

# COMBAT: Conditional World Models for Behavioral Agent Training

Anonymous ICCV submission

Paper ID \*\*\*\*\*

## Abstract

Recent advances in video generation have spurred the development of world models capable of simulating 3D-consistent environments and interactions with static objects. However, a significant limitation remains in their ability to model dynamic, reactive agents that can intelligently influence and interact with the world. To address this gap, we introduce COMBAT, a real-time, action-controlled world model trained on the complex 1v1 fighting game Tekken 3. Our work demonstrates that diffusion models can successfully simulate a dynamic opponent that reacts to player actions, learning its behavior implicitly.

Our approach utilizes a 1.2 billion parameter Diffusion Transformer, conditioned on latent representations from a deep compression autoencoder. We employ state-of-the-art techniques, including causal distillation and diffusion forcing, to achieve real-time inference. Crucially, we observe the emergence of sophisticated agent behavior by training the model solely on single-player inputs, without any explicit supervision for the opponent's policy. Unlike traditional imitation learning methods, which require complete action labels, COMBAT learns effectively from partially observed data to generate responsive behaviors for a controllable Player 1. We present an extensive study and introduce novel evaluation methods to benchmark this emergent agent behavior, establishing a strong foundation for training interactive agents within diffusion-based world models.

## 1. Introduction

As the fidelity of video generation methods improves with increased understanding of real-world phenomena, interactive world models trained on gameplay and real-world data have emerged [4, 5, 23]. The focus of these works remains on generating spatially and temporally consistent world simulations. Yet, in real-world scenarios, the most unpredictable components are reactive agents that can observe, plan, and influence their environment, such as in autonomous driving, navigation, and combat scenarios.

Recent works demonstrate that autoregressive diffusion

models are effective at world simulation. Several advances make these models real-time through distribution matching distillation (DMD) [26–28] and diffusion forcing [14] to overcome autoregressive drift. These engineering advances have enabled neural game simulations for first-person games such as Minecraft and CS:GO [17], showcasing excellent causal understanding of actions and their effects on generated frames.

However, real-world and game environments also contain rich information about how agents (humans, NPCs, and autonomous systems) respond to environmental dynamics. Current methods could greatly benefit from learning agent behavior from this observational data, but the partial observability and unstructured nature poses significant challenges. For example, while we might observe a pedestrian changing trajectory to avoid a vehicle, the exact observations and decision processes of the human agent remain hidden.

We present COMBAT (Conditional world Model for Behavioral Agent Training), an interactive world model that learns underlying agent behavior and movement dynamics directly from partially observed multi-agent systems. By training a world model on Tekken 3 gameplay with conditioning only on Player 1's input, we observe emergent tactical behavior in Player 2 without explicit behavioral supervision. We select Tekken 3 as it provides an ideal controlled environment with clear visual feedback, deterministic game mechanics, diverse movesets, and frame-precise timing requirements.

Our approach uses a 1.2B parameter diffusion transformer trained on 1.2M frames across 1,000 gameplay rounds. We first train a Deep Compression AutoEncoder (DCAE) to obtain highly compressed latent representations, then train the world model to generate temporally consistent gameplay sequences. COMBAT successfully learns to control Player 1 from conditioning signals, while Player 2 emerges with realistic combat behaviors including blocking, counterattacking, and combo execution. Through decoder distillation and CausVid DMD techniques, we achieve real-time generation at interactive frame rates.

We introduce novel benchmarking methods to evaluate emergent agent behavior, measuring behavioral diversity,

038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078



Figure 1. An overview of the COMBAT world model. (Top) The model is conditioned on the current state (visual frames and poses) and Player 1’s control inputs to autoregressively predict subsequent frames. (Bottom) Three distinct generated trajectories showcase the model’s ability to produce plausible, strategic counter-attacks from Player 2 as an emergent response to Player 1’s actions, without direct supervision of the opponent’s policy.

079 and tactical understanding. Our extensive analysis demonstrates that world models can serve as a new paradigm for  
080 learning agent behaviors from observational data, with implications for multi-agent AI systems beyond gaming.  
081  
082

## 083 2. Related Work

084 Our work is positioned at the intersection of generative  
085 world models, video diffusion architectures, and behavioral  
086 modeling. We review key advancements in these areas to  
087 contextualize our contribution.

### 088 2.1. Video Diffusion Models

089 The remarkable success of diffusion models in image synthesis [19, 20] has naturally inspired their extension to video  
090 generation. Early approaches adapted U-Net architectures  
091 from image models, achieving results in short-form video  
092 synthesis [3, 9]. However, the convolutional nature of U-  
093 Net presents challenges for video: it struggles to capture  
094 long-range temporal dependencies and scales poorly with  
095 sequence length, often leading to temporal incoherence.  
096

097 To address these limitations, Transformer-based video  
098 models have emerged. Following Peebles et al. [18], which  
099 demonstrated that Diffusion Transformers (DiT) could sur-  
100 pass U-Nets in image generation with superior scaling prop-  
101 erties, subsequent work has applied this architecture to

102 video. Models such as W.A.L.T [10] and CogVideoX [25]  
103 show that DiT self-attention mechanisms effectively model  
104 complex spatiotemporal relationships in video data, en-  
105 abling longer, more coherent sequences. Our work builds  
106 on this foundation, employing a DiT backbone tailored for  
107 action-conditioned dynamics in interactive environments.

### 108 2.2. Neural Game Engines and World Models

109 Recent advances demonstrate that generative models can  
110 serve as neural game engines, replacing traditional render-  
111 ing and state update logic. GameGAN, Kim et al. learns to  
112 imitate 2D games from raw pixels and actions using GANs  
113 with explicit memory modules [16]. More recently, diffu-  
114 sion transformers have become dominant for this task.

115 Valevski et al. introduce GameNGen, a fully neural  
116 DOOM engine that generates frames conditioned on past  
117 frames and actions, enabling real-time simulation [23].  
118 Alonso et al.’s DIAMOND trains diffusion-based world  
119 models achieving state-of-the-art RL performance while  
120 producing playable Counter-Strike simulations [1]. Che  
121 et al. extend this with GameGen-X, training on million-  
122 clip datasets to enable long-horizon, interactive open-world  
123 gameplay [6].

124 These methods validate that neural models can learn  
125 complex game dynamics from observational data. Our work

126 adopts similar architectural foundations but introduces a  
 127 novel objective: modeling emergent behavior of uncontrolled  
 128 opponents that arises solely from conditioning on  
 129 controllable player actions.

### 130 2.3. Multi-Modal and Behavioral World Models

131 While traditional world models focus on visual prediction,  
 132 recent work has pushed towards greater fidelity and behavioral  
 133 learning. Our work adopts joint RGB-pose representation  
 134 to enforce structural consistency in character movements.  
 135

136 In parallel, learning agent behavior within world models  
 137 has predominantly followed two paths. The first is **model-  
 138 based reinforcement learning**, where an agent’s policy  
 139 is trained using a learned dynamics model and an extrinsic  
 140 reward signal. Works like **DreamerV3** exemplify this,  
 141 achieving mastery in diverse domains by learning behaviors  
 142 entirely within the latent space of a world model [11].  
 143 The second path is **imitation learning**, which learns policies  
 144 from expert demonstrations. Methods like **Generative  
 145 Adversarial Imitation Learning (GAIL)** require explicit  
 146 state-action supervision for all agents to mimic expert be-  
 147 havior [13].

148 Our approach diverges from both paradigms. We demon-  
 149 strate that complex, reactive multi-agent behaviors can  
 150 emerge implicitly as a property of world modeling itself,  
 151 without engineered rewards and using only partially ob-  
 152 served data where just one agent’s actions are provided as a  
 153 condition.

### 154 2.4. Optimization Techniques for Interactive Gen- 155 eration

156 Real-time interactive generation requires addressing both  
 157 architectural efficiency and sampling speed. Recent ad-  
 158 vances in attention mechanisms include FlexAttention [8],  
 159 which enables flexible attention patterns, and Longformer  
 160 [2], which combines local sliding-window attention with  
 161 global context. We incorporate local-global attention pat-  
 162 terns inspired by these works to balance efficiency with tem-  
 163 poral coverage.

164 For sampling efficiency, Distribution Matching Distilla-  
 165 tion (DMD) [26, 28] and diffusion forcing [14] have proven  
 166 effective at reducing sampling steps while mitigating au-  
 167 toregressive drift. These techniques enable real-time neu-  
 168 ral simulation for complex games [5, 23]. We adapt DMD  
 169 through CausVid distillation to achieve interactive frame  
 170 rates while preserving behavioral quality.

171 The Muon optimizer [15] introduces orthogonalization  
 172 into momentum-based updates, improving conditioning of  
 173 weight updates and outperforming AdamW in training  
 174 speed benchmarks. We incorporate Muon optimization  
 175 to enhance training efficiency of our large-scale diffusion  
 176 transformers.

## 3. Method

177 Our approach, **COMBAT**, learns to simulate a complex,  
 178 multi-agent environment by training a generative world  
 179 model on video observations. World models have shown  
 180 promise in mastering diverse domains [11] and creating in-  
 181 teractive environments [5, 24]. We extend this paradigm to  
 182 a competitive fighting game, where the model must learn  
 183 the opponent’s behavior without explicit action labels.  
 184

### 185 3.1. Problem Formulation

186 We frame our task as learning a conditional video gener-  
 187 ation model that implicitly captures an opponent’s policy.  
 188 We select the fighting game *Tekken 3* as our environment  
 189 for three key reasons:

- 190 **Bounded Temporal Dependency:** The game state is  
 191 largely Markovian, where

$$P(s_{t+1} | s_{\leq t}) \approx P(s_{t+1} | s_{t-k:t}),$$

192 for a small history window  $k$ , since all relevant informa-  
 193 tion is contained within recent frames.

- 194 **Rich Action Space:** Characters possess diverse  
 195 movesets, with over 40 unique actions and complex  
 196 combos, providing a challenging domain for behavior  
 197 modeling.
- 198 **Strategic Depth:** Success requires a blend of rapid re-  
 199 actions and long-term tactical planning.

200 **Formal Problem Statement:** Given a dataset of par-  
 201 tially observed multi-agent trajectories

$$D = \{(s_t, a_t^{(1)}, s_{t+1})\}_{t=1}^T,$$

202 where  $s_t \in \mathbb{R}^{H \times W \times 3}$  is a game frame and  $a_t^{(1)} \in \{0, 1\}^8$   
 203 is the observed multi-hot input for Player 1. The actions of  
 204 Player 2,  $a_t^{(2)}$ , remain unobserved. Our objective is to learn  
 205 a conditional world model

$$P_\theta(s_{t+1} | s_{t-k:t}, a_{t-k:t}^{(1)})$$

206 that can accurately predict subsequent frames.

207 **Key Innovation:** Unlike traditional imitation learning  
 208 methods that require explicit action supervision for all  
 209 agents [13], COMBAT is trained without Player 2’s action  
 210 labels. The model must infer Player 2’s policy,  
 211

$$\pi^{(2)}(a_t^{(2)} | s_t, a_t^{(1)}),$$

212 as an emergent property of generating temporally consis-  
 213 tent and plausible multi-agent interactions. This forces the  
 214 world model to learn reactive and strategic opponent behav-  
 215 ior implicitly.

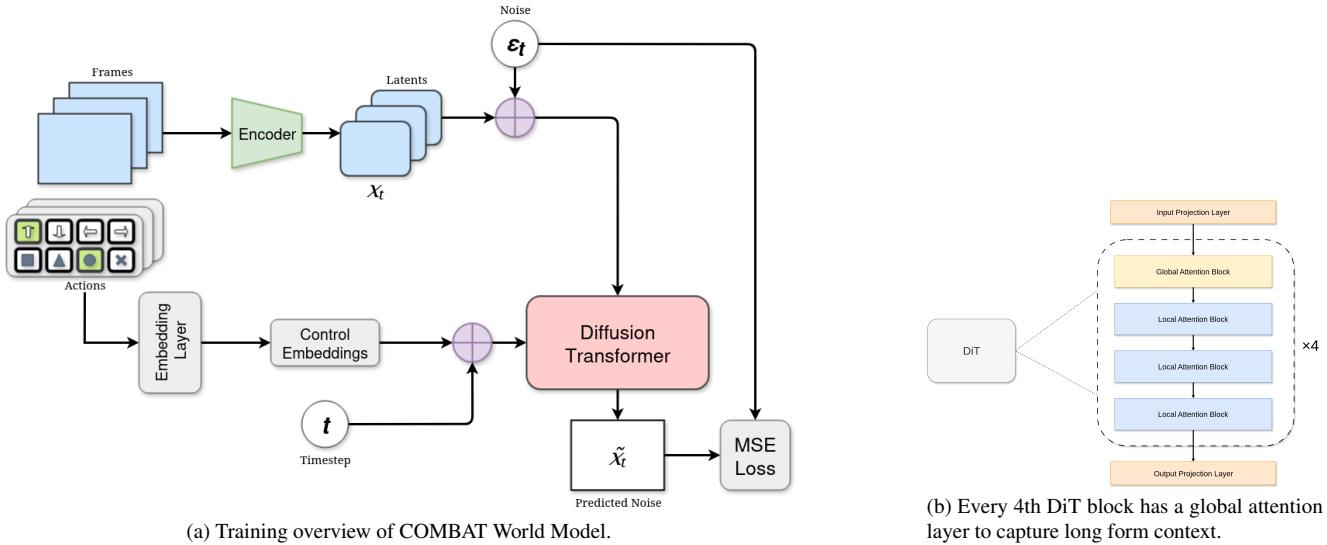


Figure 2. Architectural diagram of the COMBAT model. (a) The end-to-end training process, where a Diffusion Transformer is conditioned on action and timestep embeddings to denoise latent frame representations. (b) The internal structure of the DiT backbone, which employs a hybrid local-global attention pattern to efficiently model long-term dependencies.

219

### 3.2. Tekken 3 Gameplay Dataset

220

To train our model, we collected a large-scale dataset of *Tekken 3* gameplay, totaling 1,000 rounds (approximately 7 hours or 1.2 million frames). The data features a variety of characters and a balanced win–loss ratio between the two players. For each frame, captured at a resolution of  $3 \times 448 \times 736$ , we provide synchronized annotations, including: action inputs for both players, health and timer status, 68-point body pose coordinates, and player segmentation masks. Our data collection and annotation pipeline will be made publicly available.

230

### 3.3. Model Architecture

231

Our world model architecture integrates three main components:

232

- (1) a multi-modal variational autoencoder for high-ratio state compression,
- (2) an embedding module for player actions and diffusion timesteps, and
- (3) a Diffusion Transformer (DiT) backbone for autoregressive prediction in the latent space.

233

We train two versions of the model: one using only RGB latents and another using a joint visual–pose latent representation.

234

#### 3.3.1. Multi-Modal Latent Encoding

235

To create an efficient latent representation, we first train a 340M-parameter joint RGB–pose variational autoencoder. This model learns a shared embedding space by compressing concatenated visual frames ( $3 \times 448 \times 736$ ) and pose key-points into a compact latent tensor of shape  $128 \times 23 \times 11$ .

Our design is inspired by recent work in high-compression autoencoders for diffusion models [7]. To optimize for real-time performance, the 340M-parameter decoder is subsequently distilled to a 44M-parameter version by reducing its upsampling block count, which maintains high reconstruction quality at a fraction of the computational cost.

Player 1’s action history, encoded as a multi-hot vector over 8 buttons, is projected into a dense embedding. This action embedding is summed with a sinusoidal time embedding for the current diffusion step,  $t_{emb}$ , to form the final conditioning vector for the DiT backbone.

#### 3.3.2. Diffusion Transformer Backbone

The core of our generative model is a 1.2B-parameter Diffusion Transformer (DiT) [18], which learns to denoise and predict future latent frames. The architecture consists of 16 transformer blocks with a model dimension  $d_{model} = 2048$  and 16 attention heads. The conditioning vector is injected into each block via an Adaptive Layer Normalization Zero (AdaLNZero) layer, and tokenization is performed using linear projection layers for spatio-temporal rasterization, bypassing conventional patch-based embeddings.

Each DiT block executes the following sequence:

$$\text{AdaLN} \rightarrow \text{Attention} \rightarrow \text{Gated Residual} \rightarrow \\ \text{AdaLN} \rightarrow \text{MLP} \rightarrow \text{Gated Residual}$$

To maintain computational tractability over long 128-frame sequences, we employ a hybrid attention strategy. Most layers use a frame-causal attention mask with a local sliding window of 16 frames, while every fourth layer applies global attention across the entire 128-frame context.

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277 This structure balances long-range dependency modeling  
 278 with computational efficiency. We apply Rotary Position  
 279 Embeddings (RoPE) [21] across both spatial and temporal  
 280 axes and utilize FlexAttention for an efficient block-sparse  
 281 masking implementation.

### 282 3.4. Accelerated Inference for Real-Time Genera- 283 tion

284 Enabling real-time interaction is critical for gaming applica-  
 285 tions, but the iterative sampling process of diffusion models  
 286 is computationally intensive. To overcome this, we signifi-  
 287 cantly accelerate inference using two key optimizations.

288 First, we distill the fully trained model into a few-  
 289 step sampler using Distribution Matching Distillation  
 290 (DMD) [26, 27]. Specifically, we adopt the CausVid DMD  
 291 framework [28] to produce a 4-step distilled model that pre-  
 292 serves high generative fidelity while drastically reducing in-  
 293 ference time.

294 Second, we further enhance speed by implementing  
 295 static key-value caching, which reuses previously computed  
 296 attention states across generation steps. These optimiza-  
 297 tions are applied to both the RGB and visual–pose world  
 298 models.

## 299 4. Experiments

300 To validate our claim that a conditional world model can  
 301 learn reactive agent behavior from partial observations, we  
 302 conduct a series of experiments on the Tekken 3 dataset.  
 303 We first detail our multi-stage training pipeline and model  
 304 architectures. We then introduce our evaluation benchmarks  
 305 and present results comparing our primary models and their  
 306 distilled variants.

### 307 4.1. Implementation Details

308 Our training process is divided into three main stages: au-  
 309 toencoder training, world model training, and distillation  
 310 for real-time inference. All models were trained on a cluster  
 311 of 8× NVIDIA H200 GPUs.

312 **Stage 1: Autoencoder Training.** We first train a  
 313 340M parameter Deep Compression AutoEncoder (DCAE)  
 314 to learn a compact latent representation of the game envi-  
 315 ronment. The autoencoder is trained for 68,000 steps (ap-  
 316 prox. 75 hours) on our 1.2 million frame Tekken dataset. It  
 317 compresses raw frames ( $3 \times 448 \times 736$ ) into a latent space  
 318 of  $23 \times 11$  with 128 channels. The training objective is a  
 319 combination of L2 reconstruction loss, perceptual simila-  
 320 rity loss, and a KL divergence term to regularize the latent  
 321 space. For our pose-augmented model, we use an identical  
 322 architecture and training setup.

323 **Stage 2: World Model Training.** We train a 1.2B  
 324 parameter autoregressive Diffusion Transformer (DiT) to  
 325 function as the world model. The DiT architecture consists  
 326 of 16 layers, 16 attention heads, and a model dimension

of  $d_{model} = 2048$ . It employs a combination of local (16  
 327 frames) and global (128 frames) attention windows to cap-  
 328 ture both short-term and long-term temporal dependencies.  
 329 The model is trained on video clips with a sequence length  
 330 of 128 frames to predict the next latent frame conditioned  
 331 on Player 1’s actions. We train two distinct world models:  
 332 one using latents from the RGB-only VAE and another us-  
 333 ing latents from the pose-augmented VAE.

334 **Stage 3: Distillation for Real-Time Inference.** To  
 335 achieve interactive frame rates, we employ two separate dis-  
 336 tillation techniques:

337 • **Decoder Distillation:** We first create a lightweight VAE  
 338 decoder for real-time rendering. Using student-teacher  
 339 distillation, we reduce the number of upsampling blocks  
 340 per stage in the decoder from four to one. This process,  
 341 which took 14 hours over 50k steps, reduces the decoder’s  
 342 parameter count from 340M to a nimble 44M.

343 • **Step Distillation:** We use CausVid, a Distribution Match-  
 344 ing Distillation (DMD) method, to drastically reduce the  
 345 number of required inference steps for the world model.  
 346 We distill the fully-trained DiT into a 4-step variant. This  
 347 distillation process converges in 2,500 steps, utilizing a  
 348 combination of a DMD loss and a critic loss. We ap-  
 349 ply this technique to both the RGB-only and the pose-  
 350 augmented world models.

### 351 4.2. Evaluation Metrics and Benchmarks

352 Evaluating emergent agent behavior presents a fundamental  
 353 challenge: how do we measure intelligence that was never  
 354 explicitly supervised? Traditional video metrics assess vi-  
 355 sual fidelity, while RL metrics assume access to ground-  
 356 truth actions or rewards. Since COMBAT learns behavioral  
 357 patterns implicitly through world modeling, we need novel  
 358 evaluation approaches that can detect tactical competence  
 359 from generated gameplay alone.

#### 360 4.2.1. Standard Perceptual Metrics

361 To assess the perceptual quality of our generated trajec-  
 362 tories, we employ a suite of standard metrics. Our evalua-  
 363 tion protocol involves conditioning the models on real  
 364 Player 1 action sequences extracted from a test set of 300  
 365 ground-truth videos(1-2 seconds) consisting mixed diffi-  
 366 culty gameplays. The generated video is then compared  
 367 directly against its corresponding ground-truth counterpart  
 368 from which the actions were sourced. This setup provides  
 369 a stringent test of the model’s ability to render determinis-  
 370 tic outcomes based on specific actions,a significantly more  
 371 challenging task than unconditional video generation.

372 We report the Fréchet Video Distance (FVD)[22] to mea-  
 373 sure temporal coherence, the Fréchet Inception Distance  
 374 (FID)[12] for per-frame visual fidelity, and LPIPS to quan-  
 375 tify perceptual similarity. Given the high-fidelity nature of  
 376 the Tekken 3 environment, characterized by rapid motion  
 377 and complex visual effects, achieving strong performance

379  
380

on these metrics against the ground truth is a robust indicator of the model's precision and world-modeling capabilities.

381

Table 1. All metrics are calculated on a held-out test set of 300 video clips each with 32 frames. Lower is better for all scores.

Model	FID ↓	FVD ↓	LPIPS ↓
COMBAT: Pose	49.7	593.4	0.05
COMBAT: Non-Pose	80.9	1156.6	0.07

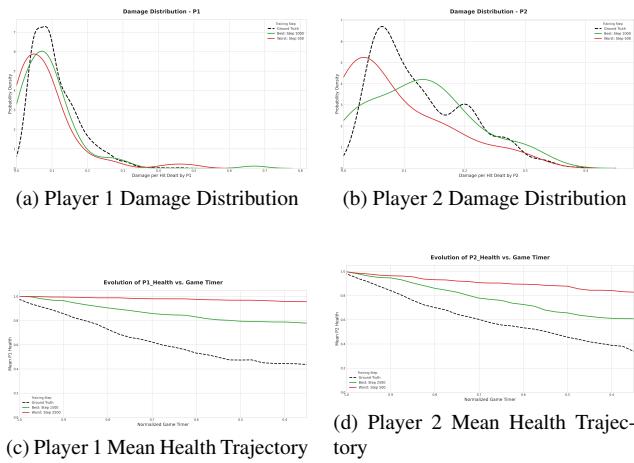


Figure 3. **Behavioral Consistency Metrics.** A comparison of generated gameplay (COMBAT) against the ground truth. **(a, b)** The per-frame damage distributions for Player 1 and Player 2, showing that our model learns a realistic mapping of actions to consequences. **(c, d)** The mean health trajectories over the course of a round, indicating that COMBAT captures the natural pacing of a match.

382

#### 4.2.2. Behavioral Consistency Metrics

383  
384  
385

To verify that our model learns the game's intrinsic rules and pacing, we propose two metrics based on in-game health data:

386

**• Damage Distribution Analysis:** This metric assesses whether the consequence of individual actions is realistic. Let  $H_i^{(t)}$  denote the health of player  $i \in \{1, 2\}$  at frame  $t$ , and define per-frame damage as  $\Delta H_i^{(t)} = \max(0, H_i^{(t-1)} - H_i^{(t)})$ . We normalize by the maximum health  $H_i^{\max}$  to obtain  $\delta_i^{(t)} = \Delta H_i^{(t)} / H_i^{\max}$ . The complete distribution of damage values from all generated sequences,  $\{\delta_{i,\text{gen}}^{(t)}\}$ , is then compared to the distribution from all ground-truth sequences,  $\{\delta_{i,\text{real}}^{(t)}\}$ , using the Wasserstein distance. A lower distance signifies that the model has learned a more accurate mapping from actions to their in-game consequences.

**• Health Trajectory Analysis:** This metric evaluates the overall temporal flow of the match. Define the normalized time  $s = t/T$ , where  $T$  is the total round duration, and let  $\bar{H}^{(s)} = \frac{1}{2} \sum_i H_i^{(t)} / H_i^{\max}$  be the average normalized health at time  $s$  for a single round.

To establish a baseline for typical match progression, we compute the **mean health trajectory** by averaging  $\bar{H}^{(s)}$  across all rounds in our ground-truth test set. We do the same for our generated rounds. The similarity between these two mean trajectories is then measured using the Mean Squared Error (MSE). A lower MSE indicates that the generated gameplay, on average, exhibits a more realistic match pace.

### 4.3. Human Evaluation of Emergent Behavior

To assess the emergent behavior of Player 2, we conduct human evaluation based on observable action patterns in gameplay. Since Player 2 is trained without explicit supervision, emergent behavior is defined as actions that react naturally to Player 1's inputs, demonstrating plausible combat strategies such as timely punches, kicks, and defensive maneuvers.

We introduce two human-interpretable metrics: **Total Action Adherence (TAA)** and **Action Ratio Consistency (ARC)**. These metrics are based on human annotations of offensive actions observed in both ground-truth and generated gameplay sequences.

#### 4.3.1. Total Action Adherence (TAA)

TAA measures whether the agent produces a comparable overall volume of offensive actions relative to human gameplay:

$$\text{TAA} = \frac{G_{\text{kicks}} + G_{\text{punch}}}{O_{\text{kicks}} + O_{\text{punch}}}$$

where  $G_{\cdot}$  denotes actions performed by the generated agent, and  $O_{\cdot}$  the actions performed in original gameplay.

A score of 1.0 indicates perfect adherence in activity level. Scores  $> 1.0$  suggest hyperactive behavior, while scores  $< 1.0$  indicate passive behavior.



Figure 4. Total Action Adherence across training checkpoints

### 434 4.3.2. Action Ratio Consistency (ARC)

435 ARC evaluates whether the stylistic balance between  
 436 punches and kicks aligns with the human player:

$$437 \quad \text{ARC} = \frac{\frac{G_{\text{punch}}}{G_{\text{kicks}}}}{\frac{O_{\text{punch}}}{O_{\text{kicks}}}}$$

438 A score of 1.0 indicates identical punch-to-kick ratio as  
 439 original gameplay. Scores above 1.0 reflect stronger pref-  
 440 erence for punches, while scores below 1.0 suggest heavier  
 441 reliance on kicks.



442 Figure 5. Action Ratio Consistency across training checkpoints

### 443 4.3.3. Results

444 We evaluated sequences at multiple training checkpoints.  
 Table 2 summarizes the results:

Training Step	TAA	ARC
Ground Truth	1.00	1.00
Step 500	3.87	1.04
Step 1000	0.88	3.90
Step 1500	1.90	1.79
Step 2000	1.79	1.47

445 Table 2. TAA and ARC scores at different training checkpoints  
 compared against human gameplay.

446 Our evaluation shows that COMBAT successfully learns  
 447 emergent Player 2 behavior through distinct phases. Initially,  
 448 the model is hyperactive, generating nearly four  
 449 times the offensive actions of human players (TAA = 3.87),  
 450 though its punch-to-kick ratio is well-aligned (ARC = 1.04).  
 451 As training progresses, the model reduces hyperactivity in  
 452 further steps. Beyond step 2000, performance declines,  
 453 with later checkpoints showing reduced adherence to origi-  
 454 nal gameplay.

455 By the final training stages, the model converges toward  
 456 stable, human-like combat patterns. It learns to regulate ac-  
 457 tivity frequency (TAA 1.8) while achieving balanced fight-  
 458 ing style (ARC 1.5). However, overall consistency de-  
 459 grades noticeably. This progression from erratic behavior  
 460 to stable patterns demonstrates that complex, emergent be-  
 461 haviors can be learned without explicit supervision.

462 The pose-augmented COMBAT model significantly out-  
 463 performs the RGB-only variant across visual quality met-  
 464 rrics, confirming that explicit pose information improves  
 generation quality.

465 **Impact of Distillation:** Our 4-step distilled models, cre-  
 466 ated using CausVid DMD, retain substantial visual qual-  
 467 ity while achieving 12.5x speedup. The pose-augmented  
 468 4-step model still outperforms the full RGB-only model,  
 469 demonstrating efficient distillation with minimal quality  
 470 trade-off.

471 Qualitatively, we observe intelligent behaviors includ-  
 472 ing combo execution, spatial awareness, and adaptation to  
 473 Player 1’s patterns. These tactical responses emerge natu-  
 474 rally from our training process without explicit behavioral  
 475 supervision.

## 476 5. Conclusion

477 In this work, we introduce COMBAT, a conditional world  
 478 model that learns complex, emergent agent behavior from  
 479 partially observed gameplay. Our key finding is that by con-  
 480 ditioning the model solely on Player 1’s actions, it success-  
 481 fully learns a reactive, tactically coherent policy for Player  
 482 2 without any direct supervision. The model correctly asso-  
 483 ciates the control inputs with the intended agent and gener-  
 484 ates plausible counter-attacks, demonstrating that intricate  
 485 behaviors can arise implicitly from the objective of tempo-  
 486 ral consistency.

487 To foster further research in this domain, we provide  
 488 an extensive analysis of emergent behavior in world mod-  
 489 els. We will also release our large-scale **Tekken 3 dataset**,  
 490 complete with synchronized pose and segmentation anno-  
 491 tations, and **open-source our pipelines** for data collection  
 492 and model training.

493 Crucially, our approach is practical for interactive applica-  
 494 tions. Through distillation, the COMBAT world model  
 495 achieves **real-time performance, operating at 85 FPS on**  
 496 **a single NVIDIA A100 GPU**. This work represents a first  
 497 step in exploring how generative world models can learn  
 498 implicit agent policies, and we hope it inspires further re-  
 499 search into multi-agent behavioral modeling in complex, in-  
 500 teractive environments.

## 501 References

- [1] Eloi Alonso, Adam Jolley, Vincent Micheli, Anssi Kanervisto, Amos Storkey, Tim Pearce, and François Fleuret. Diffusion for world modeling: Visual details matter in atari. In *Thirty-eighth Conference on Neural Information Processing Systems*, 2024. 2
- [2] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv:2004.05150*, 2020. 3
- [3] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, 502 503 504 505 506 507 508 509 510 511 512

- 513 and Robin Rombach. Stable video diffusion: Scaling latent  
514 video diffusion models to large datasets, 2023. 2
- 515 [4] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue,  
516 Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luh-  
517 man, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya  
518 Ramesh. Video generation models as world simulators.  
519 2024. 1
- 520 [5] Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-  
521 Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi  
522 Mavalankar, Richie Steigerwald, Chris Apps, Yusuf Aytar,  
523 Sarah Bechtle, Feryal M. P. Behbahani, Stephanie Chan,  
524 Nicolas Manfred Otto Heess, Lucy Gonzalez, Simon Osin-  
525 dero, Sherjil Ozair, Scott Reed, Jingwei Zhang, Konrad  
526 Zolna, Jeff Clune, Nando de Freitas, Satinder Singh, and Tim  
527 Rocktaschel. Genie: Generative interactive environments.  
528 *ArXiv*, abs/2402.15391, 2024. 1, 3
- 529 [6] Haoxuan Che, Xuanhua He, Quande Liu, Cheng Jin, and  
530 Hao Chen. Gamegen-x: Interactive open-world game video  
531 generation. In *The Thirteenth International Conference on*  
532 *Learning Representations*, 2025. 2
- 533 [7] Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang,  
534 Haotian Tang, Muyang Li, and Song Han. Deep compres-  
535 sion autoencoder for efficient high-resolution diffusion mod-  
536 els. In *The Thirteenth International Conference on Learning*  
537 *Representations*, 2025. 4
- 538 [8] Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang,  
539 and Horace He. Flex attention: A programming model for  
540 generating optimized attention kernels, 2024. 3
- 541 [9] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang,  
542 Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin,  
543 and Bo Dai. Animatediff: Animate your personalized  
544 text-to-image diffusion models without specific tuning. In *The*  
545 *Twelfth International Conference on Learning Representa-*  
546 *tions*, 2024. 2
- 547 [10] Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera  
548 Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama.  
549 Photorealistic video generation with diffusion models, 2023.  
550 2
- 551 [11] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy  
552 Lillicrap. Mastering diverse domains through world models,  
553 2024. 3
- 554 [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner,  
555 Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter.  
556 Gans trained by a two time-scale update rule converge to a  
557 nash equilibrium. *CoRR*, abs/1706.08500, 2017. 5
- 558 [13] Jonathan Ho and Stefano Ermon. Generative adversarial im-  
559 itation learning. In *Advances in Neural Information Process-  
560 ing Systems*. Curran Associates, Inc., 2016. 3
- 561 [14] Xun Huang, Zhengqi Li, Guande He, Mingyuan Zhou,  
562 and Eli Shechtman. Self forcing: Bridging the train-  
563 test gap in autoregressive video diffusion. *arXiv preprint*  
564 *arXiv:2506.08009*, 2025. 1, 3
- 565 [15] Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz  
566 Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An  
567 optimizer for hidden layers in neural networks, 2024. 3
- 568 [16] Seung Wook Kim, Yuhao Zhou, Jonah Philion, Antonio Tor-  
569 rralba, and Sanja Fidler. Learning to Simulate Dynamic En-  
570 vironments with GameGAN. In *IEEE Conference on Com-  
571 puter Vision and Pattern Recognition (CVPR)*, 2020. 2
- 572 [17] Jiaqi Li, Junshu Tang, Zhiyong Xu, Longhuang Wu, Yuan  
573 Zhou, Shuai Shao, Tianbao Yu, Zhiguo Cao, and Qinglin Lu.  
574 Hunyuan-gamecraft: High-dynamic interactive game video  
575 generation with hybrid history condition, 2025. 1
- 576 [18] William Peebles and Saining Xie. Scalable diffusion models  
577 with transformers. *arXiv preprint arXiv:2212.09748*, 2022.  
578 2, 4
- 579 [19] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray,  
580 Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever.  
581 Zero-shot text-to-image generation, 2021. 2
- 582 [20] Robin Rombach, Andreas Blattmann, Dominik Lorenz,  
583 Patrick Esser, and Björn Ommer. High-resolution image syn-  
584 thesis with latent diffusion models, 2022. 2
- 585 [21] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo  
586 Wen, and Yunfeng Liu. Roformer: Enhanced transformer  
587 with rotary position embedding, 2023. 5
- 588 [22] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach,  
589 Raphaël Marinier, Marcin Michalski, and Sylvain Gelly.  
590 FVD: A new metric for video generation, 2019. 5
- 591 [23] Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi  
592 Fruchter. Diffusion models are real-time game engines. In  
593 *International Conference on Representation Learning*, pages  
594 73754–73776, 2025. 1, 2, 3
- 595 [24] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan  
596 Tompson, Dale Schuurmans, and Pieter Abbeel. Learn-  
597 ing interactive real-world simulators. *arXiv preprint*  
598 *arXiv:2310.06114*, 2023. 3
- 599 [25] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu  
600 Huang, Jiazhen Xu, Yuanming Yang, Wenyi Hong, Xiao-  
601 han Zhang, Guanyu Feng, Da Yin, Yuxuan.Zhang, Weihan  
602 Wang, Yean Cheng, Bin Xu, Xiaotao Gu, Yuxiao Dong, and  
603 Jie Tang. Cogvideox: Text-to-video diffusion models with  
604 an expert transformer. In *The Thirteenth International Con-  
605 ference on Learning Representations*, 2025. 2
- 606 [26] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang,  
607 Eli Shechtman, Fredo Durand, and William T Freeman. Im-  
608 proved distribution matching distillation for fast image syn-  
609 thesis. In *NeurIPS*, 2024. 1, 3, 5
- 610 [27] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shecht-  
611 man, Frédo Durand, William T Freeman, and Taesung Park.  
612 One-step diffusion with distribution matching distillation. In  
613 *CVPR*, 2024. 5
- 614 [28] Tianwei Yin, Qiang Zhang, Richard Zhang, William T Free-  
615 man, Fredo Durand, Eli Shechtman, and Xun Huang. From  
616 slow bidirectional to fast autoregressive video diffusion mod-  
617 els. 2025. 1, 3, 5

# COMBAT: Conditional World Models for Behavioral Agent Training

## Supplementary Material

618

### 6. Rationale

619

Having the supplementary compiled together with the main  
620 paper means that:

621

- The supplementary can back-reference sections of the  
622 main paper, for example, we can refer to Sec. 1;
- The main paper can forward reference sub-sections  
623 within the supplementary explicitly (e.g. referring to a  
624 particular experiment);
- When submitted to arXiv, the supplementary will already  
625 included at the end of the paper.

626

To split the supplementary pages from the main paper, you  
627 can use [Preview \(on macOS\)](#), [Adobe Acrobat](#) (on all OSs),  
628 as well as [command line tools](#).