

CALM: Conditional Adversarial Latent Models for Directable Virtual Characters

Chen Tessler¹
Shie Mannor^{1,2}

Yoni Kasten¹
Gal Chechik^{1,3}

Yunrong Guo¹
Xue Bin Peng^{1,4}

¹NVIDIA

²Technion - Israel Institute of Technology

³Bar-Ilan University

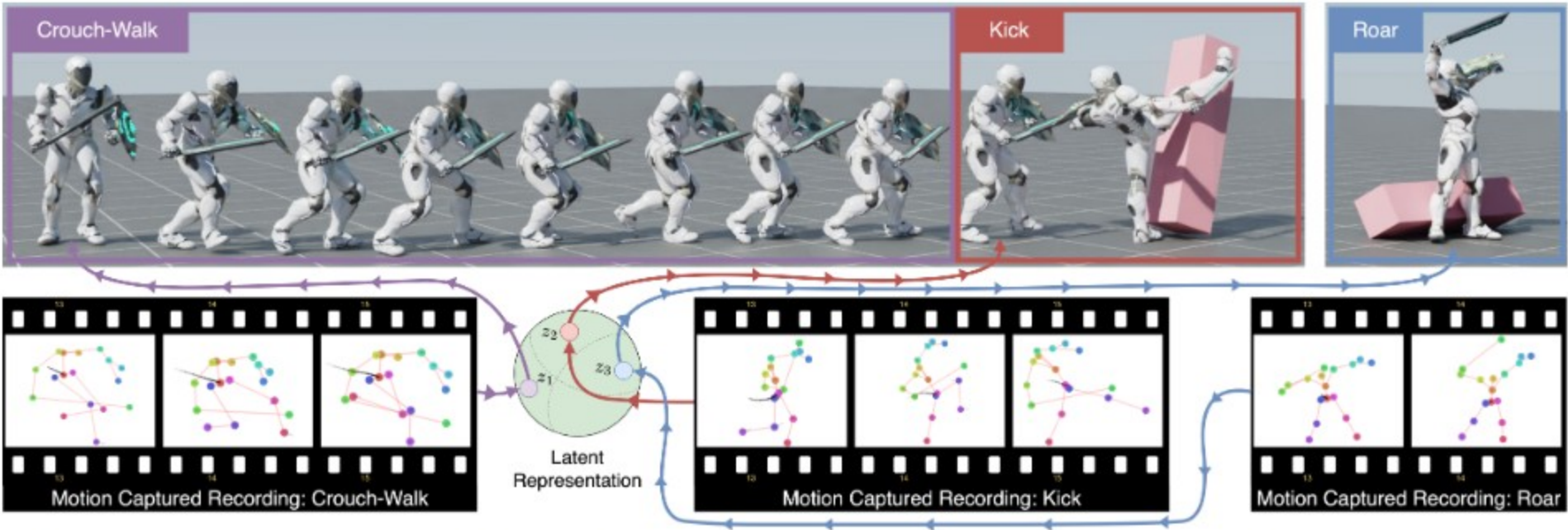
⁴Simon Fraser University

SIGGRAPH 2023

[Paper](#)

[BibTeX](#)

[Code](#)



Our framework enables users to direct the behavior of a physically simulated character using demonstrations encoded in the form of low-dimensional latent embeddings of motion capture data. In this example, the character is instructed to crouch-walk towards a target, kick when within range, and finally raise its arms and celebrate.

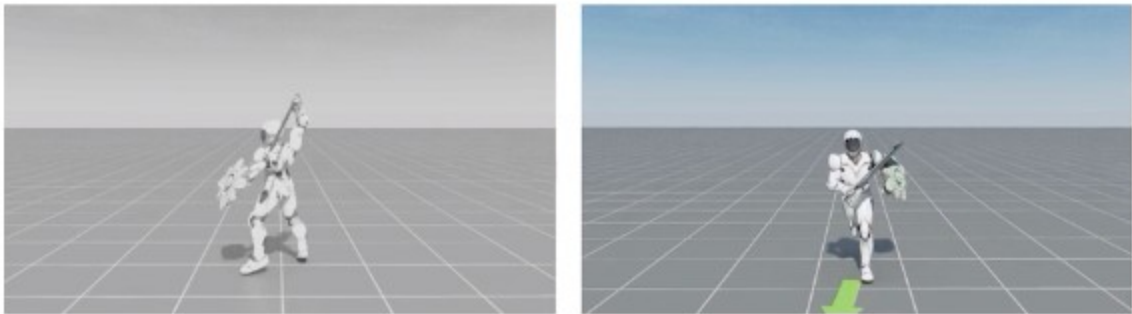
Abstract

In this work, we present Conditional Adversarial Latent Models (CALM), an approach for generating diverse and directable behaviors for user-controlled interactive virtual characters. Using imitation learning, CALM learns a representation of movement that captures the complexity and diversity of human motion, and enables direct control over character movements. The approach jointly learns a control policy and a motion encoder that reconstructs key characteristics of a given motion without merely replicating it. The results show that CALM learns a semantic motion representation, enabling control over the generated motions and style-conditioning for higher-level task training. Once trained, the character can be controlled using intuitive interfaces, akin to those found in video games.

Overview

CALM consists of 3 phases

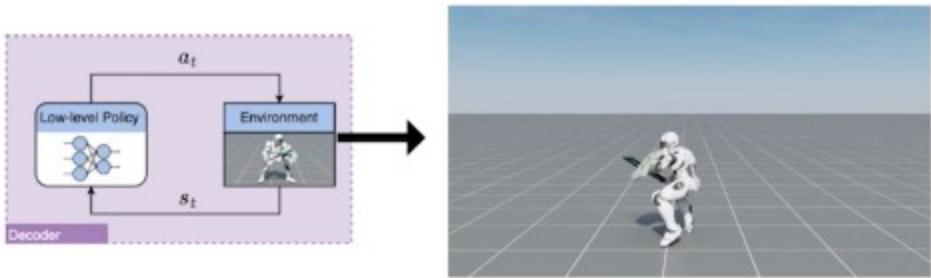
- (1) Train a low-level policy that generates diverse motions on-demand
- (2) Train a high-level policy to control the directionality of motions



To achieve zero-shot task solutions, CALM consists of 3 phases. **(1)** A motion encoder and a low-level policy (decoder) are jointly trained to map from a motion capture sequence into actions controlling the simulated character. **(2)** A high-level policy is trained using latent space conditioning, to enable control over the direction in which a motion is performed, while retaining the requested style. **(3)** Steps 1 and 2 are combined using a simple finite-state-machine in order to solve tasks without further training and without meticulous reward/termination design.

Phase 1: Low-level Training

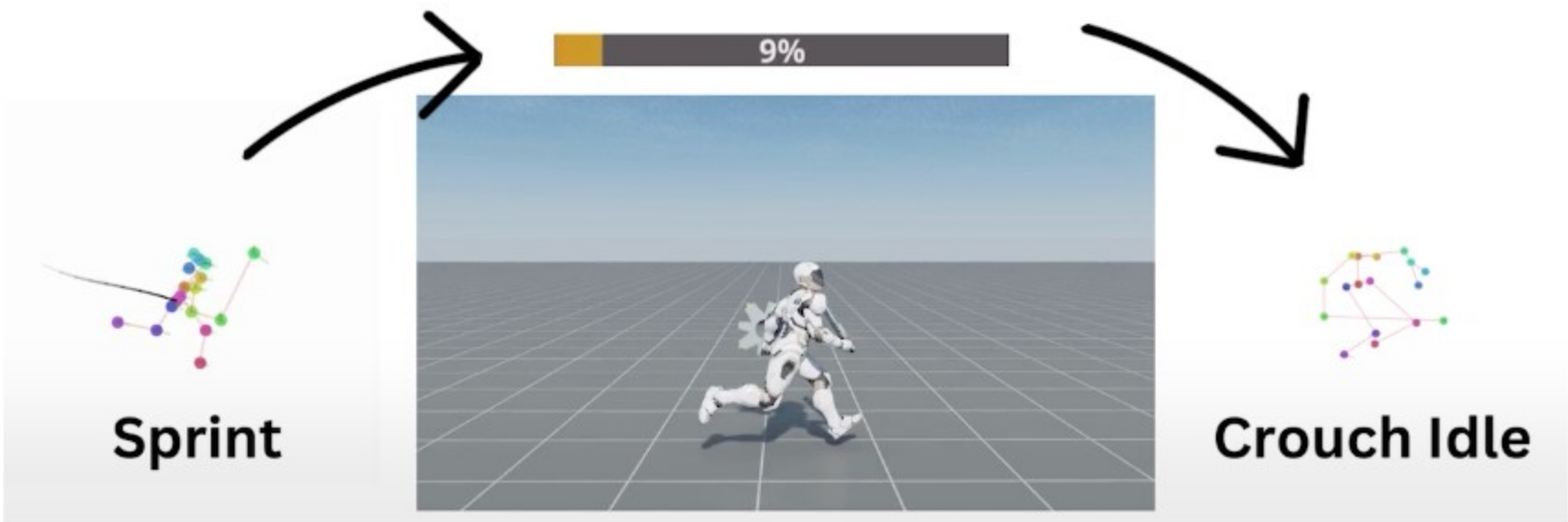
(1) Conditioned on the latent variable, a low-level controller generates motion with similar characteristics.



During low-level training, CALM learns an encoder and a decoder. The encoder takes a motion from a reference dataset of motions, a time-series of joint locations, and maps it into to a low-dimensional latent representation. Additionally, CALM also jointly learns a decoder. The decoder is a low-level policy that interacts with the simulator and generates motions similar to the reference dataset. This policy produces a variety of behaviors on demand, but is not conditioned on the directionality of the motion. For example, it can be instructed to walk, but does not enable intuitive control over the direction of walking.

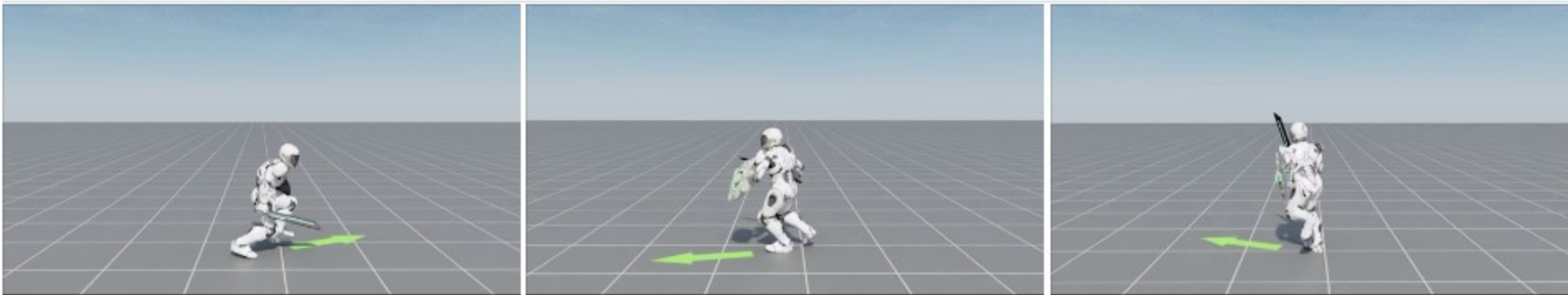
Meaningful Motion Representations

Latent space interpolation, over time, between sprint and crouch-idle



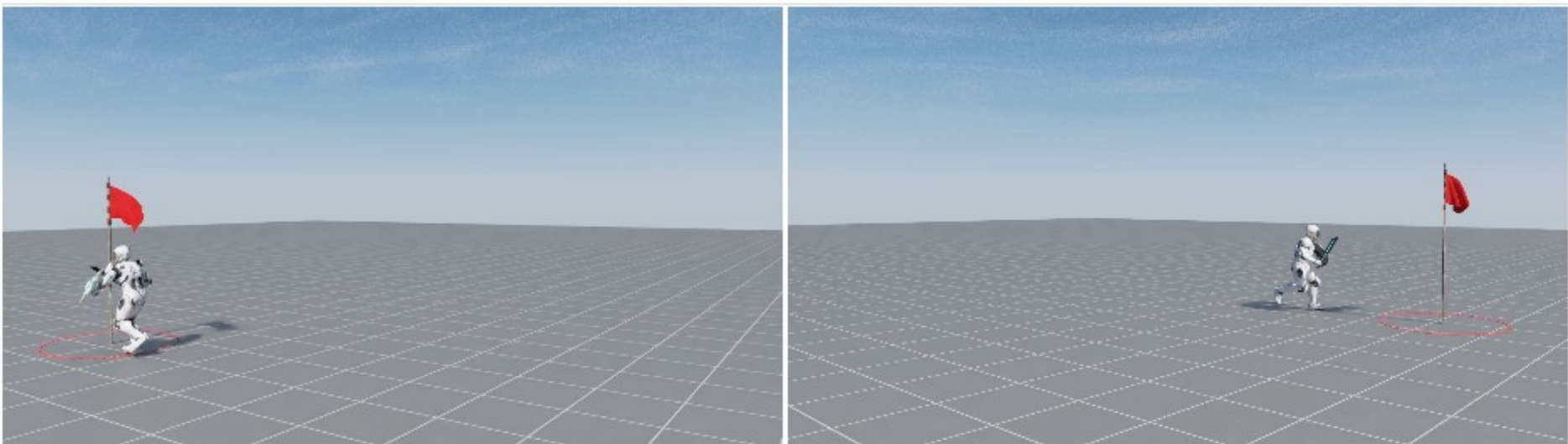
To evaluate the learned motion representation, we test the ability to interpolate between motions in the latent space. Here, the initial latent is the latent representation for sprint. The final latent is that of crouching idle. Throughout the episode, the latent is linearly interpolated over time, going from spring towards crouch-idle. The character smoothly transitions through semantically meaningful transitions, gradually reducing speed and tilting the upper body.

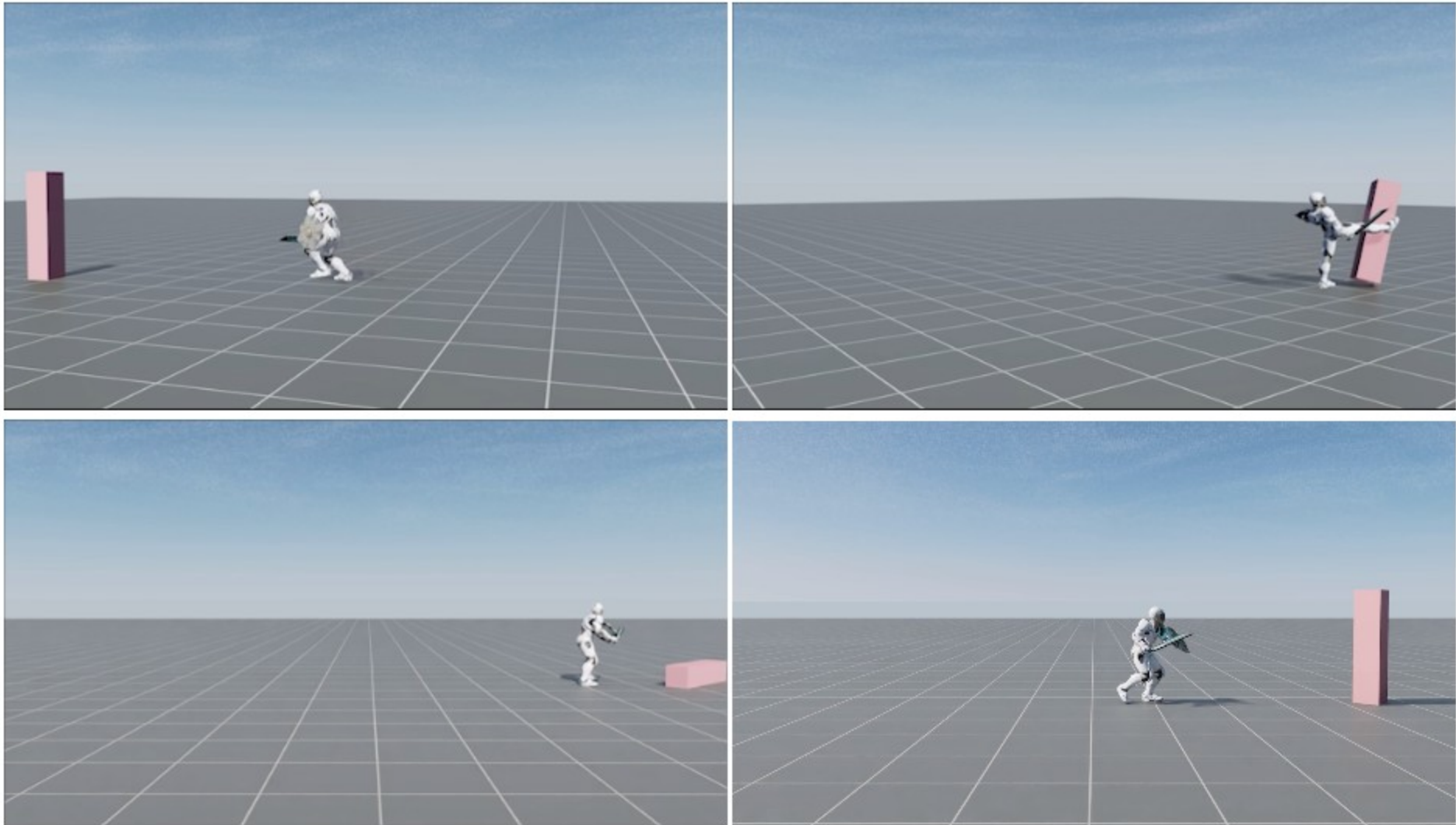
Phase 2: Directionality Control



To control motion direction, we train a high-level task-driven policy to select latent variables. These latents are provided to the low-level policy which generates the requested motion. Here, the learned motion representation enables a form of **style-conditioning**. To achieve this, the motion encoder is used to obtain the latent representation of the requested motion. The high-level policy is then provided an additional reward proportional to the cosine distance between the selected latents and the latent representing the requested style, thus guiding the high-level policy to adopt a desired behavioral style. For example, here a directionality-controller is trained to enable control over the form of loco-motion performed and the direction in which the character performs it -- crouch-walk, walk shield-up, and run.

Phase 3: Inference





Finally, the previously trained models (low-level policy and directional controller) are combined to compose complex movements without additional training. To do so, the user produces a finite-state machine (FSM) containing standard rules and commands. These determine which motion to perform, similar to how a user controls a video game character. For example, they determine whether the character should perform a simple motion, performed directly using the low-level policy, or a directed motion requiring high-level control. As an example, one may construct an FSM like (a) "crouch-walk towards the target, until distance < 1m", then (b) "kick", and finally (c) "celebrate".

Citation

```
@inproceedings{tessler2023calm,
  author = {
    Tessler, Chen
    and Kasten, Yoni
    and Guo, Yunrong
    and Mannor, Shie
    and Chechik, Gal
    and Peng, Xue Bin},
  title = {CALM: Conditional Adversarial Latent Models for Directable Virtual Characters},
  year = {2023},
  isbn = {9798400701597},
  publisher = {Association for Computing Machinery},
  address = {New York, NY, USA},
  url = {https://doi.org/10.1145/3588432.3591541},
  doi = {10.1145/3588432.3591541},
  booktitle = {ACM SIGGRAPH 2023 Conference Proceedings},
  keywords = {
    reinforcement learning,
    animated character control,
    adversarial training,
    motion capture data
  },
  location = {Los Angeles, CA, USA},
  series = {SIGGRAPH '23}
}
```

Paper

CALM: Conditional Adversarial Latent Models for Directable Virtual Characters

Chen Tessler, Yoni Kasten, Yunrong Guo, Shie Mannor, Gal Chechik, and Xue Bin Peng

- arXiv version
- Video
- BibTeX

CALM: Conditional Adversarial Latent Models for Directable Virtual Characters

CHEN TESSLER, NVIDIA, Israel
YONI KASTEN, NVIDIA, Israel
YUNRONG GUO, NVIDIA, Canada
SHIE MANNOR, NVIDIA, Israel and Technion Institute of Technology, Israel
GAL CHECHIK, NVIDIA, Israel and Bar-Ilan University, Israel
XUE BIN PENG, NVIDIA, Canada and Simon Fraser University, Canada



Fig. 1. Our framework enables users to direct the behavior of a physically-aware character using observational feedback. In the bottom row, observational inputs including a target region, task, or this example, the character is instructed to reach such location, target, but when within range, and finally when the user and character.

In this work, we present Conditional Adversarial Latent Models (CALM) as a general framework for generating and controlling behavior in a virtual environment. CALM is a generative model that takes as input a sequence of observational inputs, including a target region, task, or this example, the character is instructed to reach such location, target, but when within range, and finally when the user and character. The approach enables users to control a virtual character in a virtual environment, such as a game, by providing a sequence of observational inputs. The model takes as input a sequence of observational inputs, including a target region, task, or this example, the character is instructed to reach such location, target, but when within range, and finally when the user and character. The model takes as input a sequence of observational inputs, including a target region, task, or this example, the character is instructed to reach such location, target, but when within range, and finally when the user and character.

ACM Reference Format:
Chen Tessler, Yoni Kasten, Yunrong Guo, Shie Mannor, Gal Chechik, and Xue Bin Peng. 2024. CALM: Conditional Adversarial Latent Models for Directable Virtual Characters. In *Proceedings of the ACM Conference on Computer Graphics and Interactive Techniques (SIGGRAPH ’24)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3656448>

1 INTRODUCTION

Virtual environments and interactive simulations have become increasingly prevalent and user-friendly, but creating realistic and diverse behaviors for these virtual agents remains a challenge due to the complexity of human motion. To create interactive and immersive experiences, virtual agents must adapt to different environments and user inputs in a lifelike manner and also require the ability to perform a wide range of behaviors on demand. To that end, we need to develop virtual models that can generate complex and realistic behaviors, while taking into account the properties of the virtual world. For example, in virtual reality games, players that interact with virtual characters and objects expect them to behave realistically. This includes responding to user commands and navigating through virtual environments. When virtual agents fail to respond naturally to user input, it can disrupt the immersive experience.