

Generating Physically Realistic and Directable Human Motions from Multi-Modal Inputs

Aayam Shrestha^{*1}, Pan Liu^{*2}, German Ros^{†2}, Kai Yuan^{‡2}, and Alan Fern¹

¹ Oregon State University
² Intel Labs

1 Appendix

1.1 Implementation Details

The physics simulation is conducted using Isaac Gym [1], where the MHC runs at 30 Hz and the simulation runs at 60 Hz. The architecture of the MHC, which consists of a controller and an ensemble of discriminators, each implemented as a neural network. The *controller* policy encodes the motion directive lookahead using a 3-layer perceptron with hidden dimensions of 1024 and 512. This encoding is then concatenated with the current pose of the humanoid and fed into another 3-layer perceptron with dimensions [1024, 1024] that serves as the policy head, outputting the action of the controller. In practice, we implement the ensemble of *discriminators* as a single discriminator with different wrappers. Each wrapper masks the observations to consider only a single set of joints. The discriminators are also 3-layer perceptrons with dimensions [1024, 512]. We use SiLU activations for all perceptrons. For *Training* the MHC is trained via the Proximal Policy Optimization (PPO) reinforcement learning algorithm [2] with fixed entropy. The training process takes 7 days on a single Nvidia A6000 GPU to obtain the final policy.

* Equal contribution

† Now at NVIDIA


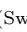
‡ Corresponding author. Email: kai.yuan@intel.com





1.2 Extended results for High level task specification

We further consider different variants of the go-to-location task and the heading task following [3] to showcase the flexibility afforded by the MHC framework. The rewards for both tasks are defined following [3].

For the **heading task**, a random heading direction and a velocity direction are sampled as targets, and the FSM should be able to follow these directives. For example, "Head east while facing west." Furthermore, we can also specify the speed and height during these motions. "Run" maps to a speed of 2.5 m/s with a root height of 0.85m, while "crouch walk" refers to a root height of 0.4m with a speed of 1m/s. We evaluate the resulting motions using task-specific rewards, which reward matching the desired direction and heading. We find that the produced FSMs can reliably generate motions that match the higher-level heading task objectives.

For the **go-to-location task**, we consider different variations in movement time and finishing motion. Unlike [3], where only a particular skill can be requested as a finishing motion, our framework allows giving any kind of full or partial directive as a finishing motion. Here, we consider sword swing and taunt motions as finishing directives. The movement directives are chosen following the heading task. It is worth noting that our framework also allows for different upper body movements throughout the movement as well. Table 1 shows that FSMs using MHC generate motions that achieve high task rewards across all go-to-location task variants.

Table 1: Quantitative evaluation of directional motion control(Heading) and zero-shot task solution. We consider two forms of locomotion Run and Crouch walk, each characterized by a different speed and style. We consider various finishing motions the location task  (Sword Swing) and  (Taunt).

Motion	Heading		Location	
	Style	Score	Ending	Score
Run	1	0.92		0.98
				0.98
Crouch Walk	0.94	0.91		0.96
				0.96

References

1. Makoviychuk, V., Wawrzyniak, L., Guo, Y., Lu, M., Storey, K., Macklin, M., Hoeller, D., Rudin, N., Allshire, A., Handa, A., State, G.: Isaac gym: High performance gpu-based physics simulation for robot learning. ArXiv [abs/2108.10470](#) (2021)
2. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. ArXiv [abs/1707.06347](#) (2017)
3. Tessler, C., Kasten, Y., Guo, Y., Mannor, S., Chechik, G., Peng, X.B.: Calm: Conditional adversarial latent models for directable virtual characters. ACM SIGGRAPH 2023 Conference Proceedings (2023)