

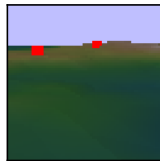
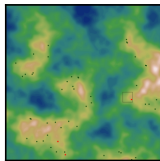
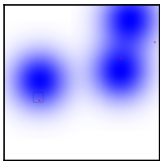
# Learning to Search for Targets

with Deep Reinforcement Learning

Oskar Lundin

Linköping University

May 19, 2022



# Outline

# Problem Description

Autonomous search for a set of targets in an scene with a camera.

- ▶ Limited region of scene visible at any given time.
- ▶ Camera can be moved to change visible region.
- ▶ Locate targets by bringing them into view and indicating that they are visible.
- ▶ Should locate all targets while minimizing the number of actions.
- ▶ Applications in search and rescue, fire detection, surveillance, etc.

- ▶ In a small or random scene with uniformly distributed targets, random or exhaustive search is sufficient.
- ▶ Most real-world search tasks are not random, but exhibit structure.
- ▶ Cues in the searched scene can be used to find targets quicker.
  - ▶ Books are in bookshelves.
  - ▶ Cars can be found on roads.
  - ▶ Some targets spread out/close together.
- ▶ Patterns and cues may be subtle and difficult to pick up.
- ▶ Manually engineering a searching system with domain knowledge be difficult and costly.
- ▶ **Can a system learn to search intelligently from a set of samples and generalize to similar search tasks?**

# Challenges

- ▶ Prioritize regions with high probability of targets based on previous experience.
- ▶ Learn correlations between scene appearance and target probability.
- ▶ Search exhaustively while avoiding searching the same region twice.
- ▶ Real-world tasks have limited number of training samples.

# Research Questions

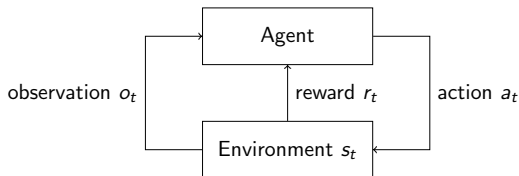
1. How can an agent that learns to intelligently search for targets be implemented with reinforcement learning?
2. How does the learning agent compare to random walk, exhaustive search, and a human searcher with prior knowledge of the characteristics of the searched scenes?
3. How does the agent's ability to generalize to unseen in-distribution scenes depend on the number of training samples?

# Markov Decision Process (MDP)

Framework for modeling decision making in partly random processes.

In our case, *partially observable* MDP [?]:

- ▶ *Agent* interacts with *environment* over discrete time steps  $t = 0, 1, 2, \dots, T$ .
- ▶ Takes *action*  $a_t$  in state  $s_t$ .
- ▶ Perceives (partial) *observation* of state  $o_t$ .
- ▶ Receives scalar reward  $r_t$  that indicates whether action is good or bad.
- ▶ New state  $s_{t+1}$  depends only on history of interactions.
- ▶ Agent usually maintains some internal state depending on history  
→ memory.



# Reinforcement Learning (RL)

Paradigm for learning from interactions how to achieve a goal.

- ▶ Tasks usually formalized as (partially observable) MDPs.
- ▶ Policy  $\pi(a|s)$  is a mapping from states to actions.
- ▶ Find  $\pi$  that maximizes cumulative reward  $\mathbb{E} \left[ \sum_{k=0}^T \gamma^{k-t-1} r_k \right]$ .
- ▶ Often involves estimating the value  $v_{\pi}(s)$  of a state under policy  $\pi$  (useful for training).

Deep RL: Approximate  $\pi$  (and  $v_{\pi}$ ) with deep neural networks. Has been used to play Atari [?], Go [?], StarCraft II [?], etc.



# Search with Reinforcement Learning

- ▶ Object localization ( $[?, ?, ?]$ ).
- ▶ Visual navigation (...).
- ▶ Todo: add more related work.

# Problem Formulation

- ▶ Agent searches scene  $S \subset \mathbb{R}^d$ .
- ▶ Scene contains set of targets  $\{t_0, \dots, t_n\}$ ,  $t_i \in S$ .
- ▶ Agent perceives view  $V \subset S$ .
- ▶ Move actions transform view to new subspace.
- ▶ Trigger action indicates that a target is in view.
- ▶ Select actions that maximize the probability of finding all targets while minimizing cost in time.
- ▶ NP complete [?], intractable to solve optimally.

# Environments

- ▶ Three environments with varying characteristics.
- ▶ Search space discretized into  $10 \times 10$  camera positions.
- ▶ Each camera position has a unique view  $V \subset S$ .
- ▶ Three targets in all scenes.
- ▶ Target probability correlated with scene appearance.
- ▶ Should be possible to do better than exhaustive search on average.
- ▶ Scenes procedurally generated:
  - ▶ Pseudorandom seed determines scene appearance and target positions.
  - ▶ Gives control over difficulty to solve.
  - ▶ Can vary training and test set sizes.

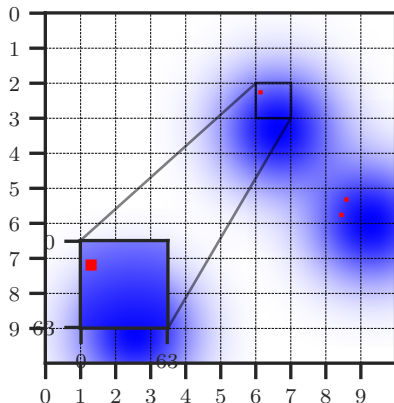
# Observation, Action and Reward

At each time step  $t$ :

- ▶ The agent receives observation  $o_t = \langle x_t, p_t \rangle$ , where
  - ▶  $x_t \in \mathbb{R}^{3 \times 64 \times 64}$  is an RGB image of current view, and
  - ▶  $p_t \in \{0, \dots, 9\} \times \{0, \dots, 9\}$  is the position of the camera.
- ▶ Takes action  $a_t \in \{\text{TRIGGER}, \text{UP}, \text{DOWN}, \text{LEFT}, \text{RIGHT}\}$ , where
  - ▶ TRIGGER indicates that a target is in view, and
  - ▶ UP, DOWN, LEFT, RIGHT move the view in each cardinal direction.
- ▶ Receives reward  $r_t = h - 0.001$  where  $h = |\mathcal{T} \cap V|$  is the number of targets in view.
  - ▶ Rewarded for finding targets.
  - ▶ Constant penalty encourages quick episode completion.

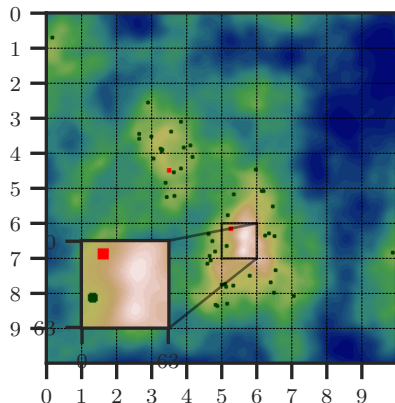
# Gaussian Environment

- ▶ 2D scene.
- ▶ Three gaussian kernels with random center.
- ▶ Sum of kernels determine appearance of scene and probability of targets.
- ▶ Clear correlation between appearance and desired behavior.



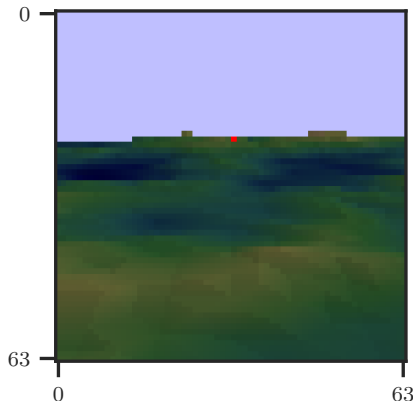
# Terrain Environment

- ▶ Similar to previous environment.
- ▶ Terrain seen from above.
- ▶ Gradient noise used to generate height map.
- ▶ Color determined by height.
- ▶ Targets placed with uniform probability across coastlines.
- ▶ More realistic, higher variance.
- ▶ Analogous to search and rescue with UAV.



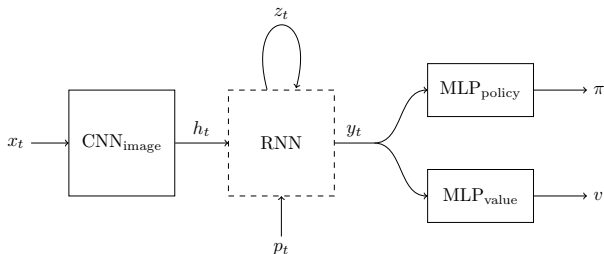
# Camera Environment

- ▶ 3D scene viewed from a perspective projection camera.
- ▶ Height map from terrain environment turned into mesh, same appearance and target probability as before.
- ▶ Camera location fixed at center of scene.
- ▶ Moving actions control pan and tilt (pitch and yaw).
- ▶ Visually complex, difficult to interpret.



# Architecture

- ▶ Actor-critic method trained with PPO [?].
- ▶ Image  $x_t$  passed through CNN.
- ▶ Latent image representation  $h_t$  and position  $p_t$  passed through RNN.  
Two variants:
  1. LSTM with input  $[h_t, p_t]$ .
  2. Spatial memory.
- ▶ Policy and value heads approximate  $\pi$  and  $v_\pi$  with MLPs.





# Recurrent Steps

## 1. LSTM:

- ▶ Proven to work for POMDPs [?, ?, ?, ?].
- ▶ May struggle with remembering over many time steps.
- ▶ Important for exhaustive search and scene understanding.

## 2. Spatial memory (inspired by [?]):

- ▶ Structured memory  $M_t \in \mathbb{R}^{C \times 10 \times 10}$  as hidden state (one slot per camera position  $p_t$  / unique view  $V$  / image  $x_t$ ).
- ▶ Read vector  $r_t = f(M_t)$ ,  $f$  is CNN.
- ▶ Write vector  $w_t = g([h_t, r_t])$ ,  $g$  is MLP.
- ▶ Action probabilities  $\pi([r_t, w_t])$  and value  $v([r_t, w_t])$ .
- ▶  $r_t$  contains information from the whole explored scene.
- ▶  $w_t$  written to index  $p_t$  of  $M_{t+1}$ .

# Experiments

- ▶ Train for 25M time steps.
- ▶ Results reported across 3 runs with different seeds.
- ▶ Interval estimates via stratified bootstrap confidence intervals.
- ▶ Separate training and test sets.
- ▶ Same hyperparameters in all runs.

# Experiment I: Search Space and Reward Signal

- ▶ Larger search spaces take longer to train:
  - ▶ Sparse reward might not be sufficient.
  - ▶ Stronger demands on memory (remember searched positions, scene understanding).
- ▶ Investigate impact by comparing agents on  $10 \times 10$ ,  $15 \times 15$ , and  $20 \times 20$  versions of gaussian environment.
- ▶ Evaluate two additional reward signals that may speed up training:
  - ▶  $r'_t = r_t + e$ , where  $e = 0.1$  if  $a_t \neq \text{TRIGGER}$  moves the view to an unexplored region and 0 otherwise.
  - ▶  $r''_t = r_t + d$ , where  $d = 0.1$  if  $a_t \neq \text{TRIGGER}$  moves the view towards the nearest target and 0 otherwise.

## Experiment II: Search Performance

- ▶ Compare to random searcher, exhaustive searcher, human searcher with prior knowledge of scenes.
- ▶ Use held out samples as test set.
- ▶ Average number of steps on test set.
- ▶ SPL metric [?], with  $N$  as the number of test samples,  $S_i$  indicating success,  $p_i$  as the number of steps and  $l_i$  as the shortest path length:

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

## Experiment III: Generalization with Limited Training Samples

- ▶ Limit number of scene samples seen during training to 100, 1000, 10 000, ....
- ▶ Use terrain environment, high appearance variance and somewhat realistic.
- ▶ Fix seed pool used to generate scenes seen during training.
- ▶ Train agents until convergence (or for a fixed number of time steps).
- ▶ Test on held out scenes from full distribution.

# Implementation

- ▶ OpenAI Gym environment interface.
- ▶ PyTorch for models and automatic differentiation.
- ▶ Intel Core i9-10900X CPU.
- ▶ NVIDIA GeForce RTX 2080 Ti GPU.

# Preliminary Results

## Status:

- ▶ Done with implementation.
- ▶ Exploratory experiments done, know what seems to work and what does not.
- ▶ Spatial memory scales to larger search spaces than LSTM.

## Problems:

- ▶ Many configurable parts (hyperparameters, environment settings, agent architectures).
- ▶ Long training times, difficult to predict hyperparameter interplay, expensive to tune.
- ▶ Difficult to design reward signal, magnitudes matter.
- ▶ Initial problems with scaling up to large search spaces where intelligent search is more important, hopefully solved now.

# Future Steps

1. Collect results across multiple seeds for all experiments.
2. Find opponent.
3. More baselines vs. ablation studies?
4. Discussion and conclusion.
5. Tidy up report.
6. Presentation preparation.



# References I