

# X-Nav: Learning End-to-End Cross-Embodiment Navigation for Mobile Robots

Haitong Wang, *Student Member, IEEE*, Aaron Hao Tan, *Student Member, IEEE*, Angus Fung, *Student Member, IEEE*, and Goldie Nejat, *Member, IEEE*

**Abstract**—Existing navigation methods are primarily designed for specific robot embodiments, limiting their generalizability across diverse robot platforms. In this paper, we introduce X-Nav, a novel framework for end-to-end cross-embodiment navigation where a single unified policy can be deployed across various embodiments for both wheeled and quadrupedal robots. X-Nav consists of two learning stages: 1) multiple expert policies are trained using deep reinforcement learning with privileged observations on a wide range of randomly generated robot embodiments; and 2) a single general policy is distilled from the expert policies via navigation action chunking with transformer (Nav-ACT). The general policy directly maps visual and proprioceptive observations to low-level control commands, enabling generalization to novel robot embodiments. Simulated experiments demonstrated that X-Nav achieved zero-shot transfer to both unseen embodiments and photorealistic environments. A scalability study showed that the performance of X-Nav improves when trained with an increasing number of randomly generated embodiments. An ablation study confirmed the design choices of X-Nav. Furthermore, real-world experiments were conducted to validate the generalizability of X-Nav in real-world environments.

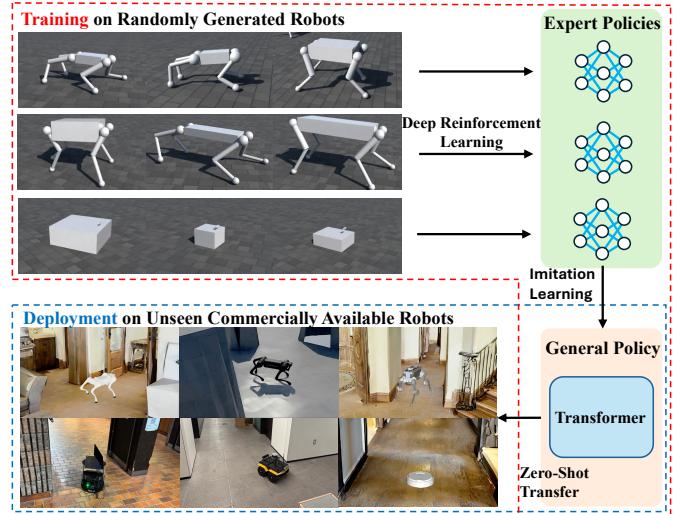
**Index Terms**—Mobile robot navigation, cross-embodiment, expert policy learning, general policy distillation, zero-shot transfer

## I. INTRODUCTION

Robot navigation in diverse and challenging environments is crucial for mobile robots to perform tasks such as person search and detection [1]-[3], exploration in unknown environments [4], [5], and robot-guided navigation [6], [7]. However, existing robot navigation methods are dependent on embodiment-specific kinematics, dynamics, and sensory configurations. Embodiment-specific design limits generalization across embodiments, as policies trained for one robot embodiment often cannot be transferred to other robots with different morphological properties [8].

To address this limitation, cross-embodiment navigation methods train a single generalized policy to be deployed on a wide range of robot embodiments of the same (e.g., wheeled robots [9]) or different robot types (e.g., wheeled and legged robots [10]). Namely, existing methods have mainly used either imitation learning (IL) [9]-[15], or deep reinforcement learning (DRL) [8], [16]-[21]. In particular, IL methods learn a navigation policy from datasets collected on heterogeneous

This work was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC). Corresponding author: Haitong Wang. The authors are with the Autonomous Systems and Biomechatronics Laboratory (ASBLab), Department of Mechanical and Engineering, University of Toronto, Toronto, ON M5S 3G8, Canada  
(e-mail: [haitong.wang@mail.utoronto.ca](mailto:haitong.wang@mail.utoronto.ca); [aaronhao.tan@utoronto.ca](mailto:aaronhao.tan@utoronto.ca); [angus.fung@mail.utoronto.ca](mailto:angus.fung@mail.utoronto.ca); [nejat@mie.utoronto.ca](mailto:nejat@mie.utoronto.ca)).



**Fig. 1.** An overview of X-Nav. X-Nav trains a single end-to-end general navigation policy using randomly generated robot embodiments, which achieves zero-shot transfer to a variety of unseen wheeled and quadrupedal robots in simulated and real-world environments.

robot platforms (e.g., wheeled and quadrupedal robots). On-the-other-hand, DRL methods learn to navigate by training on diverse existing or randomly generated wheeled or quadrupedal robots. However, these methods rely on embodiment-specific modules such as waypoint tracking controller [14], or dynamics models to predict future robot poses [16]. Therefore, the parameters of these embodiment-specific modules need to be tuned for each robot embodiment. Moreover, for DRL methods that focus on locomotion, their policies *only learn to track* given velocity commands and they are unable to *plan velocities*. As a result, these methods require human teleoperation [19] or a separate navigation planner to plan robot velocities [22].

In this paper, we propose X-Nav, a novel two-stage learning framework for end-to-end cross-embodiment navigation, where a *single* generalized navigation policy is developed that can be deployed on a wide range of mobile robots, Fig. 1. In particular, we focus on wheeled and quadrupedal robots as they are the most widely used robot mobility types for navigation. In particular, wheeled robots have energy-efficient motion on flat terrains, whereas quadrupeds can traverse across diverse and challenging terrains [23]. Our navigation policy directly maps visual and proprioceptive observations to executable control commands without relying on embodiment-specific models.

The key contributions of this work are: 1) the development of the *first* end-to-end cross-embodiment navigation approach, which can be deployed on different types of mobile robots including wheeled and quadrupedal robots; 2) the introduction of a novel two-stage learning framework that integrates both expert policy learning using DRL with privileged observations,

and general policy distillation using IL with a transformer model. The general policy implicitly infers embodiment information from proprioception, which enables zero-shot generalization to unseen robot embodiments.

## II. RELATED WORKS

We discuss the existing literature on cross-embodiment navigation that have used: 1) IL [9]-[15], or 2) DRL [8], [16]-[21].

### A. Imitation Learning-based Methods

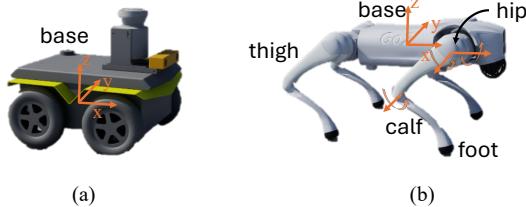
IL-based methods learn a cross-embodiment navigation policy by training on datasets collected from heterogeneous robot platforms in real-world environments. These methods use robot visual observations (i.e., RGB images) and goal location images as inputs. They then extract spatial features from these input images using visual encoders such as EfficientNet [24] and ResNet [25]. The extracted image features are used to generate relative waypoints [10]-[13] or velocities [9], [15]. This is achieved through the use of fully connected layers (FCLs) [9], [11], [15], transformer blocks, [10], [12]-[14] or diffusion action heads [10], [13]. FCLs directly map the image features to robot actions. Transformer blocks use self-attention layers to account for the spatial-temporal features of robot observations [26]. Diffusion head generates robot actions by progressively denoising action vectors with image features as conditioning [27]. Then, embodiment-specific controllers are used to track the generated waypoints or velocities.

IL methods have been trained on datasets containing robot navigation trajectories collected on various robot platforms in real-world environments. These datasets typically include RGB image observations and robot poses collected in indoor environments (e.g., SACSoN [28]), or outdoor off-road and sidewalk environments (e.g., GNM [11]). Evaluation of the IL methods were conducted in both indoor [9]-[15] and outdoor [9]-[13] real-world environments using different wheeled robots [9]-[14], and/or quadruped robots [10], [15].

### B. Deep Reinforcement Learning-based Methods

DRL-based methods consist of: 1) hierarchical [16], [17], and 2) monolithic methods [8], [18]-[21].

Hierarchical methods integrate high-level planning and low-level control modules. For example, in [16], a model predictive control (MPC) framework consisting of a dynamics module and a perception module was used for wheeled robot navigation in outdoor environments. In [17], a quadrupedal robot navigation system was developed using the SAC+AE [29] DRL method. It consisted of a general high-level and a low-level policy. The high-level policy used a robot embedding network to generate robot-specific embeddings and a multi-layer perceptron (MLP)



**Fig. 2.** The base frame of a: (a) wheeled robot, and (b) quadrupedal robot.

to generate robot base velocities. The low-level policy used an MPC policy to track robot velocities. It was evaluated in indoor house environments using different quadrupeds.

Monolithic methods have been used to map observations such as proprioception (e.g., joint positions, joint velocities) [8], and joint descriptions (e.g., torque limits, velocity limits) [18] directly to robot actions. The generated actions are either high-level discrete actions (move forward, turn left) for wheeled robots [21] or low-level continuous actions (joint positions) for quadrupeds [8], [18]-[20]. To learn policies that generalize across different embodiments, these methods were trained using randomly generated [8], [19]-[21] or existing robot embodiments [18]. Monolithic methods have mainly been trained using Proximal Policy Optimization (PPO) [30]. They have been evaluated in both indoor [8], [20], [21] and outdoor [18], [19] environments using various wheeled robots [21] or quadrupedal robots [8], [18]-[20].

### C. Summary of Limitations

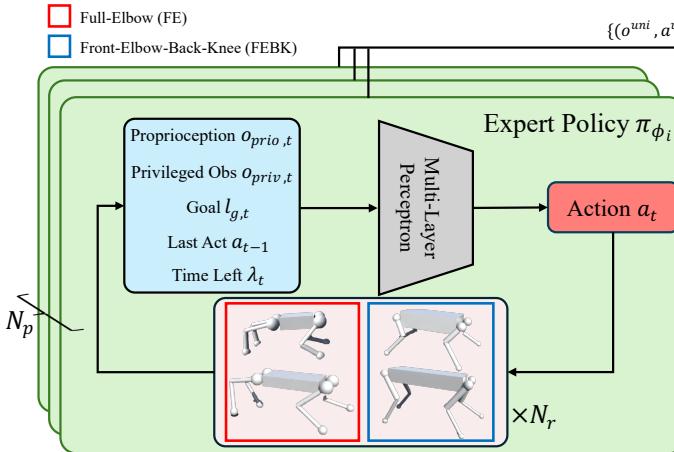
The aforementioned IL-based methods depend on embodiment-specific controllers to track the predicted waypoints [10], [11] or velocities [9], [15]. Each robot embodiment requires individualized tuning of these controllers to ensure accurate execution of navigation commands. DRL-based methods also require embodiment-specific models to either generate low-level control commands [16], predict future robot poses [17], or execute high-level discrete actions [21], which requires embodiment-specific tuning when deployed on new embodiments. Furthermore, DRL methods that have focused only on quadruped locomotion primarily learn to track velocity commands but lack the ability to plan robot velocities [8], [18]-[20]. Therefore, they require manual teleoperation [19] or external velocity planner to perform navigation [22].

To address the above limitations of relying on embodiment-specific modules and the lack of velocity planning, X-Nav provides an end-to-end navigation approach that directly maps robot observations to executable low-level commands, removing the need for an embodiment-specific module to plan velocities or track waypoints. This is achieved through a novel two-stage training framework, where multiple expert policies are trained on a diverse set of randomly generated robot embodiments, and distilled into a general navigation policy, enabling zero-shot generalization to unseen embodiments.

## III. PROBLEM FORMULATION

The cross-embodiment navigation problem requires a mobile robot  $r$  to navigate from a starting position  $l_s \in \mathbb{R}^2$  to a given goal position  $l_g \in \mathbb{R}^2$  in an unknown cluttered environment. The robot navigation is guided by: 1) visual observations which are represented by depth images  $o_{depth} \in \mathbb{R}^{H \times W}$  from an onboard depth camera, 2) the 2D goal position  $l_g$ , and 3) the proprioceptive observations  $o_{prio}$  obtained from an inertial measurement unit (IMU) and motor encoders. Based on these observations, the robot executes an action  $a$  at each timestep  $t$ . The objective of the robot  $r$  is to minimize its travel distance between  $l_s$  and  $l_g$ :

## Stage 1: Expert Policy Learning



**Fig. 3.** The proposed X-Nav framework consists of two stages: 1) Expert Policy Learning, and 2) General Policy Distillation. In Stage 1, multiple expert policies are trained on randomly generated robot embodiments using DRL with privileged observations. In Stage 2, the knowledge of expert policies is distilled into a single general policy using IL.

$$\min \mathbb{E}[d_{nav}(l_s, l_g)], \quad (1)$$

where  $d_{nav}(\cdot, \cdot)$  is a function representing travel distance.

We consider both wheeled and quadruped robots. For a wheeled robot  $r \in \mathcal{R}_{wheel}$ , the proprioception  $o_{prio}^{wheel} \in \mathbb{R}^2$  is a 2D vector representing robot linear and angular velocity in the robot base frame (Fig. 2 (a)), and action  $a^{wheel} \in \mathbb{R}^2$  is a 2D vector representing the desired linear and angular velocity. For a quadruped  $r \in \mathcal{R}_{quad}$ , the proprioception  $o_{prio}^{quad} \in \mathbb{R}^{30}$  is a 30D vector representing the concatenation of the robot base velocity  $o_{prio,vel}^{quad} \in \mathbb{R}^3$ , the gravity projected in the robot base frame  $o_{prio,g}^{quad} \in \mathbb{R}^3$ , the positions of the 12 joints (i.e., hip, thigh and knee joints)  $o_{prio,jp}^{quad} \in \mathbb{R}^{12}$  and the corresponding joint velocities  $o_{prio,jv}^{quad} \in \mathbb{R}^{12}$ . The robot base frame is defined with its origin at the center of the robot trunk, Fig. 2 (b). The action  $a^{quad} \in \mathbb{R}^{12}$  represent the 12 desired joint positions. A Proportional-Derivative (PD) controller is used to track the desired position for each joint of a quadruped robot [31]:

$$\tau = K_p(a^{quad} - q^{quad}) - K_d\dot{q}^{quad}, \quad (2)$$

where  $\tau$  is the motor torque,  $K_p$ ,  $K_d$  are the PD gains, and  $q^{quad}$  represents the actual joint position.

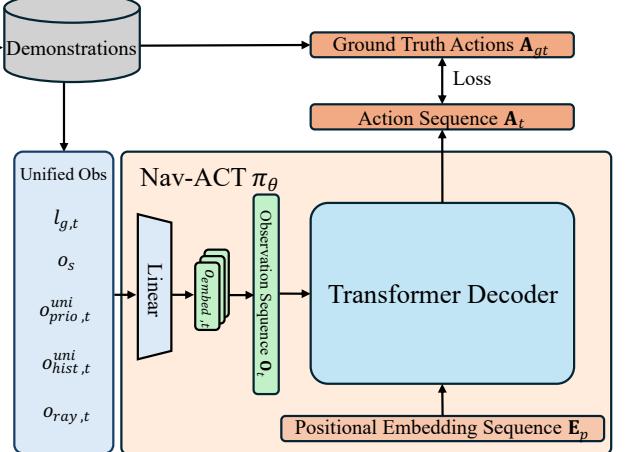
## IV. X-NAV ARCHITECTURE

The proposed X-Nav architecture, shown in Fig. 3, consists of two stages: 1) Expert Policy Learning, and 2) General Policy Distillation. In Stage 1, randomized robot embodiments are used to train expert policies using DRL with privileged observations. In Stage 2, demonstrations generated by the expert policies are collected and utilized to train a single general navigation policy using a transformer model.

### A. Stage 1: Expert Policy Learning

In Stage 1, we train  $N_p$  expert policies  $\pi_{\phi_i}$  ( $i = 1, \dots, N_p$ ) in simulation using DRL. Each expert policy is trained on  $N_r$  different robot embodiments of a same type.

## Stage 2: General Policy Distillation



### 1) Observation and Action Space

At each timestep  $t$  ( $0 \leq t < T$ ), the policy observation is  $o_t = [a_{t-1}, o_{prio,t}, l_{g,t}, \lambda_t, o_{priv,t}]$ , where  $a_{t-1}$  denotes the last robot action at timestep  $t-1$ ,  $\lambda_t = 1 - \frac{t}{T}$  denotes the normalized time remaining in the episode,  $T$  is the maximum timesteps of the episode.  $o_{priv,t} = [o_{embod}, o_{scan,t}]$  denotes the privileged observations which include embodiment parameters  $o_{embod}$  and ground-truth terrain height scans  $o_{scan,t}$ . For wheeled robots, embodiment parameters include the robot mass and its 3D size. For quadrupeds, embodiment parameters include base size, base mass, thigh size, thigh mass, calf size, calf mass, and motor PD gains  $K_{p/d}$ .  $o_{scan}$  represents the terrain height of a local region surrounding the robot. The actions of the expert policies (i.e.,  $a^{wheel}$  or  $a^{quad}$ ) are defined above in Section III.

### 2) Embodiment Randomization

Embodiment randomization is implemented to ensure that real-world robot embodiments are within the distribution of the randomly generated robots used during training [8]. We generate random wheeled and quadrupedal robots by sampling  $o_{embod}$ . To ensure that the randomly generated quadrupeds can generate sufficient motor torques, we design a robot template and utilize it to compute motor PD gains [8]:

$$K_{p/d} = v_{temp} \times \frac{m_{gen}}{m_{temp}} \times v, \quad (3)$$

where  $v$  denotes the randomly sampled value,  $v_{temp}$  denotes the PD gains of the robot template,  $m_{gen}$  denotes the total mass of the generated robot, and  $m_{temp}$  denotes the total mass of the robot template. Furthermore, we consider two leg configurations for the generated quadruped embodiments, Fig. 3: Full-Elbow (FE), and Front-Elbow Back-Knee (FEBK) as they are the most commonly used configurations [8].

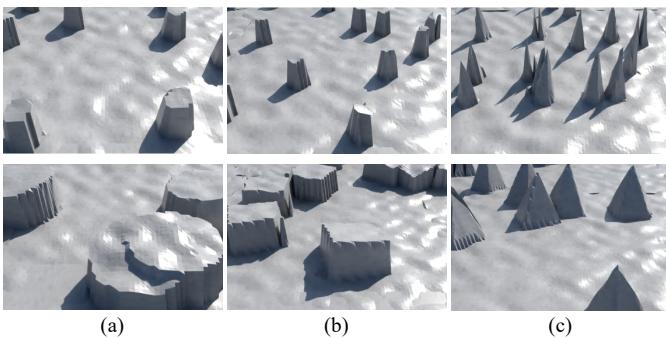
### 3) Deep Reinforcement Learning

We utilized the PPO DRL algorithm [30] to train the expert policies due to its stability and efficiency in learning robot

TABLE I: REWARD FUNCTION

Reward	Expression	Explanation
$r_{pos,(\text{soft/hard})}$	$r_{pos,(\text{soft/hard})} = \frac{c_{1,(\text{soft/hard})}}{1 + \left\  \frac{d_{goal}}{c_{2,(\text{soft/hard})}} \right\ ^2} \cdot \mathbb{1}(t > T - T_r)$	Encourages the robot to move to the goal position [32]. $T_r$ denotes the duration of timesteps for which $r_{pos}$ is activated. $r_{pos,\text{soft}}$ is a dense reward used to encourage exploration using a larger $c_1, c_{2,\text{soft}}$ , while $r_{pos,\text{hard}}$ is a sparse reward for accurate position tracking using a smaller $c_1, c_{2,\text{hard}}$ .
$r_{task}$	$r_{fwd} = c_{fwd} \cdot \text{CLIP}\left(\text{ReLU}\left(\frac{v_x}{v_{max}}\right), 1\right) \cdot \mathbb{1}(\delta_{goal} < \sigma_{direct})$ $r_{stop} = \frac{c_{1,stop}}{1 + \left\  \frac{o_{prio,vel}}{c_{2,stop}} \right\ ^2} \cdot \mathbb{1}(t > T - T_r) \cdot \mathbb{1}(d_{goal} < \sigma_{hard})$	Encourages the robot to move forward [33]. $\delta_{goal}$ is the heading error relative to the goal, and $\sigma_{direct}$ is an angle threshold for the robot heading in the correct direction.
	$r_{collide} = c_{collide} \cdot \mathbb{1}_{collide}$	Ensures the robot stops at its goal position. $\sigma_{hard}$ is a distance threshold used to determine when the robot has reached its goal position.
$r_{wheel\_reg}$	$r_{wheel\_reg} = c_a \ \dot{a}\ ^2$ $r_{v_z} = c_{v_z} v_z^2$ $r_{\omega} = c_{\omega} (\omega_{xy}^2)$ $r_{\tau} = c_{\tau} \ \tau\ ^2$ $r_{\ddot{q}} = c_{\ddot{q}} \ \ddot{q}\ ^2$ $r_{flat} = c_{flat} \ o_{prio,g,xy}\ ^2$ $r_{air} = c_{air} \sum(t_{air,f} - 0.5), f \in \{1,2,3,4\}$	Penalizes collisions. $\mathbb{1}_{collide}$ indicates whether a collision has happened between the robot and obstacles. $\dot{a}$ denotes the action rate. $v_z$ is the vertical velocity of the robot base, and $\omega_{xy}$ is the angular velocity of the base around $x$ and $y$ axes. $\tau$ is the joint torque. $\dot{a}$ is the action rate. $\ddot{q}$ is the joint acceleration. $o_{prio,g,xy}$ is the $x, y$ component of the projected gravity. $t_{air,f}$ is the time duration when the robot foot $f$ is in the air.

\*  $c_{1,(\text{soft/hard})}, c_{2,(\text{soft/hard})}, c_{fwd}, c_{1,stop}, c_{2,stop}, c_{collide}, c_a, c_{v_z}, c_{\omega}, c_{\tau}, c_{\ddot{q}}, c_{flat}, c_{air}$  are all constants.



**Fig. 4.** The obstacles used for expert policy learning. (a) cylinders, (b) boxes, (c) pyramids.

locomotion and navigation [22], [31]. Each policy network  $\pi_{\phi_i}$  was implemented as an MLP. We perform massively parallel training [31] by simultaneously training on  $N_r$  robots.

#### 4) Rewards

The reward function,  $r$ , has been designed to encourage mobile robots to perform obstacle-free navigation and promote smooth and efficient movements.  $r$  includes task rewards  $r_{task}$ , and regularization rewards  $r_{reg}$ :  $r = r_{task} + r_{reg}$ . Wheeled and quadrupedal robots share the same task rewards but different regularization rewards  $r_{wheel\_reg}$ , and  $r_{quad\_reg}$ . Table I presents the definition of all rewards.

#### 5) Curriculum Learning

We use a game-inspired curriculum [31] to progressively train the navigation policy from simple to complex environments. The overall environment consists of  $n_x \times n_y$  subfields. Each subfield contains randomly generated obstacles of one of the three types: box, pyramid, and cylinder, Fig. 4. These fundamental geometric primitives are used as approximations of common real-world structures such as walls, beds, desks. Each obstacle type has difficulty levels represented by its density and the obstacle size. Terrain roughness is also varied across the difficulty levels. Training starts at the lowest difficulty level. The progression of robots through the difficulty levels is governed by two distance thresholds  $\sigma_{close}$  and  $\sigma_{far}$  ( $\sigma_{close} < \sigma_{far}$ ). When a robot reaches its goal position at the end of an episode ( $d_{goal} < \sigma_{close}$ ), the robot is promoted to a higher

difficulty level. Conversely, if the goal distance exceeds  $\sigma_{far}$  ( $d_{goal} > \sigma_{far}$ ), the robot is demoted to a lower difficulty level.

#### 6) Domain Randomization

We use domain randomization [34] to improve robustness of the expert policies for sim-to-real transfer. First, friction coefficients and robot base mass are sampled using a uniform distribution from predefined ranges at the beginning of each episode. Second, push disturbances are simulated for quadruped training by adding a randomly sampled velocity to the robot base every  $T_{push}$  seconds. Third, random noise is added to the proprioceptive observations  $o_{prio}$  and actions  $a$ .

#### B. Stage 2: General Policy Distillation

We distill the trained expert policies into a single general policy  $\pi_\theta$  using imitation learning. We first define a unified observation and action space, then collect demonstrations using the expert policies, and train  $\pi_\theta$  on these demonstrations.

##### 1) Unified Observation and Action Space

The unified observation  $o_t^{uni}$  at timestep  $t$  is:

$$o_t^{uni} = [l_{g,t}, o_s, o_{prio,t}^{uni}, o_{hist,t}^{uni}, o_{ray,t}], \quad (4)$$

where  $o_s \in \mathbb{R}^2$  is the robot base length and width, and  $o_{prio,t}^{uni}$  is the unified proprioception. For wheeled robots, where joint position and velocity data are absent from proprioception  $o_{prio}^{wheel}$ , we pad the dimensions with zeros to match  $o_{prio}^{quad}$ .  $o_{hist,t}^{uni}$  is the concatenation of the last  $N_{hist}$  frames of  $a^{uni}$  and  $o_{prio}^{uni}$ .  $o_{ray,t}$  is the unified laser rays for distance measurement derived from the raw depth image  $o_{depth}$ . We project depth pixels of  $o_{depth}$  to a horizontal plane and extract distances at evenly spaced angles. The resulting laser scan is then interpolated into a unified format with a field of view (FOV) of  $\theta_{fov}$ . This ensures consistency across different cameras.

The unified action,  $a^{uni} \in \mathbb{R}^{14}$ , is a 14D vector with the first two dimensions representing the linear and angular velocity of wheeled robots, and the last 12 dimensions representing the target joint positions for quadrupeds.

##### 2) Demonstration Collection

$N_d$  demonstrations are obtained for each of the expert policies  $\pi_{\phi_i}$  trained in Stage 1.  $N_d$  is identical to  $N_r$  to ensure

TABLE II: PARAMETERS FOR EMBODIMENT RANDOMIZATION

Type	Parameters	Range	Parameters	Range
Small-Sized Quadrupeds	Base length	[0.24, 0.91] m	Thigh mass	[0.56, 1.69] kg
	Base width	[0.16, 0.39] m	Calf radius	[0.02, 0.05] m
	Base height	[0.06, 0.21] m	Calf length	[0.12, 0.39] m
	Base mass	[4.8, 19.5] kg	Calf mass	[0.12, 0.39] kg
	Thigh radius	[0.02, 0.05] m	Motor P gain	[0.7, 1.3]
	Thigh length	[0.16, 0.46] m	Motor D gain	[0.7, 1.3]
Large-Sized Quadrupeds	Base length	[0.56, 1.04] m	Thigh mass	[2, 5.2] kg
	Base width	[0.28, 0.52] m	Calf radius	[0.02, 0.04] m
	Base height	[0.14, 0.26] m	Calf length	[0.24, 0.36] m
	Base mass	[24, 39] kg	Calf mass	[0.4, 0.6] kg
	Thigh radius	[0.03, 0.05] m	Motor P gain	[0.5, 1.3]
	Thigh length	[0.24, 0.39] m	Motor D gain	[0.5, 1.3]
Wheeled Robots	Base length	[0.3, 0.8] m	Base height	[0.15, 0.3] m
	Base width	[0.2, 0.65] m	Base mass	[5, 20] kg

the dataset covers the full diversity of robot embodiments. Each demonstration contains a sequence of unified observations  $o_t^{uni}$  and actions  $a_t^{uni}$  collected from a successful episode in a randomly generated environment.

### 3) Navigation Action Chunking with Transformer

We introduce Navigation Action Chunking with Transformer (Nav-ACT), a transformer model to generate  $a_t^{uni}$  from  $o_t^{uni}$  for cross-embodiment navigation, Fig. 3. Nav-ACT consists of: 1) an encoder to convert the unified observations  $o_t^{uni}$  into embeddings, and 2) a transformer decoder to generate an action chunk conditioned on the observation sequence.

At each timestep  $t$ , the unified observation  $o_t^{uni}$  is converted into an embedding  $o_{embed,t}$  using a linear encoder.  $S_o$  consecutive steps of observation embeddings  $o_{embed}$  are stacked to construct an observation sequence  $\mathbf{O}_t = [o_{embed,t-S_o+1}, \dots, o_{embed,t}]$ . The transformer decoder takes as input a sequence of  $S_a$  learnable positional embeddings  $E_p$ . It then uses the observation embedding sequence  $\mathbf{O}_t$  for cross-attention computation to generate an action sequence  $\mathbf{A}_t$ ,  $\mathbf{A}_t = [a_{t-S_a+1}, \dots, a_t]$ , where  $S_a$  denotes the sequence length.  $\mathbf{A}_t$  is generated in a single forward pass for computation efficiency. Nav-ACT is trained using the mean squared error (MSE) loss:

$$\mathcal{L}_{Nav-ACT} = \frac{1}{S_a} \sum_{i=1}^{S_a} \left\| \mathbf{A}_{pred}^{(i)} - \mathbf{A}_{gt}^{(i)} \right\|^2, \quad (5)$$

where  $\mathbf{A}_{pred}$  and  $\mathbf{A}_{gt}$  are the predicted and ground truth actions.

### 4) Inference

At inference time, for: 1) wheeled robots, we use temporal ensemble (TE) [35] to improve smoothness and avoid jerky movements, and 2) for quadrupeds, we disable TE and only take the first action from the action sequence  $\mathbf{A}_t$  to ensure real-time adaptation to dynamic changes during navigation. Namely, for wheeled robots, at each timestep  $t$ , Nav-ACT generates action predictions  $\mathbf{A}_t$ , then the last  $S_a$  predictions  $\mathbf{A}_i$ ,  $i = (t - S_a + 1, \dots, t)$  are taken to compute the action prediction for the current timestep  $t$  from each  $\mathbf{A}_i$ . Weighted average is applied to generate the action:

$$a_t = \frac{\sum_{i=t-S_a+1}^t w_i \mathbf{A}_i}{\sum_{i=t-S_a+1}^t w_i}, \quad (6)$$

TABLE III: PARAMETERS AND VALUE RANGES FOR DOMAIN RANDOMIZATION

Parameters	Ranges	Parameters	Ranges
Static friction	[0.7, 1.1]	Push disturbance	[-0.5, 0.5] m/s
Dynamic friction	[0.6, 1.0]	Push Interval	[4, 8] s
Added mass	[0.0, 2.0] kg	Ray distance noise	[-0.1, 0.1] m
Linear velocity noise	[-0.1, 0.1] m/s	Angular velocity noise	[-0.1, 0.1] rad/s
Projected gravity noise	[-0.05, 0.05] m/s <sup>2</sup>	Joint velocity noise	[-1.0, 1.0] rad/s
Joint position noise	[-0.01, 0.01] rad		

TABLE IV: HYPERPARAMETERS AND VALUES

Parameter	Value	Parameter	Value	Parameter	Value
$x_{scan}$	2.5 m	$y_{scan}$	2.5 m	$c_{collide}$	-40
$c_{1,soft}$	10	$c_{2,soft}$	5	$T_r$	1.5 s
$c_{1,hard}$	15	$c_{2,hard}$	0.5	$c_{air}$	0.5
$v_{max}$	1.0 m/s	$c_\omega$	-0.05	$T$	8 s
$\sigma_{direct}$	1.75 rad	$c_\tau$	-0.0002	$c_{v_z}$	-2
$\sigma_{hard}$	0.5 m	$c_a$	-0.01	$\sigma_{close}$	0.5 m
$c_{1,stop}$	10	$c_q$	-2.5e-7	$\sigma_{far}$	3.0 m
$c_{2,stop}$	0.2	$c_{flat}$	-5	$c_{fwd}$	2

where  $w_i = \exp(-k * (i - t + S_a - 1))$  and  $k$  is a positive constant. For quadrupeds, at each timestep, we take the first action from the latest action sequence,  $a_t = \mathbf{A}_t[0]$ .

## V. POLICY TRAINING

### A. Expert Policy Learning

1) *Setup*: We trained  $N_p = 3$  expert policies using DRL, with one policy for each of the following types: small-sized quadrupeds (i.e., mass < 30kg), large-sized quadrupeds (i.e., mass > 30kg) and wheeled robots. Quadrupeds were split into small and large categories due to the significant differences in their dynamics (e.g., inertia, torque) to reduce variability and promote stable training. For each robot type,  $N_r = 4096$  robot embodiments were generated by sampling parameters from Table II. Domain randomization followed the ranges defined in Table III. For curriculum learning, the overall training environment consisted of 384 subfields arranged in 6 rows and 64 columns, with each subfield measuring 10 × 10 m. Terrain roughness was increased from 0 to 8 cm across the difficulty levels. Rough terrain was applied to quadruped training to increase its generalizability and robustness. The other hyperparameters are defined in Table IV. The values of the hyperparameters were empirically determined during training. Each policy network has (1024, 512, 256) units. The Isaac Sim simulator and Isaac Lab framework [36] were used for training.

2) *Training*: All training was done on a workstation with an NVIDIA RTX 4090 GPU, an Intel Core i9-13900KF CPU and 32GB RAM. Each expert policy was trained with a batch size of 24576 for 4000 epochs. The Adam optimizer [37] was used with a learning rate of 0.001.

### B. General Policy Distillation

1) *Setup*: For observation, we used  $N_{hist}$  of 5,  $\theta_{fov}$  of 90°, the minimum and maximum ray distance of 0.2 m and 8 m, and number of laser rays of 128. For demonstration collection,  $N_d = 4096$  demonstrations were collected for each expert policy. We set  $S_a$  as 6, and  $S_o$  as 4. Nav-ACT has 4 transformer layers and 4 heads with an embedding size of 256, totaling 4.85M parameters.

TABLE V: COMPARISON STUDY WITH UNSEEN ROBOT EMBODIMENTS

Method	Go2		A1		ANYmal B		Jackal		Dingo		Create3	
	SR	SPL										
BC	63.6	0.55	70.2	0.60	45.1	0.40	57.5	0.51	31.2	0.24	46.7	0.36
BCT	66.8	0.58	70.2	0.59	45.6	0.43	41.6	0.34	16.0	0.11	16.8	0.12
DP	64.7	0.56	71.1	0.59	36.2	0.32	66.6	0.60	18.6	0.14	20.2	0.16
CP	27.1	0.24	40.2	0.36	33.0	0.30	51.2	0.45	15.4	0.11	20.9	0.16
<b>X-Nav</b>	<b>69.1</b>	<b>0.60</b>	<b>71.6</b>	<b>0.61</b>	<b>59.6</b>	<b>0.52</b>	<b>90.4</b>	<b>0.84</b>	<b>83.5</b>	<b>0.76</b>	<b>93.1</b>	<b>0.85</b>

2) *Training*: Nav-ACT was trained with a batch size of 256, learning rate of 0.0001 for 100 epochs. The Adam optimizer [37] was used with weight decay of 0.001.

## VI. SIMULATED EXPERIMENTS

We evaluated the performance of X-Nav by conducting: 1) a comparison study with state-of-the-art (SOTA) learning methods to assess the generalizability of X-Nav on unseen robot embodiments, 2) a scalability study to explore X-Nav’s performance when trained with increasing numbers of random robot embodiments, and 3) an ablation study to investigate the design choices of X-Nav.

### A. Comparison Study with SOTA

We used commercially available robot platforms that were unseen during the training of X-Nav. For wheeled robots, the Clearpath Jackal, Clearpath Dingo, and iRobot Create3 were used. For the quadrupeds, the Unitree A1, Unitree Go2, and ANYmal B were used. Each robot was equipped with a front-facing depth camera with a FOV of 90°. The performance metrics utilized were the: 1) success rate (SR), and 2) the success weighted by the normalized inverse path length (SPL) [38] for measuring the trajectory efficiency. A hundred unseen in-distribution environments were randomly generated using the same obstacles in Section IV.A.5. Quadrupeds were tested on rough terrain and wheeled robots were tested on flat terrain. 1) *Comparison Methods*: We compared our X-Nav method with the following four SOTA IL learning methods. For comparison purposes, we extended *these* methods for the end-to-end cross-embodiment navigation problem by incorporating the same observation and action representation as X-Nav. All methods were trained using the same demonstrations as X-Nav.

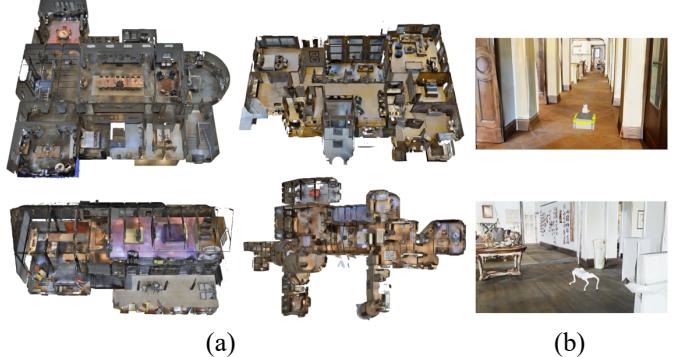
**Behavior Cloning (BC)** [39]: The BC method uses an MLP to generate a robot action based on the current observation. BC was selected as a standard baseline method.

**Behavior Cloning with Transformer (BCT)** [26]: The BCT method uses a transformer decoder to generate robot actions. The transformer decoder takes the last 6 observations and actions as input to generate the next action. BCT was selected as a representative transformer-based method.

**Diffusion Policy (DP)** [27]: The DP method uses the EDM scheduler [40] and a transformer decoder to generate actions. DP used 80 steps for training and 10 steps for inference. DP was selected for its ability to model multi-modal action distributions.

**Consistency Policy (CP)** [41]: The CP method uses a transformer decoder to generate actions. It uses DP as the teacher model. CP was selected for its computational efficiency during inference by using one-step denoising.

2) *Procedure*: We conducted 3000 trials for each method and robot embodiment. At the beginning of each trial, both robot starting and goal positions were randomly placed. A robot



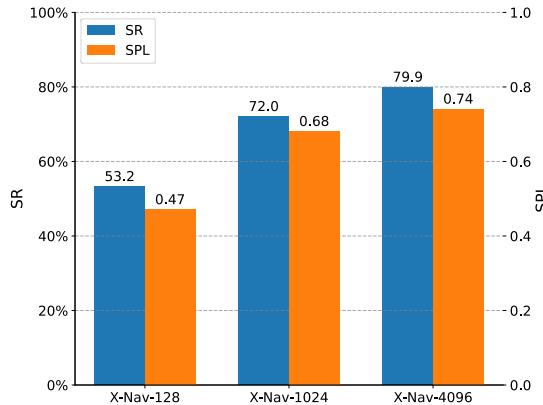
**Fig. 5.** (a) Top view of photorealistic indoor environments from MP3D dataset, (b) Jackal and Go2 deployed inside the environments.

successfully achieved the goal when the distance between its current and goal position was within 0.5m at the end of a trial. Each timestep was 0.02s. A trial terminated when the total timesteps exceeded 750.

3) *Results*: The SR, SPL of X-Nav and comparison methods are presented in Table V. X-Nav achieved the highest SR and SPL compared to four SOTA methods. BC achieved lower SR and SPL than X-Nav across all embodiments. This is due to BC generating the next robot action solely based on the latest observation without accounting for history of observations, which is essential for spatial understanding to avoid getting stuck in local minima such as dead ends. X-Nav implicitly predicts future states by generating an action sequence, while BCT only predicts the next immediate action which can lead to suboptimal trajectory for navigation. The inference of DP is computationally intensive due to the iterative denoising process [27]. Though DP has been used for quadruped locomotion task [42], its deployment requires external desktop-level GPU, which limits its integration on resource-constrained robot platforms for general navigation applications. CP, on the other hand, has faster inference speed than DP by directly generating the denoised action through one single forward pass. However, both DP and CP performed worse than X-Nav. We postulate that this is due to the diffusion-based methods not being able to accurately model the two distinct action distributions of wheeled and quadrupedal robots at the same time. Specifically, the actions of quadrupeds require faster phase changes to switch between swing and stance phases to maintain stability, while wheeled robots require smoother and more consistent actions.

### B. Scalability Study

The objective of this study is to investigate if X-Nav’s performance improves with an increasing number of randomized robot embodiments used during training. Three configurations of X-Nav were tested, denoted as X-Nav-128, X-Nav-1024, and X-Nav-4096, where the number indicates the total number of randomly generated robot embodiments used by each expert policy for training and demonstration collection. For each configuration, we trained 3 expert policies and collected 4096 demonstrations for each expert policy to maintain a consistent total dataset size. For evaluation environments, we used the Matterport3D (MP3D) dataset [43] which contains 3D meshes of real-world indoor home environments. This dataset was selected as it is widely used for robot indoor navigation [44]. In particular, four houses with an



**Fig. 6.** SR and SPL of three X-Nav configurations, using 128, 1024, and 4096 robot embodiments for expert policy training and data collection.

Variants	Go2		Jackal	
	SR	SPL	SR	SPL
X-Nav w L1	55.3	0.44	85.4	0.78
X-Nav w EC	41.0	0.36	89.1	0.82
X-Nav w/o TE	69.1	0.60	74.9	0.65
X-Nav w TE	42.7	0.38	90.4	0.84
<b>X-Nav</b>	<b>69.1</b>	<b>0.60</b>	<b>90.4</b>	<b>0.84</b>

average size of  $22 \times 25$  m from MP3D were imported into Isaac Sim, Fig. 5. In each house, we randomly generated 20 pairs of start and goal positions and conducted 500 trials for each of the six robot embodiments used in Section VI.A. SR and SPL were computed and averaged across all test robots.

*1) Results:* Fig. 6 presents the SR and SPL for all configurations. Overall, X-Nav’s SR and SPL increased with the number of random embodiments used for training. This is due to X-Nav being exposed to a broader range of embodiment parameters during training, which allowed it to implicitly infer the embodiment parameters of unseen robots with higher accuracy. As a result, X-Nav was able to generate adaptive navigation actions that accounted for the robot’s morphology and dynamics, leading to improved navigation performance. Furthermore, these results demonstrate X-Nav’s capability to achieve zero-shot transfer to out-of-distribution photorealistic environments. The average SR and SPL of X-Nav-4096 are comparable to its performance in in-distribution environments (SR 77.9%, SPL 0.7) from Section VI.A. This is due to the unified ray-based distance representation as input, which allows X-Nav to generalize across different camera configurations. X-Nav was trained with obstacles of varying placement, density, and size to expose the policy to diverse obstacle interactions, enabling generalization to out-of-distribution environments.

### C. Ablation Study

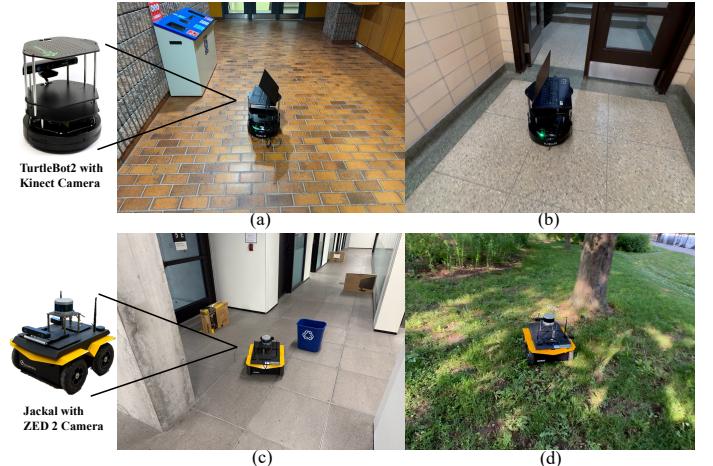
We conducted an ablation study with different variants of X-Nav to evaluate the impact of our design choices on the loss function and inference strategy. These included:

**X-Nav with L1 Loss (w L1):** L1 loss was used for training.

**X-Nav with Executing Chunk (w EC):** It executed the entire action chunk  $\mathbf{A}_t$  before predicting the next action chunk  $\mathbf{A}_{t+S_a}$ .

**X-Nav without Temporal Ensemble (w/o TE):** It generated an action sequence  $\mathbf{A}_t$  at each timestep and always executed the first action from the current predicted sequence without TE.

**X-Nav with Temporal Ensemble (w TE):** This variant applied TE to both wheeled and quadrupedal robots at inference time.



**Fig. 7.** X-Nav was deployed on the TurtleBot2 and the Jackal robots in both indoor and outdoor environments, including: (a) a hallway, (b) a doorway, (c) a corridor, and (d) a park.

We conducted 3000 trials for each of the above variants using the Unitree Go2 and Clearpath Jackal platforms. SR and SPL were used as performance metrics. We used the same test environments as in Section VI.A.

*1) Results:* The results are shown in Table VI. Overall, X-Nav achieved the highest SR and SPL among all the variants. X-Nav w L1 achieved lower SR and SPL than X-Nav, as the L1 loss penalizes less severely for large action prediction errors compared to the L2 loss, which led to less precise navigation actions. X-Nav w EC and X-Nav w TE had the lowest SR and SPL values for quadrupeds, due to the quadruped navigation requiring fast action changes to maintain stability. Namely, X-Nav w EC delayed robot adaptation to dynamic changes of robot pose and terrain as it required the execution of an entire action chunk. X-Nav w TE reduced robot responsiveness as it smoothed actions by averaging past predictions, which led to navigation failures. X-Nav w/o TE achieved the lowest SR and SPL for wheeled robots, highlighting the importance of TE for wheeled robot navigation. With TE, robots can generate smoother and more consistent motions, which helps reduce abrupt velocity and orientation changes.

## VII. REAL-WORLD EXPERIMENTS

We conducted real-world experiments in both indoor environments including a hallway, a doorway, a corridor, and outdoor environments including a park with uneven terrain, Fig. 7. To validate X-Nav’s effectiveness across different robot embodiments, we used two distinct wheeled robots: the TurtleBot2 and Clearpath Jackal. They differ in physical size, mass, kinematics, and sensor configurations. The TurtleBot2 was equipped with a Kinect sensor and the Jackal robot was equipped with a ZED 2 stereo camera. Both robots used the Robot Operating System (ROS) Noetic with X-Nav running at 50 Hz. No additional training was done for the real-world deployment. The TurtleBot2 robot navigated a  $21 \times 8$  m indoor environment with a hallway with obstacles such as door mullions and garbage bins, Fig. 7 (a), and a  $19 \times 5$  m indoor environment that included a narrow doorway, Fig. 7 (b). The Jackal robot navigated in a  $23 \times 7$  m indoor narrow corridor with obstacles, Fig. 7 (c), and a  $20 \times 10$  m outdoor park with

uneven terrain, Fig. 7 (d). A total of twenty trials were conducted using random start and goal locations.

X-Nav achieved an SR of 85% and SPL of 0.79, respectively. These results demonstrate the zero-shot transfer capabilities of X-Nav by successfully adapting to distinct robot embodiments, camera configurations, and environments. A video of X-Nav performing cross-embodiment navigation is provided at our project webpage: <https://cross-embodiment-nav.github.io/> and on the ASBLab YouTube channel: <https://youtu.be/vGwR1Rqujik>.

### VIII. CONCLUSION

In this paper, we present a novel two-stage learning framework, X-Nav, to address the problem of cross-embodiment navigation for both wheeled and quadrupedal robots. The first stage utilized deep reinforcement learning with privileged observations to train multiple expert policies tailored to specific groups of embodiments. The second stage distilled these expert policies into a single general policy using a transformer model, Nav-ACT. Through extensive simulated experiments, we demonstrated the effectiveness of X-Nav in zero-shot transfer to unseen robot embodiments and photorealistic environments. A scalability study showed that X-Nav's performance scales with increasing number of random embodiments used during training. An ablation study validated our design choices and inference strategy. Real-world experiments validated the generalizability of X-Nav in both indoor and outdoor environments. Future work will extend X-Nav to more robot types such as humanoid robots, and object-goal navigation to expand its applicability.

### ACKNOWLEDGMENT

The authors would like to thank the following members of the ASBLab at UoT for their assistance with the experimental setup: Sourabh Prasad, Matthew Lisondra, and Daniel Choi.

### REFERENCES

- [1] A. Fung, B. Benhabib, and G. Nejat, "Robots Autonomous Detecting People: A Multimodal Deep Contrastive Learning Method Robust to Intra-class Variations," *IEEE Robot. Autom. Lett.*, vol. 8, no. 6, pp. 3550–3557, Jun. 2023.
- [2] A. Fung, B. Benhabib, and G. Nejat, "LDTrack: Dynamic People Tracking by Service Robots Using Diffusion Models," *Int. J. Comput. Vis.*, Jan. 2025.
- [3] H. Wang, A. H. Tan, and G. Nejat, "NavFormer: A Transformer Architecture for Robot Target-Driven Navigation in Unknown and Dynamic Environments," *IEEE Robot. Autom. Lett.*, vol. 9, no. 8, pp. 6808–6815, Aug. 2024.
- [4] A. H. Tan, F. P. Bejarano, Y. Zhu, R. Ren, and G. Nejat, "Deep Reinforcement Learning for Decentralized Multi-Robot Exploration With Macro Actions," *IEEE Robot. Autom. Lett.*, vol. 8, no. 1, pp. 272–279, Jan. 2023.
- [5] A. H. Tan, S. Narasimhan, and G. Nejat, "4CNet: A Diffusion Approach to Map Prediction for Decentralized Multi-robot Exploration," 2024, *arXiv:2402.17904*.
- [6] N. Tyagi, D. Sharma, J. Singh, B. Sharma, and S. Narang, "Assistive Navigation System for Visually Impaired and Blind People: A Review," in *Int. Conf. Artif. Intell. Mach. Vis.*, Sep. 2021, pp. 1–5.
- [7] J. Guerreiro, D. Sato, S. Asakawa, H. Dong, K. M. Kitani, and C. Asakawa, "CaBot: Designing and Evaluating an Autonomous Navigation Robot for Blind People," in *Int. Conf. Comput. Accessibility*, 2019, pp. 68–82.
- [8] G. Feng *et al.*, "GenLoco: Generalized Locomotion Controllers for Quadrupedal Robots", 2022, *arXiv:2209.05309*.
- [9] N. Hirose, D. Shah, A. Sridhar, and S. Levine, "ExAug: Robot-Conditioned Navigation Policies via Geometric Experience Augmentation," 2022, *arXiv:2210.07450*.
- [10] J. Yang *et al.*, "Pushing the Limits of Cross-Embodiment Learning for Manipulation and Navigation," 2024, *arXiv: 2402.19432*.
- [11] D. Shah, A. Sridhar, A. Bhorkar, N. Hirose, and S. Levine, "GNM: A General Navigation Model to Drive Any Robot," 2023, *arXiv:2210.03370*.
- [12] D. Shah *et al.*, "ViNT: A Foundation Model for Visual Navigation," 2023, *arXiv:2306.14846*.
- [13] A. Sridhar, D. Shah, C. Glossop, and S. Levine, "NoMaD: Goal Masked Diffusion Policies for Navigation and Exploration," 2023, *arXiv:2310.07896*.
- [14] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine, "Scaling Cross-Embodied Learning: One Policy for Manipulation, Navigation, Locomotion and Aviation," 2024, *arXiv: 2408.11812*.
- [15] W. Liu *et al.*, "X-MOBILITY: End-To-End Generalizable Navigation via World Modeling," 2024, *arXiv:2410.17491*.
- [16] K. Kang, G. Kahn, and S. Levine, "Hierarchically Integrated Models: Learning to Navigate from Heterogeneous Robots," 2021, *arXiv:2106.13280*.
- [17] J. Truong *et al.*, "Learning Navigation Skills for Legged Robots with Learned Robot Embeddings," in *IEEE Int. Conf. Intell. Robot Syst.*, Sep. 2021, pp. 484–491.
- [18] N. Bohlinger *et al.*, "One Policy to Run Them All: an End-to-end Learning Approach to Multi-Embodiment Locomotion," 2024, *arXiv:2409.06366*.
- [19] Z. Luo *et al.*, "MorAL: Learning Morphologically Adaptive Locomotion Controller for Quadrupedal Robots on Challenging Terrains," *IEEE Robot. Autom. Lett.*, vol. 9, no. 5, pp. 4019–4026, May 2024.
- [20] F. Di Giuro, F. Zargarbashi, J. Cheng, D. Kang, B. Sukhija, and S. Coros, "Meta-Reinforcement Learning for Universal Quadrupedal Locomotion Control," 2024, *arXiv:2407.17502v1*.
- [21] P. Putta *et al.*, "Embodiment Randomization for Cross Embodiment Navigation," in *IEEE Int. Conf. Intell. Robot Syst.*, Oct. 2024, pp. 5527–5534.
- [22] D. Hoeller, N. Rudin, D. Sako, and M. Hutter, "ANYmal parkour: Learning agile navigation for quadrupedal robots," *Sci. Robot.*, vol. 9, no. 88, p. eadl7566, Mar. 2024.
- [23] F. Rubio, F. Valero, and C. Llopis-Albert, "A review of mobile robots: Concepts, methods, theoretical framework, and applications," *Int. J. Adv. Robot. Syst.*, vol. 16, no. 2, p. 1729881419839596, Mar. 2019.
- [24] M. Tan and Q. V. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," 2020, *arXiv:1905.11946*.
- [25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2016, pp. 770–778.
- [26] A. Vaswani *et al.*, "Attention is all you need," *Adv. Neural Inf. Process. Syst.*, vol. 2017-Decem, pp. 5999–6009, 2017.
- [27] C. Chi *et al.*, "Diffusion Policy: Visuomotor Policy Learning via Action Diffusion," 2024, *arXiv:2303.04137*.
- [28] N. Hirose, D. Shah, A. Sridhar, and S. Levine, "SACSoN: Scalable Autonomous Control for Social Navigation," 2023, *arXiv:2306.01874*.
- [29] D. Yarats, A. Zhang, I. Kostrikov, B. Amos, J. Pineau, and R. Fergus, "Improving Sample Efficiency in Model-Free Reinforcement Learning from Images," 2020, *arXiv:1910.01741*.
- [30] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," 2017, *arXiv:1707.06347*.
- [31] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to Walk in Minutes Using Massively Parallel Deep Reinforcement Learning," 2022, *arXiv:2109.11978*.
- [32] J. Jin *et al.*, "Resilient Legged Local Navigation: Learning to Traverse with Compromised Perception End-to-End," 2023, *arXiv:2310.03581*.
- [33] T. He, C. Zhang, W. Xiao, G. He, C. Liu, and G. Shi, "Agile But Safe: Learning Collision-Free High-Speed Legged Locomotion," 2024, *arXiv:2401.17583*.
- [34] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *IEEE Int. Conf. Intell. Robot Syst.*, Dec. 2017, pp. 23–30.
- [35] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn, "Learning Fine-Grained Bimanual Manipulation with Low-Cost Hardware," 2023, *arXiv:2304.13705*.

- [36] M. Mittal *et al.*, “Orbit: A Unified Simulation Framework for Interactive Robot Learning Environments,” 2024, *arXiv:2301.04195*.
- [37] D. P. Kingma and J. Lei, “Adam: A Method for Stochastic Optimization,” in *Proc. Int. Conf. Learn. Representations*, 2015.
- [38] P. Anderson *et al.*, “On Evaluation of Embodied Navigation Agents,” 2018, *arXiv:1807.06757*.
- [39] D. A. Pomerleau, “ALVINN: An Autonomous Land Vehicle in a Neural Network,” in *Adv. Neural Inf. Process. Syst.*, 1988.
- [40] T. Karras, M. Aittala, T. Aila, and S. Laine, “Elucidating the Design Space of Diffusion-Based Generative Models,” 2022, *arXiv:2206.00364*.
- [41] A. Prasad, K. Lin, J. Wu, L. Zhou, and J. Bohg, “Consistency Policy: Accelerated Visuomotor Policies via Consistency Distillation,” 2024, *arXiv:2405.07503*.
- [42] X. Huang *et al.*, “DiffuseLoco: Real-Time Legged Locomotion Control with Diffusion from Offline Datasets,” 2024, *arXiv:2404.19264*.
- [43] A. Chang *et al.*, “Matterport3D: Learning from RGB-D Data in Indoor Environments,” 2017, *arXiv:1709.06158*.
- [44] M. Savva *et al.*, “Habitat: A Platform for Embodied AI Research,” in *IEEE Int. Conf. Comput. Vis.*, 2019, pp. 9338–9346.