

# Learning Human Search Behavior from Egocentric Visual Inputs

Maks Sorokin<sup>1</sup>  Wenhao Yu<sup>1,2</sup>  Sehoon Ha<sup>1,2</sup>  C. Karen Liu<sup>3</sup> 

{maks,wenhaoyu,sehoonha}@gatech.edu, karenliu@cs.stanford.edu

<sup>1</sup> Georgia Institute of Technology, Atlanta, GA, 30308, USA

<sup>2</sup> Robotics at Google, Mountain View, CA, 94043, USA

<sup>3</sup> Stanford University, Stanford, CA, 94305, USA



**Figure 1:** A humanoid character learns to navigate and search for a target object (the mustard bottle) in a photorealistic 3D scene using its own egocentric vision and locomotion capability. Top: third-person view. Bottom: first-person view.

## Abstract

“Looking for things” is a mundane but critical task we repeatedly carry on in our daily life. We introduce a method to develop a human character capable of searching for a randomly located target object in a detailed 3D scene using its locomotion capability and egocentric vision perception represented as RGBD images. By depriving the privileged 3D information from the human character, it is forced to move and look around simultaneously to account for the restricted sensing capability, resulting in natural navigation and search behaviors. Our method consists of two components: 1) a search control policy based on an abstract character model, and 2) an online replanning control module for synthesizing detailed kinematic motion based on the trajectories planned by the search policy. We demonstrate that the combined techniques enable the character to effectively find often occluded household items in indoor environments. The same search policy can be applied to different full body characters without the need of retraining. We evaluate our method quantitatively by testing it on randomly generated scenarios. Our work is a first step toward creating intelligent virtual agents with humanlike behaviors driven by onboard sensors, paving the road toward future robotic applications.

## CCS Concepts

• *Computing methodologies* → *Procedural animation; Motion processing;*

## 1. Introduction

“We spend about 5000 hours of our lives looking for things around the home” [IKE19]. Indeed, searching for objects in complex indoor environments is a frequent event in our daily life—we look for ingredients in the kitchen for a recipe, we locate grocery items on the shelves in a supermarket, and we seem to always be in search of phones, keys or glasses many times a day. The goal of this paper

is to model such important and ubiquitous behaviors by developing a virtual human capable of using its egocentric vision perception and locomotion capability to search for any randomly placed target object in a complex 3D scene.

Search behaviors depend on simultaneous locomotion and survey of the environment, requiring modeling not only the physical motor skills, but also human sensory and decision making. Con-

ventional character animation assumes full knowledge of the 3D environment and utilizes it to achieve optimal movements. While optimality is indeed observed in many human locomotion and manipulation tasks, it is also at odds with the stochastic nature of human sensing and decision making capabilities. Having the “oracle ability” to know exactly the 3D position of every vertex in the scene will likely result in an “optimal but unnatural” search behaviors. Alternatively, imitating or tracking motion capture trajectories directly can potentially lead to natural human behaviors. However, pre-scripting or pre-planning a reference trajectory can be challenging for a search controller tasked to find objects placed at random locations in random scenes.

This paper builds on the hypothesis that equipping the virtual character with human-like sensing capabilities can lead to more natural behaviors, as demonstrated by previous work [YLN<sup>12</sup>, NZC<sup>18</sup>, EHSN<sup>20</sup>]. In particular, we limit the virtual character to egocentric vision perception from RGBD images when performing the search task. Without the full state information about the environment and the global position and orientation of itself, the character is forced to coordinate its motor capabilities to navigate and scan the scene to find the target object, naturally inducing humanlike decision-making behaviors under partial observation. Our method consists of two components: 1) a search control policy that determines where the character should move to and look at based on an abstract model, and 2) an online replanning control module for synthesizing the detailed kinematic motion based on the planned trajectories from the search policy. This decoupling of the task provides several benefits: first, training the search policy with an abstract model that only includes a torso and a head/camera is more computationally efficient for high-dimensional and large observation space, and second, a trained search policy can be re-used for multiple characters that share the same abstract model without retraining.

We use deep reinforcement learning (DRL) to train the search policy that takes as input the visual perception from the abstract model and predicts where it should move to and look at. Our training framework shares some similarities with existing works in visual navigation while having a few key distinctions. Unlike most vision-guided navigation applications, our agent can actively choose where to look independent of the body moving direction, enabling more plausible head movements and effective search strategies for the character. However, decoupling the body movement with the gaze direction makes the learning problem more challenging due to the larger observation space and the higher dimensional action space. We show that by combining Soft Actor Critic (SAC) with a contrastive proposed by Srinivas *et al.* [SLA<sup>20</sup>], a vision-guided policy can effectively learn the features from high dimensional pixels for our search problem. In addition, transferring the policy trained for an abstract model to a full-body character presents many challenges. Drawing analogy from sim-to-real transfer learning, we propose a zero-shot online replanning method to transfer a model-agnostic policy to the biped human model and a wheel-based robot model. Combining offline visuomotor policy learning with online trajectory planning results in a virtual human capable of making motion plans using egocentric vision perception.

We demonstrate our method on a human and a robotic character searching household items in realistic indoor scenes. We show that the character is able to find a small object, such as a pair of glasses in a large space including an open kitchen and a living room, populated with furniture and other objects. To get the overall performance of our policy, we report the success rate of search policy tested on randomly created scenarios. We further demonstrate the importance of enabling the head movement of the character for both better learning performance and better search behavior. Our work is a first step toward creating intelligent virtual agents with human-like behaviors driven by onboard sensors, paving the road toward future robotic applications.

## 2. Related Work

Our research is inspired by prior work in visual navigation, deep reinforcement learning, and data-driven kinematic animation. Below we will review each of them in turn.

### 2.1. Visual Navigation

Training an autonomous agent to navigate complex environments from visual inputs has been an important topic in computer graphics, robotics, and machine learning [ACC<sup>18</sup>, KWR<sup>16</sup>, MPV<sup>16</sup>, ZMK<sup>17</sup>, PZI<sup>19</sup>].

Some of the work by Kuffner *et al.* [KJL<sup>99</sup>] tackles the problem using path-planning and path-following algorithms that utilize the privileged information about the environment (e.g., floor layouts) and aims to generate collision free paths. Shim *et al.* [SYT<sup>17</sup>] and Wang *et al.* [WSY<sup>18</sup>] avoid the use of the privileged information and learn to approach the goal object which perform feature-based goal identification while tackling the searching via random exploration. In this work, the character is deprived of privileged information and perceives the world only via visual observations, which leads to a learned searching behaviour.

Enabling these agents to work with real images is essential for applying them to real-world applications. However, directly training visual navigation agents in real environments is expensive, especially since we usually want to train agents for a large variety of environments. To this end, researchers have developed simulated environments that leverage modern 3D scanning techniques to reproduce real-world scenarios and allow agents to observe photo-realistic visual inputs in a scalable way [XSL<sup>20</sup>, MAO<sup>19</sup>, KMH<sup>17</sup>]. These tools enable rapid advancements in learning algorithms and neural network structures in training visual navigation agents [GDL<sup>17</sup>, PS<sup>17</sup>, ZTB<sup>17</sup>, FTFFS<sup>19</sup>, WKM<sup>20</sup>]. For example, Fang *et al.* [FTFFS<sup>19</sup>] proposed a scene memory transformer architecture that saves history of observations into the memory and extracts relevant information using the Transformer architecture [VSP<sup>17</sup>]. Wijmans *et al.* developed a decentralized distributed proximal policy optimization (DD-PPO) that allows large scale training of visual navigation agents and demonstrated that with large scale training one can obtain agents that generalize to novel scenarios [WKM<sup>20</sup>].

Our method also utilizes the simulation tools developed by other researchers to train our virtual human in a realistic environment.

Specifically, we used iGibson, which provides a suite of realistic indoor environments [XSL\*20]. Unlike prior work in visual navigation that focus on agents moving in 2D or 2.5D spaces, e.g. mobile robots, our work solves visual navigation tasks with an additional challenge of controlling egocentric perspective.

## 2.2. Deep Reinforcement Learning

Deep reinforcement learning (DRL) provides a general framework for automatic design of controllers for complex motor skills from simple reward functions. Within the graphics community, researchers have applied DRL algorithms on a variety of physics-based control problems, such as locomotion [YTL18, PBVDP17, PALvdP18], manipulation [CYT\*18], aerial behaviors [WPKL17], and soft-body motion [MWL\*19]. However, most of these methods use low-dimensional character states or exploit privileged 3D information in the input space to the policy in order to simplify the learning problem. Directly learning a controller from egocentric vision inputs remains a challenging problem. Recent advancements in image-based deep reinforcement learning have shown promising progress in addressing this challenge [SLA20, HFW\*19, OLV18]. For example, Srinivas *et al.* proposed to learn an embedding of the visual input by minimizing a contrastive loss between randomly cropped input images from the replay buffer [SLA20]. Their method demonstrated superior performance on a set of image-based robotic control problems. In our work, we apply the method by Srinivas *et al.* [SLA20] for training egocentric vision-based policies to accomplish the searching task.

Similar to our work, Merel *et al.* investigated the problem of creating full-body human motions with egocentric vision-based control policies [MAP\*19, MTA\*20]. In particular, they developed a hierarchical control scheme that exploits egocentric vision to coordinate low-level motor skill modules derived from motion capture demonstrations. Their learning approach is able to train the character locomotion and whole-body manipulation using visual inputs. Our work also takes 2D images as input, but our 3D scenes contain detailed geometry with photo-realistic appearance, resulting in a much more complex observation space than those used in the previous work. In addition, the vision inputs are more critical to the search task and require careful coordination between the locomotion and the gaze direction to enable the character to navigate in a cluttered environment while thoroughly surveying the environment.

## 2.3. Data-driven Kinematic Animation

Data-driven kinematic animation has been an effective approach for generating realistic human animations from example motion trajectories. Early work constructs graph-based structures to automatically transition between recorded motion clips [LCR\*02, AF02, KGP02]. Although these methods can successfully generate whole-body motions, the output motions are limited to motion clips in the database. To overcome this limitation, interpolation techniques based on linear bases [SHP04, CH05] or statistical transition models [CH07, LWB\*10] are adopted for predicting more expressive motions from a smaller number of examples. These methods are further extended by exploiting neural networks, such as conditional Restricted Boltzmann Machine (cRBM) [TH09] or an

Encoder-Recurrent-Decoder (ERD) network [FLFM15]. Recently, many researchers have demonstrated that deep neural network can successfully learn human motion manifolds for bipedal locomotion [HSK16, HKS17], quadrupedal locomotion [ZSKS18], object interactions [SZKS19], and motion retargetting [AWL\*19]. We utilize previous work by Holden *et al.* [HKS17] and Starke *et al.* [SZKS19] to generate detailed full-body animations from an abstract trajectory planned by the search policy.

## 3. Overview

Given a realistic 3D indoor scene populated with furniture and household items, we develop a virtual human capable of using its own egocentric vision and locomotion capability to search for any randomly placed target object in the scene, including those occluded by furniture or due to the layout of the room. We take a hierarchical approach which consists of two components, a *search policy* operates on an abstract agent to determine the movement and gaze direction at every time step, and a *motion synthesis* module that synthesizes the kinematic motion on a full-body character to realize the actions determined by the search policy.

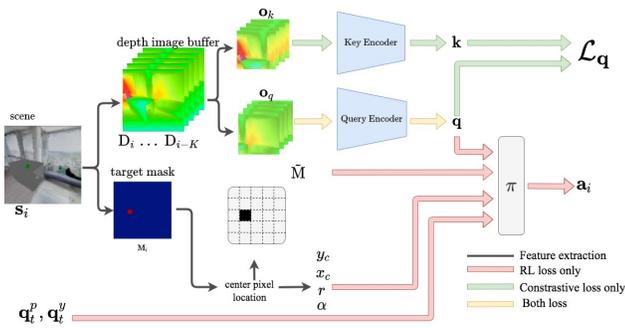
We made a few assumptions in our framework. The vision input includes RGBD images and a mask channel in which the target object has been segmented and labeled. Automatic segmentation and object recognition for general scenes and objects from raw images is a challenging computer vision problem, beyond the scope of this work. In addition, we assume that the 3D environment has furniture and partial divisions that block the line of sight, but it is roughly one connected space, as we do not attempt to solve a maze navigation problem.

## 4. Search Control Policy

While end-to-end DRL approaches have demonstrated success in learning motor skills, solving a visuomotor policy with a detailed full-body character in a large and highly textured environment remains challenging due to the large and complex observation space of the agent. In this work, we propose to first use the learning approach for training an agent-agnostic search policy that has a vision-based observation space but an abstract action space. Once trained, the search policy can be applied to characters with various kinematics to synthesize search behaviors with full-body movements. To this end, we define an abstract model that consists of a cylinder-shape main body and a camera connected to the main body via a universal joint with two degrees of freedom (dofs). The abstract model can move around in the 3D space while the camera can point at different directions independently from the body movement. This additional head movement allows a character to simultaneously navigate and look around, which is essential to model human-like search behaviors.

### 4.1. Problem Formulation

We formulate the vision-based search task as Partially Observable Markov Decision Processes (POMDPs),  $(\mathcal{S}, \mathcal{O}, \mathcal{A}, \mathcal{T}, r, p_0, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{O}$  is the observation space,  $\mathcal{A}$  is the action space,  $\mathcal{T}$  is the transition function,  $r$  is the reward function,  $p_0$



**Figure 2:** Overview of the learning pipeline for training the search policy.

is the initial state distribution and  $\gamma$  is a discount factor. We take the approach of model-free reinforcement learning to find a policy  $\pi$ , such that it maximizes the accumulated reward:

$$J(\pi) = \mathbb{E}_{\mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_T} \sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t), \quad (1)$$

where  $\mathbf{s}_0 \sim p_0$ ,  $\mathbf{a}_t \sim \pi(\mathbf{o}_t)$ ,  $\mathbf{o}_t \sim c(\mathbf{s}_t)$  and  $\mathbf{s}_{t+1} = \mathcal{T}(\mathbf{s}_t, \mathbf{a}_t)$ . The state space contains 3D information of the environment and the global position and orientation of the agent for reward evaluation during training, but it is not available to the policy during testing. Instead, the policy can only access limited information observable to the onboard sensors. Our POMDP is defined as follows:

**Observation space.** The observation space consists of two sensing modalities: vision and proprioception (Figure 2). The proprioception for the abstract model only contains the joint position between the camera and the cylinder: the pitch angle  $q^p$  and the yaw angle  $q^y$ . The agent does not have an access to its global position and orientation.

The vision perception is represented as 2D images observable by standard RGBD cameras, augmented by mask images that provide the segmentation of the target object. For the searching tasks, we exclude the color information (RGB) because we found that the depth and mask images alone contain sufficient information for finding collision-free paths to complete the task when navigating in a cluttered 3D scene. The depth image  $\mathbf{D}$  ( $84 \times 84$ ), obtained with the field of view (FOV) angle of  $90^\circ$ , has the maximal depth at 5 m, normalized to the range of  $[0, 1]$ . This setting has been proven effective for visual navigation tasks, such as Point-Goal navigation [SKM\*19].

The mask image  $\mathbf{M}$  is a binary image which contains 1s at the target object’s pixel locations and 0s otherwise. When the target object is not in the field of view,  $\mathbf{M}$  provides no information. We process the raw mask image  $\mathbf{M}$  to obtain a feature vector  $\mathbf{m} = [x_c, y_c, r, \alpha, \tilde{\mathbf{M}}]$ , where  $x_c$  and  $y_c$  are the average coordinates of the pixels with value 1,  $r = \sqrt{x_c^2 + y_c^2}$  and  $\alpha = \arctan(y_c, x_c)$  are their polar coordinates, and  $\tilde{\mathbf{M}}$  is the downsampled mask image of size  $5 \times 5$ . Using our mask feature vector instead of the raw mask image reduces the dimensionality of the state space and forces the agent to learn a policy agnostic to the object shapes. Note that  $x_c, y_c, r$  and  $\alpha$  are all defined in the image frame. The observation at every time

step  $t$  is defined as  $\mathbf{o}_t = [\mathbf{D}_t, \mathbf{D}_{t-1}, \dots, \mathbf{D}_{t-K+1}, \mathbf{m}_t, q_t^p, q_t^y]$  where we concatenate the  $K$  history of the depth images so the policy has some “short-term memory” of the environment. We set  $K = 5$  for all the experiments.

**Action space.** We define a compact action space that only determines the agent’s 2D global movement and the camera direction. Specifically, the action vector is defined as  $\mathbf{a} = [\Delta x, \Delta y, \Delta \theta, \Delta q^p, \Delta q^y]$ ,

which are the relative movements in the forward direction, lateral direction, yaw orientation, camera pitch angle, and camera yaw angle, respectively. We use the action at the current time step to modify the target global configuration and camera angles for the next time step. The next state of the abstract model is simulated by tracking the target global configuration and the camera pose using position control in a physics simulator:  $\mathbf{s}_{t+1} = \mathcal{T}(\mathbf{s}_t, \mathbf{a})$ .

**Reward function.** Unlike policy execution, when we evaluate the reward function during training, we utilize all 3D information relevant to the reward calculation, such as global position and orientation of the agent, the 3D coordinate of the target object and the 3D meshes of the environment. Inspired by the work of Savva *et al.* [SKM\*19], we define the following reward function:

$$r(\mathbf{s}_t) = w_1 r_s(\mathbf{s}_t) + w_2 r_d(\mathbf{s}_t) + w_3 r_l(\mathbf{s}_t) + w_4 r_c(\mathbf{s}_t). \quad (2)$$

$r_s$  is the success reward of 10 when the agent successfully finds the object, which is only awarded once the success-checking terminal condition is invoked.  $r_d$  measures the negative distance to the goal while returns 0 if the goal is not visible, which encourages the agent to move toward and look at the target.  $r_l$  is the live penalty of the value  $-0.1$ . Finally,  $r_c$  checks the collision between the cylinder (the main body) and the environment mesh and penalizes collision by  $r_c(\mathbf{s}_t) = \text{clip}(-0.1n_{col}, -3.0, 0.0)$ , where  $n_{col}$  is the number of collision in the history. For all experiments, we use the same weight vector  $[1.0, 1.0, 1.0, 0.1]$ .

**Initial state distribution.** For training a robust searching behavior, we randomize both the agent’s initial location and the target object’s location at each episode. We collect the candidate locations for both initial positions by sampling random places and filter out the invalid candidates that collide with any other objects. Note that the target object is not necessarily always blocked from the agent’s line of sight and might be visible immediately, which successfully facilitates the learning process.

**Termination conditions.** There are two occasions at which the trajectory rollout can be terminated. First is the time based condition, which caps the total number of actions in the environment at  $T_{max}$  steps, which is set to 100 for all the experiments. Second condition control is a success check, which is triggered once the object is in the nearby proximity of  $0.5m$  and the agent is observing the object, i.e. mask image  $\mathbf{M}$  contains some 1s.

**Domain randomization.** Similar to many sim-to-real transfer learning problems, the success of the search policy depends on whether it can be transferred to the target character with different state and action space. We apply the approach of domain randomization to increase robustness of the search policy when transferred to a different character. We identify that the global vertical position of the character can be drastically different from that of the abstract

model due to its designed body height. We therefore randomize the height of the abstract model during training for the duration of a trajectory rollout in the range of  $[1.0, 1.8]m$ . In addition, legged characters may exhibit natural vertical oscillation during locomotion, while the abstract model moves at a constant height. Therefore, we inject white noises with the range of  $[-0.1, 0.1]m$  to the vertical position of the abstract model at each time step. The randomized vertical movement will affect visual observations by changing the camera position, essential to success transfer of search policy to different characters.

## 4.2. Policy Training Process

Training control policies for simulated character with visual perception input has several challenges. First, these policies usually have complex structures and a large number of parameters, making it computationally expensive to obtain reliable gradients for updating the policy. In addition, learning to extract useful features from images depending solely on the task reward might lead to suboptimal features that do not generalize well to new scenarios.

We train the policy using Soft Actor Critic (SAC) and Contrastive Unsupervised Representations for Reinforcement Learning (CURL), as proposed in the work of Srinivas *et al.* [SLA20]. SAC [HHZ\*18] is an off-policy model-free reinforcement learning algorithm, which has been applied to challenging robotic control problems with desirable sample efficiency [HHZ\*18]. CURL augments the SAC algorithm for learning effective features from high dimensional pixels by jointly optimizing the DRL loss and a contrastive loss to learn a compact latent space.

Specifically, from each input image, CURL will randomly crop two sets of smaller images named queries and keys. These cropped images are then passed through a query encoder and a key encoder to obtain a low-dimensional latent representation of the image  $\mathbf{q}$  and  $\mathbf{k}$ . CURL formulates a contrastive loss:

$$\mathcal{L}_{\mathbf{q}} = \log \frac{\exp(\mathbf{q}^T W \mathbf{k}_+)}{\exp(\mathbf{q}^T W \mathbf{k}_+) + \sum_{i=0}^{K-1} \exp(\mathbf{q}^T W \mathbf{k}_i)}, \quad (3)$$

where  $W$  is a learn-able weight matrix and  $\mathbf{k}_+$  are keys that are from the same time instance of  $\mathbf{q}$ . The contrastive loss encourages the encoded latent features of queries and keys from the same time instance to be close to each other while being far away from the latent features from different time instances under a bilinear product distance. We optimize this contrastive loss jointly with the RL objective (Equation 1). Figure 2 illustrates the data flow of our learning algorithm and indicates the paths where the gradient information is propagated through. We refer readers to the original paper of SAC and CURL for more details about the algorithms.

We represent our search policy as a convolutional neural network. The history of depth image is passed through a Pixel CNN into an embedding of 128 dimensions [SLA20]. The depth image embedding is then concatenated with the mask feature and the agent state, which are then fed into 3 fully connected layers with 1024 neurons to obtain the final action output. We train the search policy using SAC with CURL for 0.75 million simulation samples.

## 5. Full-body Motion Synthesis

The search policy generates global/root configuration and head/camera movements for the abstract model. However, this trajectory cannot be transferred directly to the actual character due to the discrepancy in the state and action spaces and the transition functions, between the actual character and the abstract model. In particular, the policy will receive different input images due to the oscillating head height of the characters and due to the character exhibiting the walking motions and not searching. Furthermore, given the same command, such as moving forward, the characters will achieve different resulting states depending on the transition function of the locomotion model. The differences in the height of the characters can be mitigated by domain randomizing and injecting noise to the height of the abstract model during training to robustify the controller. While, the difference in the head motion will be resolved by querying the controller as if it was to follow the trajectory generated by actual character.

To overcome the discrepancy in the transition function of locomotion model, we employ an online replanning scheme to generate character motions that best match the planned trajectory from the abstract model. Our method is analogous to the model predictive control (MPC) framework in that every time step we replan a long trajectory using the abstract model and execute only a small portion of the trajectory under full-body dynamics. Please note that the trajectories must be collision-free for both abstract and full-body dynamics.

At current time step during testing, we first rollout the search policy (with abstract model) for  $T$  time steps from the current state and observation,  $\mathbf{s}_0$  and  $\mathbf{o}_0$ , by repeatedly querying the policy,  $\mathbf{a}_t = \pi(\mathbf{o}_t)$ , stepping forward in the environment,  $\mathbf{s}_{t+1} = \mathcal{T}(\mathbf{s}_t, \mathbf{a}_t)$ , and observing the environment,  $\mathbf{o}_{t+1} = \text{render}(\mathbf{s}_{t+1})$ , i.e. rendering the environment to images. Note that the state of the abstract model  $\mathbf{s}$  includes the root configuration and the head pose. We pass the planned state trajectory  $\mathbf{s}_0, \dots, \mathbf{s}_T$  to the locomotion generator Phase-Functioned Neural Network (PFNN) [HKS17] or Neural State Machine (NSM) [SZKS19] which generates legged motion on a human character to match the plan. Our method is agnostic to the full-body motion generator if it can synthesize a reasonable full-body motion trajectory for the given abstract plan: we will generally refer them to as “MG” thereafter.

However, the root and body/head states along the full-body trajectory are likely to deviate from the planned state trajectory due to the difference in the transition function between MG and our abstract model. Applying the strategy of online replanning, we only consider the first  $M$  (where  $M \ll T$ ) steps of the full-body trajectory  $\mathbf{q}_0, \dots, \mathbf{q}_{M-1}$  to be valid, and replan the abstract model at the  $M^{\text{th}}$  step.

One issue with our online replanning scheme is that the locomotion generator, MG, does not generate vision perception during motion synthesis, but replanning at the  $M^{\text{th}}$  step requires a short history of depth images as the short-term memory. Furthermore, the planned head motion becomes suboptimal for the searching task since the traversed trajectory deviates from the plan. As such, the vision observations generated by the head poses in  $\mathbf{q}_0, \dots, \mathbf{q}_{M-1}$  can also result in suboptimal next plan. To overcome this issue, we iteratively update the history of abstract model’s state  $\mathbf{s}_t$  from  $t = 0$

to  $t = M - 1$  using the root configuration from the full-body pose  $\mathbf{q}_t$  and the head pose from the search policy. We re-generate observations  $\mathbf{o}_t$  using the modified abstract state and store the depth images in the buffer to recover the “memory lapse” from  $t = 0$  to  $t = M - 1$ . Finally, the abstract model makes a new plan from  $t = M$  to  $t = M + T$  with optimal head movements and restored memory (history of depth images). Our algorithm applied at every  $M$  time steps can be summarized in Algorithm 1.

**Algorithm 1** Online Motion Replanning and Synthesis

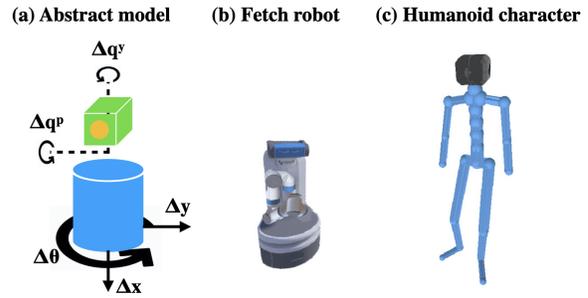
- 1: Input:  $\mathbf{s}_0, \mathbf{o}_0, \mathbf{q}_0$
- 2:  $\mathcal{B} = \{\mathbf{s}_0\}$  ▷ initialize plan buffer
- 3: **for**  $t = 0 : T - 1$  **do** ▷ generate an abstract trajectory
- 4:      $\mathbf{a} = \pi(\mathbf{o}_t)$  ▷ query policy
- 5:      $\mathbf{s}_{t+1} = \mathcal{T}(\mathbf{s}_t, \mathbf{a})$  ▷ advance in environment
- 6:      $\mathbf{o}_{t+1} = \text{render}(\mathbf{s}_{t+1})$  ▷ render observation
- 7:     Store  $\mathbf{s}_{t+1}$  in  $\mathcal{B}$
- 8:  $\mathbf{q}_0, \dots, \mathbf{q}_{M-1} = \text{MG}(\mathcal{B})$  ▷ generate  $M$  human motion poses
- 9: **for**  $i = 0 : M - 1$  **do** ▷ regenerate head orientations & history
- 10:      $\mathbf{s}_i.\text{root} = \mathbf{q}_i.\text{root}$  ▷ update the root position
- 11:      $\mathbf{s}_i = \text{set\_root\_height}(\mathbf{q}_i)$  ▷ update the camera height
- 12:      $\mathbf{o}_i = \text{render}(\mathbf{s}_i)$
- 13:     Update depth image buffer with  $\mathbf{o}_i$
- 14:      $\mathbf{a} = \pi(\mathbf{o}_i)$  ▷ query policy with new observation
- 15:      $\mathbf{s}_{i+1} = \mathcal{T}(\mathbf{s}_i, \mathbf{a})$  ▷ modify camera pose according to  $\pi$
- 16:      $\mathbf{q}_{i+1}.\text{head} = \mathbf{s}_{i+1}.\text{head}$  ▷ overwrite head orientation
- 17:  $\mathbf{o}_M = \text{render}(\mathbf{s}_M)$  ▷ render at last camera pose
- 18: Return  $\mathbf{s}_M, \mathbf{o}_M, \mathbf{q}_0, \dots, \mathbf{q}_{M-1}$

**6. Evaluation**

We evaluate our method on a humanoid character and a wheel-based robot character, Fetch Robot [WFK\*16]. For the humanoid character, we applied two motion generators to synthesize the full-body motion. Figure 3 shows the humanoid character and the robot we use in our experiments. We test our search controller in a modern home scene with an open kitchen and a living area separated by a countertop, as well as a master bedroom connected to a bathroom (Figure 4). To increase complexity of the scene, we add a few open cabinets in which the target objects can be placed. We use iGibson [XSL\*20] environment that provides 3D scans reconstructed from realistic indoor environments and a photorealistic renderer to generate vision inputs to the character. For physics simulation, we use PyBullet [CB17] to simulate the motion of abstract model and to check collision with the environment.

**6.1. Evaluation of Search Policy**

We generate 100 random scenarios to evaluate the success rate of the search policy for an abstract model. At the beginning of each test, the agent is randomly assigned to a collision-free location in the scene with random orientation. Similarly, the target object is placed randomly on any surface in the scene, including the interior of cabinets (Figure 5). If the agent can get close to the target object within 100 steps (~15 seconds), we consider it a successful trial.



**Figure 3:** An abstract model (left) for policy search and two animated characters, the Fetch robot (middle) and a humanoid character (right), for full-body animation.



**Figure 4:** The home scene used in our experiments. (a) a kitchen and a living area separated by a countertop. (b) a bedroom connected to a bathroom.



**Figure 5:** Examples of object placements in our experiments.

Since the search policy will be used by different agents with specific body heights, we evaluate the performance of the search policy by setting the abstract model to three different heights: 1.65m, 1.05m, and 0.45m, where the first two correspond to the height of the human character and the Fetch robot, and 0.45m is selected for further comparison. We compute the success rate of the policy with those three different heights and show the results in Table 1. In general, the advantage of height gives the tall characters a better view

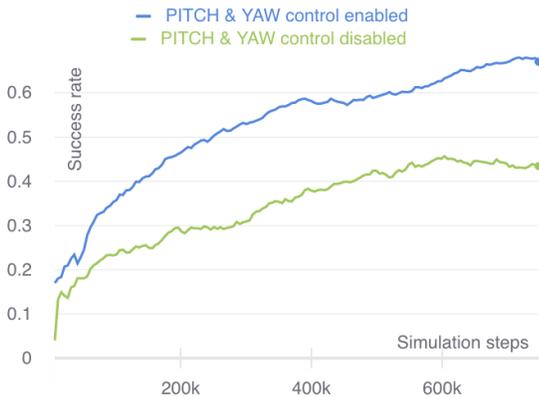
of surfaces while shorter characters struggle to see objects placed on the surface above their heights. As such, there is a near 20% drop in success rate for the shortest character.

We also test the policy on two sampling schemes of the target locations of objects: 1) sampling everywhere except for the inside area of the low cabinets on the floor, and 2) sampling everywhere. The results show that the tall characters (1.05m and 1.65m) perform worse on those challenging cases where objects are inside the low cabinets, while the success rate of the shorter character (0.45m) increases when we allow objects to be sampled inside the cabinet.

Head height	Excluding cabinets	Everywhere
0.45m	50%	58%
1.05m	92%	75%
1.65m	90%	76%

**Table 1:** Performance of the search policy with different camera heights and object locations. In general, abstract models with higher camera locations show better success rates. However, lower cameras are beneficial if we include the challenging case of the object is hidden in the cabinets on the floor.

## 6.2. Comparison to Search Policy without Head Movements



**Figure 6:** Training curves comparison of abstract agents with and without head movements. The result indicates that the agent with head movements shows a 20% higher success rate than the agent without head movements.

Our key hypothesis is that active head movements lead to more effective search behavior by allowing the character to look at different parts of the scene. To evaluate this, we train a baseline search policy for an abstract character without the degrees of freedom to move its head relative to the body. To look around the environment, the agent needs to rotate its entire body around. The result shows that the head movement is crucial to achieve 20% higher success rate in the searching task (Figure 6).

## 6.3. Evaluation of Full-body Characters

To evaluate the performance of our algorithm we use human character with the height of 1.65m. To animate the character as a kine-

matic motion generation model we use two distinct Motion Generators: Phase-Functioned Neural Network (PFNN) [HKS17] and Neural State Machine (NSM) [SZKS19]. We also utilize both to generate legged locomotion for the character and incorporate the discrepancy using the online replanning control scheme described in Section 5. Additionally, we apply the searching policy on another drastically different model, a Fetch Robot with height of 1.05m. Fetch is a wheel-based robot with telescopic degree of freedom to adjust its height. Due to similarity between our control model and Fetch’s differential drive we directly apply the actions produced by the search policy on Fetch.

We show that the search policy can be successfully realized by the full body characters even for challenging cases in which the target objects are placed inside of cabinets, on the other side of room, or occluded by furniture. Different search strategies emerge around different locations in the scene. For example, the character tends to move slower around cabinets and look carefully at the interior part where objects are likely to be placed. The character also utilizes backward steps to have a better view of the surface in front of it. The results are best viewed in the supplemental video where we show the full-body motion in the first-person-view as well as multiple third-person views.

## 6.4. Performance Analysis

To further analyze the performance of our method and understand the choices of different hyperparameters, we construct a few baselines for comparison:

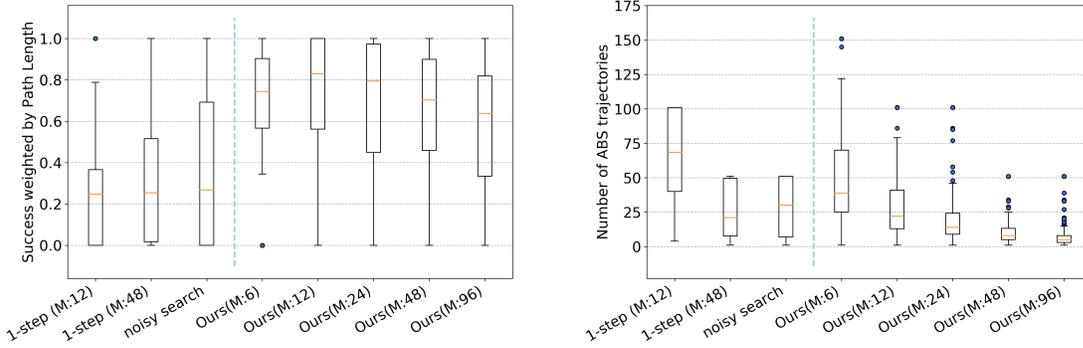
- **1-step:** To evaluate the importance of long-term planning for the searching task, we design a baseline to control the agent by querying the searching policy for one single action and passing it to MG. This baseline does not utilize the long term trajectory planning but instead extrapolates a straight line path in the action direction, with respective orientations. We evaluate two settings with longer and shorter horizons for MG.
- **Noisy Search:** To evaluate the importance of the searching policy before the target object is seen, we create a simple control scheme which in the presence of the object will query the trained searching policy (to use it as an approach mechanism with obstacle avoidance), and when there is no object apply actions sample uniformly at random to move and look around the scene.
- **Ours with different Ms:** To understand the effects of the hyperparameter choices of motion synthesis algorithm, we evaluate multiple values controlling the parameter  $M$  which specifies the number of steps MG returns when attempting to follow the trajectory.

The above baselines and our method are compared using the following metrics:

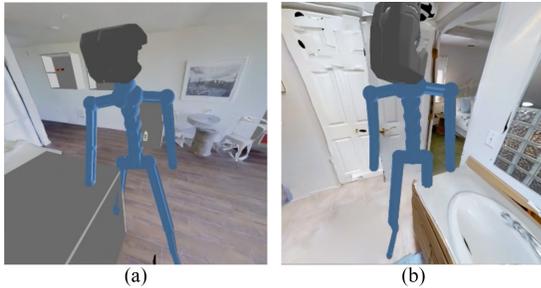
- **Number of attempts:** Number of trajectories generated by the abstract searching policy before the character reached the goal.
- **SPL:** Success weighted by the path length, represented as:

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

where  $\ell_i$  is the shortest distance path length,  $p_i$  is a traversed path length, and  $S_i$  is a binary indicator if the rollout success.



**Figure 7:** Performance of baselines (to the left of dashed line) and our online motion planning with different hyperparameter choices (to the right of dashed line). **left** - Success weighted by the path length - showing the importance of long-term trajectory planning and that simple noise based exploration is not sufficient for finding the object in the designed environment. **right** - shows that a larger number of abstract trajectories generated by abstract model is required for smaller values of  $M$ , which is a trade-off for higher SPL.



**Figure 8:** Failure cases where the character penetrates furniture.

Figure 7 left shows that **1-step** performs poorly for both MG horizons compared to our method with long-term planning. On average, our method performs 40% better than **1-step** on SPL metric. The performance of **noisy search** is also worse than our method by 40%, which shows that structured exploration is required to efficiently find the target object. The performance of our method is sensitive to the choice of the hyperparameter  $M$  (the number of executed frames from the motion plan). In our experiments, smaller values of  $M$  perform better in terms of the SPL because the algorithm replans more frequently, but it also requires more computational costs to generate a larger number of generated abstract trajectories (Figure 7 right).

Most common failure modes for the baselines are related to collisions with the environment obstacles, which causes the character to get stuck in the furniture or the walls. Collisions can also happen to the full-body character if it significantly deviates from the plan. However, we found that collisions of the full-body character occur less frequently because the agent is trained with penalty to avoid the obstacles. Some example scenarios are illustrated in Figure 8.

## 7. Discussion and Conclusion

In this work, we have developed a virtual character exhibiting searching behaviors inspired by the human. We have demonstrated that by training an agent-agnostic search policy and using a re-planning algorithm for transferring the planned abstract motion to actual characters, we can obtain successful and plausible searching behaviors in complex environments. The decomposition of the search task enabled us to reuse the single trained search policy for

characters with different shapes and motor capabilities. Our key insight is that by depriving the privilege 3D information from the character, humanlike behaviors emerge because the character is forced to rely on its own egocentric vision perception and locomotion to complete the task. Furthermore, allowing the head of the character to move independently from the rest of the body leads to more natural searching behaviors, while facilitating the learning of a more effective policy.

One limitation of our current system is the dependence on the mask channel to recognize and identify the target object. This assumption can be lifted by incorporating the state-of-the-art object recognition methods, such as [HGDG17, XGD\*17]. In our current implementation, we chose to not incorporate memory structure in our policy beyond a very short history of the vision observations. We find that reasonable search behavior can be obtained without using memory for the set of environments we are working with. On the other hand, when working with more complex scenes, such as navigating in an entire building, memory becomes essential for localizing the agent and recognizing places that have been explored in the past. On the locomotion side, we notice that the collision checking during search policy learning sometimes is not sufficient when the policy is applied to the full character (Figure 8). A finer resolution collision checking might be needed to further improve the motion quality. Lastly, our scheme performs reasonably well for characters with relatively small variations in height during locomotion such as the ones presented here, however, the characters with more complex dynamics in the head motion or tasks that require active control of the character height might be of a challenge and further research on the topic is necessary.

There are a few promising future avenues for our work. First, is enabled interactions between the character and the environment, such as opening the fridge or drawers. This will allow the emergence of more intricate searching behaviors. Furthermore, our algorithm takes the character from a random location in the room to be in front of the object of interest. This provides an ideal initial pose for the character to perform downstream manipulation tasks such as pouring water into a cup, or pick up a phone. How to incorporate manipulation into our system and achieve more complex human behaviors is thus an important future work direction.

**Appendix A: Hyperparameters**

In the Table. 2 we provide a complete list of the hyperparameters used for training the CURL.

Parameter	Setting
Image Size	100x100
Augmentation	Random Crop (84x84)
Image History Buffer Size	5
Head History Buffer Size	1
Replay buffer	100000
Discount rate $\gamma$	0.99
Number of training steps	0.75M
Batch-size	32
Alpha (SAC) : initial temperature	0.1
Alpha (SAC) : learning rate	0.0001
Alpha (SAC) : optimizer $\beta_1$	0.5
Alpha (SAC) : optimizer $\beta_2$	0.999
Actor : learning rate	0.001
Actor : optimizer $\beta_1$	0.9
Actor : optimizer $\beta_2$	0.999
Actor : Number of layers	4
Actor : Hidden dim	1024
Actor : Activation function	ReLU
Critic : learning rate	0.001
Critic : $\tau$ (polyak averaging)	0.01
Critic : optimizer $\beta_1$	0.9
Critic : optimizer $\beta_2$	0.999
Critic : Number of layers	4
Critic : Hidden dim	1024
Critic : Activation function	ReLU
Encoder (CNN) : learning rate	0.001
Encoder (CNN) : $\tau$ (polyak averaging)	0.05
Encoder (CNN) : Number of layers	4
Encoder (CNN) : Number of filters	32
Encoder (CNN) : Latent dimension	128
Encoder (CNN) : Activation function	ReLU

**Table 2:** A complete set of CURL hyperparameters used to conduct all of the training experiments.

**References**

[ACC\*18] ANDERSON P., CHANG A., CHAPLOT D. S., DOSOVITSKIY A., GUPTA S., KOLTUN V., KOSECKA J., MALIK J., MOTTAGHI R., SAVVA M., ET AL.: On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757* (2018). 2

[AF02] ARIKAN O., FORSYTH D. A.: Interactive motion generation from examples. *ACM Transactions on Graphics (TOG)* 21, 3 (2002), 483–490. 3

[AWL\*19] ABERMAN K., WU R., LISCHINSKI D., CHEN B., COHEN-OR D.: Learning character-agnostic motion for motion retargeting in 2d. *arXiv preprint arXiv:1905.01680* (2019). 3

[CB17] COUMANS E., BAI Y.: Pybullet, a python module for physics simulation in robotics, games and machine learning., 2016–2017. URL: <http://pybullet.org>. 6

[CH05] CHAI J., HODGINS J. K.: Performance animation from low-dimensional control signals. In *ACM SIGGRAPH 2005 Papers*. ACM New York, NY, USA, 2005, pp. 686–696. 3

[CH07] CHAI J., HODGINS J. K.: Constraint-based motion optimization using a statistical dynamic model. In *ACM SIGGRAPH 2007 papers*. ACM New York, NY, USA, 2007, pp. 8–es. 3

[CYT\*18] CLEGG A., YU W., TAN J., LIU C. K., TURK G.: Learning to dress: Synthesizing human dressing motion via deep reinforcement learning. *ACM Transactions on Graphics (TOG)* 37, 6 (2018), 1–10. 3

[EHSN20] EOM H., HAN D., SHIN J. S., NOH J.: Model predictive control with a visuomotor system for physics-based character animation. *ACM Trans. Graph.* 39 (July 2020). 2

[FLFM15] FRAGKIADAKI K., LEVINE S., FELSEN P., MALIK J.: Recurrent network models for human dynamics. In *Proceedings of the IEEE International Conference on Computer Vision* (2015), pp. 4346–4354. 3

[FTFFS19] FANG K., TOSHEV A., FEI-FEI L., SAVARESE S.: Scene memory transformer for embodied agents in long-horizon tasks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2019), pp. 538–547. 2

[GDL\*17] GUPTA S., DAVIDSON J., LEVINE S., SUKTHANKAR R., MALIK J.: Cognitive mapping and planning for visual navigation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017), pp. 2616–2625. 2

[HFV\*19] HE K., FAN H., WU Y., XIE S., GIRSHICK R.: Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722* (2019). 3

[HGDG17] HE K., GKIOXARI G., DOLLÁR P., GIRSHICK R.: Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision* (2017), pp. 2961–2969. 8

[HHZ\*18] HAARNOJA T., HA S., ZHOU A., TAN J., TUCKER G., LEVINE S.: Learning to walk via deep reinforcement learning. *arXiv preprint arXiv:1812.11103* (2018). 5

[HKS17] HOLDEN D., KOMURA T., SAITO J.: Phase-functioned neural networks for character control. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–13. 3, 5, 7

[HSK16] HOLDEN D., SAITO J., KOMURA T.: A deep learning framework for character motion synthesis and editing. *ACM Transactions on Graphics (TOG)* 35, 4 (2016), 1–11. 3

[HZH\*18] HAARNOJA T., ZHOU A., HARTIKAINEN K., TUCKER G., HA S., TAN J., KUMAR V., ZHU H., GUPTA A., ABBEEL P., ET AL.: Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905* (2018). 5

[IKE19] IKEA: How much time do we spend searching for things around the home?, 2019. URL: <https://www.ikea.com/es/en/ideas/how-much-time-do-we-spend-searching-for-things-around-the-home-pubec2a8ae0>. 1

[KGP02] KOVAR L., GLEICHER M., PIGHIN F.: Motion graphs. In *ACM Transactions on Graphics*. ACM New York, NY, USA, 2002, pp. 1–10. 3

[KJL99] KUFFNER JR J. J., LATOMBE J.-C.: Perception-based navigation for animated characters in real-time virtual environments. 2

[KMH\*17] KOLVE E., MOTTAGHI R., HAN W., VANDERBILT E., WEIHS L., HERRASTI A., GORDON D., ZHU Y., GUPTA A., FARHADI A.: AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv* (2017). 2

[KWR\*16] KEMPKA M., WYDMUCH M., RUNC G., TOCZEK J., JAŚKOWSKI W.: Vizdoom: A doom-based ai research platform for visual reinforcement learning. In *2016 IEEE Conference on Computational Intelligence and Games (CIG)* (2016), IEEE, pp. 1–8. 2

[LCR\*02] LEE J., CHAI J., REITSMA P. S., HODGINS J. K., POLLARD N. S.: Interactive control of avatars animated with human motion data. In *Proceedings of the 29th annual conference on Computer graphics and interactive techniques*. ACM, 2002, pp. 491–500. 3

[LWB\*10] LEE Y., WAMPLER K., BERNSTEIN G., POPOVIĆ J., POPOVIĆ Z.: Motion fields for interactive character locomotion. In *ACM SIGGRAPH Asia 2010 papers*. ACM, 2010, pp. 1–8. 3

- [MAO\*19] MANOLIS SAVVA\*, ABHISHEK KADIAN\*, OLEKSANDR MAKSYMETS\*, ZHAO Y., WIJMANS E., JAIN B., STRAUB J., LIU J., KOLTUN V., MALIK J., PARIKH D., BATRA D.: Habitat: A Platform for Embodied AI Research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019). 2
- [MAP\*19] MEREL J., AHUJA A., PHAM V., TUNYASUVUNAKOOL S., LIU S., TIRUMALA D., HEESS N., WAYNE G.: Hierarchical visuomotor control of humanoids. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019* (2019). 3
- [MPV\*16] MIROWSKI P., PASCANU R., VIOLA F., SOYER H., BALLARD A. J., BANINO A., DENIL M., GOROSHIN R., SIFRE L., KAVUKCUOGLU K., ET AL.: Learning to navigate in complex environments. *arXiv preprint arXiv:1611.03673* (2016). 2
- [MTA\*20] MEREL J., TUNYASUVUNAKOOL S., AHUJA A., TASSA Y., HASENCLEVER L., PHAM V., EREZ T., WAYNE G., HEESS N.: Catch and carry: Reusable neural controllers for vision-guided whole-body tasks. *ACM Trans. Graph.* 39 (July 2020). 3
- [MWL\*19] MIN S., WON J., LEE S., PARK J., LEE J.: Softcon: simulation and control of soft-bodied animals with biomimetic actuators. *ACM Transactions on Graphics (TOG)* 38, 6 (2019), 1–12. 3
- [NZC\*18] NAKADA M., ZHOU T., CHEN H., WEISS T., TERZOPOULOS D.: Deep learning of biomimetic sensorimotor control for biomechanical human animation. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–15. 2
- [OLV18] OORD A. V. D., LI Y., VINYALS O.: Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018). 3
- [PALvdP18] PENG X. B., ABBEEL P., LEVINE S., VAN DE PANNE M.: Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–14. 3
- [PBVYDP17] PENG X. B., BERSETH G., YIN K., VAN DE PANNE M.: Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–13. 3
- [PS17] PARISOTTO E., SALAKHUTDINOV R.: Neural map: Structured memory for deep reinforcement learning. *arXiv preprint arXiv:1702.08360* (2017). 2
- [PZI\*19] PAN X., ZHANG T., ICHTER B., FAUST A., TAN J., HA S.: Zero-shot imitation learning from demonstrations for legged robot visual navigation. *arXiv preprint arXiv:1909.12971* (2019). 2
- [SHP04] SAFONOVA A., HODGINS J. K., POLLARD N. S.: Synthesizing physically realistic human motion in low-dimensional, behavior-specific spaces. *ACM Transactions on Graphics (ToG)* 23, 3 (2004), 514–521. 3
- [SKM\*19] SAVVA M., KADIAN A., MAKSYMETS O., ZHAO Y., WIJMANS E., JAIN B., STRAUB J., LIU J., KOLTUN V., MALIK J., PARIKH D., BATRA D.: Habitat: A platform for embodied ai research. In *The IEEE International Conference on Computer Vision (ICCV)* (October 2019). 4
- [SLA20] SRINIVAS A., LASKIN M., ABBEEL P.: Curl: Contrastive unsupervised representations for reinforcement learning. *arXiv preprint arXiv:2004.04136* (2020). 2, 3, 5
- [SYT17] SHIM V. A., YUAN M., TAN B. H.: Automatic object searching by a mobile robot with single rgb-d camera. In *2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)* (2017), IEEE, pp. 056–062. 2
- [SZKS19] STARKE S., ZHANG H., KOMURA T., SAITO J.: Neural state machine for character-scene interactions. *ACM Transactions on Graphics (TOG)* 38, 6 (2019), 1–14. 3, 5, 7
- [TH09] TAYLOR G. W., HINTON G. E.: Factored conditional restricted boltzmann machines for modeling motion style. In *Proceedings of the 26th annual international conference on machine learning* (2009), pp. 1025–1032. 3
- [VSP\*17] VASWANI A., SHAZEER N., PARMAR N., USZKOREIT J., JONES L., GOMEZ A. N., KAISER Ł., POLOSUKHIN I.: Attention is all you need. In *Advances in neural information processing systems* (2017), pp. 5998–6008. 2
- [WFK\*16] WISE M., FERGUSON M., KING D., DIEHR E., DYMESICH D.: Fetch and freight: Standard platforms for service robot applications. In *Workshop on autonomous mobile service robots* (2016). 6
- [WKM\*20] WIJMANS E., KADIAN A., MORCOS A., LEE S., ESSA I., PARIKH D., SAVVA M., BATRA D.: DD-PPO: Learning near-perfect pointgoal navigators from 2.5 billion frames. *International Conference on Learning Representations (ICLR)* (2020). 2
- [WPKL17] WON J., PARK J., KIM K., LEE J.: How to train your dragon: example-guided control of flapping flight. *ACM Transactions on Graphics (TOG)* 36, 6 (2017), 1–13. 3
- [WSY\*18] WANG J., SHIM V. A., YAN R., TANG H., SUN F.: Automatic object searching and behavior learning for mobile robots in unstructured environment by deep belief networks. *IEEE Transactions on Cognitive and Developmental Systems* 11, 3 (2018), 395–404. 2
- [XGD\*17] XIE S., GIRSHICK R., DOLLÁR P., TU Z., HE K.: Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 1492–1500. 8
- [XSL\*20] XIA F., SHEN W. B., LI C., KASIMBEG P., TCHAPMI M. E., TOSHEV A., MARTÍN-MARTÍN R., SAVARESE S.: Interactive gibson benchmark: A benchmark for interactive navigation in cluttered environments. *IEEE Robotics and Automation Letters* 5, 2 (2020), 713–720. 2, 3, 6
- [YLNPI2] YEO S. H., LESMANA M., NEOG D. R., PAI D. K.: Eyecatch: Simulating visuomotor coordination for object interception. *ACM Trans. Graph.* 31, 4 (July 2012). URL: <https://doi.org/10.1145/2185520.2185538>, doi:10.1145/2185520.2185538. 2
- [YTL18] YU W., TURK G., LIU C. K.: Learning symmetric and low-energy locomotion. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–12. 3
- [ZMK\*17] ZHU Y., MOTTAGHI R., KOLVE E., LIM J. J., GUPTA A., FEI-FEI L., FARHADI A.: Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *2017 IEEE international conference on robotics and automation (ICRA)* (2017), IEEE, pp. 3357–3364. 2
- [ZSKS18] ZHANG H., STARKE S., KOMURA T., SAITO J.: Mode-adaptive neural networks for quadruped motion control. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–11. 3
- [ZTB\*17] ZHANG J., TAI L., BOEDECKER J., BURGARD W., LIU M.: Neural slam: Learning to explore with external memory. *arXiv preprint arXiv:1706.09520* (2017). 2