

# Genie: Generative Interactive Environments



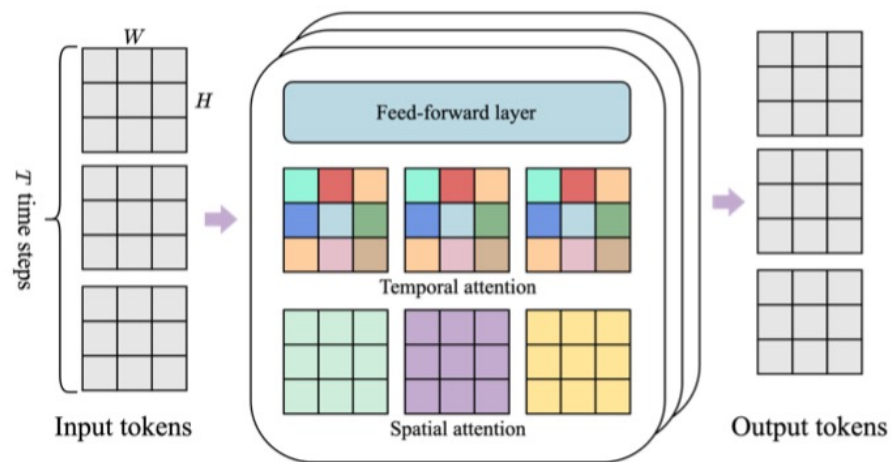
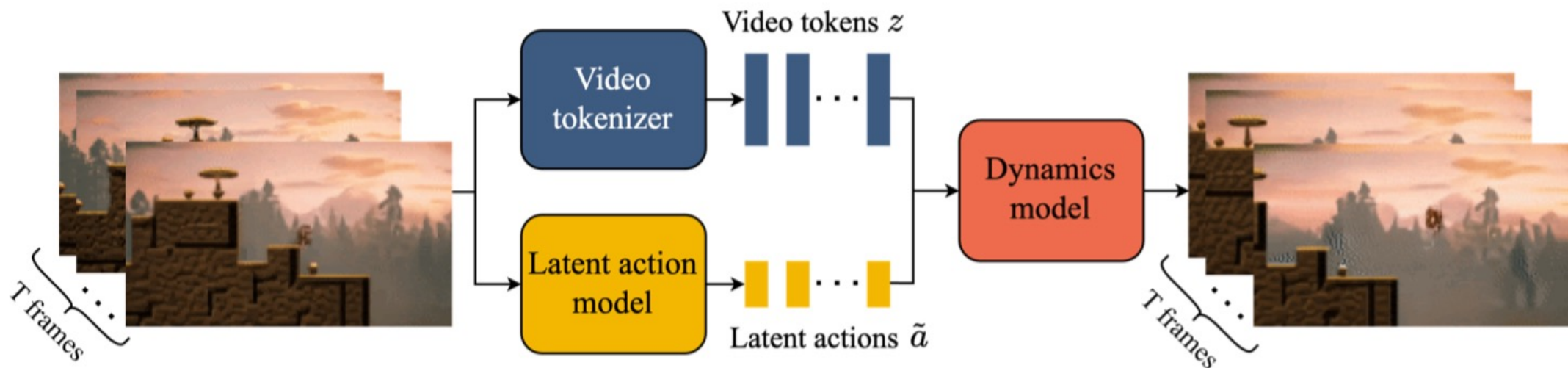
# Genie: Generative Interactive Environments

Jake Bruce<sup>\*,1</sup>, Michael Dennis<sup>\*,1</sup>, Ashley Edwards<sup>\*,1</sup>, Jack Parker-Holder<sup>\*,1</sup>, Yuge (Jimmy) Shi<sup>\*,1</sup>, Edward Hughes<sup>1</sup>, Matthew Lai<sup>1</sup>, Aditi Mavalankar<sup>1</sup>, Richie Steigerwald<sup>1</sup>, Chris Apps<sup>1</sup>, Yusuf Aytar<sup>1</sup>, Sarah Bechtle<sup>1</sup>, Feryal Behbahani<sup>1</sup>, Stephanie Chan<sup>1</sup>, Nicolas Heess<sup>1</sup>, Lucy Gonzalez<sup>1</sup>, Simon Osindero<sup>1</sup>, Sherjil Ozair<sup>1</sup>, Scott Reed<sup>1</sup>, Jingwei Zhang<sup>1</sup>, Konrad Zolna<sup>1</sup>, Jeff Clune<sup>1,2</sup>, Nando de Freitas<sup>1</sup>, Satinder Singh<sup>1</sup> and Tim Rocktäschel<sup>\*,1</sup>

<sup>\*</sup>Equal contributions, <sup>1</sup>Google DeepMind, <sup>2</sup>University of British Columbia

ICML 2024 Best Paper

# Training Pipeline



# Model

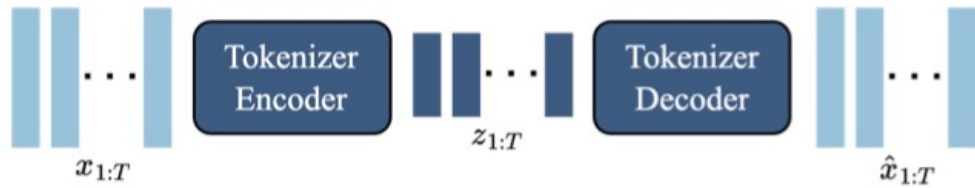


Figure 6 | **Video tokenizer:** a VQ-VAE with ST-transformer.

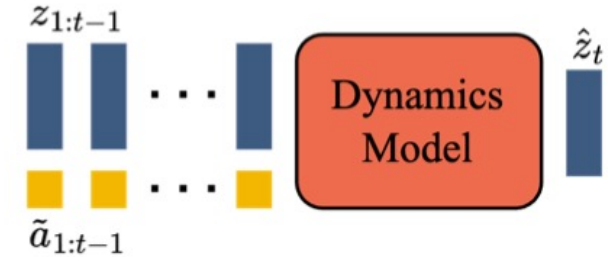


Figure 7 | **Dynamics model:** takes in video tokens and action embeddings, and predicts future masked video tokens.

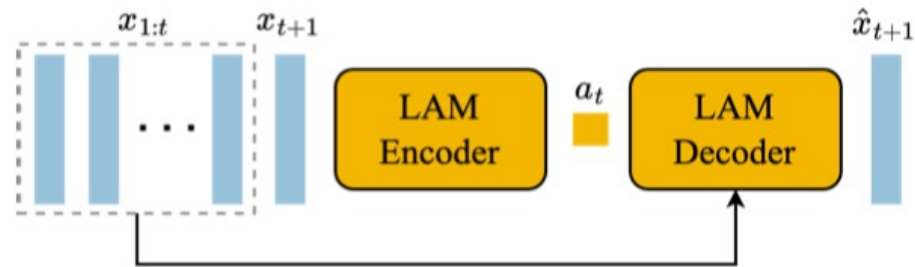
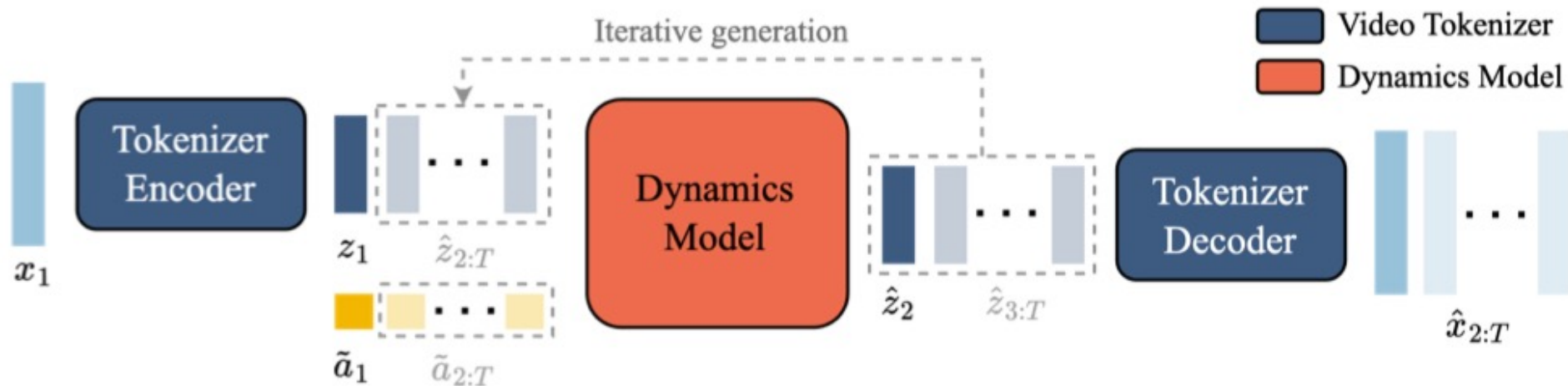
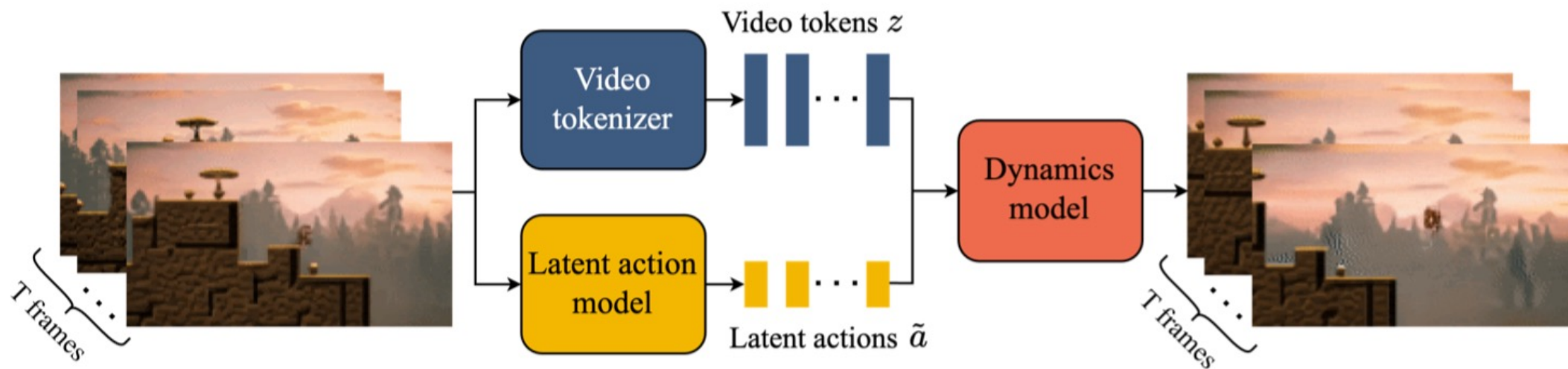


Figure 5 | **Latent action model:** learns actions  $a_t$  unsupervised from unlabelled video frames.

# Model



# Evaluation

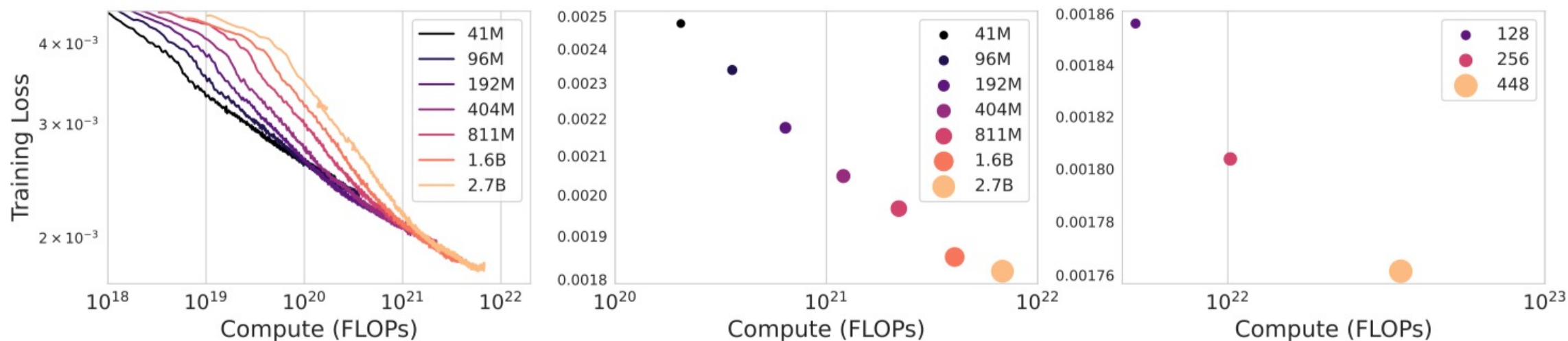


Figure 9 | **Scaling results.** **Left:** Training curves for different model sizes, **Middle:** Final training loss for each model size, averaged over the last 300 updates, **Right:** Final training loss for a 2.3B model with different batch sizes.





Figure 10 | **Playing from Image Prompts:** We can prompt Genie with images generated by text-to-image models, hand-drawn sketches or real-world photos. In each case we show the prompt frame and a second frame after taking one of the latent actions four consecutive times. In each case we see clear character movement, despite some of the images being visually distinct from the dataset.

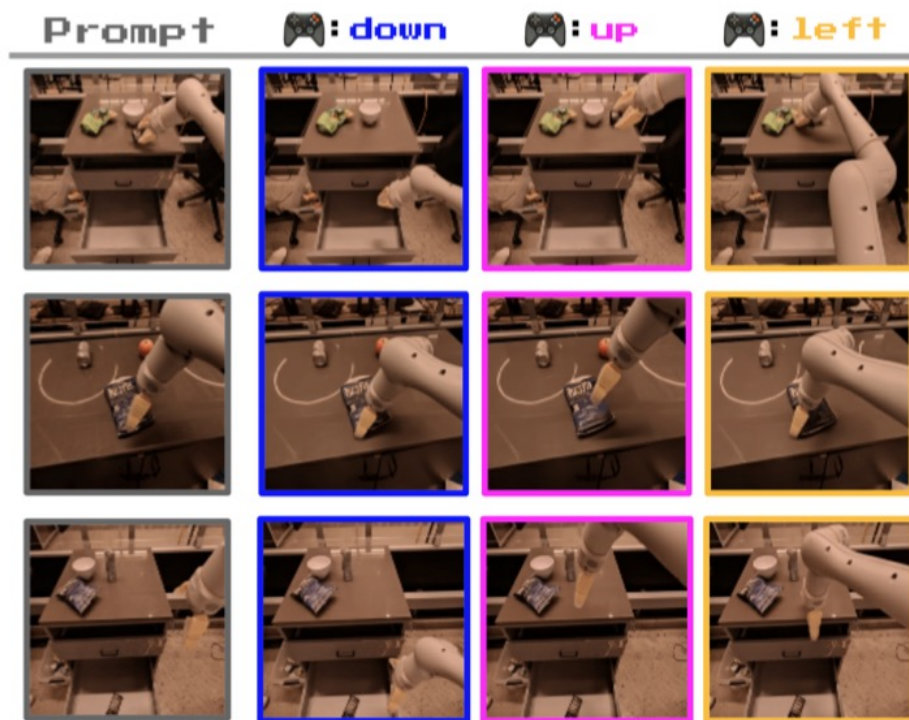


Figure 13 | **Controllable, consistent latent actions in Robotics:** trajectories beginning from three different starting frames from our Robotics dataset. Each column shows the resulting frame from taking the same latent action five times. Despite training without action labels, the same actions are consistent across varied prompt frames and have semantic meaning: *down*, *up* and *left*.



Figure 14 | **Playing from RL environments:** Genie can generate diverse trajectories given an image of an unseen RL environment.

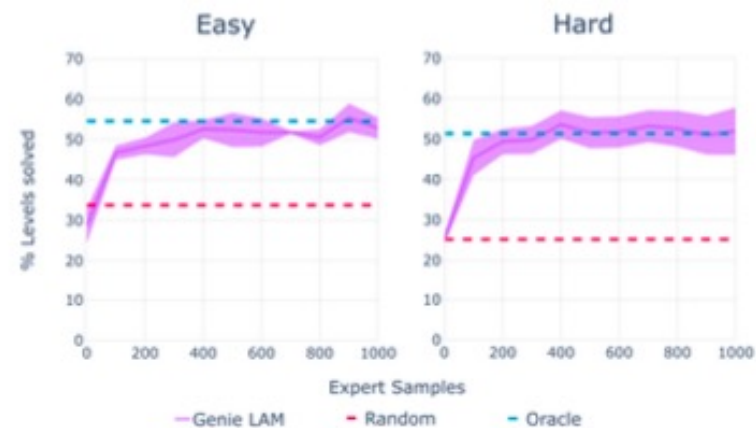


Figure 15 | **BC results.** Mean percentage of levels solved out of 100 samples, averaged over 5 seeds with 95% confidence intervals.



**Thank You**