

Supplementary Materials for Learning Diverse Natural Behaviors for Enhancing the Agility of Quadrupedal Robots

S1. Nomenclature

\mathbf{s}	state	1197
\mathbf{o}	observation	1198
\mathbf{a}	action	1199
\mathbf{a}^D	discrete action	1200
\mathbf{a}^C	continuous action	1201
r	reward	1202
π	policy	1203
H	observation horizon	1204
D	discriminator	1205
Q	predictor	1206
\mathbf{c}	latent skill variable	1207
ϵ	latent shifting variable	1208
y	behavior label	1209
\mathcal{H}	entropy	1210
$I(\cdot)$	mutual information	1211
$(\cdot)^I$	imitation quantity	1212
$(\cdot)^T$	task quantity	1213
$(\cdot)^S$	simulator quantity	1214
(\cdot)	real-world quantity	1215
(\cdot)	predicted quantity	1216
$(\cdot)^{SS}$	semi-supervised quantity	1217
$(\cdot)^{US}$	unsupervised quantity	1218
\mathbb{I}	indicator function	1219
$\ \cdot\ $	l_2 norm	1220

S2. Training details

The training of the integrated controller is divided into three stages: BBC pretrain-
ing, BBC fine-tuning, and TSC training. All training processes are conducted in the
IsaacGym physics simulation environment.

In the first stage, the simulator’s physical parameters are set empirically. We use
semi-supervised InfoGAIL to train a multimodal BBC based on motion capture data
of common dog behaviors, with less than 5% of the data labeled. Aside from motion
retargeting, the data requires no additional preprocessing, such as cropping, label-
ing, or alignment. To address category imbalance, labeled data is simply repeated to
rebalance minority categories. The training pseudocode is presented in Algorithm S1.

In each training iteration i , the latent skill variable \mathbf{c} is sampled from a categorical
distribution $p_i(\mathbf{c})$, and the latent shifting variable ϵ is drawn from a uniform distribu-
tion $p(\epsilon)$. The policy interacts with the environment based on observations and latent
variables, generating a trajectory distribution d^π . The discriminator D and predictor

1243 Q are updated via gradient descent on Equation (10). The policy π is then optimized
 1244 using the PPO algorithm [58] to maximize the reward in Equation (11), allowing it to
 1245 replicate the multimodal behaviors in the motion capture data. The value function V
 1246 is updated using $TD(\lambda)$.

1247 To align with the imbalanced motion capture data, $p(\mathbf{c})$ is estimated by Q^c and
 1248 updated using the exponential moving average:

$$1249$$

$$1250 \quad p_{i+1}(\mathbf{c}) = \beta p_i(\mathbf{c}) + (1 - \beta) \frac{1}{N} \sum_i Q^c(\hat{y} | \mathbf{o}_i^I), \quad (\text{S1})$$

$$1251$$

$$1252$$

1253 where β is the weight coefficient. The hyperparameters for semi-supervised InfoGAIL
 1254 are shown in Table S1.

1255 In addition to imitating motion capture data, the BBC is also trained to follow
 1256 task-related command inputs to enhance its controllability. The command ranges for
 1257 each behavior mode are shown in Table S2. During training, task commands and
 1258 latent variables are sampled every six seconds. Additionally, randomly generated rough
 1259 terrains and random disturbances applied to the robot are introduced to improve the
 1260 robustness of the BBC.

1261 In the second stage, we use EASI to learn the simulator’s parameter distribution
 1262 from a small set of real-world samples, allowing the simulator to better approximate
 1263 the real world. The BBC is then fine-tuned in this enhanced simulator. The real-
 1264 world samples, collected using the pre-trained BBC (or any other controller), consist
 1265 of 4,000 steps at a control frequency of 50 Hz, providing about 80 seconds of data.
 1266 The simulator parameter optimization pseudocode is presented in Algorithm S2.

1267 In each generation, N parameter sets are sampled from the distribution $\Xi^{(i)}$, includ-
 1268 ing the PD parameters for the hip, thigh, and calf joints of the quadrupedal robot,
 1269 while the known component masses remain fixed. The policy π then collects trajec-
 1270 tories d^π in N parallel simulations with the sampled parameters. A batch of state
 1271 transitions is subsequently used to update the discriminator \tilde{D}_k to \tilde{D}_{k+1} following
 1272 Equation (13). Finally, the updated discriminator calculates the rewards $r^{(i)}$, which are
 1273 used by ES to generate the next parameter distribution $\Xi^{(i+1)}$. The hyperparameters
 1274 for EASI are shown in Table S3.

1275 In the last stage, the TSC is trained using privileged learning to perform specific
 1276 tasks. The pseudocode is presented in Algorithm S3. To address the issue of sparse
 1277 rewards, the teacher policy π^{TSC} is first trained to maximize task reward and imitation
 1278 reward, with the latter calculated by the discriminator from the BBC. Hybrid-PPO [47]
 1279 is employed to handle the hybrid action space. Subsequently, the student policy $\hat{\pi}^{\text{TSC}}$
 1280 takes depth images as input and outputs actions to mimic the teacher’s behavior.
 1281 The positions of obstacles and the robot are randomly initialized at the beginning
 1282 of each episode. Additionally, to enhance the robustness of the student policy in real
 1283 world, we incorporate the BYOL loss [48]. We apply random augmentations to the
 1284 input depth images, such as rectangular noise, white noise, and Gaussian blur. The
 1285 hyperparameters for privileged learning are shown in Table S4.

1286 The networks are implemented using Pytorch [62], and are trained on a computer
 1287 with an i5-12600KF CPU and a Nvidia GTX 4080 GPU. We use the parallelized
 1288

physics simulator IsaacGym [51] for training. The BBC is trained using real dog motion capture data from [57], which includes approximately 20 minutes of motion data covering common dog behaviors. Training the BBC takes about 4 days on a single GPU, while the TSC-teacher and TSC-student require approximately 10 hours and 14 hours, respectively.

S3. Deployment details

We deploy our integrated controller on the Unitree Go2 robot, which is equipped with a Jetson Orin NX, an Intel RealSense D435i depth camera [63] and features four limbs comprising 12 joints.

In our real-world experiments, we deploy our integrated controller that incorporates both BBC and TSC on the Jetson Orin NX embedded in the Unitree Go2 robot, achieving stable control frequency at 50 Hz while performing deep image processing.

To ensure that image processing and policy inference are completed within each 20ms cycle, all computations are accelerated using CUDA. Intercommunication among the robot’s various sensors and low-level motors is facilitated via the Cyclone DDS and LCM platforms, with the motors executing torque at a rate of 1000 Hz based on Unitree Go2 integrated PD controllers. Depth images are acquired at a sensor resolution of 848×480, and then these images are subsequently downsampled and cropped to yield a final resolution of 87×58. After normalization in which each pixel value is scaled to the range from -0.5 to 0.5, the processed depth image serves as the input to the visual encoder network, which operates concurrently with the motion policy at a frequency of 50 Hz. Moreover, owing to the inherent imaging principles of the RealSense D435i, the raw depth images often contain considerable invalid regions and noise. To address these limitations, RGB-Depth Processing module in OpenCV is employed to process the depth map data, effectively enhancing the image quality.

S4. Additional Rewards for BBC

To achieve a smoother and more stable controller, we introduce the following regularization rewards for the BBC:

- Joint torque reward: $r_{\text{torque}}^T = -10^{-5} \cdot \sum \|\tau\|^2$ penalizes excessive joint torques.
- Torque change reward: $r_{\text{delta_torque}}^T = -10^{-7} \cdot \sum \|\Delta_\tau\|^2$ penalizes rapid changes in joint torque, where Δ_τ is the difference between consecutive steps.
- Joint acceleration reward: $r_{\text{dof_acc}}^T = -2.5 \times 10^{-7} \cdot \sum \|\ddot{\mathbf{q}}\|^2$ penalizes excessive joint acceleration.
- Joint position limitation reward: $r_{\text{joint_pos_limit}}^T = -0.1 \cdot \sum (\max(0, \mathbf{q}_{\min} - \mathbf{q}) + \max(0, \mathbf{q} - \mathbf{q}_{\max}))$ penalizes joint positions that exceed their defined limits.
- Joint velocity limitation reward: $r_{\text{joint_vel_limit}}^T = -0.1 \cdot \sum (\min(1, \max(0, |\dot{\mathbf{q}}| - \mathbf{q}_{\text{vel_limit}})))$ penalizes joint velocities exceeding the velocity limit.
- Joint torque limitation reward: $r_{\text{joint_torque_limit}}^T = -0.03 \cdot \sum (\max(0, |\tau| - \tau_{\text{limit}}))$ penalizes joint torques that approach or exceed the torque limit.

- 1335 • Collision reward: $r_{\text{collision}}^T = -10 \cdot \mathbb{I}_{\text{collision}}$ penalizes collisions involving the thigh
1336 and calf.
- 1337 • Action smoothness reward: $r_{\text{delta_action}}^T = -0.1 \cdot \|\Delta_a\|^2$ penalizes abrupt changes in
1338 action commands.

1339

1340 S5. Additional rewards for TSC

1341 To achieve a smoother and safer TSC, we introduce the following regularization
1342 rewards:
1343

- 1344 • Action smoothness reward: $r_{\text{delta_action}}^{\text{TSC}} = -0.1 \cdot \|\Delta_{a^T}\|^2$ penalizes abrupt changes
1345 in action commands.
- 1346 • Collision reward: $r_{\text{collision}}^{\text{TSC}} = -20 \cdot \mathbb{I}_{\text{collision}}$ penalizes collisions involving the thigh
1347 and calf.
- 1348 • Edge contact reward: $r_{\text{feet_edge}}^{\text{TSC}} = -\sum_{i=0}^4 c_i^{\text{foot}} M(p_i)$ penalizes each foot contact
1349 with the edge of the obstacles. $c_i^{\text{foot}} = 1$ if foot i is in contact. M is a boolean
1350 function which is 1 if the foot position p_i within the obstacle edge.
1351

1352 S6. Neural network architectures

1353
1354 The BBC’s policy network outputs target positions for the quadrupedal robot’s 12
1355 joints. This network is an MLP with 3 hidden layers {512, 256, 128}. The input includes
1356 commands, measured and estimated proprioception, and latent physical parameters.
1357 Estimated proprioception is generated by a two-layer MLP (sizes {128, 64}), based on
1358 measured proprioception. Latent physical parameters are inferred from proprioception
1359 history using two 1D convolutional layers with kernel sizes {4, 2} and strides {2, 1}.

1360 The TSC-teacher’s policy network outputs both discrete and continuous com-
1361 mands. Discrete commands are probabilities for 5 behaviors, while continuous
1362 commands consist of 5×6 values corresponding to the behaviors. The TSC policy net-
1363 work is an MLP with 3 hidden layers {512, 256, 128}. Its input consists of scandots,
1364 privileged information, and proprioception. Scandots are an 11×12 matrix, which is
1365 flattened and input into a three-layer MLP encoder (sizes {128, 64, 32}), with the 32-
1366 dimensional output used as input to the policy network. Estimated proprioception is
1367 predicted in the same way as in the BBC.

1368 The TSC-student’s policy network is an MLP with 3 hidden layers {512, 256, 128}.
1369 It’s outputs are identical to the TSC-teacher. Instead of scandots, it uses the current
1370 depth map, sized 87×58 as exteroceptive perception. This depth map is processed
1371 with 2 convolutional layers (kernel sizes {5, 3}), one max pooling layer (kernel size 2,
1372 stride 2), and two fully connected layers to produce a 32-dimensional depth feature.
1373 This feature is concatenated with the measured proprioception and passed into a GRU
1374 with a 512-unit hidden size to infer the privileged information from historical data.

1375

1376 S7. Human-robot collaboration

1377 In the quadrupedal agility challenge, the robot follows a specific route through ran-
1378 domly placed obstacles. However, in situations where the route is ambiguous, human
1379 collaboration is needed to complete the task, similar to the role of handlers in dog
1380

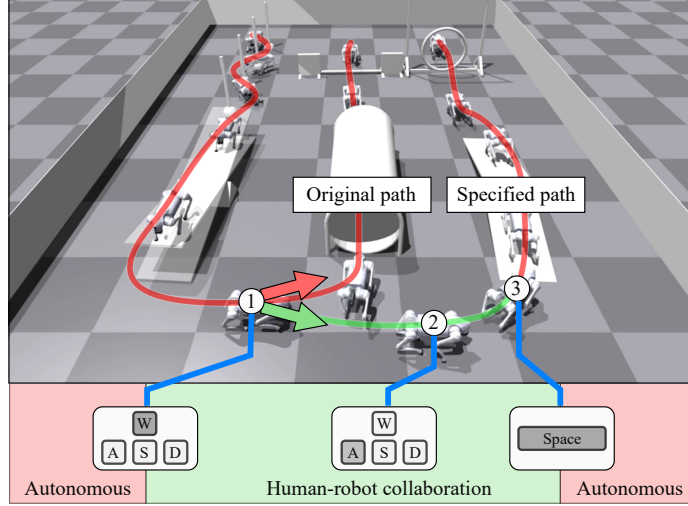


Fig. S1 Human-robot collaboration. The robot receives the operator’s input at point 1, where it switches from an immediate left turn to moving straight. The operator only adjusts the linear and angular velocities, while the robot’s behavior style remains controlled by the TSC. At point 3, the operator presses the Space key, and the robot transitions from human-robot collaboration to autonomous movement.

agility, where handlers are informed of the route in advance. We simulate this collaborative mode in a scenario with bidirectional routes (Fig. S1). The robot receives the operator’s input at point 1, where it switches from an immediate left turn to moving straight. The operator only adjusts the linear and angular velocities, while the robot’s behavior style remains controlled by the TSC. At point 3, the operator presses the Space key, and the robot transitions from human-robot collaboration to autonomous movement. We use a remote control for this task, though other methods could also achieve the same goal, such as voice commands, gestures, or behavioral intentions, which we plan to explore in future work. In addition to altering the route, the handler can also modify the skill variable inputs to switch between different behaviors.

S8. Ablation of the BYOL loss for student training

To intuitively demonstrate the improvement in the robustness of the student policy by the BYOL loss, we use Gradient-weighted Class Activation Mapping (Grad-CAM) [64] to visualize the areas of focus in the depth image encoder, as shown in Fig. S2. The first row shows the original depth image containing the bar-jump, as well as the depth images with added random rectangles, white noise, and Gaussian blur. The second row displays the areas of focus for our method’s encoder. The third row shows the areas of focus for the encoder without BYOL. With the use of BYOL, the encoder consistently focuses on important areas under various noise disturbances, ensuring the policy’s robustness in real-world scenarios.

Table S1 Hyperparameters for semi-supervised InfoGAIL.

Optimizer	Adam
Policy & value learning rate	1.0 E-3
Discriminator learning rate	5.0 E-4
Encoder learning rate	1.0 E-3
Observation horizon	2
Discount factor	0.99
TD(λ)	0.95
PPO clip threshold	0.2
Moving average weight	1.0 E-3
Imitation reward weight	0.2
Task reward weight	0.2

Table S2 Task command ranges for quadrupedal agility challenge.

	Walk	Pace	Trot	Canter	Jump
v_x (m/s)	[0.0, 0.6]	[0.5, 1.5]	[0.5, 1.5]	[0.8, 2.0]	[0.8, 1.8]
v_y (m/s)	[-0.15, 0.15]	[-0.3, 0.3]	[-0.3, 0.3]	[-0.5, 0.5]	[-0.3, 0.3]
ω_z (rad/s)	[-1.0, 1.0]	[-1.57, 1.57]	[-1.57, 1.57]	[-0.5, 0.5]	[-0.5, 0.5]
jump height (m)	-	-	-	-	[0.45, 0.55]
locomotion height (m)	[0.25, 0.34]	[0.25, 0.34]	[0.25, 0.34]	[0.25, 0.34]	-

Table S3 Hyperparameters for EASI.

Optimizer	RMSprop
Discriminator learning rate	3.0 E-4
Discriminator training epoch	10
Number of env.	300
Evolution survival rate	0.5
Trajectory length	500

Table S4 Hyperparameters for privileged learning.

Optimizer	Adam
Policy & value learning rate	5.0 E-4
Discount factor	0.99
TD(λ)	0.95
PPO clip threshold	0.2
Imitation reward weight	0.2
Task reward weight	2.0

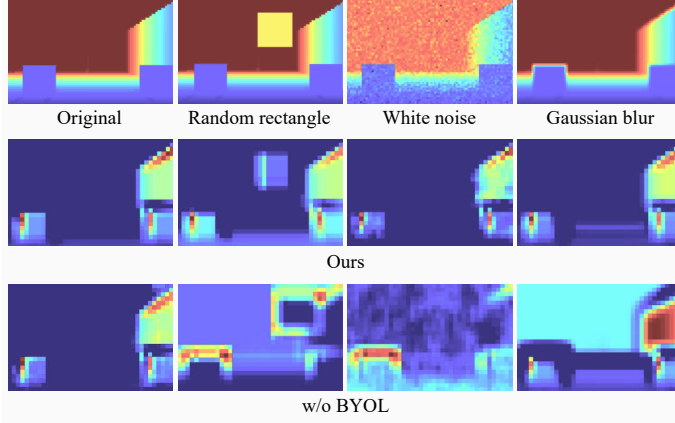


Fig. S2 Visualization of the depth image encoder’s focused areas. The first row shows the original depth image, followed by depth images with added random rectangles, white noise, and Gaussian blur. The second row displays the areas of focus for our method’s encoder. The third row shows the areas of focus for the encoder without BYOL.

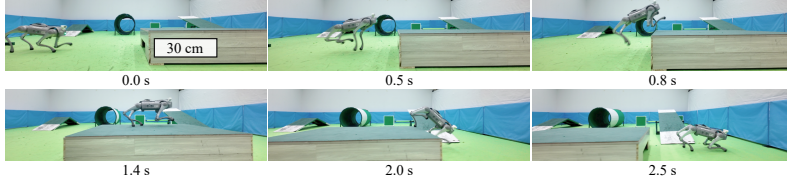


Fig. S3 Box jumping process. The box is 30 cm high, approximately the same as the robot’s normal standing height. With our controller, the quadrupedal robot can quickly jump onto the obstacle without relying on its front legs for support.

Algorithm S1 BBC training with semi-supervised InfoGAIL

- 1: Initialize policy π_0 , value function V_0 , discriminator D_0 , predictor Q_0 , and latent skill distribution $p_0(\mathbf{c})$. Prepare unlabeled and labeled demonstrations d^E and d^{EL}
 - 2: **for** $i = 0, 1, 2, \dots N$ **do**
 - 3: Sample a batch of latent variables $\mathbf{c} \sim p_i(\mathbf{c})$, $\epsilon \sim p(\epsilon)$
 - 4: Interact with the environment using π_i , \mathbf{c} and ϵ , and obtain d^π
 - 5: Sample a batch of observations $b^E \sim d^E$, $b^{EL} \sim d^{EL}$, $b^\pi \sim d^\pi$
 - 6: Update D_i to D_{i+1} and Q_i to Q_{i+1} according to Equation (10)
 - 7: Update π_i to π_{i+1} using PPO [58] with reward according to Equation (11)
 - 8: Update V_i to V_{i+1} using $TD(\lambda)$
 - 9: Update $p_i(\mathbf{c})$ to $p_{i+1}(\mathbf{c})$ with Q_{i+1} and b^E
 - 10: **end for**
-

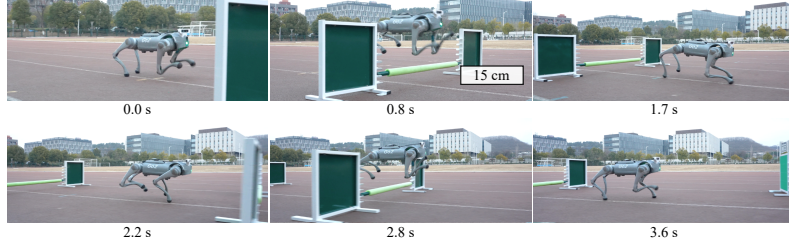


Fig. S4 Hurdling process. We showcase two jumping phases in the hurdling task. The robot utilizes depth vision to naturally switch behaviors, completing the task like a real dog. The whole process can be found in Supplementary Video 3.

Algorithm S2 Simulator parameter optimization with EASI

```

1: Initialize parameter distribution  $\Xi^{(0)}$ , discriminator  $\tilde{D}_0$ , policy  $\pi$ , real-world
   demonstration  $\tilde{d}^E$ 
2: for generation  $i = 0, 1, 2, \dots, G$  do
3:   Sample individuals by  $\xi_j^{(i)} \sim \Xi^{(i)}$ ,  $j = 1, 2, \dots, N$ 
4:   for simulation environment  $j = 1, 2, \dots, N$  do
5:     Set simulator parameters to  $\xi_j^{(i)}$ 
6:     Use  $\pi$  to sample trajectory  $\tau_j$ 
7:     Store  $\tau_j$  in  $d^\pi$ 
8:   end for
9:   for update step  $k = 1, 2, \dots, K$  do
10:    Sample a batch of transitions  $\tilde{b}^E \sim \tilde{d}^E$ ,  $b^\pi \sim d^\pi$ 
11:    Update  $\tilde{D}_k$  to  $\tilde{D}_{k+1}$  according to Equation (13)
12:   end for
13:   Calculate discriminator reward  $r^{(i)}$  according to Equation (2)
14:   Using ES to find next generation distribution  $\Xi^{(i+1)} = \text{ES}(\xi^{(i)}, r^{(i)})$ 
15: end for

```

Algorithm S3 TSC training with privileged learning	1565
1: Initialize teacher policy π_0^{TSC} , student policy $\hat{\pi}_0^{\text{TSC}}$, value function V_0	1566
2: for $i = 0, 1, 2, \dots N$ do	1567
3: Initialize obstacles and robot at a random position	1568
4: Interact with the environment using π_i^{TSC} , and obtain d^π	1569
5: Update π_i^{TSC} to π_{i+1}^{TSC} using Hybrid-PPO [47] with reward according to Equation (14)	1570
6: Update V_i to V_{i+1} using $TD(\lambda)$	1571
7: end for	1572
8: for $j = 0, 1, 2, \dots M$ do	1573
9: Initialize obstacles and robot at a random position	1574
10: Interact with the environment using $\hat{\pi}_j^{\text{TSC}}$, and obtain \hat{d}^π	1575
11: Update $\hat{\pi}_j^{\text{TSC}}$ to $\hat{\pi}_{j+1}^{\text{TSC}}$ by descending with gradients according to Equation (15)	1576
12: end for	1577
	1578
	1579
	1580
	1581
	1582
	1583
	1584
	1585
	1586
	1587
	1588
	1589
	1590
	1591
	1592
	1593
	1594
	1595
	1596
	1597
	1598
	1599
	1600
	1601
	1602
	1603
	1604
	1605
	1606
	1607
	1608
	1609
	1610