# Target Driven Visual Navigation with Hybrid Asynchronous Universal Successor Representations

**Shamane Siriwardhana**
Augmented Human Lab
The University of Auckland
New Zealand
shamane@ahlab.org

**Rivindu Weerasekera**
Department of Electrical and Computer Engineering
The University of Auckland
New Zealand
rwee015@aucklanduni.ac.nz

**Suranga Nanayakkara**
Augmented Human Lab
The University of Auckland
New Zealand
s.nanayakkara@auckland.ac.nz

## Abstract

Being able to navigate to a target with minimal supervision and prior knowledge is critical to creating human-like assistive agents. Prior work on map-based and map-less approaches have limited generalizability. In this paper, we present a novel approach, Hybrid Asynchronous Universal Successor Representations (HAUSR), which overcomes the problem of generalizability to new goals by adapting recent work on Universal Successor Representations with Asynchronous Actor-Critic Agents. We show that the agent was able to successfully reach novel goals and we were able to quickly fine-tune the network for adapting to new scenes. This opens up novel application scenarios where intelligent agents could learn from and adapt to a wide range of environments with minimal human input.

## 1 Introduction

Visual navigation is a core problem in the fields of robotics and machine vision. Previous research had used map-based, map-building or map-less approaches [Bonin-Font et al., 2008, Oriolo et al., 1995, Borenstein and Koren, 1991]. The first two approaches had been favoured in the past, however, they depended on an accurate mapping of the environment and a careful human-guided training phase that limited its generalizability [Filliat and Meyer, 2003].

With recent advances in Deep Reinforcement Learning [Mnih et al., 2015, Silver et al., 2016, 2017], map-less navigation [Zhu et al., 2017, Zhang et al., 2017] has seen major advancements. Reinfocement Learning (RL) systems are composed of an agent that learns by trial & error. The agent can learn from what worked best for each situation in the past and apply that knowledge when presented with another similar situation. Ideally, the agent transfers its learning among navigation for various destinations or goals within an environment. However, classical RL algorithms struggle to generalize to changing tasks or goals because rewards and value functions, are generally defined in terms of just one goal or target [Sutton and Barto, 2011]. For example, the value function for a robot navigating to the microwave will be quite different from the one used to navigate to the sofa. Previous successful attempts at target-driven navigation using RL have a marked drop in performance when adapting to new goals even though they use models which encourage generalizability [Zhu et al., 2017, Zhang et al., 2017].

To overcome challenges in generalizability and performance, in this paper, we present a new approach using Hybrid Asynchronous Universal Successor Representations (HAUSR). This is based on the concept of Universal Successor Representations (USR) which are able to learn representations of the environment dynamics that are transferable between different goals.

We contribute with showing how HAUSR can be adapted to the problem of target driven visual navigation in a complex photo-realistic environmend inside AI2THOR [Kolve et al., 2017]. With evaluation of HAUSR in AI2THOR, we demonstrate that an agent was able to successfully reach unseen goals.

## 2 Related Work

### 2.1 General Value Functions

One question central to RL is how to learn a feature representation that is scalable and re-uses learned information between tasks. Designing a value function that is capable of adapting to different tasks, would greatly help generalizability. This line of research builds on a concept called General Value Functions (GVF) [Sutton et al., 2011]. GVFs generalize the concept of the value function to capture not only the goal-oriented semantics of a task but also attempt to capture a more general form of knowledge about the world. An extension to this idea, called Universal Value Function Approximators (UVFAs), was introduced by Schaul et al. [2015]. The main idea of UVFAs is to represent a large set of optimal value functions by a single, unified function approximator that generalises over both states and goals. Although theoretically sound, learning a good UVFA is a challenging task [Ma et al., 2018].

### 2.2 Successor Representations

Successor Representation (SR) [Dayan, 1993] emerged from the field of cognitive science and is modeled on how the brain is able to create a reusable predictive map. SR was combined with Deep Learning to create Deep Successor Reinforcement Learning (DSR) by Kulkarni et al. [2016] which decoupled reward and environment dynamics. Based on Deep Q-Network (DQN) fundamentals, they were able to learn task-specific features that were able to quickly adapt to distal changes in the reward by fine-tuning only the reward prediction feature vector. Transfer in RL was also evaluated on multiple similar tasks by Barreto et al. [2017] who introduced Successor Features (SF). They adapted SR to the continuous domain and were able to show how classic policy improvement can be extended to multiple policies. Extending Barreto et al. [2017] to Deep Learning, Barreto et al. [2018] showed how these models could be trained in a stable way. For the problem of visual navigation, a SR-based DRL architecture similar to Kulkarni et al. [2016] was used by Zhang et al. [2017]. Unlike our approach, they showcase their solution in simple maze-like environments using DQN as the baseline method, while we use actor-critic methods in a photorealistic simulation environment. DQN-based techniques frequently suffer from stability issues when applied to complex problems like large-scale navigation [Barreto et al., 2018].

Universal Successor Representations (USR) [Ma et al., 2018] are a recent extension to SR. Unlike previous methods which were based on DQN, USR learns a policy directly by modeling it with actor-critic methods. Similar to SR, USR modifies the policy with successor features. DQN-based approaches learn an optimal action-value function indirectly. USR attempts to obtain a General Value Function which can be directly used to obtain an optimal policy. It can be seen as a combination of the SF and UVFA methods as discussed earlier. Unlike methods based on DQN, USR is able to directly optimize an agent to learn multiple tasks simultaneously. However, USR has not been tested on high-dimensional complex problems. This paper shows how USR could be adapted to the problem of target driven visual navigation in a photorealistic environment.

### 2.3 Target-Driven Navigation

There were several attempts at solving the visual navigation with DRL. The most relevant to our approach used a deep siamese actor critic agent that shares parameters across multiple goals [Zhu et al., 2017]. They used the asynchronous advantage actor critic (A3C) [Mnih et al., 2016] algorithm and were able to learn multiple goals simultaneously by feeding in both state and goal features as inputs to train across multiple goals simultaneously. Although their approach intuitively resembles

General Value Function approximation, their results show a marked decrease in performance when navigating to targets which the agent has not seen during training. Unlike our approach, they did not use any form of successor representations for generalizations in their architecture. Many RL algorithms are trained in simulated environments that are usually simplified lower dimensional versions of reality [Beattie et al., 2016]. As a result, the agent does not use rich visual information to make decisions. However with the development of photo-realistic simulated environments like AI2THOR [Kolve et al., 2017], recent work have been able to use rich visual information and successfully transfer knowledge from a simulated environment to the real world [Zhu et al., 2017].

## 3 Background

### 3.1 Reinforcement Learning

We formalize the goal-directed navigation task as a Markov Decision Process (MDP). The transition probability $p(s_{t+1}|s,a)$ defines the probability of reaching the next state $s_{t+1}$ when action $a_t \in \mathcal{A}$ is taken in state $s_t \in \mathcal{S}$. For any goal $g \in \mathcal{G}$ (very often $\mathcal{G} \subseteq \mathcal{S}$), we define a pseudo-reward function $r_g(s_t, a_t, s_{t+1}) \in \mathbb{R}$ and a pseudo-discount function $\gamma_g(s) \in [0,1]$ (for terminal state, $\gamma_g = 0$). For any policy $\pi(a_t|s_t)$, the General Value Function [Sutton et al., 2011] can be defined as:

$$V_g^\pi(s) = \mathbb{E}^\pi \left[ \sum_{t=0}^{\infty} r_g(s_t, a_t, s_{t+1}) \prod_{k=0}^{t} \gamma_g(s_k) \middle| s_0 = s \right] \tag{1}$$

We assume for any goal $g$, there exists an optimal value function $V_g^*(s) = V_g^{\pi_g^*}(s)$ evaluated according to a goal oriented optimal policy $\pi_g^*$. The aim of policy learning is to find the optimal policy $\pi^*$ that maximizes the future discounted rewards starting from $s_0$ and following $\pi^*$.

To generalize over the goal space $\mathcal{G}$, the agent needs to learn multiple optimal policy and optimal value functions for navigating to each goal. Each goal is considered a new task for the agent and it should be able to quickly adapt to these new tasks and learn $V_g^{\pi_g^*}(s)$ and $\pi_g^*$.

### 3.2 Asynchronous Advantage Actor Critic (A3C)

Asynchronous advantage actor-critic (A3C) is an on-policy policy search method introduced by Mnih et al. [2016] based on actor-critic methods [Witten, 1977, Barto et al., 1983]. A3C is a fast, robust, scalable and has achieved state of the art results in the video game domain [Mnih et al., 2016].

A3C maintains estimates of both a policy $\pi(a_t|s_t)$ and the value function $V(s_t)$. The agent uses the value function (the critic) to update it's policy (the actor) by training multiple threads in parallel and asynchronously updating a shared set of model parameters. It has been shown that the parallel threads are able to stabilize each other. In practice we can learn both policy and value functions in a single neural network with multiple heads.

An update is performed by calculating $\nabla_{\theta_\pi} \log \pi(a_t|s_t) A(s_t, a_t)$ where $A$ is the advantage term. The advantage term can be thought of as quantifying how much better an action turned out to be than expected. For a state, the advantage can be estimated by the difference between the discounted rewards of an episode and the value of that state. The full update can therefore be written as:

$$\nabla_{\theta_\pi} \log \pi(a_t|s_t; \theta_\pi)(R_t - V(s_t; \theta_V)) + \beta \nabla_{\theta_\pi} H(\pi(s_t; \theta_\pi)) \tag{2}$$

$H$ is an entropy terms which encourages exploration of the agent controlled by hyperparameter $\beta$.

### 3.3 Universal Successor Representations (USR)

A key idea in applying SR to deep architectures is being able to decouple and approximate the reward function $r_g(s_t, a_t, s_{t+1})$ as a linear combination of learned state features $\phi(s_t, a_t, s_{t+1})$ and a reward weight vector $\omega(g)$ [Kulkarni et al., 2016, Barreto et al., 2017].

$$r_g \approx \phi(s_t, a_t, s_{t+1}; \theta_\phi)^\top \omega(g_t; \theta_\omega) \approx \phi(s_{t+1}; \theta_\phi)^\top \omega(g_t; \theta_\omega) \tag{3}$$

where $\theta_\phi$ and $\theta_\omega$ are learnt parameters of a function approximator such as a neural network. In this case, equation 1 can be rewritten as

(a) State Representation and Reward Prediction Networks

(b) Reinforcement Learning Network

Figure 1: Proposed Network Architecture for Hybrid Asynchronous Universal Successor Representations (HAUSR). Network shown in Figure 1a is trained first with Algortihm 1 and the reinforcement learning network in Figure 1b is trained with Algorithm 2.

$$V_g^\pi(s) = \mathbb{E}^\pi \left[ \sum_{t=0}^\infty \phi(s_t, a_t, s_{t+1}; \theta_\phi) \prod_{k=0}^t \gamma_g(s_k) \middle| s_0 = s \right]^\top \omega(g_t; \theta_\omega) = \psi^\pi(s_t, g_t; \theta_\psi)^\top \omega(g_t; \theta_\omega) \tag{4}$$

where $\psi^\pi(s_t, g_t)$ is defined as the USR of state $s_t$ [Ma et al., 2018]. Intuitively, $\psi^\pi(s_t, g_t)$ can be thought of as the expected future state occupancy.

USR makes it easier for us to transfer knowledge between goals. If $\omega(g_t; \theta_\omega)$ can be effectively computed for any $g_t$, we can compute a value function and then an optimal policy for any goal. Learning the USR can then be done in the same way as the value function update using the following Bellman update:

$$\psi^\pi(s_t, g_t; \theta_\psi) = \mathbb{E}^\pi[\phi(s_t, a_t, s_{t+1}) + \gamma_g(s)\psi^\pi(s_{t+1}, g_t; \theta_\psi)] \tag{5}$$

## 4 Problem Formulation

The objective of our target driven navigation agent is to learn a stochastic policy function $\pi(s_t, g_t)$ where $s_t$ is the representation of the current state and $g_t$ is a representation of the target state. The output of policy $\pi$ would be a probability distribution over actions $\mathcal{A}$. In our formulation, the agent has four discrete actions: 1) to move forward 0.5m, 2) move backward 0.5m, 3) turn left or 4) turn right. We use the "bathroom 02" which contains 180 unique states in the photo realistic 3D AI2THOR [Kolve et al., 2017] environment to train our policy and test the transferability of the learned policy to other novel goals.

Each scene is divided into $0.5m \times 0.5m$ grids similar to a grid world environment. The state representation coming from the AI2THOR environment is a $300 \times 400$ raw colour



Figure 2: States coming from the AI2THOR environment are **300 × 400** RGB images

pixel image 2. We scale this down to a $110 \times 110$ gray-scale image and concatenate four consecutive frames to create our state vector $s_t$.

The agent receives a reward of $+1$ for reaching the goal state and $-0.01$ for each time-step. The objective of our agent is to navigate to a goal location in a minimal number of steps. After training, the agent should be able to generalize to novel goal locations which were not used to train our model.

### 4.1 Challenges in merging USR with A3C

USR and actor-critics methods have previously only been applied to simple maze-like environments. We found applying USR under the baseline method of A3C directly to the complex problem of visual navigation to be unstable due to a few reasons.

First, learning a good state embedding $\phi$ (Equation 5) using just an auto-encoder as proposed by Ma et al. [2018] in a complex environment (like AI2THOR) was a challenge. The state representation is a crucial part required for the agent to gain an understanding of it's environment. Mistakes in the state representation caused degraded performance of the agent.

Second, we found that it was not easy to accurately predict the reward prediction vector ($\omega(g_t)$ in Equation 3). This vector and state representations are in turn are used in Equation 4 and help us approximate the advantage function and scalar reward. Unlike ad-hoc rewards, the predicted scalar rewards using Equation 5 have wide variations. This is especially the case when training for multiple goals simultaneously. In the next section we describe how we address these challenges.

In the next section we will go over our Hybrid Asynchronous Universal Successor Representations (HAUSR) architecture. To our knowledge, this is the first time USR has been combined with A3C for target driven visual navigation.

## 5 Hybrid Asynchronous Universal Successor Representations (HAUSR)

### 5.1 Network Architecture

Different from previous methods which used DQN as the baseline method with SR to improve the generalizability and transfer learning capabilities, we used a deep actor critic method and extend USR to work with A3C. This has the advantage of opening up our algorithm to be used in real world sceneraios due to A3C's scalability [Mnih et al., 2016]. Our architecture is also more stable and scalable than the vanilla USR network as proposed by Ma et al. [2018]. The three main networks of our architecture are shown in Figure 1. The state representation and reward prediction networks (Figure 1a) were trained before the reinforcement learning network (Figure 1b) was trained.

#### 5.1.1 State Representation Network

The goal of the state representation network was to generate $\phi$ which encodes useful dynamics of the state that would later be used to approximate the successor feature representations. Instead of training $\phi$ with just an auto-encoder loss, we used a new architecture with auxiliary losses to make a more robust state representation. We added auxiliary branches in the form of a foreword dynamics loss and an inverse dynamics loss, in addition to the autoencoder reconstruction loss (as shown in Figure 1a). Forward and inverse dynamics have been widely discussed in literature and have shown to improve state representations for DRL tasks [Christiano et al., 2016, Pathak et al., 2018]. Both autoencoder reconstruction and forward dynamics losses use a mean squared error loss function, while the inverse dynamics loss uses a cross-entropy loss.



Figure 3: The top row shows some of the states in the simulator. The bottom row shows the reconstruction of those states generated by feeding the learnt state representation $\phi$ through the auto-encoder decoder

In the early stages of training, states were sampled from exploration of the agent with a randomly initialized policy. In each interaction with the environment

we collected roll-outs consisting of $s_t, a_t, s_{t+1}....$ Using these rollouts, we were able to train an informative bottleneck layer using Algorithm 1 and generate $\phi$. Figure 3 shows some of the reconstructed images of different states using $\phi$.

### 5.1.2 Reward Prediction Network

The goal of the reward prediction network was to train the reward prediction vector $\omega_g$. This vector should be able to transform the state representation $\phi$ into a scalar reward as

$$r_g \approx \omega \cdot \phi^\top \tag{6}$$

We used a three layer Convolutional Neural Network (CNN) which took in the target image as input as generated $\omega$ as shown in Figure 1a. The predicted rewards should be either $-0.01$ if the agent had not reached the goal or $+1$ if the agent had reached the goal. Generating $\omega$ with this simple network gave us the ability to quickly obtain the reward prediction vector for different goals. First, we pre-trained this network with the collected roll-outs and then trained simultaneously with the A3C agent.

### 5.1.3 Reinforcement Learning Network

The final network was our Actor-Critic agent network (Figure 1b) which implements our A3C agent with successor features.

We created a new advantage function in order to achieve stability during training. The advantage function is normally the driving force for policy improvement and is based on ad-hoc rewards (like a small penalty for every step and +1 for reaching the goal). Advantage is calculated as the difference between the discounted reward and the value of that state. An ad-hoc scaler rewards are very useful for optimizing agents in complex environments, especially with parallel architectures like A3C. Unlike DQN, A3C is stable out of the box and generally require less hyperparameter tuning.

In our work, we combined the normal A3C architecture with USR by calculating two advantage functions. One advantage function was calculated with ad-hoc scaler rewards (like in traditional A3C), while the other was calculated with the output from the reward prediction network (which took the successor features into account). We modified the conventional advantage and return functions in the A3C agent with USR as follows:

$$\hat{A}_t^\psi = [\phi(s_t) + \gamma\psi(s_{t+1}, g_t) - \psi(s_t, g_t)]^\top \omega_g \tag{7}$$

$$\hat{A}_t^V = \sum_{k=0}^\infty \gamma^k r_{t+k} - V_t \tag{8}$$

We took a weighted sum of the two advantages with the hyperparameter $\lambda$, determining the contribution from successor feature rewards. As a result, the agent was able to update its policy parameters not only with respect to ad-hoc rewards but also took into account the state and environment dynamics. We believe this is key to learning better feature representations that are able to generalize across goals.

In the RL network, we sent both the goal image and current state through a shared three layer siamese CNN. We then concatanated the two output vectors from the CNN and passed it through fully connected layers. The output of the network was a value function $V$, policy $\pi$ and a Universal Successor Representation Approximator $\psi$. It was trained as described in Algorithm 2. The model can be easily extended to multiple environments by creating branches after the convolutional layers of the network.

### 5.2 Training Protocol

Constructing a stable training architecture was very important to applying Deep Reinforcement Learning to high dimensional problems, such as visual navigation.

Both the return and the advantage functions of our network depend on the reward prediction vector ($\omega_g$) and the state representation vector ($\phi$). Conventional A3C agents rely on a stable reward

generating mechanism. Obtaining perfect state representation and reward prediction vectors was highly unlikely due the noise and complexities in the real world. During training, we initially set $\lambda$ to a very low value to prioritize learning of the more stable value function. After $5k$ iterations with 32 agents, we increased $\lambda$ to $0.001$. After we increased $\lambda$, the value loss also increased, providing evidence that the SR loss had an effect on training. The combined loss function consisted of the policy loss, value function loss, successor feature approximation loss and an entropy loss [Williams and Peng, 1991].

$$\mathcal{L}_{total} = \lambda \mathcal{L}_\psi + \mathcal{L}_V + \mathcal{L}_\pi + \mathcal{L}_H \tag{9}$$

It is our understanding that this is the first time that the conventional value loss was used in conjunction with successor features. The addition of auxiliary losses to help training has been done before [Jaderberg et al., 2016] and we believe this concept will be key to creating successful real world RL agents in the future.

## 6  Evaluation

Figure 4 shows the n-step temporal difference error with respect to the conventional Value function as well as the Successor Feature approximator. We can see how these stabilize over time. As evident in the graph, the SR loss does not reach zero but this is sufficient to give generalization abilities to our agent. It also provides evidence for why we train in a hybrid manner and still use the traditional Value loss in our training. We didn't compare our model with a baseline model of A3C since we were more concerned about



Figure 4: Shows how our SR loss and Value Loss stabilize over time

creating a stable training architecture for combining A3C with USR in a high dimensional complex problem.

We tested the capability of our model for transfer learning and generalization by conducting two experiments. One to test the generalizability of the agent without retraining and the other to test how quickly we can fine tune the agent for novel goals. For both experiments, we trained our model until convergence in the bathroom environment of AI2THOR with five goal locations.

### 6.1  Generalization to novel goals

We checked the generalizability of our agent to new tasks *without* fine-tuning. When generalizing to new goals, the most important factor to consider was if the agent was successful in reaching the new goal. Therefore, following [Zhu et al., 2017], we used the success rate of reaching the goal as a metric. For each novel goal, we ran 100 episodes of the agent with each starting from a random initial state. If the agent reached the goal within 500 steps, we considered it to be a successful episode. We tested our agent's ability to navigate to 50 goal states in the same environment. Figure 6 shows the success rate of our agent in reaching these goals. Note that we do not fine tune our network for these new goals.



Figure 5: Starting from random states, our agent was trained to navigate to only five trained images. When directed to navigate to the over 19 other targets, our agent was successful in completing these new tasks without any additional training (only some successful targets are shown above). Even though some states had only minor differences, the agent was successful in differentiating and reaching both goals.

Figure 6: Success rate of agent navigating to goal for all tested states. The red highlighted states are the targets obtained from the AI2THOR simulator and used for training. The black states were not used for training. The agent was able to successfully generalize it's learnings and navigate to a large number of new goals. Note: the states are not necessarily in order of closeness

## 6.2 Transfer Learning Ability

Transfer learning has been extensively studied in
the supervised learning domain and leads to fast training on new tasks [Sharif Razavian et al., 2014]. If a model is able learn a rich representation of a problem, that knowledge should be useful when learning new tasks. This enables us to have an agent that can quickly adapt to new tasks. We tested our network's transfer learning ability by fine-tuning the network for new goals. Figure 7 shows that our agent is able to very quickly learn to navigate to goals that it previously not able to reach. The number of steps taken to reach a new goal dramatically decreases and the agent is able learn in under $15 minutes$.

## 6.3 Future Work

In the recent past, DRL has been very effective at solving complex tasks [Silver et al., 2016, Mnih et al., 2015]. Real world navigation is still an open problem in DRL. New breakthrough algorithms are frequently only tested in simple toy environments. Scalability is another factor that needs to be considered when applying DRL to complex domains. We identified that merging new concepts into complex domains is challenging and a key research question is how these ideas can be implemented with stable architectures. We hope that our initial attempt at merging two novel concepts leads to more research in creating fine-tuned architectures that are able to perform reliably in real-world applications.



Figure 7: Transfer Learning: the agent is able to learn to navigate to novel goals it was previously unable to reach in under 15 minutes of fine-tuning

## 7 Conclusion

We show the applicability of Universal Sucessor Representations to the complex domain of target driven visual navigation
in a photorealistic environment. We present a new architecture (HAUSR) that is able to create rich successor representations and train asynchronously with A3C. Despite encouraging results, there are several opportunities to exploit hybrid Successor Representation based approaches to achieve higher generalizability and transferrability of RL agents.

## References

André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. In *Advances in neural information processing systems*, pages 4055–4065, 2017.

Andre Barreto, Diana Borsa, John Quan, Tom Schaul, David Silver, Matteo Hessel, Daniel Mankowitz, Augustin Zidek, and Remi Munos. Transfer in deep reinforcement learning using successor features

and generalised policy improvement. In *International Conference on Machine Learning*, pages 510–519, 2018.

Andrew G Barto, Richard S Sutton, and Charles W Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE transactions on systems, man, and cybernetics*, (5): 834–846, 1983.

Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab. *arXiv preprint arXiv:1612.03801*, 2016.

Francisco Bonin-Font, Alberto Ortiz, and Gabriel Oliver. Visual navigation for mobile robots: A survey. *Journal of intelligent and robotic systems*, 53(3):263–296, 2008.

Johann Borenstein and Yoram Koren. The vector field histogram-fast obstacle avoidance for mobile robots. *IEEE transactions on robotics and automation*, 7(3):278–288, 1991.

Paul Christiano, Zain Shah, Igor Mordatch, Jonas Schneider, Trevor Blackwell, Joshua Tobin, Pieter Abbeel, and Wojciech Zaremba. Transfer from simulation to real world through learning deep inverse dynamics model. *arXiv preprint arXiv:1610.03518*, 2016.

Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5(4):613–624, 1993.

David Filliat and Jean-Arcady Meyer. Map-based navigation in mobile robots:: I. a review of localization strategies. *Cognitive Systems Research*, 4(4):243–282, 2003.

Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397*, 2016.

Eric Kolve, Roozbeh Mottaghi, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.

Tejas D Kulkarni, Ardavan Saeedi, Simanta Gautam, and Samuel J Gershman. Deep successor reinforcement learning. *arXiv preprint arXiv:1606.02396*, 2016.

Chen Ma, Junfeng Wen, and Yoshua Bengio. Universal successor representations for transfer reinforcement learning. *arXiv preprint arXiv:1804.03758*, 2018.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.

Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937, 2016.

Giuseppe Oriolo, Marilena Vendittelli, and Giovanni Ulivi. On-line map building and navigation for autonomous mobile robots. In *Robotics and Automation, 1995. Proceedings., 1995 IEEE International Conference on*, volume 3, pages 2900–2906. IEEE, 1995.

Deepak Pathak, Parsa Mahmoudieh, Guanghao Luo, Pulkit Agrawal, Dian Chen, Yide Shentu, Evan Shelhamer, Jitendra Malik, Alexei A Efros, and Trevor Darrell. Zero-shot visual imitation. In *International Conference on Learning Representations*, 2018.

Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In *International Conference on Machine Learning*, pages 1312–1320, 2015.

Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 806–813, 2014.

David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.

David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.

Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. 2011.

Richard S Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M Pilarski, Adam White, and Doina Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pages 761–768. International Foundation for Autonomous Agents and Multiagent Systems, 2011.

Ronald J Williams and Jing Peng. Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268, 1991.

Ian H Witten. An adaptive optimal controller for discrete-time markov environments. *Information and control*, 34(4):286–295, 1977.

Jingwei Zhang, Jost Tobias Springenberg, Joschka Boedecker, and Wolfram Burgard. Deep reinforcement learning with successor features for navigation across similar environments. In *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*, pages 2371–2378. IEEE, 2017.

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 *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 3357–3364. IEEE, 2017.

# Appendix

---

**Algorithm 1** State Representation and Reward Prediction Learning

---

1: Initialize transition history $\mathcal{H}$
2: **for** $g$ training goals **do**
3:     Collect random rollouts $\mathcal{H} \leftarrow \mathcal{H} \cup (g, s_1, a_1, r_1, s_2, ..., s_N, a_N, r_N)$
4: **for** time step $t$ in $\mathcal{H}$ **do**
5:     Pick random transition $(g, s_t, a_t, r_t, s_{t+1}, a_{t+1}, r_{t+1})$ from $\mathcal{H}$
6:     Calculate autoencoder, forward and inverse dynamics losses as $\mathcal{L}_\phi$
7:     Perform gradient decent on $\mathcal{L}_\phi$ to update state representation network parameters
8:     Perform gradient decent on $\mathcal{L}_\omega = [r_t - \phi(s_{t+1})^\top \omega(g; \theta_\omega)]^2$ w.r.t. $\theta_\omega$

---

---

**Algorithm 2** Async-USR

---

1: **for** $\mathcal{M}$ agents simultaneuously **do**
2: **for** $ns$ steps **do**
3:     Obtain rollout $(g, s_1, a_1, r_1, s_2, ..., s_N, a_N, r_N)$ by following $\pi(s_t)$
4:     Compute $L_\psi = \|\phi(s_t) + \gamma_t \psi(s_{t+1}, g; \theta_\psi) - \psi(s_t, g; \theta_\psi)\|^2$
5:     Compute $L_V = \|r(s_t) + \gamma_t V(s_{t+1}, g; \theta_V) - V(s_t, g; \theta_V)\|^2$
6:     Compute $\hat{A}_t^\psi = [\phi(s_t) + \gamma_t \psi(s_{t+1}, g) - \psi(s_t, g)]^\top \omega_g$
7:     Compute $\hat{A}_t^V = \sum_{k=0}^{\infty} \gamma^k r_{t+k} - V_t$
8:     Perform gradient descent on loss $(\lambda \hat{A}_t^\psi + \hat{A}_t^V) log(\pi(s_t, g; \theta_\pi) + 0.5 L_V + 0.1 L_\psi - \beta H(\pi)$ w.r.t. $\theta_\pi, \theta_V$ where $H(\pi)$ is the entropy loss and $\lambda \& \beta$ are hyperparamters

---