
3.4	Experiments . . . . .	23
3.4.1	Quality of Search Behavior . . . . .	23
3.4.2	Size of Search Space . . . . .	24
3.4.3	Number of Training Samples . . . . .	25
3.4.4	Ablations . . . . .	25
3.5	Implementation . . . . .	25
<b>4</b>	<b>Results</b>	<b>26</b>
4.1	Quality of Search Behavior . . . . .	26
4.2	Size of Search Space . . . . .	26
4.3	Number of Training Samples . . . . .	26
<b>5</b>	<b>Discussion</b>	<b>30</b>
5.1	Results . . . . .	30
5.2	Method . . . . .	30
5.3	The work in a wider context . . . . .	31
<b>6</b>	<b>Conclusion</b>	<b>32</b>
<b>A</b>	<b>Training</b>	<b>33</b>
	<b>Bibliography</b>	<b>34</b>

# List of Figures

2.1	Partially observable Markov decision process . . . . .	8
3.1	Gaussian environment . . . . .	19
3.2	Terrain environment . . . . .	20
3.3	Camera environment . . . . .	21
3.4	Network architecture . . . . .	22
4.1	Search space size learning curve. . . . .	27
4.2	Generalization results. . . . .	29

# List of Tables

3.1	PPO hyperparameters . . . . .	23
4.1	Quality results. . . . .	28



# Todo list

Illustrate map updates. . . . .	23
---------------------------------	----



# 1 Introduction

In this thesis project, the problem of searching for target objects in unknown but familiar environments is addressed. This chapter presents the motivation behind the project, the research questions that are addressed, and the delimitations.

## 1.1 Motivation

The ability to visually search for things in an environment is fundamental to intelligent behavior. We humans are constantly looking for things, be it be it the right book in the bookshelf, a certain keyword in an article or blueberries in the forest. In many cases, it is important that this search is strategic, efficient, and fast. Animals need to quickly identify predators, and drivers need to be able to search for pedestrians crossing the road they are driving on.

Automating visual search is of great interest for several reasons. Visual search is crucial for applications such as search and rescue, surveillance, fire detection, etc. Autonomous vehicles can both reduce risk and potentially exhibit more intelligent searching behavior than human-controlled ones.

However, while visual search is often seemingly effortless to us humans, it is a complex process. Attempts to understand and recreate human visual search in machines has been a big challenge [23]. At the root of the visual search problem is partial observability. A searcher can only perceive, or pay attention to, a limited region of the searched environment at once. Therefore which regions to observe and in what order becomes an important decision.

How humans and animals search for things has been studied extensively in neuroscience [23, 59, 46]. When we search, we use features of the environment to guide our attention [60, 23]. For example, we know to look for berries at the forest floor, and not to look for boats on land. Furthermore, search is not purely reactive but involves the use of memory. We also use memory to take the history of the search into account when deciding where to move our attention [60].

Such features can in some cases be quite subtle and difficult to pick up, even for humans. Manually engineering guidance in accordance with these features can be expensive, especially if a searching system should be deployed in many different environments. If one could instead learn a good searching strategy from a limited set of sample environments this would be circumvented. Such a system could be taught to search in arbitrary environments without the use of manually encoded environment-specific rules.

Reinforcement learning [56] (RL) is a paradigm that is suited for learning mappings from sensor values to actions. In recent years, RL has been combined with deep learning [27] with tremendous success. It has been used to master arcade games [45], board games [54], and even complex real-time strategy games [58]. Several works have also applied RL to tasks involving embodied agents with visual input [40, 43, 66, 41]. This makes it interesting to see if RL can also be applied to visual search.

## 1.2 Aim

The aim of this thesis is to investigate how an intelligent agent that learns to search for targets using visual input can be implemented with deep reinforcement learning. Such an agent should learn the characteristics of the environments it is trained on and utilize this knowledge to search strategically in unseen environments. Specifically, we consider scenarios where the agent can only observe a small portion of its environment at any given time through a camera whose movements are unrestricted. The agent has to actively choose where to look in order to gain new information about the environment.

We postulate that an effective searcher learns how the distribution of targets and prioritizes regions where they are more likely according to previous experience. The distribution of targets may be correlated with the appearance of the searched scene. A good searcher should integrate information over time to build an internal representation of the environment and use it to make informed decisions. The agent should be able to search the environment exhaustively while avoiding visiting the same region twice. Finally, the agent should be able to locate multiple targets while minimizing the number of steps taken.

If such a system is to be trained and deployed for a real-world task, there is likely a limited set of samples to learn from. Therefore it is also of interest to investigate how many samples are required to infer how to effectively search in similar environments. While similar problems have been addressed in the past, both with learning agents [40, 41, 26, 14] and non-learning agents [53, 24], our impression is that this is the first work to address learning to search in unseen environments where emphasis is placed on how arbitrary visual cues can guide search. Our contributions are as follows:

- We provide a set of environments to train and evaluate visual search agents.
- We propose a method for solving the visual search task with reinforcement learning that can be trained end-to-end.
- We compare the method to a set of common baseline agents.
- We investigate how well each agent is able to generalize to unseen environments.

## 1.3 Research questions

This thesis will address the following questions:

1. How can an agent that learns to intelligently search for targets be implemented with reinforcement learning?
2. How does the learning agent compare to random walk, greedy search, and a human searcher with prior knowledge of the searched environment?
3. How many training samples are needed for the agent to generalize to unseen indistribution environments?

## 1.4 Delimitations

We focus on the behavioral and decision-making aspects of the presented problem, and delimit ourselves from difficult detection problems. For this reason, targets will deliberately be made easy to detect once visible. Furthermore, we make the assumption that the searched environment is static. The appearance of the environment and the location of the targets does not change from one observation to the next. Finally, we are specifically interested in deep reinforcement learning approaches.



## 2 Theory

This chapter introduces background and related work.

### 2.1 Background

In this section, we give background to the tackled problem and outline the theory behind our approach. Sections 2.1.1 and 2.1.2 provide perspective to the problem from neuroscience and machine perception respectively. Section 2.1.3 overviews the theory behind function approximation with neural networks, and some common neural network architectures. Section 2.1.4 summarizes the foundations of reinforcement learning and deep reinforcement learning.

#### 2.1.1 Visual Search and Attention

The perceptual task of searching for something in a visual environment is usually referred to as *visual search* [59]. The object or feature that is being searched for is referred to as the *target*, and the other objects or features in the environment as *distractors*. This task has been studied extensively in psychology and neuroscience.

Two big limitations of performance in visual search are processing power, and limited observability. Processing high-dimensional sensory input like images is expensive, and environments are often too large and complex to view all at once. An observer that searches a scene for targets has to direct its *visual attention* to one region at a time. Humans scan environments by directing their gaze (*overt attention*) and shifts of attention in the current visible region (*covert attention*) [34]. While covert attention tends to be reactionary, overt attention usually integrates more features over time. In this work we focus specifically on the former, moving the gaze to bring targets into view.

Humans control overt attention through both eye movements and head movements. Eye movement is almost instant between locations, and the cost of directing attention with eye movements is constant regardless of distance. This is not true for head movements, and in general not true for robotic systems: movements induced by motors tend to come at a cost that is proportional to the distance moved, both in time and energy. The cost of directing overt attention means that there is a need to do so strategically.

If the searched environment is a random field, visual search is akin to a random process [46]. Experiments in humans have shown that search in featureless environments is con-

sistent with random walk with no memory. Optimal search algorithms in such environments would be exhaustive, such as those found in coverage path planning [25]. Humans have also been shown to use the history of observations to improve visual search efficiency. Simple memory mechanisms, like some inhibition-of return mechanism [34] that prevents searching visited locations twice, improve search time in random fields considerably.

Natural environments are in fact seldom completely random, but instead tend to exhibit some structure. When finding targets, there is usually something about this structure that can be utilized to improve search performance. Such regularities can be learned from past experience and used during future searches. Improvements in human search performance have been measured when certain locations have higher probabilities of containing targets [23, 60]. Furthermore, if targets usually co-occur with other visible elements they tend to be easier to find [23, 60].

### 2.1.2 Active Vision and Object Search

Much of past and present research in machine perception involves a passive observer which samples and perceives images from a fixed distribution. Animal perception, in contrast, is active - we do not only see, but also decide where to look. In the *active vision* paradigm, an observer has some control of its sensory input. [2]

Active vision, and *active perception* in general, is a problem of intelligent data acquisition. An active observer must control its sensory inputs to constrain the interpretation of its environment. One of the difficulties of active perception problems is that they are scene and context dependent. A thorough understanding of the data acquisition parameters and the goal of the visual processing is needed. [9]

In active vision, searching for objects with a camera in unknown environments is referred to as *object search* or *active object localization*. A searching observer has to both recognize and localize its targets, while controlling its sensory input. To search effectively, it must also model its environment and perform path planning. [16]

Solving the object search problem optimally involves determining a sequence of camera controlling actions that maximizes the probability of finding the target while satisfying a cost constraint. The state of the searcher is uniquely determined by the control parameters of the camera. Actions that adjust these parameters and adjust the observed region come at some cost in time or energy. Finally, the agent has some prior knowledge of the probability of targets which it updates after each observation. Finding an optimal solution in three dimensional search spaces has been shown to be NP-complete, necessitating approximate solutions. [63, 4]

### 2.1.3 Deep Learning

Deep learning is a family of techniques in which hypothesis are represented as computation graphs with tunable weights. The computation graphs are inspired by biological neurons in the brain and are referred to as *neural networks*. Deep neural networks consist of *nodes* arranged in *layers*: one input layer, zero or more hidden layers and one output layer. Each layer receives an input *representation* [11] from the previous layer and outputs a transformed representation to the next layer. Given some input, a neural network optimizes its output representation with regard to some criterion. Usually, a loss function  $\mathcal{L}$  is minimized by updating the weights  $\mathbf{w}$  of the network with some variant of *gradient descent* with learning rate  $\alpha$ :

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \quad (2.1)$$

The only requirement on the functions computed by each node is that it is differentiable. As long as this holds, layers can be stacked arbitrarily and the gradients can be computed

with the chain rule. This way, errors in the output can be passed back through the network (*back-propagation*) and used to update the weights. [51, 27]

The architecture of a neural network imposes some bias onto the learning that its expected to be useful for generalizing to unseen samples. We now describe three neural network architectures that will be used in this work.

### 2.1.3.1 Feedforward Neural Network

A feed-forward neural network, also known as a multi-layer perceptron (MLP) [27], only has connections in one direction. Each node in the network receives inputs from its predecessors and outputs the result of a function of those inputs. The output  $y$  of each node is usually computed by taking the weighted sum of its inputs  $x$  and applying some non-linear function

$$y_j = g_j(\mathbf{w}_j^T \mathbf{x}), \quad (2.2)$$

where  $y_j$  is the output of node  $j$ ,  $g_j$  is a non-linear *activation function*,  $\mathbf{w}_j$  is the vector of weights leading into node  $j$ , and  $\mathbf{x}$  is the vector of inputs to the node. By convention, each layer also has some *bias* that allows the total weighted input to  $g_j$  to be non-zero even when the outputs from the previous layer are zero. The bias is included as an extra input  $x_0$  fixed to 1, and an extra tunable weight  $w_{0,j}$ . The non-linearity ensures that a network with at least two layers can approximate any continuous function. [51]

### 2.1.3.2 Convolutional Neural Network

Convolutional neural networks (CNNs) contain spatially local connections. They have patterns of weights, called *kernels*, that are replicated across units in each layer. With some input vector  $\mathbf{x}$  of size  $n$  and a vector kernel  $\mathbf{k}$  of size  $l$ , the (discrete) convolution operation  $\mathbf{z} = \mathbf{x} * \mathbf{k}$  is defined as

$$z_i = \sum_{j=1}^l k_j x_{j+1-\frac{l+1}{s}}, \quad (2.3)$$

where  $s$  is the *stride*. This operations can be generalized up to more than one dimension, such as 2 dimensions for images and 3 dimensions for volumes. With multiple input channels, kernels are stacked into a *filter*. The outputs of each kernel are then summed over, giving one output channel per filter.

There are several advantages to using CNNs for structured input data where neighboring values are correlated. Kernels are smaller than the input, which means that fewer parameters have to be stored. These *sparse interactions* give CNNs reduced memory requirements, as well as improved statistical and computational efficiency.

Furthermore, the same parameters are also used for more than one function in the CNN. *Parameter sharing* across input locations mean that layers in a CNN have *equivariance* to translation. The output of one kernel is the same regardless of the input location. This property of CNNs is useful for images where similar features may be useful regardless of their location in the input. [27]

### 2.1.3.3 Recurrent Neural Network

Recurrent neural networks (RNNs) extend feed-forward networks by allowing cycles in the computation graph. Each cycle has a delay so that some *hidden state* from the previous computation is used as input to the current computation. A recurrent layer with input  $\mathbf{x}_t$ , output  $\mathbf{y}_t$  and hidden state  $\mathbf{z}_t$  is defined by

$$\begin{aligned} \mathbf{z}_t &= f_{\mathbf{w}}(\mathbf{z}_{t-1}, \mathbf{x}_t) \\ \mathbf{y}_t &= g_y(\mathbf{W}_{z,y}, \mathbf{z}_t), \end{aligned} \tag{2.4}$$

where  $f_{\mathbf{w}}$  is the update process for the hidden state and  $g_y$  is the activation function for the hidden layer. This model can be turned into a feed-forward network over a sequence of input vectors  $\mathbf{x}_1, \dots, \mathbf{x}_T$  and observed outputs  $\mathbf{y}_1, \dots, \mathbf{y}_T$  by *unrolling* it for  $T$  steps. The weights are shared across all time steps. This means that RNNs can operate on inputs of arbitrary lengths. The hidden state is used as a summary of all previous items in the sequence. Thus, RNNs make a Markov assumption. [51]

In practice, conventional RNNs struggle with learning long-term dependencies. During back-propagation, gradients can tend to zero for long sequences, something known as the vanishing gradient problem [27]. An architecture that addresses this issue is long short-term memory (LSTM) [33]. LSTMs include a *memory cell*  $c$  in the hidden state that is copied from time step to time step, and three soft *gating units* that govern the information flow in the hidden state update process  $f$ . This makes LSTMs particularly useful for learning over long sequences.

## 2.1.4 Reinforcement Learning

Reinforcement learning (RL) [56] is a subfield of machine learning concerned with learning from interaction how to achieve a goal. This section introduces the fundamental concepts of RL.

### 2.1.4.1 Partially Observable Markov Decision Processes

The problem of learning from interaction to achieve some goal is often framed as a Markov decision process (MDP). A learning *agent* interacts continually with its *environment*. The agent takes the *state* of the environment as input, and select an *action* to take. This action updates the state of the environment and gives the agent a scalar *reward*. It is assumed that the next state and reward depend only on the previous state and the action taken. This is referred to as the *Markov* property. [35]

In an MDP, the agent can perceive the state of the environment with full certainty. For many problems, including the one we consider here, this is not the case. The agent can only perceive a partial representation of the environment's state. Such a process is referred to as a partially observable Markov decision process (POMDP). A POMDP is formally defined as a 7-tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O}, \gamma \rangle$ , where

- $\mathcal{S}$  is a finite set of states,
- $\mathcal{A}$  is a finite set of actions,
- $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \Pi(\mathcal{S})$  is a state-transition function,
- $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is a reward function,
- $\Omega$  is a finite set of observations,
- $\mathcal{O} : \mathcal{S} \times \mathcal{A} \rightarrow \Pi(\Omega)$  is an observation function, and
- $\gamma \in [0, 1]$  is a discount factor.

Assume that the environment is in state  $s_t \in \mathcal{S}$ , and the agent selects action  $a_t \in \mathcal{A}$ . Then,  $T(s_t, a_t, s_{t+1})$  is the probability of ending in state  $s_{t+1}$  and  $r_t = R(s_t, a_t)$  is the expected reward gained by the agent. The agent also receives an observation  $o_t \in \Omega$  with probability  $\mathcal{O}(s_{t+1}, a_t, o_t)$ . [35] Figure 2.1.4.1 illustrates the interaction between agent and environment.



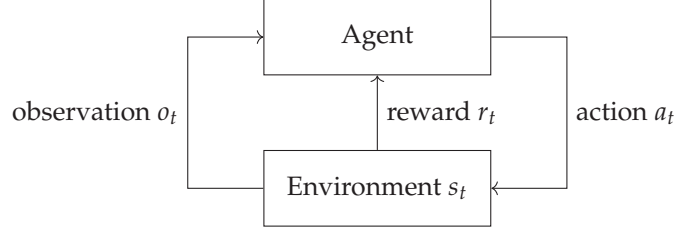


Figure 2.1: Interaction between agent and environment in a partially observable Markov decision process.

The agent and environment interact over a sequence of discrete time steps  $t = 0, 1, \dots, T$ , giving rise to an *episode* of length  $T$ . At each time step  $t$ , the goal of the agent is to select the action that maximizes the expected *discounted return*:

$$\mathbb{E} \left[ \sum_{k=0}^T \gamma^{k-t-1} r_k \right]$$

Since the agent receives partial observations of the environment’s state, it has to act under uncertainty. Planning in a POMDP is undecidable, and solving one is often computationally intractable. Approximate solutions are more common, where the agent usually maintains an internal *belief state* [35] which it acts on. The belief state summarizes the agent’s previous experience and is therefore dependent on the full *history* of actions and observations. It does not need to summarize the whole history, but generally only the information that helps the agent maximize the expected reward. From here on we will use the belief state and the environment state  $s$  interchangeably.

#### 2.1.4.2 Policies and Value Functions

The behavior of the agent is described by its *policy*. A policy  $\pi$  is a mapping from perceived environment states to actions. Policies are often stochastic and specify probabilities for each action, with  $\pi(a|s)$  denoting the probability of taking action  $a$  in state  $s$ . [56]

Most RL solutions methods also approximate a *value function*. A value function  $V_\pi$  estimates how good it is to be in a state. The value function  $V_\pi(s)$  is the expected (discounted) return when starting at state  $s$  and following policy  $\pi$  until the end of the episode. There are two common alternative value functions: The *quality function*  $Q_\pi(s, a)$  gives the value of state  $s$  under policy  $\pi$  where  $a$  is the first action taken. Given a quality function, *action-value* methods choose the action greedily at every state as  $\arg \max Q_\pi(s, a)$ . The *advantage function*  $A_\pi(s, a)$  instead represents the relative advantage of actions,  $A_\pi = Q_\pi - V_\pi$ . [56]

For problems with large state and action spaces, it is common to represent value functions with *function approximation*. In such cases, it is common to encounter states that have never been encountered before. This makes it important that the estimated value function can generalize from seen to unseen states. With examples from the true value function, an approximation can be made with supervised learning methods. We write  $\hat{V}(s, \mathbf{w}) \approx V_\pi(s)$  for the approximate value of state  $s$  with some weight vector  $\mathbf{w} \in \mathbb{R}^d$ . [56]

An alternative to action-value methods is to approximate the policy itself. *Policy gradient* [55] methods learn a parametrized policy that select actions without a value function. We denote a parametrized policy as  $\pi(a|s, \theta)$  with  $\theta \in \mathbb{R}^{d'}$  as the parameters to the policy. The policy parameters are usually learned based on the gradient of some performance measure  $L(\theta)$ . As long as  $\pi(a|s, \theta)$  is differentiable with respect to its parameters, the parameters can be updated so as to optimize for the objective.

Advantages of policy parametrization over action-value methods include stronger convergence guarantees [55] and more flexibility in parametrization [56]. In practice, value func-

tions are often still used to learn the policy parameter, but they are not needed for action selection. Such methods are called *actor-critic* methods, with actor referring to the learned policy and critic referring to the learned value function. In these cases, there might also be some overlap between the weights  $\mathbf{w}$  of the value function estimate and  $\theta$  of the policy estimate.

One important aspect of function approximation is its interplay with partial observability. If there is a state variable that is not observable, as for partially observable environments, then the parametrization can be chosen such that the approximate value does not depend on that state variable. Because of this, function approximation is applicable to the partially observable case. [56]

### 2.1.4.3 Design Challenges

One of the challenges that arises in reinforcement learning is the exploration-exploitation trade-off. An RL agent should *exploit* knowledge gained from previous experiences and prefer actions that has yielded reward in the past. It should also *explore* in order to learn better actions to take in the future. Agents that fail to both exploit and explore will lead to failure at the task, and striking a good balance between the two is non-trivial. [56]

Another challenge is the design of reward signals. For some tasks, like certain video games, the objective is simply to maximize the score obtained. In this case there is an inherent reward signal and the agent achieves its task simply by maximizing this inherent signal. Other times, we have a task we want the agent to solve and have to design a reward signal around that task. Designing rewards is not straight-forward and can often have unintended effects [56]. Special care has to be taken to ensure that the reward encourages the desired behavior.

Here, the problem of *sparse rewards* also comes into play. The agent has to be reward frequently enough to allow it to achieve its goal once. Often it has to incentivize it to achieve its goal efficiently, with multiple different starting conditions. If rewards are too sparse, the agent may explore aimlessly and take too long to find achieve its goal. If the received reward is temporally distant from the action that caused it, the agent may have difficulty connecting the two. This is known as the *credit assignment problem* [39].

In practice, rewards are often designed through trial-and-error. Through *reward shaping* [38], the reward is designed so as to guide the agent towards achieving its goal by giving additional rewards along the way. Several iterations of a reward signal are tried until one yields expected and sufficient results.

### 2.1.4.4 Deep Reinforcement Learning

As mentioned in Section 2.1.4.2, policies and value functions are often approximated. Neural networks have good properties for function approximation and have been used for RL with success. One early example is TD-Gammon [57], a neural network trained with RL that reached expert Backgammon performance in 1995.

More recently, the successes of deep learning have bled over into the field of RL. In 2015, Mnih et al. [45] extend [44] and introduce DQN, which combines deep neural networks with RL. DQN successfully plays Atari games using only visual input. It approximates the quality function  $q(s, a)$  with a CNN architecture, and selects actions greedily. To incorporate some memory, images from the 4 previous time steps are stacked and used as input to the neural network. The input is fed through three convolutional layers and a hidden fully connected layer, all with ReLU activation functions. The output layer has one output for each valid action, representing Atari controller buttons.

DQN inspired many several follow-up works which use deep neural networks to approximate policy and/or value functions. Such methods are often referred to as deep reinforcement learning (deep RL) methods.

### 2.1.4.5 Proximal Policy Optimization

Proximal policy optimization (PPO) [52] is a family of policy-gradient algorithms that have turned out to strike a good balance between simplicity, stability and ease of tuning while achieving strong results on several tasks [52, 30, 18, 58, 5].

PPO alternates between sampling agent-environment interactions over  $T$  time steps and optimizing the policy using those interactions. The parameters  $\theta$  of the policy  $\pi$  are updated with the objective

$$\mathcal{L}_{\text{CLIP}}(\theta) = \hat{\mathbb{E}}_t [\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t)].$$

The expectation  $\hat{\mathbb{E}}$  indicates the empirical average over a finite batch of samples collected under the old policy with parameters  $\theta_{\text{old}}$ . Here,  $r_t(\theta) = \frac{\pi(a_t|s_t, \theta)}{\pi(a_t|s_t, \theta_{\text{old}})}$  and  $\hat{A}_t$  is an estimate of the advantage function. The clip range  $\epsilon$  is used to clip the surrogate objective, putting a pessimistic bound on the product of the probability ratio and the advantage estimate. This in turn ensures that the size of the policy updates is limited. In an actor-critic approach, the advantage is estimated over using a learned state value function  $V(s)$  as

$$\hat{A}_t = \delta_t + (\gamma\lambda)\delta_{t+1} + \dots + (\gamma\lambda)^{T-t+1}\delta_{T-1}$$

where  $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ . If the policy and value functions estimations share parameters, for example in a multi-headed neural network architecture, a loss function that combines the policy loss and the value function error term must be used. It is also common to introduce an entropy bonus to ensure sufficient exploration. The full PPO objective for an actor-critic approach is

$$\mathcal{L}_t(\theta) = \hat{\mathbb{E}}_t [\mathcal{L}_{\text{CLIP}}(\theta) - c_1 \mathcal{L}_{\text{VF}}(\theta) + c_2 S[\pi](s_t, \theta)]$$

where  $c_1$  and  $c_2$  value loss is the entropy bonus coefficients,  $S$  is an entropy bonus and  $\mathcal{L}_{\text{VF}}$  is a squared-error loss in the value output. Algorithm 2.1.4.5 shows the steps of the PPO algorithm.

---

**Algorithm 1** Proximal Policy Optimization

---

```

for iteration = 1, 2, ... do
  for actor = 1, 2, ..., N do
    Run policy  $\pi_{\theta_{\text{old}}}$  for  $T$  time steps
    Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
  end for
  Optimize  $\mathcal{L}$  wrt  $\theta$ , with  $K$  epochs and mini-batch size  $M \leq NT$ 
   $\theta_{\text{old}} \leftarrow \theta$ 
end for

```

---

## 2.2 Related Work

Works that consider tasks that are related to searching for targets in unknown environments can be found in several fields under many names. In this section we survey some of the more relevant ones, focusing on those that employ RL methods.

### 2.2.1 Search with Reinforcement Learning

There are several rigorous attempts to implement robotic object search systems using non-learning methods. Forssén et al. [24] implement a mobile robot that searches for specific objects within a cluttered environment. Their robot explores by moving towards unexplored regions and looks around to identify potential objects. Each objects is examined from several

perspectives and ranked by probability of being the queried object. Shubina and Tsotsos [53] model the search problem as in [63], and solve with a greedy action selection strategy. They do not take environment appearance into account, and rely on prior knowledge of target distribution being provided as a spatial probability map. Some works also take appearance and semantics of the environment into account to search more efficiently [6, 7].

These systems solve more complex tasks than the one we consider, such as handling multiple target object categories and navigating in obstructed environments. However, they are also restricted by their reliance on human knowledge to act well in their environments. While reinforcement learning is not among the more traditional solution methods for active vision tasks [16], learning systems that are less dependent on expert knowledge have the potential to be more generally applicable.

### 2.2.1.1 Sequential Visual Attention

Minut and Mahadevan[40] propose an sequential model of selective visual attention for visual search tasks. An agent is tasked with finding a particular object in a scene. A policy for controlling a fixed pan-tilt-zoom camera is learned using reinforcement learning. The goal of the agent is to aim the camera at the region where a target is most likely to be found. It has to decide where to fixate next based on visual information only. Despite being limited to a single environment, this is an early example of visual search modelled as an RL problem.

Mnih et al. [43] take inspiration from visual attention and foveated vision found in humans and propose to use a similar mechanism for computer vision tasks. Applying a CNN to a large images can be expensive, as the complexity scales linearly with the number of images pixels. They propose a recurrent model that extracts information from images by adaptively selecting a sequence of smaller regions to process at high resolution. At each time step, the agent receives an image observation of the environment. Through actions, the agent selects a limited region of this observation to view in high resolution. It is given a reward that is dependent on the task the agent should perform. It can access this image via a bandwidth-limited sensor which it focuses on a limited region. The agent maintains an internal state using an LSTM layer. Though the model is not differentiable, it is trained using RL with a policy gradient method. The authors evaluate the agent for image classification and a game-like task. They find that the model outperforms a similar convolutional architecture for cluttered object classification tasks. This is attributed to its ability to focus its attention on important regions.

### 2.2.1.2 Active Object Detection

In computer vision, *object detection* is the task of detecting semantic objects of a certain class in images. Detecting an object entails *recognizing* that it is present in an image, and *localizing* it by determining its bounding box. State of the art object detection use deep learning techniques and usually involve a deep convolutional architecture to recognize objects [65]. For localization, a region proposal process considers the whole image and returns a set of bounding box candidates. The object recognizer is then run on each region proposal. Region proposals are akin to visual attention, as they limit the region of the image that is processed by the recognizer.

Although we do not focus on difficult recognition problems in this work, object localization is highly relevant. The difference is that images in object detection are passively sampled - they are drawn from some distribution and all are independent. Furthermore, the whole scene that is searched for objects is visible in the image. In *active object detection*, the analyzed images are instead chosen so as to help the detector perform its task better.

Caicedo and Lazebnik [14] propose to use deep reinforcement learning for active object localization in images where the object to be localized is fully visible. An agent is trained to successively improve a bounding box using translating and scaling transformations. They use a reward signal that is proportional to how well the current box covers the target object.

An action that improves the region is rewarded with +1, and given a punishment of -1 otherwise. They find that this reward communicates more clearly which transformations keep the object inside the box and which take the box away from the target. When there is no action that improves the bounding box, the agent may select an indication action (which would be the only action that does not give a negative reward) which resets the box. This way the agent may select additional bounding boxes. Each indication action modifies the environment by marking it so that the agent may learn to not select the same region twice.

A similar work by Ghesu et al. [26] present an agent for anatomical landmark detection trained with DRL. Different from [14] is that the entire scene is not visible at once. The agent sees a limited region of interest in an image, with its center representing the current position of the agent. The actions available to the agent translate the view up, down, left and right. A reward is given to the agent that is equal to the supervised relative distance-change to the landmark after each action. Three datasets of 891 anatomical images are used. The agent starts at random positions in the image close to the target landmark and is tasked with moving to the target location. While achieving strong results (90% success rate), the scenes and targets are all drawn from a distribution with low variance. Most real-world search tasks exhibit larger variance than a population of anatomical images of the human body.

Chen and Gupta [17] use a spatial memory for context reasoning in object detection. They argue that object detection systems require memory to perform well. Furthermore, they pose that this memory should capture the spatial layout of the scene in order to model object-object relationships. Their proposed memory remembers features and locations of past detections in an image and uses this to improve future ones. A CNN is used to extract features from the memory and used together with the original image as input to a standard region proposal network. The approach gives a small improvement over the standalone region proposal network.

### 2.2.1.3 Visual Navigation

Visual navigation is the process of determining a suitable path between a start point and destination point using visual input [12]. Visual navigation has been studied for use in both autonomous ground vehicles, where obstacles have to be avoided, as well as for unmanned aerial vehicles which typically don't have the problem of avoiding obstacles. The visual search problem can be considered analogous to a visual navigation in unrestricted environments. In recent years, several works have used deep RL to solve various visual navigation tasks.

Zhu et al. [66] create a model for target-driven visual navigation in indoor scenes with DRL. An observer is given a partial image of its scene as well as an image of the target object, and is tasked with navigating to the object in the scene with a minimal number of steps. The agent moves forwards, backwards, and turns left and right at constant step lengths. They use a reward signal with a small time penalty to incentivize task completion in few steps. They compare their approach to random walk and the shortest path and achieve promising results. This setup is quite similar to the one considered in this report, but the authors make a few assumptions that we do not. They use a set of 32 scenes, each of which contain a fixed number of object instances. They focus on learning spatial relationships between objects in these specific scenes, and have scene-specific layers to achieve this. Thus, while they show that they can adapt a trained network to a new scene, their approach is unable to zero-shot generalize to new scenes.

Ye et al. [62] propose an alternative architecture to [66] which, in addition to an image of the camera view and the target object, uses the output of an object recognition module to select actions. Specifically, the output of the localization component is fed as input. This allows the agent to know if the target object is in view, and where. The object recognition module is trained separately from the policy. While requires a set of images of potential target objects labelled with their class and location, it also illustrates how object recognition



can be offloaded to a separate module in reinforcement learning agents that benefit from such functionality.

Yang et al. [61] propose to use semantic scene priors to improve navigation towards objects in unseen scenes. The agent observes its environment through a camera. Prior knowledge of spatial relations between objects types is encoded as a graph, which is encoded with a graph convolutional network. Their agent is also told what object to look for in the form of a word embedding. The encoded object graph and word embedding is used as input together with the camera image to an actor critic network that selects actions that navigate towards the target object. As the agent collects experiences, it also updates its prior knowledge of spatial relationships between objects. It is shown that prior knowledge encoded in such a way improves the agents ability to navigate to objects.

Mnih et al. [42] use a recurrent policy with only RGB images to navigate in a labyrinth. 3D labyrinths are randomly generated, and an agent is tasked with finding objects in them. The same architecture as in [45] is used, but with 256 LSTM cells after the final hidden layer. At each episode, a maze is randomly generated with objects that give rewards then traversed scattered around. They train the agent using an actor-critic approach and manage to get the agent to successfully navigate in unseen but similar environments.

Mirowski et al. [41] propose a new approach to tackle a similar maze navigation problem as [42]. They introduce a new architecture, that in addition to an image, also observes the reward, velocity and the action from the previous time step. Observing the reward may help the agent learn features of desirable frames, but also limits the set of possible reward signals as the reward has to be present during testing. As in [42], the agent has a recurrent LSTM layer. The architecture is trained using an auxiliary depth prediction objective. The authors argue that understanding the depth of the searched environment helps the agent navigate it better. This agent is evaluated against three baselines on both static and randomly generated maze environments. It is found that an agent that only receives image observations performs worse than one that also has an LSTM layer. Adding the previous reward, action and velocity to the observations bring modest performance improvements. The depth prediction is found to improve performance in all environments, indicating that the correct auxiliary task can help learning.

Dhiman et al. [22] critically investigate deep RL for navigation, using a CNN and RNN architecture as in [29] and [41]. They ask whether DRL algorithms are inherently able to gather and exploit environmental information for during navigation. Experimentally, they find that an agent is able to exploit environment information when trained and tested on the same map. When trained and tested on different maps, the agent does not seem to be able to exploit environment information to navigate more effectively. The authors use relatively small training set of less than 1000 samples, but do not consider that their agents could have overfit this set. Finally, they find that the agent is not able to consistently find optimal navigation paths in seen environments when its starting position is randomized.

### 2.2.2 Memory Architectures

A recurring theme in the related work above is the use of memory architectures. In most cases where active sensing is involved, memory is a requirement for good performance. As mentioned in Section 2.1.4.1, effective behavior in partially observable processes usually requires some form of memory. Agents need to summarize the history of interactions with their environment for good action selection. Works that don't use memory rely on either full observability, like [14], or low variance so that search can be done reactively, like [26]. Memory is intuitively important for the search task we consider here, as integrating features over time should lead to more efficient search strategies.

Anderson et al. [3] emphasize the importance of memory mechanisms that support the construction of rich internal representations of environments in navigation agents. They argue that the nature of the internal representation is central to the study of embodied naviga-

tion. Simple agents that are purely reactive and act on the sensory input at the current time step only work for simple tasks. Simple memory mechanisms, like the frame stacking used in [45] can model velocity and work well in certain environments where reactive action is sufficient. Once a task requires integration of features over longer time, more advanced memory is required.

Augmentations like recurrent update mechanisms add more potential. Hausknecht and Stone [29] investigate the effects of adding memory in the form of a recurrent step to DQN in order to tackle POMDPs. They use the same convolutional network as [45], but only use the most recent frame as input and replace the hidden layer with a recurrent LSTM layer. It is found that the agent is able to integrate information over time and achieves comparable performance to the original DQN agent. Several other works use a recurrent step to solve tasks that require memory, like [43] and [42]. General-purpose recurrent neural networks can in theory remember anything, but have practical issues like the ones described in Section 2.1.3.3. Furthermore, they might not be easily trained for all tasks.

Specialized memory architectures can provide better performance for their intended tasks. Oh et al. [47] propose one such memory that can retain information over a fixed limited number of time steps. During each time step, the agent encodes its observation and shifts it into the memory. It also reads the memory using a soft attention [8] mechanism whose weights are determined by the encoded observation. This way, the agent can recall information from the past that is useful to the present. The authors evaluate this architecture in several visual navigation tasks, where image observations are encoded with a CNN. In their tests, their memory architecture performs better than LSTM layers, and generalizes better to unseen environments.

Several works propose the use of spatial memories rather than temporal ones, which the agent can address with its location and use to store structured information of its environment. Parisotto and Salakhutdinov [49] introduce a general-purpose spatial memory architecture. They assume that the environment can be discretized into a set of positions. The agent retains a feature map with one feature vector per position it can occupy in the environment. At each time step, the agent reads the feature map from the previous time step and writes information to the slot for its current position. This architecture is shown to outperform LSTM and the architecture in [47] in several navigation tasks which require remembering over many time steps. Several similar architectures have been proposed [31, 28, 15].

### 2.2.3 Generalization and Inductive Bias

In RL, generalization refers to an agents ability to act well in environments that have not been seen during training. This is vital if RL is to be applied to real-world problems where conditions are diverse and unpredictable. Agents have to be robust to such variations without being trained on them (*zero-shot* policy transfer). In the case of visual search, we want a system that can search in novel environments by generalizing from those seen during training.

While deep neural networks have proved to be effective function approximators for RL, they are also prone to *overfitting*. High-capacity models trained over a long time may memorize the distribution seen during training rather than general patterns. While studied in supervised learning, overfitting and generalization has generally been neglected in deep RL [37]. Training and evaluation stages are typically not separated. Instead, the final return on the training environments is used as a measure of agent performance. This results in agents that perform badly on environments that are only slightly different from those seen during training.

Zhang et al. [64] study overfitting and generalization in deep RL. With experiments, they show that RL agents are capable of memorizing training data, even when completely random. When the number of training samples exceeds the capacity of the agent, they overfit to them. When exposed to new but statistically similar environments during testing, test performance could vary significantly despite consistent training performance. The authors argue

that good generalization requires the *inductive bias* of the algorithms to be compatible with the bias of the problems. The inductive bias refers to a priori algorithmic preferences, like neural network architecture. When comparing MLPs with CNNs, they find that MLPs tend to be better at fitting the training data but worse at generalizing. When rewards are spatially invariant, CNNs generalize much better than MLPs. The authors advocate for carefully designed testing protocols for detecting overfitting. The effectiveness of stochastic-based evaluation depends on the properties of the task. Agents could still learn to overfit to random training data. For this reason, they recommend isolation of statistically tied training and test sets.

In a similar spirit, Cobbe et al. [19] construct distinct training and test sets to measure generalization in RL. They find that agents can overfit to surprisingly large training sets, and that deep convolutional architectures can improve generalization. Methods from supervised learning, like L2 regularization, dropout, data augmentation and batch normalization are also shown to aid with generalization.

Many current deep RL agents do not optimize the true objective that they are evaluated against, but rather a handcrafted objective that incorporates biases to simplify learning. Stronger biases can lead to faster learning, while weaker biases potentially lead to more general agents. Hessel et al. [32] investigate the trade-off between generality and performance from the perspective of inductive biases. Through experimentation with common reward sculpting techniques, they find that learned solutions are competitive with domain heuristics like handcrafted objectives. Learned solutions also seem to be better at generalizing to unseen domains. For this reason, they argue for avoiding biases determined with domain knowledge in future research.

Cobbe et al. [18] introduce a benchmark for sample efficiency and generalization in RL. They make use of procedural generalization, dependent on a random seed, to decide many parameters of the initial state of the environment. This forces agents to learn policies that are robust to variation and avoid overfitting. To evaluate sample efficiency of agents in the benchmark, they train and test on the full distribution of states. To evaluate generalization, they limit the number of training samples and then test on held out levels. When an episode ends, a new sample is drawn from the training set. Agents may train for arbitrarily many time steps. The number of training samples required to generalize is dependent on the particulars and difficulty of the environment. The authors choose the training set size to be near the region when generalization begins to take effect. Empirically they find that larger model architectures improve both sample efficiency and generalization. Agents strongly overfit to small training sets and need many samples to generalize. Interestingly, training performance improves as the training set grows past a certain threshold. The authors attribute this to the implicit curriculum of the distribution of levels.

Kirk et al. [37] survey generalization in deep RL. They distinguish between methods that try to generalize from training to testing sets by increasing their similarity, and those that attempt to handle their differences. To handle differences between training and test sets, an agent’s policy should rely only on features which will behave similarly in the training and testing context. If the similarities between training and testing contexts are known, they can be encoded as inductive bias for stronger generalization. When there is not specific inductive bias to encode, standard regularization can be used. Finally, when we can not rely on specific inductive biases or standard regularization, we have to rely on learning the invariances from training and hope that these generalize to testing.

They further propose a formalism for collections of problems called contextual CMDPs. A CMDP which is an MDP  $\mathcal{M}$  (or POMDP) where the state can be decomposed into  $s = \langle c, s' \rangle$ , where  $s \in \mathcal{S}$  is the underlying state and  $c \in \mathcal{C}$  is the *context*. The context takes the role of a seed and determines the sample drawn from the underlying distribution of task instances. In a CMDP, separate train and test tasks can be defined by creating  $\mathcal{C}_{\text{train}} \subseteq \mathcal{C}$  and  $\mathcal{C}_{\text{test}} \subseteq \mathcal{C}$  such that  $\mathcal{C}_{\text{train}} \cap \mathcal{C}_{\text{test}} = \emptyset$ . It is argued that the only choice when evaluating generalization with procedurally generated environments is the training set size. Evaluation has to be performed on the full distribution of (held out) contexts.



### 2.2.4 Evaluation of Agents

A problem in state-of-the-art RL is reproducibility. There is often non-determinism, both in the methods and environments used. Furthermore, many methods have intrinsic variance which can make published results difficult to interpret. In some cases, training an algorithm in an environment twice can give learning curves that do not fall within the same distribution [30]. This has meant that reproducing state-of-the-art deep RL results is difficult. A common solution is to report average results and variance across multiple different runs.

However, the sample inefficiency of current deep RL algorithms together with the rise in popularity of more challenging benchmarks has led to long training times. This has made it less feasible to measure performance over many training runs, which in turn has led to a shift to only evaluating a small number of runs ( $N \leq 5$ ) per task. [1] This is problematic, as several analyses indicate that as many as  $N = 20$  runs are required for statistically significant results [21, 20].

Henderson et al. [30] investigate reproducibility methods in deep RL empirically, focusing on policy gradient methods. They try to reproduce results in works with published code bases, and investigate how various modifications affect results.

They find that hyperparameters and the choice of network architecture for policy and value function approximation can affect performance significantly. Furthermore, rescaling rewards can have a large effect, although it is difficult to predict how. It is found that ReLU activations tend to perform best across environments and algorithms. For PPO, the use of large networks may require changing other hyperparameters like learning rate. Confidence bounds with sample bootstrapping is used to show that PPO is among the more stable algorithms.

For certain environments, learning curves can indicate successful optimization but the learned behavior may not be satisfactory. It is therefore important to not only show returns, but also demonstrations of the learned policy in action. Interestingly, implementation differences that are not reflected in publications can have a dramatic impact on performance. It is therefore necessary to enumerate implementation details and package code bases with publications. Performance of baseline experiments should also match original baseline publication code.

Finally, [30] make the point that more emphasis should be placed on applying RL algorithms to real-world tasks. It could be more useful to propose a set of tasks that an algorithm could be used for than to show performance on fictional tasks.

Anderson et al. [3] discuss problem statements and evaluation measures for embodied navigation agents, and make a set of recommendations. A navigation agent should be equipped with a special action that indicates that it has concluded the episode. The agent should be evaluated at the time this action is made, and not at some more favorable time step. Proximity to a goal should be measured using geodesic distance, the shortest distance in the environment. They recommend success weighted by (normalized inverse) path length (SPL) as the primary measure of navigation performance. With  $N$  test episodes, SPL is computed as

$$\frac{1}{N} \sum_{i=1}^N S_i \frac{l_i}{\max(p_i, l_i)} \quad (2.5)$$

where  $S_i$  is a binary indicator of success,  $l_i$  is the shortest path distance from the agent's starting position to the goal, and  $p_i$  is the length of the path actually taken in the episode. SPL takes the quality of the solution into account. If 50% of test episode are successful and the agent takes the optimal path in all of them, its SPL is 0.5. By measuring SPL of human subjects, what is a good score can be calibrated.

## 3 Method

In this chapter, the method used is described. Section 3.1 formalizes the problem solved. Section 3.2 details the environment used to evaluate solutions. Section ?? describes the baseline learning method. Section 3.3 describes the approach used to solve the problem with a learning agent. Section 3.4 describes the experiments conducted to answer research questions 2 and 3.

### 3.1 Problem Statement

We can now formally define the problem of searching for targets in unknown environments, adopting the CMDP formalism [37]. Let the task be a contextual POMDP  $\mathcal{M}$  where the state includes the context  $c$ , which we refer to as the *seed*. The state includes a space  $S \subset \mathbb{R}^d$  which we refer to as the *scene*. In the scene, there is a set of  $N$  targets  $T = \{t_0, t_1, \dots, t_N\}$  such that  $t_i \in S$ . At each time step, the agent perceives a subspace  $V \subset S$  of the environment which we refer to as the *view*. If  $T \cap V \neq \emptyset$  there are  $|T \cap V|$  targets in view.

The actions  $a \in \mathcal{A}$  transform the view into a different subset of the scene. With a final action, the agent can indicate that there is one or more target in the view. The observations  $o \in \Omega$  are tuples  $o = \langle x, p \rangle$ . Here,  $x \in \mathbb{R}^{3 \times w \times h}$  is an RGB image representing the current view, and  $p \in S$  is the position of the agent which uniquely identifies the view. The goal of the agent is to select actions that bring each target into view and indicate that they are visible, while minimizing the number of total steps. There is no inherent reward  $\mathcal{R}$  for this problem, so it has to be designed. Finally, seed  $c$  determines the initial view  $V$ , the location of the targets  $T$ , the initial position  $p_0$  as well as the image observations  $x_t$  at each position  $p_t$ .

### 3.2 Environments

To train and test an agent for the problem, we use three different environments with similar observation spaces, action spaces and reward signal. Each environment has different characteristics that test the applicability of the evaluated approaches to different types of search problems. In all environments, the appearance of the scenes and the location of targets are correlated to some degree. This means that the agent should be able to learn common characteristics of each environment and use those to search more efficiently.

As [18] and [42], we leverage procedural generation in all environments. This gives us control over the difficulty of the environments as well as the number of training and test samples the agent is exposed to. A seed determines the appearance of the scene, the location of the targets and the initial position of the agent.

Each episode is terminated when all targets have been found, or after 1000 time steps. Terminating episodes early this way is common to speed up training [48].

### 3.2.1 Observations, Actions and Reward

All three environments use the same observation space, action space and reward signal. The position and image observations are

$$o_t = \langle x_t, p_t \rangle, \text{ where} \quad (3.1)$$

$$x_t \in \mathbb{R}^{3 \times 64 \times 64}, \text{ and} \quad (3.2)$$

$$p_t \in \{0, \dots, H-1\} \times \{0, \dots, W-1\} \quad (3.3)$$

All image observations are  $64 \times 64$  RGB images. The agent moves in a  $H \times W$  grid, and we assume the presence of some oracle that provides the agent with its position. In many realistic scenarios it is possible to determine the global position of an agent (GPS, pan/tilt, etc.). If how each action moves the agent is well-defined, the relative position can be used instead of the absolute one.

The action space is the same in all environments:

$$a_t \in \{\text{UP, DOWN, LEFT, RIGHT, INDICATE}\}, \quad (3.4)$$

where UP, DOWN, LEFT, and RIGHT move the view one step in each direction. This action space is realistic for many real-world search tasks, such as search and rescue with a UAV, where the actions correspond to translations in each cardinal direction, and surveillance with a pan-tilt-camera, where the actions correspond to pitch and yaw rotations. The final action, INDICATE, is used to indicate that a target is in view.

The reward signal should be designed so that the agent learns a policy that achieves the goal of finding all targets with a minimal number of actions. We use the following reward signal:

$$r_t = h - 0.01 + 0.005d + 0.005e \quad (3.5)$$

Here,  $h = |T \cap V|$  if  $a_t = \text{INDICATE}$ , and  $h = 0$  otherwise. This term is equal to the number of targets that were found at this time step. The constant penalty of  $-0.01$  ensures that the agent is rewarded for quick episode completion. The term  $d = 1$  if  $p_t$  is closer to the nearest target than  $p_{t-1}$ , otherwise  $d = 0$ . Similarly,  $e = 1$  if  $p_t$  has not been explored previously and  $e = 0$  otherwise.

Through experimentation, we find that a larger time penalty than  $-0.01$  dominates the reward for finding targets. This is potentially related to the episode time – with a time penalty that is too large, the reward  $h$  for finding targets has a relatively small impact on episode return for long episodes. The two next terms are bonus rewards, intended to speed up learning by further encouraging desired behavior. Actions that move the agent towards the nearest target and move the view to previously unseen regions are desirable. Importantly, the sum of these bonuses is not larger than the magnitude of the time penalty. This is to ensure that exploration and moving towards targets does not seem more important to the agent than finishing the episode quickly. Therefore, the bonuses may guide the agent towards finding targets but do not cause it to steer away from the underlying goal.

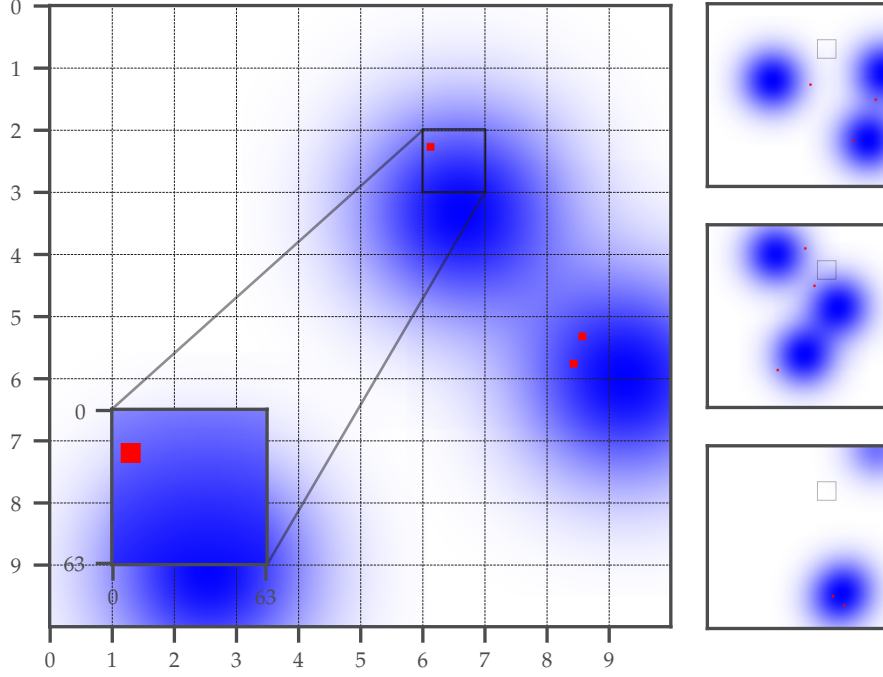


Figure 3.1: Four samples of the first environment. There are three blue bumps and three red targets in each scene. Targets are more likely where the intensity of the bumps is high.

### 3.2.2 Gaussian Environment

The first environment is the simplest environment, where the correlation between scene appearance and target probability is clear. During each episode reset, a  $1024 \times 1024$  RGB image is generated conditioned on the seed. The image contains three blue bumps, whose intensity is highest towards their center and wears off radially outwards. The intensity wears off as an (approximate) Gaussian function. The sum of the intensity in the blue channel in the image is used as a probability density function to select the location of three targets. Targets are characterized by red  $8 \times 8$  squares. The higher the intensity in the blue channel, the more likely that there is a target there.

The image is divided into a grid of  $10 \times 10$  steps, one per position  $p_t$ . The image  $x_t$  is the  $64 \times 64$  sub-image at the current position in the grid. Each moving action translates the agent one step in corresponding cardinal direction in the grid.

The idea with this environment is to test whether an agent is able to utilize scene characteristics to search quicker. An efficient searcher in this environment should prioritize locations where the intensity in the blue channel is high. It should be able to use the gradient of the underlying Gaussian function to move towards locations with higher target probabilities, while avoiding revisiting locations and intelligently planning its search path.

### 3.2.3 Terrain Environment

The second environment is similar to the first one, but intended to simulate a search scenario in realistic terrain. Actions and observations behave the same as in the first environment. At the start of each episode, a  $1024 \times 1024$  height map is generated using gradient noise. The height map is used to determine the color of an image of the same size. Lower heights are filled with blue ocean, and higher areas with green grass and brown mountains. Three red targets are located with uniform probability along the shores, between mountains and water.

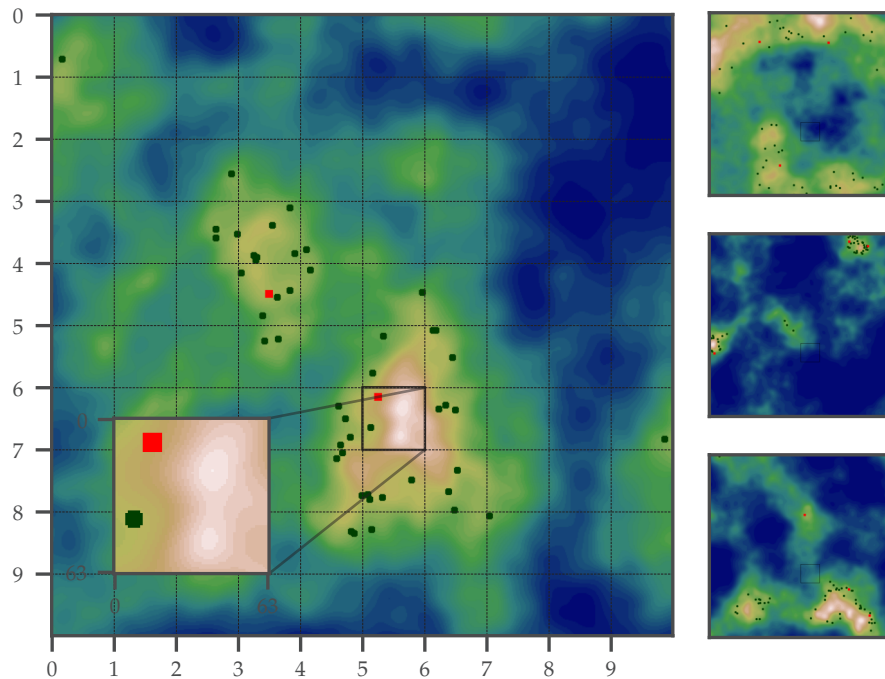


Figure 3.2: Four samples of the second environment. Terrain seen from above with red targets scattered and green trees scattered along shores.

Green trees are also scattered around each scene, whose positions are sampled from the same distribution as that of the targets.

While this environment is similar to the first environment, it is less clear how to search efficiently in it. There is higher variance between scene samples. It is also less clear how the scene appearance is correlated to the probability of targets. It is desirable that a searching agent should learn to not search oceans and mountains. Instead it should prioritize searching along the edges of land masses. For some scene samples, there are multiple islands with small patches of land between them. These patches could be used to quickly prioritize land while avoiding water. One can draw parallels to search-and-rescue scenarios with UAVs or fire detection.

### 3.2.4 Camera Environment

The third environment is a three-dimensional version of the second one. The height map is turned into a three-dimensional mesh, and the agent is placed at its center. The agent observes the scene through a pan-tilt perspective camera. Targets are, as before, placed along island edges.

The `LEFT` and `RIGHT` actions control the yaw of the camera, while `DOWN` and `UP` control its pitch. The yaw is divided into 20 steps between 0 and 360 degrees. The yaw angle wraps around, so that the agent can look around freely. Similarly, the pitch is divided into 10 steps between 0 (straight forward) and -90 degrees (straight down), but without wrapping around. This means that the camera can take 200 different positions.

Targets are always visible from at least one camera position. They may be visible from multiple positions, and the agent is expected to use the `INDICATE` action only when a target is as close to the center of the view image as possible. This environment is intended to model more realistic scenarios where the image is more difficult to interpret. In some scenarios it can be important that objects are not only localized, but localized accurately.

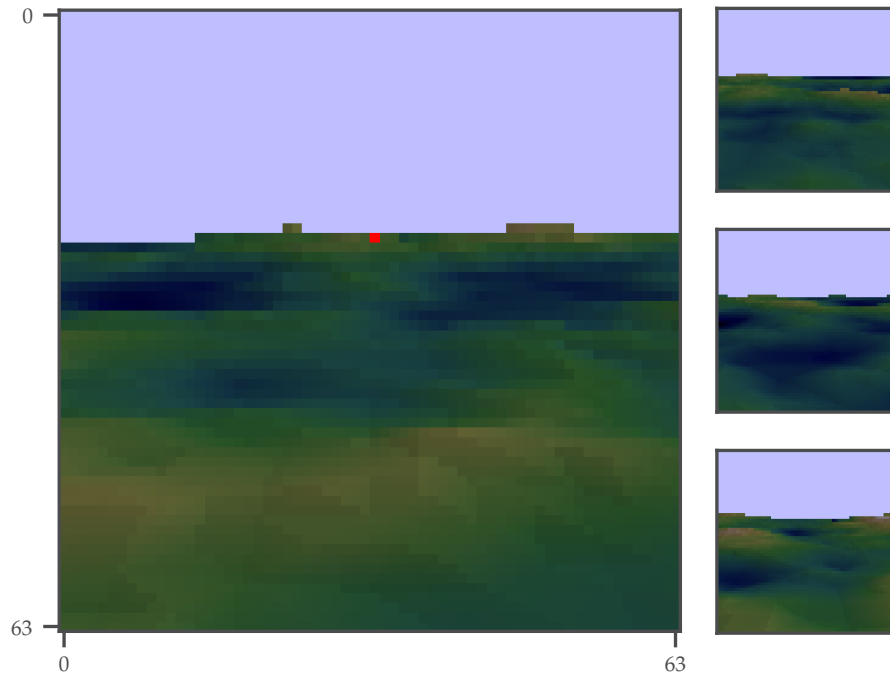


Figure 3.3: Four samples from the camera environment. Terrain seen from a pan-tilt camera. The pan and tilt of the camera can be adjusted to move the view around.

### 3.3 Approach

To design an agent that effectively solves the task, we draw inspiration from several previous works and adapt them to better suit this particular task. The final agent should be able to recognize targets, regardless of where they appear in view. It is likely important for the agent to have access to its location, especially for large search spaces. As environments are procedurally generated, a partial image observation of the scene is not sufficient to ground the agent in it. It should also be able to integrate features over time in order to remember which locations have been visited. Remembering visual features of these locations may also be of importance, as it can provide clues for what is in their proximity.

Due to the advantages described in Section 2.1.4.2, we limit our approaches to policy gradient methods. Specifically, we employ an actor-critic approach that estimates a policy and value function with a multi-headed neural network. The neural network architecture should reflect the aforementioned requirements.

#### 3.3.1 Architecture

We design our neural network architecture as follows: A CNN takes the observed image  $x_t$  and encodes it into a latent representation  $h_t$ . This allows the agent to extract translation invariant features from observed images. The latent representation, as well as the current position of the agent  $p_t$  is used as input to an RNN. Feeding both an image representation and the position of the agent to a recurrent step lets the agent remember visited locations and their appearance. The output of the recurrent step is in turn connected to an actor MLP head and a critic MLP head, which approximate the policy  $\pi$  and value function  $v$  respectively. The architecture of the neural network is presented in Figure 3.3.1.

For the CNN, we use the same architecture as [45]. The input image  $x_t \in \mathbb{R}^{3 \times 64 \times 64}$  is fed through three convolutional layers: the first layer convolves 32 filters of size  $8 \times 8$  with stride

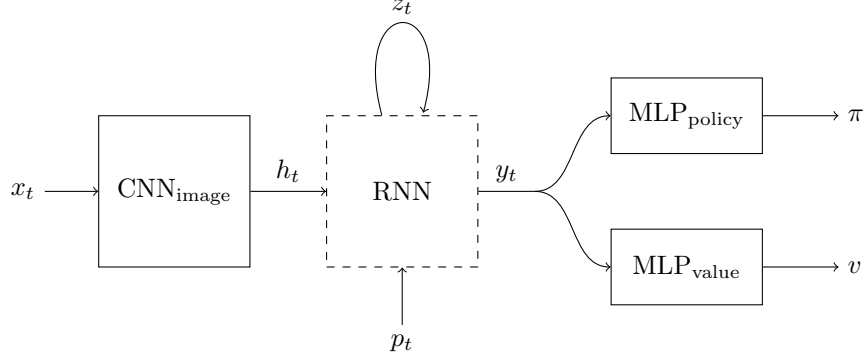


Figure 3.4: The neural network architecture used for estimating the policy and value functions from image and position observations.

4, the second convolves 64 filters of size  $4 \times 4$  with stride 2, and the third layer convolves 64 filters of size  $3 \times 3$  with stride 1. This is followed by a final fully connected hidden layer with 512 outputs for the latent representation  $h_t \in \mathbb{R}^{512}$ .

Two RNN architectures are compared: one temporal, and one spatial. While a temporal memory can retain location and appearance information over time, it is not cleared how this information is stored and whether it can be utilized properly. We hypothesize that agents using temporal memories may struggle with learning good policies in large search spaces and those that require scene understanding, such as reasoning over previously visited locations. Several results from works in embodied visual navigation indicate that spatial memories can give better results than temporal ones [49, 31, 28, 15]. Taking inspiration from this, we investigate whether a spatial memory can be more useful than a temporal one when searching for targets.

The temporal memory is a single LSTM [33] layer, as proposed by [29] and used in [42], [41], and [28]. As input to the LSTM, we use the latent image representation  $h_t$  concatenated with one-hot encodings of each dimension of  $p_t$ . The LSTM layer has 128 hidden cells so that  $y_t$  has 128 dimensions. Finally,  $z_t$  contains the hidden and cell states.

The spatial memory is composed of a readable and writable feature map, similar to those in [49], [31], [28] and [15]. It is defined by the following set of operations:

$$\rho_t = \text{CNN}_{\text{read}}(z_t) \quad (3.6)$$

$$\omega_t = \text{MLP}_{\text{write}}([h_t, \rho_t, z_t^{(p_t)}]) \quad (3.7)$$

$$q_t = \text{MLP}_{\text{position}}(\text{onehot}(p_t)) \quad (3.8)$$

$$y_t = [\rho_t, \omega_t, q_t] \quad (3.9)$$

$$z_{t+1}^{(p')} = \begin{cases} \omega_t & \text{if } p' = p_t \\ z_t^{(p')} & \text{if } p' \neq p_t \end{cases} \quad (3.10)$$

Here,  $z_t \in \mathbb{R}^{64 \times 10 \times 10}$  is a feature map with one 64-dimensional feature vector per possible position  $p_t$ . At each time step, the agent reads from this map using a CNN with three convolutional layers. Each layer has 32 filters of size  $3 \times 3$ , and uses a padding of 1. A final fully connected layer produces a 64-dimensional read vector  $\rho_t$ . An MLP takes the read vector concatenated with  $h_t$  and the feature vector stored in the map at the current position  $z_t^{(p_t)}$ , and produces a 64-dimensional write vector  $\omega_t$ . Another MLP takes a one-hot encoding of



Table 3.1: PPO hyperparameters used during training.

Parameter	Value
total time steps	$25 \times 10^6$
discount factor $\gamma$	.999
advantage factor $\lambda$	.95
parallel environments $N$	64
time steps per rollout $T$	256
epochs per rollout $K$	3
mini-batches size $M$	2048
value loss coefficient $c_1$	0.5
entropy bonus coefficient $c_2$	0.01
clip range $\epsilon$	.2
learning rate $\alpha$	$5 \times 10^{-4}$
reward normalization	yes
rate schedule	linear

the position as input and outputs a 64-dimensional vector  $q_t$ . The read vector  $\rho_t$ , the write vector  $\omega_t$ , and the position vector  $q_t$  are concatenated and used as input  $y_t$  to the value and policy networks. This means that they can make use of visual features of previously explored locations, spatial relationships between these features, and the current position of the agent. Finally, the feature map is updated so that  $z_{t+1}$  contains the write vector  $\omega_t$  at position  $p_t$ .

Illustrate map updates.

The hyperparameters for the two RNN architectures are chosen so that they have a comparable number of trainable parameters. Both the policy and value network are fully connected networks. The value network has one output for the value estimate. The policy network has 5 outputs, the logits for each action. Applying the softmax operation on this output gives the final action probabilities. All network layers have ReLU activation functions, as suggested by [30].

### 3.3.2 Training

We train both agents with PPO [52], as described in Section 2.1.4.5. Early experiments show that PPO gives good results, stable learning curves and good sample efficiency, which is in line with results reported by [5]. Furthermore, we use similar hyperparameters to [18]. These are presented in Table 3.1. Many parallel environments seem to both speed up and stabilize training. As suggested by [5], we initialize the policy output weights so that their mean is 0 and their standard deviation is low ( $10^{-2}$ ).

We normalize the reward using a moving average. This is done to limit the scale of the error derivatives. Early experiments show that not normalizing the reward destabilized learning. Similar results have been reported by [5] and [44]. The Adam optimizer [36] is used in all experiments. Finally, we decay the learning rate linearly so that it is zero at the final time step. We find that this helps the agent with finding a better local optimum.

## 3.4 Experiments

We conduct four different experiments to evaluate our approaches and answer the research questions in 1.3. Following the recommendations of [30], [20] and [1], we report mean and standard deviation across a handful of seeds.

### 3.4.1 Quality of Search Behavior

To evaluate the quality of the learned policy, we compare our approaches to a set of baselines. The first baseline acts randomly, and the second one greedily. The random baseline is detailed



in Algorithm 3.4.1, and the greedy baseline in Algorithm 3.4.1. Both automatically indicate when a target is visible. Without offloading the detection from the baselines, they would not find reasonable solutions. Furthermore, we feel that it is still an interesting comparison as it does not benefit our approach. We also investigate how the search performance of our approach compares to that of a human. Human searchers are given the same observations as the other agents.

---

**Algorithm 2** Random Baseline Policy

---

```

if there is a target at  $p_t$  then
   $a_t \leftarrow \text{INDICATE}$ 
else
   $\mathcal{A}_{\text{move}} = \{\text{LEFT}, \text{UP}, \text{RIGHT}, \text{DOWN}\}$ 
  sample  $a_t$  from  $\mathcal{A}_{\text{move}}$ 
end if

```

---



---

**Algorithm 3** Greedy Baseline Policy

---

```

if there is a target at  $p_t$  then
   $a_t \leftarrow \text{INDICATE}$ 
else
   $\mathcal{A}_{\text{move}} = \{\text{LEFT}, \text{UP}, \text{RIGHT}, \text{DOWN}\}$ 
   $\mathcal{A}_{\text{explore}} \leftarrow \{a | a \in \mathcal{A}_{\text{move}} \text{ moves the agent to an unvisited location}\}$ 
  if  $\mathcal{A}_{\text{explore}} \neq \emptyset$  then
    sample  $a_t$  from  $\mathcal{A}_{\text{explore}}$ 
  else
    sample  $a_t$  from  $\mathcal{A}_{\text{move}}$ 
  end if
end if

```

---

For each environment, approach and baseline, we measure the search path length and success rate after training. For this, we use a fixed set of 100 test samples from each environment. We report both mean and standard deviation for the 5 human searchers and three runs of our two approaches. We also report the SPL metric [3] for all agents, where the shortest path length is the optimal travel distance between the targets and the initial position of the agent. The distance between two points is computed as the minimal number of actions to transform the view between the two. Although finding the optimal path is not realistic in partially observable environments, SPL can still be useful when compared to that of a human.

### 3.4.2 Size of Search Space

The number of steps required to locate all targets in a scene is dependent on the size of the search space. In small search spaces, the difference in performance between an intelligent searcher and a less intelligent one may be difficult to quantify. Furthermore, a large scene is more difficult to search intelligently - more capacity is needed to remember visited locations and plan future steps.

To investigate how well our approach scales to different search space sizes, we train our approaches on three variants of the gaussian environment. The scene image is scaled up so that there are  $10 \times 10$ ,  $15 \times 15$  and  $20 \times 20$  possible camera positions in each environment. We train and test both approaches on the full distribution of samples from all three search space sizes respectively. By varying the search space size in each environment, we can get a feel for how the size affects training times and the quality of the learned policies. For each agent and search space size, we report how the average length and success rate changes while learning.

### 3.4.3 Number of Training Samples

In realistic scenarios, training agents on an unlimited number of samples is not possible. It is more likely that there is a limited set of samples for an agent to learn from. For this reason, it is interesting to quantify how the number of training samples affects generalization to the full distribution.

We do this by training our two agents on 500, 1000, 5000, 10000 and an unlimited number samples of the terrain environment. Both agents are tested on held out samples from the full distribution of samples, as suggested by [18]. This should illustrate how many samples are needed to generalize, and if overfitting is an issue. For each agent and training set size, we report how the average length and success rate changes while learning.

### 3.4.4 Ablations

Finally, we investigate which aspects of our architecture are important for good search behavior through a set of ablation studies. In the gaussian environment, we train three additional variants of our agent architecture. One without the recurrent step, one without image observations, and one without position observations. We compare the achieved performance to that of our full architectures.

## 3.5 Implementation

The environment is implemented with using the Gym [13] interface. The agent is implemented and RL algorithms are implemented with PyTorch [50] for automatic differentiation of the computation graphs. PPO [52] was implemented following the official implementation by OpenAI, and verified by testing on the benchmarks used in the original paper. All experiments are conducted on an Intel Core i9-10900X CPU and an NVIDIA GeForce RTX 2080 Ti GPU. The source code for environments, models and algorithms is available at <https://gitlab.liu.se/osklu414/tqdt33-masters-thesis>.



## 4 Results

This chapter presents the results for each of the experiments described in Section 3.4.

### 4.1 Quality of Search Behavior

Table 4.1 shows the average search path length, success rate and SPL metric on a fixed set of 100 levels from each environment. These metrics are presented for our approaches trained on the full distribution of environments, as well as baselines. Human results are collected from 5 individuals.

The random and greedy baselines are consistently worse than the other agents. The temporal memory agent seems to...

### 4.2 Size of Search Space

The results of the search space experiments in the gaussian environment are presented in Figure 4.2. Results were collected across four different runs, and the plots show the mean and standard deviation. For the search space of  $10 \times 10$ , both architectures initially improve their policy quickly. Past a certain time step, they keep improving at a reduced pace. At the end of training, the spatial memory architecture has reached a policy that seems to find targets quicker than the temporal memory. Both seem to be able to find the all three targets in every episode.

For the larger search space sizes with  $15 \times 15$  and  $20 \times 20$  the difference between the two architectures is greater. While the spatial memory seems to consistently find targets in a number of steps that is comparable to the number of positions in the search space, the search paths of the agent with the temporal memory are substantially longer.

### 4.3 Number of Training Samples

Figure 4.3 shows how the average length and success rate in the terrain environment is affected by the number of samples seen during training. These metrics are presented for the limited training set and unlimited testing set respectively.

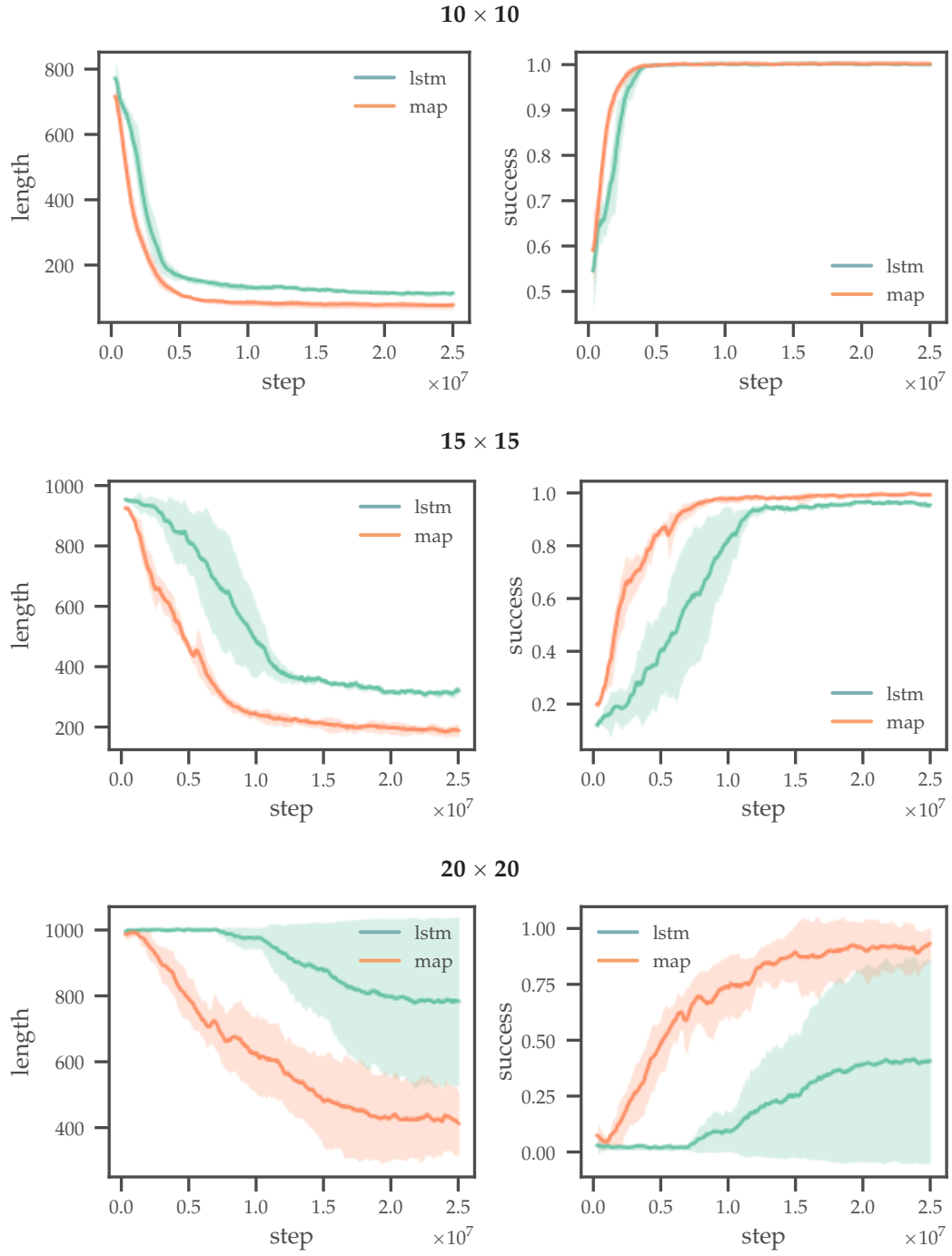


Figure 4.1: Reward and episode length curves during training for three different search space sizes. Mean and standard deviation across 4 runs.

Table 4.1: Average search path length, success rate and SPL metric on a fixed set of a 100 samples from each environment. Metrics for spatial memory, temporal memory, random baseline, greedy baseline and human searchers.

Gaussian			
Agent	SPL	Length	Success
random	$0.07 \pm 0.01$	$417.69 \pm 30.87$	$0.93 \pm 0.02$
greedy	$0.15 \pm 0.00$	$148.69 \pm 3.41$	$1.00 \pm 0.01$
exhaustive	$0.21 \pm 0.00$	$82.10 \pm 3.64$	$1.00 \pm 0.00$
human	$0.23 \pm 0.03$	$80.97 \pm 13.49$	$1.00 \pm 0.00$
lstm	$0.25 \pm 0.02$	$108.26 \pm 12.67$	$0.99 \pm 0.01$
map	$0.30 \pm 0.02$	$74.25 \pm 12.55$	$1.00 \pm 0.01$

Terrain			
Agent	SPL	Length	Success
random	$0.06 \pm 0.01$	$434.67 \pm 51.75$	$0.91 \pm 0.02$
greedy	$0.18 \pm 0.00$	$136.30 \pm 3.85$	$1.00 \pm 0.00$
exhaustive	$0.22 \pm 0.00$	$82.47 \pm 2.63$	$1.00 \pm 0.00$
human	$0.26 \pm 0.02$	$76.73 \pm 5.33$	$1.00 \pm 0.00$
lstm	$0.25 \pm 0.01$	$103.09 \pm 8.58$	$1.00 \pm 0.01$
map	$0.27 \pm 0.00$	$86.04 \pm 10.83$	$1.00 \pm 0.01$

Camera			
Agent	SPL	Length	Success

It seems like both architectures can overfit to training sets of as many as 1000 samples, For smaller training set sizes, search path lengths on the test get gradually worse past a certain time step. Interestingly the LSTM architecture seems to overfit more severely to small training sets. Its performance on the test set decreases substantially past a certain time step. The spatial memory architecture does not As many as 10000 training samples are needed to generalize to the full distribution of scenes and achieve equal training and test performance.

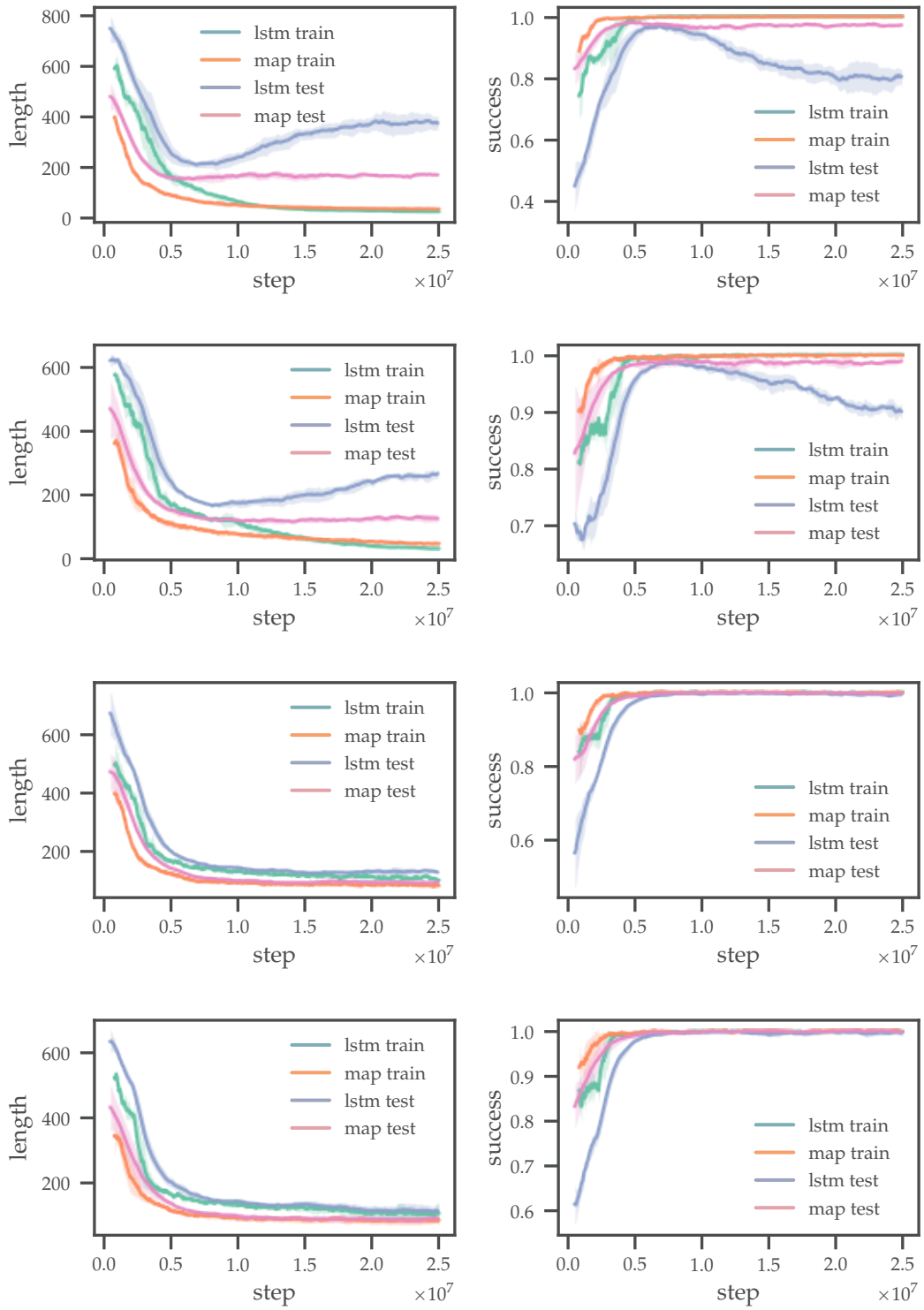


Figure 4.2: Reward and episode length curves during training for three training set sizes. Mean and standard deviation across 3 seeds.



## 5 Discussion

This chapter contains the following sub-headings.

### 5.1 Results

### 5.2 Method

- Compare and analyze results.
- Evaluate generalization and viability for real world training sets.
- Discuss advantages and disadvantages of RL for this task.
- The less bias we introduce, the more general the method has the potential to be.
- However, we could make the observation space clearer (for example, give the agent visited position directly).
- Search space size and its impact on performance.
- Larger search space could correspond to larger area or higher granularity.
- Is it better than exhaustive search?
- Is it better than a human?
- Not just visual search – also a practical example of RL for exploration, generalization, ...
- Stacked LSTM was unstable, even with dropout. Could it have remembered more?
- Differences in number of weights for different environment sizes. ...
- Comparison to exploration problems, other solutions methods?
- Camera movements in three dimensions: method could be expanded with higher-dimensional convolutions.

- Discuss if indication is necessary
- Look into epsilon greedy
- Can gamma be used for quick exploration
- Rephrase target?
- Could give bad vibes...
- Greedy and random baselines have unfair advantage, but this is fine since we are not putting our method at an advantage.
- The map approach seems to scale better - with more time we could have trained for a larger search space.
- Neither scales well (weights in map scale increase quadratically with search space size – fine for fixed positions, bad for movable cameras).
- If we find that bonuses speed up training but converge to poor solutions, we could discuss reward shaping. Remove bonus after convergence.
- The scene characteristics matter - in some cases, amortized probability of targets is not uniform, and patterns are more complex. We need more types of environments.
- Discuss the generality of the proposed approach - what types of cues can it pick up? Spatial and semantic relationships, non-uniform probabilities, etc.
- An ego-centric architecture might be more general, does not need the position and can work over larger territories (although not take whole territory into account).
- Show some search paths and discuss them?
- The more specialized the reward, the higher the risk that it is good for some environments and bad for some. We saw with the unspecialized reward that it did not work well for the terrain environment. If we end up using only a specialized reward, discuss cases when it would not be able to find an optimal solution.
- Discuss lack of completion action (as suggested in related work).
- Talk about optimal paths and SPL (not realistic).
- Should have looked at more interesting environments (at least for optimal paths).
- Hopefully the camera environment is interesting.
- Early experiments with stacking multiple LSTM layers lead to heavy overfitting and unstable learning curves.
- The difficulty of each environment is related to the potential region targets can appear in. We can measure this in the environments themselves.
- We can visualize attention as in "A critical investigation of deep reinforcement learning for navigation"
- We have illustrated that reinforcement learning agents overfit to training samples and must be tested on separate sets.

### 5.3 The work in a wider context





## 6 Conclusion

In this work, we have presented...

Bajcsy, Aloimonos and Tsotsos [10] connect past work in active vision with recent advances in robotics, artificial intelligence and computer vision. They argue that a complete artificial agent must include active perception. The ultimate goal of artificial intelligence research is the computational generation of intelligent behavior. Agents that choose their behavior based on their context and know why they behave as they do would certainly seem to embody this.

An agent is an active perceiver if it knows why it wishes to sense, and then chooses what to perceive, and determines how, when and where to achieve that perception.



# A Training



## Bibliography

- [1] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G. Bellemare. “Deep Reinforcement Learning at the Edge of the Statistical Precipice”. In: *arXiv:2108.13264 [cs, stat]* (Jan. 2022). arXiv: 2108.13264. URL: <http://arxiv.org/abs/2108.13264> (Cited on pages 16, 23).
- [2] John Aloimonos, Isaac Weiss, and Amit Bandyopadhyay. “Active vision”. In: *International Journal of Computer Vision* 1.4 (Jan. 1988). 705 citations (Crossref) [2022-02-07], pp. 333–356. ISSN: 0920-5691, 1573-1405. DOI: 10/cn4mdc. URL: <http://link.springer.com/10.1007/BF00133571> (visited on 02/07/2022) (Cited on page 5).
- [3] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, and Amir R. Zamir. “On Evaluation of Embodied Navigation Agents”. In: *arXiv:1807.06757 [cs]* (July 2018). arXiv: 1807.06757. URL: <http://arxiv.org/abs/1807.06757> (Cited on pages 13, 16, 24).
- [4] Alexander Andreopoulos and John K Tsotsos. “A Theory of Active Object Localization”. en. In: (), p. 8 (Cited on page 5).
- [5] Marcin Andrychowicz, Anton Raichuk, Piotr Stańczyk, Manu Orsini, Sertan Girgin, Raphael Marinier, Léonard Hussenot, Matthieu Geist, Olivier Pietquin, Marcin Michalski, Sylvain Gelly, and Olivier Bachem. “What Matters In On-Policy Reinforcement Learning? A Large-Scale Empirical Study”. In: *arXiv:2006.05990 [cs, stat]* (June 2020). arXiv: 2006.05990. URL: <http://arxiv.org/abs/2006.05990> (Cited on pages 10, 23).
- [6] A. Aydemir, K. Sjøo, J. Folkesson, A. Pronobis, and P. Jensfelt. “Search in the real world: Active visual object search based on spatial relations”. en. In: *2011 IEEE International Conference on Robotics and Automation*. 48 citations (Crossref) [2022-02-28]. Shanghai, China: IEEE, May 2011, pp. 2818–2824. ISBN: 978-1-61284-386-5. DOI: 10.1109/ICRA.2011.5980495. URL: <http://ieeexplore.ieee.org/document/5980495/> (Cited on page 11).
- [7] Alper Aydemir, Andrzej Pronobis, Moritz Göbelbecker, and Patric Jensfelt. “Active Visual Object Search in Unknown Environments Using Uncertain Semantics”. In: *IEEE Transactions on Robotics* 29.4 (Aug. 2013). 65 citations (Crossref) [2022-02-28], pp. 986–1002. ISSN: 1941-0468. DOI: 10.1109/TRO.2013.2256686 (Cited on page 11).

- [8] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. “Neural Machine Translation by Jointly Learning to Align and Translate”. In: *arXiv:1409.0473 [cs, stat]* (May 2016). arXiv: 1409.0473. URL: <http://arxiv.org/abs/1409.0473> (Cited on page 14).
- [9] R. Bajcsy. “Active perception”. In: *Proceedings of the IEEE* 76.8 (Aug. 1988). 646 citations (Crossref) [2022-02-28] Conference Name: Proceedings of the IEEE, pp. 966–1005. ISSN: 1558-2256. DOI: 10.1109/5.5968 (Cited on page 5).
- [10] Ruzena Bajcsy, Yiannis Aloimonos, and John K. Tsotsos. “Revisiting active perception”. In: *Autonomous Robots* 42.2 (Feb. 1, 2018). 86 citations (Crossref) [2022-03-13], pp. 177–196. ISSN: 1573-7527. DOI: 10.1007/s10514-017-9615-3. URL: <https://doi.org/10.1007/s10514-017-9615-3> (visited on 03/13/2022) (Cited on page 32).
- [11] Yoshua Bengio, Aaron Courville, and Pascal Vincent. “Representation Learning: A Review and New Perspectives”. In: *arXiv:1206.5538 [cs]* (Apr. 2014). arXiv: 1206.5538. URL: <http://arxiv.org/abs/1206.5538> (Cited on page 5).
- [12] Francisco Bonin-Font, Alberto Ortiz, and Gabriel Oliver. “Visual Navigation for Mobile Robots: A Survey”. In: *Journal of Intelligent and Robotic Systems* 53.3 (Nov. 2008), pp. 263–296. ISSN: 0921-0296, 1573-0409. DOI: 10.1007/s10846-008-9235-4 (Cited on page 12).
- [13] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. “OpenAI Gym”. In: *arXiv:1606.01540 [cs]* (June 2016). arXiv: 1606.01540. URL: <http://arxiv.org/abs/1606.01540> (Cited on page 25).
- [14] Juan C. Caicedo and Svetlana Lazebnik. “Active Object Localization with Deep Reinforcement Learning”. In: *arXiv:1511.06015 [cs]* (Nov. 18, 2015). arXiv: 1511.06015. URL: <http://arxiv.org/abs/1511.06015> (visited on 02/03/2022) (Cited on pages 2, 11–13).
- [15] Devendra Singh Chaplot, Dhiraj Gandhi, Abhinav Gupta, and Ruslan Salakhutdinov. “Object Goal Navigation using Goal-Oriented Semantic Exploration”. In: *arXiv:2007.00643 [cs]* (July 2020). arXiv: 2007.00643. URL: <http://arxiv.org/abs/2007.00643> (Cited on pages 14, 22).
- [16] Shengyong Chen, Youfu Li, and Ngai Ming Kwok. “Active vision in robotic systems: A survey of recent developments”. In: *The International Journal of Robotics Research* 30.11 (Sept. 2011), pp. 1343–1377. ISSN: 0278-3649. DOI: 10.1177/0278364911410755 (Cited on pages 5, 11).
- [17] Xinlei Chen and Abhinav Gupta. “Spatial Memory for Context Reasoning in Object Detection”. In: *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2017, pp. 4106–4116. ISBN: 978-1-5386-1032-9. DOI: 10.1109/ICCV.2017.440. URL: <http://ieeexplore.ieee.org/document/8237702/> (Cited on page 12).
- [18] Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. “Leveraging Procedural Generation to Benchmark Reinforcement Learning”. In: *arXiv:1912.01588 [cs, stat]* (July 2020). arXiv: 1912.01588. URL: <http://arxiv.org/abs/1912.01588> (Cited on pages 10, 15, 18, 23, 25).
- [19] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. “Quantifying Generalization in Reinforcement Learning”. In: *arXiv:1812.02341 [cs, stat]* (July 2019). arXiv: 1812.02341. URL: <http://arxiv.org/abs/1812.02341> (Cited on page 15).
- [20] Cédric Colas, Olivier Sigaud, and Pierre-Yves Oudeyer. “A Hitchhiker’s Guide to Statistical Comparisons of Reinforcement Learning Algorithms”. en. In: *arXiv:1904.06979 [cs, stat]* (Apr. 2019). arXiv: 1904.06979. URL: <http://arxiv.org/abs/1904.06979> (Cited on pages 16, 23).

- [21] Cédric Colas, Olivier Sigaud, and Pierre-Yves Oudeyer. “How Many Random Seeds? Statistical Power Analysis in Deep Reinforcement Learning Experiments”. en. In: *arXiv:1806.08295 [cs, stat]* (July 2018). arXiv: 1806.08295. URL: <http://arxiv.org/abs/1806.08295> (Cited on page 16).
- [22] Vikas Dhiman, Shurjo Banerjee, Brent Griffin, Jeffrey M. Siskind, and Jason J. Corso. “A Critical Investigation of Deep Reinforcement Learning for Navigation”. In: *arXiv:1802.02274 [cs]* (Jan. 2019). arXiv: 1802.02274. URL: <http://arxiv.org/abs/1802.02274> (Cited on page 13).
- [23] M. P. Eckstein. “Visual search: A retrospective”. In: *Journal of Vision* 11.5 (Dec. 30, 2011). 207 citations (Crossref) [2022-02-28], pp. 14–14. ISSN: 1534-7362. DOI: 10.1167/11.5.14. URL: <http://jov.arvojournals.org/Article.aspx?doi=10.1167/11.5.14> (visited on 02/22/2022) (Cited on pages 1, 5).
- [24] Per-Erik Forssen, David Meger, Kevin Lai, Scott Helmer, James J. Little, and David G. Lowe. “Informed visual search: Combining attention and object recognition”. In: *2008 IEEE International Conference on Robotics and Automation*. May 2008, pp. 935–942. DOI: 10.1109/ROBOT.2008.4543325 (Cited on pages 2, 10).
- [25] Enric Galceran and Marc Carreras. “A survey on coverage path planning for robotics”. In: *Robotics and Autonomous Systems* 61.12 (Dec. 2013), pp. 1258–1276. ISSN: 09218890. DOI: 10/f5j2n5 (Cited on page 5).
- [26] Florin C. Ghesu, Bogdan Georgescu, Tommaso Mansi, Dominik Neumann, Joachim Hornegger, and Dorin Comaniciu. “An Artificial Agent for Anatomical Landmark Detection in Medical Images”. In: *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2016*. Ed. by Sebastien Ourselin, Leo Joskowicz, Mert R. Sabuncu, Gozde Unal, and William Wells. Vol. 9902. Lecture Notes in Computer Science. Springer International Publishing, 2016, pp. 229–237. ISBN: 978-3-319-46725-2. DOI: 10.1007/978-3-319-46726-9\_27. URL: [https://link.springer.com/10.1007/978-3-319-46726-9\\_27](https://link.springer.com/10.1007/978-3-319-46726-9_27) (Cited on pages 2, 12, 13).
- [27] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. Adaptive Computation and Machine Learning series. MIT Press, Nov. 2016. ISBN: 978-0-262-03561-3 (Cited on pages 2, 6, 7).
- [28] Saurabh Gupta, Varun Tolani, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. “Cognitive Mapping and Planning for Visual Navigation”. In: *arXiv:1702.03920 [cs]* (Feb. 2019). arXiv: 1702.03920. URL: <http://arxiv.org/abs/1702.03920> (Cited on pages 14, 22).
- [29] Matthew Hausknecht and Peter Stone. “Deep Recurrent Q-Learning for Partially Observable MDPs”. In: *arXiv:1507.06527 [cs]* (Jan. 2017). arXiv: 1507.06527. URL: <http://arxiv.org/abs/1507.06527> (Cited on pages 13, 14, 22).
- [30] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. “Deep Reinforcement Learning That Matters”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 32.11 (Apr. 2018). ISSN: 2374-3468. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/11694> (Cited on pages 10, 16, 23).
- [31] Joao F. Henriques and Andrea Vedaldi. “MapNet: An Allocentric Spatial Memory for Mapping Environments”. In: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, June 2018, pp. 8476–8484. ISBN: 978-1-5386-6420-9. DOI: 10.1109/CVPR.2018.00884. URL: <https://ieeexplore.ieee.org/document/8578982/> (Cited on pages 14, 22).
- [32] Matteo Hessel, Hado van Hasselt, Joseph Modayil, and David Silver. “On Inductive Biases in Deep Reinforcement Learning”. In: *arXiv:1907.02908 [cs, stat]* (July 2019). arXiv: 1907.02908. URL: <http://arxiv.org/abs/1907.02908> (Cited on page 15).

- [33] Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-Term Memory”. In: *Neural Computation* 9.8 (Nov. 1997), pp. 1735–1780. ISSN: 0899-7667. DOI: 10.1162/neco.1997.9.8.1735 (Cited on pages 7, 22).
- [34] Laurent Itti and Christof Koch. “Computational modelling of visual attention”. In: *Nature Reviews Neuroscience* 2.33 (Mar. 2001), pp. 194–203. ISSN: 1471-0048. DOI: 10.1038/35058500 (Cited on pages 4, 5).
- [35] Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. “Planning and acting in partially observable stochastic domains”. In: *Artificial Intelligence* 101.1–2 (May 1998), pp. 99–134. ISSN: 00043702. DOI: 10.1016/S0004-3702(98)00023-X (Cited on pages 7, 8).
- [36] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *arXiv:1412.6980 [cs]* (Jan. 2017). arXiv: 1412.6980. URL: <http://arxiv.org/abs/1412.6980> (Cited on page 23).
- [37] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. “A Survey of Generalisation in Deep Reinforcement Learning”. In: *arXiv:2111.09794 [cs]* (Jan. 2022). arXiv: 2111.09794. URL: <http://arxiv.org/abs/2111.09794> (Cited on pages 14, 15, 17).
- [38] Maja J Mataric. “Reward Functions for Accelerated Learning”. In: *Machine Learning Proceedings 1994*. Ed. by William W. Cohen and Haym Hirsh. Morgan Kaufmann, Jan. 1994, pp. 181–189. ISBN: 978-1-55860-335-6. DOI: 10.1016/B978-1-55860-335-6.50030-1. URL: <https://www.sciencedirect.com/science/article/pii/B9781558603356500301> (Cited on page 9).
- [39] Marvin Minsky. “Steps toward Artificial Intelligence”. In: *Proceedings of the IRE* 49.1 (Jan. 1961), pp. 8–30. ISSN: 2162-6634. DOI: 10.1109/JRPROC.1961.287775 (Cited on page 9).
- [40] Silviu Minut and Sridhar Mahadevan. “A reinforcement learning model of selective visual attention”. In: *Proceedings of the fifth international conference on Autonomous agents - AGENTS '01*. ACM Press, 2001, pp. 457–464. ISBN: 978-1-58113-326-4. DOI: 10/dbwckq. URL: <http://portal.acm.org/citation.cfm?doid=375735.376414> (Cited on pages 2, 11).
- [41] Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andrew J. Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, Dharmashan Kumaran, and Raia Hadsell. “Learning to Navigate in Complex Environments”. In: *arXiv:1611.03673 [cs]* (Jan. 2017). arXiv: 1611.03673. URL: <http://arxiv.org/abs/1611.03673> (Cited on pages 2, 13, 22).
- [42] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. “Asynchronous Methods for Deep Reinforcement Learning”. In: *arXiv:1602.01783 [cs]* (June 2016). arXiv: 1602.01783. URL: <http://arxiv.org/abs/1602.01783> (Cited on pages 13, 14, 18, 22).
- [43] Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. “Recurrent Models of Visual Attention”. In: *arXiv:1406.6247 [cs, stat]* (June 2014). arXiv: 1406.6247. URL: <http://arxiv.org/abs/1406.6247> (Cited on pages 2, 11, 14).
- [44] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. “Playing Atari with Deep Reinforcement Learning”. In: *arXiv:1312.5602 [cs]* (Dec. 2013). arXiv: 1312.5602. URL: <http://arxiv.org/abs/1312.5602> (Cited on pages 9, 23).



- [45] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. "Human-level control through deep reinforcement learning". In: *Nature* 518.75407540 (Feb. 2015), pp. 529–533. ISSN: 1476-4687. DOI: 10.1038/nature14236 (Cited on pages 2, 9, 13, 14, 21).
- [46] Ken Nakayama and Paolo Martini. "Situating visual search". In: *Vision Research*. Vision Research 50th Anniversary Issue: Part 2 51.13 (July 2011), pp. 1526–1537. ISSN: 0042-6989. DOI: 10.1016/j.visres.2010.09.003 (Cited on pages 1, 4).
- [47] Junhyuk Oh, Valliappa Chockalingam, Satinder Singh, and Honglak Lee. "Control of Memory, Active Perception, and Action in Minecraft". In: *arXiv:1605.09128 [cs]* (May 2016). arXiv: 1605.09128. URL: <http://arxiv.org/abs/1605.09128> (Cited on page 14).
- [48] Fabio Pardo, Arash Tavakoli, Vitaly Levnik, and Petar Kormushev. "Time Limits in Reinforcement Learning". In: *arXiv:1712.00378 [cs]* (Jan. 2022). arXiv: 1712.00378. URL: <http://arxiv.org/abs/1712.00378> (Cited on page 18).
- [49] Emilio Parisotto and Ruslan Salakhutdinov. "Neural Map: Structured Memory for Deep Reinforcement Learning". In: *arXiv:1702.08360 [cs]* (Feb. 2017). arXiv: 1702.08360. URL: <http://arxiv.org/abs/1702.08360> (Cited on pages 14, 22).
- [50] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. "PyTorch: An Imperative Style, High-Performance Deep Learning Library". In: (), p. 12 (Cited on page 25).
- [51] Stuart J. Russell and Peter Norvig. *Artificial intelligence: a modern approach*. In collab. with Ming-wei Chang, Jacob Devlin, Anca Dragan, David Forsyth, Ian Goodfellow, Jitendra Malik, Vikash Mansinghka, Judea Pearl, and Michael Woolridge. Fourth Edition. Pearson Series in Artificial Intelligence. Hoboken, NJ: Pearson, 2021. 1115 pp. ISBN: 978-0-13-461099-3 (Cited on pages 6, 7).
- [52] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. "Proximal Policy Optimization Algorithms". In: *arXiv:1707.06347 [cs]* (Aug. 2017). arXiv: 1707.06347. URL: <http://arxiv.org/abs/1707.06347> (Cited on pages 10, 23, 25).
- [53] Ksenia Shubina and John K. Tsotsos. "Visual search for an object in a 3D environment using a mobile robot". In: *Computer Vision and Image Understanding*. Special issue on Intelligent Vision Systems 114.5 (May 2010), pp. 535–547. ISSN: 1077-3142. DOI: 10.1016/j.cviu.2009.06.010 (Cited on pages 2, 11).
- [54] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. "Mastering the game of Go with deep neural networks and tree search". In: *Nature* 529.75877587 (Jan. 2016), pp. 484–489. ISSN: 1476-4687. DOI: 10.1038/nature16961 (Cited on page 2).
- [55] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. "Policy Gradient Methods for Reinforcement Learning with Function Approximation". In: *Advances in Neural Information Processing Systems*. Ed. by S. Solla, T. Leen, and K. Müller. Vol. 12. MIT Press, 1999. URL: <https://proceedings.neurips.cc/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf> (Cited on page 8).

- [56] Richard S. Sutton and Andrew G. Barto. *Reinforcement learning: an introduction*. Second edition. Adaptive computation and machine learning series. Cambridge, Massachusetts: The MIT Press, 2018. 526 pp. ISBN: 978-0-262-03924-6 (Cited on pages 2, 7–9).
- [57] Gerald Tesauro et al. “Temporal difference learning and TD-Gammon”. In: *Communications of the ACM* 38.3 (1995), pp. 58–68 (Cited on page 9).
- [58] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. In: *Nature* 575.77827782 (Nov. 2019), pp. 350–354. ISSN: 1476-4687. DOI: 10.1038/s41586-019-1724-z (Cited on pages 2, 10).
- [59] Jeremy M. Wolfe. “Visual search”. In: *Current biology : CB* 20.8 (Apr. 27, 2010). 64 citations (Crossref) [2022-03-02] Publisher: NIH Public Access, R346. DOI: 10.1016/j.cub.2010.02.016. URL: <https://www.ncbi.nlm.nih.gov/labs/pmc/articles/PMC5678963/> (visited on 03/02/2022) (Cited on pages 1, 4).
- [60] Jeremy M. Wolfe and Todd S. Horowitz. “Five factors that guide attention in visual search”. In: *Nature Human Behaviour* 1.3 (Mar. 8, 2017). 300 citations (Crossref) [2022-02-28] Number: 3 Publisher: Nature Publishing Group, pp. 1–8. ISSN: 2397-3374. DOI: 10.1038/s41562-017-0058. URL: <https://www.nature.com/articles/s41562-017-0058> (visited on 02/28/2022) (Cited on pages 1, 5).
- [61] Wei Yang, Xiaolong Wang, Ali Farhadi, Abhinav Gupta, and Roozbeh Mottaghi. “Visual Semantic Navigation using Scene Priors”. In: *arXiv:1810.06543 [cs]* (Oct. 2018). arXiv: 1810.06543. URL: <http://arxiv.org/abs/1810.06543> (Cited on page 13).
- [62] Xin Ye, Zhe Lin, Haoxiang Li, Shibin Zheng, and Yezhou Yang. “Active Object Perceiver: Recognition-Guided Policy Learning for Object Searching on Mobile Robots”. In: *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). ISSN: 2153-0866. Oct. 2018, pp. 6857–6863. DOI: 10.1109/IROS.2018.8593720 (Cited on page 12).
- [63] Yiming Ye and John K. Tsotsos. “A Complexity-Level Analysis of the Sensor Planning Task for Object Search”. en. In: *Computational Intelligence* 17.4 (Nov. 2001). 13 citations (Crossref) [2022-02-28], pp. 605–620. ISSN: 0824-7935, 1467-8640. DOI: 10.1111/0824-7935.00166 (Cited on pages 5, 11).
- [64] Chiyuan Zhang, Oriol Vinyals, Remi Munos, and Samy Bengio. “A Study on Overfitting in Deep Reinforcement Learning”. In: *arXiv:1804.06893 [cs, stat]* (Apr. 2018). arXiv: 1804.06893. URL: <http://arxiv.org/abs/1804.06893> (Cited on page 14).
- [65] Zhong-Qiu Zhao, Peng Zheng, Shou-Tao Xu, and Xindong Wu. “Object Detection With Deep Learning: A Review”. In: *IEEE Transactions on Neural Networks and Learning Systems* 30.11 (Nov. 2019), pp. 3212–3232. ISSN: 2162-237X, 2162-2388. DOI: 10.1109/TNNLS.2018.2876865 (Cited on page 11).



- [66] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J. Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. "Target-driven Visual Navigation in Indoor Scenes using Deep Reinforcement Learning". In: *arXiv:1609.05143 [cs]* (Sept. 16, 2016). arXiv: 1609.05143. URL: <http://arxiv.org/abs/1609.05143> (visited on 03/14/2022) (Cited on pages 2, 12).