

# ADAPTING VISION-LANGUAGE MODELS FOR EVALUATING WORLD MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

World models – generative models that simulate environment dynamics conditioned on past observations and actions – are increasingly central to planning, simulation, and embodied AI. However, evaluating their rollouts remains challenging: existing metrics provide coarse, semantics-agnostic signals, while human evaluation is costly and hard to scale. Effective evaluation requires fine-grained, temporally grounded assessment of action alignment and semantic consistency – capabilities that vision-language models (VLMs) possess but have not been systematically explored for this purpose. We present a case study investigating whether lightweight VLM adaptation can provide reliable semantic evaluation of world model rollouts. We introduce a semantic evaluation protocol targeting two core recognition tasks – action recognition and character recognition – assessed across binary, multiple-choice, and open-ended question formats. To support this protocol, we develop UNIVERSE (*UN*ified *V*ision-language *E*valuator for *R*ollouts in *S*imulated *E*nvironments), a parameter-efficient VLM adaptation method tailored to rollout evaluation. Through extensive experiments totaling over 5,154 GPU-days, we explore full, partial, and parameter-efficient adaptation methods across various task formats, context lengths, sampling methods, and data compositions. We demonstrate that UNIVERSE matches the performance of task-specific checkpoints while using significantly less training data and parameters. The results demonstrate that VLMs can serve as lightweight, semantics-aware evaluators of world models, and highlight promising directions for extending such evaluators to more complex environments.

## 1 INTRODUCTION

World models – generative models trained to predict future observations conditioned on past observations and actions (Ha & Schmidhuber, 2018; Hafner et al., 2025; Alonso et al., 2024) – are rapidly becoming central to interactive AI. They provide a powerful abstraction for learning, reasoning, and planning in complex interactive environments, and underpin advances in neural game engines (Kanervisto et al., 2025; Guo et al., 2025; Gao et al., 2025; Chen et al., 2025), embodied AI (Du et al.; Yang et al., 2024), and autonomous driving (Russell et al., 2025; Hu et al., 2023a; Ni et al., 2025).

Yet, evaluating world models remains a bottleneck. Rollouts are semantically rich and temporally grounded, requiring metrics that assess (i) alignment between generated frames and action sequences at the timestamp level (Yang et al., 2024), and (ii) consistent entity tracking over time (Kanervisto et al., 2025). Existing approaches fall short: (i) early distributional metrics focus on images and are sensitive to low-level variations (Salimans et al., 2016; Heusel et al., 2017; Binkowski et al., 2018), (ii) motion-aware metrics like FVD (Unterthiner et al., 2018) lack semantic grounding, and (iii) multimodal metrics ignore timestamp-level action conditioning (Jayasumana et al., 2024). Emerging text-to-video benchmarks (Liu et al., 2024b; Huang et al., 2024; Liao et al., 2024) focus on open-ended generation but neglect the fine-grained control central to world model evaluation. Even cutting-edge LLMs fail in this setting (Appendix G.1, Figure 15). Human evaluation remains the gold standard (Agarwal et al., 2025; Analysis, 2024), however, it remains costly and hard to scale.

To address this gap, we propose a novel evaluation protocol targeting two important dimensions of rollout quality: action alignment and character consistency, formalized as recognition tasks – Action Recognition (AR), Character Recognition (CR) – across formats of varying complexity. The protocol

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  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

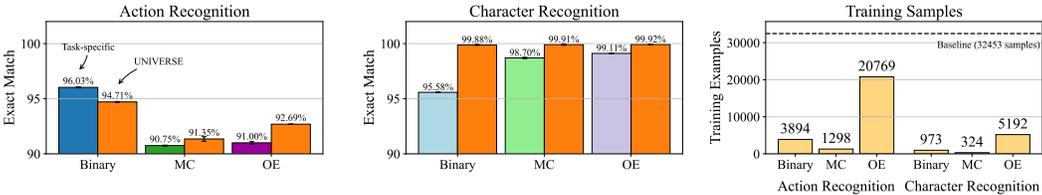


Figure 1: Performance and efficiency of UNIVERSE (orange bars throughout) compared to task-specific baselines (multiple colours), all models trained for 10 epochs. **Left and Center:** Action recognition and Character Recognition accuracy across binary, multiple-choice, and open-ended settings. **Right:** Sample efficiency – our adaptation recipe achieves strong performance with substantially fewer training samples per epoch.

provides a foundation for semantic evaluation of rollouts, and extends upon existing evaluation method by providing insight into semantic quality of generated rollouts. Inspired by the LLM-as-a-Judge research direction Zheng et al. (2023), we explore whether Vision-Language Models (VLMs) can serve as effective evaluators in this setting. VLMs have demonstrated strong generalization across multimodal tasks (Li et al., 2023; Driess et al., 2023; Chen et al., 2023; Wang et al., 2024; Abdin et al., 2024; Liu et al., 2024a; Deitke et al., 2024; McKinzie et al., 2024), and show promise as automatic judges of generative models (Lee et al., 2024; Mañas et al., 2024; Lin et al., 2024; Chen et al., 2024). Yet, off-the-shelf VLMs struggle with the temporal grounding and domain-specific knowledge (see Sec. 4, Zero-Shot Evaluation). Limited resources and sparse text supervision introduce another layer of complexity to the setting.

We therefore conduct a systematic study of adaptation strategies under realistic data and compute constraints. Throughout the study totaling over 5,154 GPU-days, we analyze the impact of supervision regime, frame sampling strategy, visual context length, and training budget. The result of the study is UNIVERSE (*UNified Vision-language Evaluator for Rollouts in Simulated Environments*), a lightweight, adaptable, and semantics-aware evaluator for world models. UNIVERSE achieves parity with six task-specific checkpoints while using a single unified model (see Figure 1. To validate reliability, we conduct a large-scale human study on rollouts spanning diverse environments, model scales, and rollout fidelities. UNIVERSE’s judgments show strong agreement with human ratings, demonstrating its effectiveness as a practical, semantics-aware evaluator—particularly valuable in settings where ground-truth annotations are unavailable or prohibitively expensive. To support reproducibility, we release our code, evaluation dataset, and human annotation dataset<sup>1</sup>.

## 2 RELATED WORK

**Challenges in Evaluating World Models.** World models are generative systems that learn predictive representations of environment dynamics (Ha & Schmidhuber, 2018), originally proposed for model-based RL (Sutton, 1991) and now central to domains such as neural game engines (Kanervisto et al., 2025; Guo et al., 2025; Gao et al., 2025; Chen et al., 2025), embodied AI (Du et al.; Yang et al., 2024), and autonomous driving (Russell et al., 2025; Hu et al., 2023a; Ni et al., 2025). Recent models such as Dreamer v1–3 (Hafner et al., 2020; Hafner et al.; 2023; 2025), MuZero (Schrittwieser et al., 2020), IRIS (Micheli et al., 2023), UniSim (Yang et al., 2024), and DIAMOND (Alonso et al., 2024) have improved rollout fidelity and controllability. Yet evaluation largely focuses on downstream success metrics—e.g., game score or goal completion (Bellemare et al., 2013; Kaiser et al., 2020; Guss et al., 2021; Baker et al., 2022; Beattie et al., 2016)—which provide only coarse, indirect signals of rollout quality. Genie (Bruce et al., 2024; Parker-Holder et al., 2024) decouples world model learning and agent training, but its evaluation still emphasizes visual quality and control, without probing semantic or causal fidelity. Cosmos (Agarwal et al., 2025) proposes a structured protocol that combines FID/FVD with structure-from-motion-based 3D consistency checks and human ratings on instruction following, object permanence, and visual verity. While insightful, this approach is tied to simulator-specific infrastructure and requires costly manual comparison. Human-in-the-loop

<sup>1</sup><https://anonymous.4open.science/r/vlms-for-wms-2651/README.md>

108 protocols such as the Video Generation Arena (Analysis, 2024) also rely on pairwise comparison to  
109 assess rollout quality. These methods, though informative, are expensive and hard to scale.

110  
111 **Evaluation Metrics and Protocols for Visual Generation.** Early evaluations of generative models  
112 relied on full-reference metrics such as PSNR and SSIM (Wang et al., 2004), which capture pixel-level  
113 and perceptual similarity but are sensitive to spatial misalignments and fail to reflect semantic fidelity.  
114 To address this, distributional metrics like Inception Score (IS) (Salimans et al., 2016), Fréchet  
115 Inception Distance (FID) (Heusel et al., 2017), and Kernel Inception Distance (KID) (Binkowski  
116 et al., 2018). Other proposals such as PPL (Karras et al., 2019), Parzen likelihoods (Goodfellow et al.,  
117 2014), and HYPE (Zhou et al., 2019) attempt to quantify perceptual smoothness or human realism,  
118 but remain focused on static images. For video generation, FVD (Unterthiner et al., 2018) generalizes  
119 FID using I3D features (Carreira & Zisserman, 2017), introducing a motion-aware distributional  
120 baseline. Yet, FVD also lacks semantic grounding and does not account for causal structure or  
121 goal alignment. To improve semantic grounding, metrics based on text-image alignment have been  
122 proposed. CLIPScore (Hessel et al., 2021) and CLIPSIM (Wu et al., 2021) compute similarity  
123 between generated visuals and textual or visual references using CLIP embeddings (Radford et al.,  
124 2021), while Jayasumana et al. (2024) extend this to distributional comparisons via MMD. However,  
125 all operate at the frame level. Structured evaluation protocols using vision-language reasoning have  
126 also emerged. VQA Accuracy (Mañas et al., 2024) uses LLMs to score answers on static image  
127 questions, and VQAScore (Lin et al., 2024) probes alignment via templated binary queries. Lee  
128 et al. (2024) propose VLM evaluator to evaluate other VLMs responses given user criteria. These  
129 approaches introduce task structure but remain limited to single-frame evaluation. Recent text-to-  
130 video (T2V) benchmarks such as EvalCrafter (Liu et al., 2024b), VBench (Huang et al., 2024), and  
131 DEVIL (Liao et al., 2024) introduce curated prompts and metrics covering text alignment, motion  
132 realism, and perceptual quality. While these protocols push forward evaluation of open-ended video  
133 generation, they lack timestamp-level action grounding.

134  
135 **Vision-Language Model Adaptation.** VLMs have emerged as powerful tools for multimodal  
136 understanding, demonstrating strong performance across tasks such as captioning, retrieval, visual  
137 question answering, and instruction following (Hendriksen et al., 2022; Li et al., 2023; Driess et al.,  
138 2023; Chen et al., 2023; Wang et al., 2024; Abdin et al., 2024; Liu et al., 2024a; Deitke et al.,  
139 2024; McKinzie et al., 2024). Adaptation approaches can be broadly categorized into prompt-level  
140 and weight-level methods. One prominent prompt-level adaptation techniques is prompt tuning,  
141 which injects task information directly into the input space (Miyai et al., 2023; Zhou et al., 2024;  
142 Wu et al., 2024a), and in-context learning (ICL), where models such as GPT-3 (Brown et al.,  
143 2020) and Flamingo (Alayrac et al., 2022) condition on task demonstrations at inference time  
144 without updating parameters. Retrieval-augmented generation (RAG) (Lewis et al., 2020) combines  
145 parametric models with non-parametric memory, and multimodal variants incorporate external visual  
146 or auditory context (Hu et al., 2023b; Chen et al., 2022a). While lightweight, these approaches are  
147 limited in their ability to model temporal dependencies or align with structured rollouts. Weight-level  
148 adaptation enables stronger domain alignment but incurs higher computational cost. Full finetuning  
149 remains effective yet costly, while partial finetuning (Ye et al., 2023) offers a trade-off by updating  
150 only selected layers. Parameter-efficient finetuning (PEFT) provides a scalable alternative and can be  
151 grouped into low-rank and adapter-based strategies (Han et al., 2024). Low-rank methods, such as  
152 LoRA (Hu et al., 2022), inject rank-constrained updates into frozen layers. Recent extensions improve  
153 upon this via weight decomposition (Liu et al.), quantization-aware adaptation (Dettmers et al., 2023;  
154 Xu et al., 2024), mixture-of-experts routing (Wu et al., 2024b), and long-context support (Chen  
155 et al.). Adapter-based methods insert lightweight modules between frozen layers to enable modular  
156 adaptation with minimal overhead (Luo et al., 2023; Zhao et al., 2024). A parallel line of work  
157 investigates multimodal few-shot learning. Frozen (Tsimpoukelli et al., 2021) was among the first  
158 to explore this setting, followed by works combining prompting and ICL for improved sample  
159 efficiency (Jin et al., 2022; Song et al., 2022), and works introducing a learnable meta-mapper to  
160 bridge frozen VLM components for few-shot meta-learning (Najdenkoska et al., 2023).

161  
**Our Focus.** While prior efforts have explored related challenges, none directly address the evaluation  
of the structured, action-conditioned fidelity and semantics of world model rollouts using adapted  
VLM. To this end, we introduce: (i) an evaluation protocol for world model rollouts, targeting  
fine-grained, temporally grounded assessment of semantic fidelity; (ii) UNIVERSE, a VLM-based  
method to support the protocol. We validate its alignment with human judgments and demonstrate its  
scalability and semantic sensitivity across rollout conditions.

### 3 METHODOLOGY

We consider the problem of evaluating rollouts generated by *world models* in interactive environments. A world model  $W$  is trained to predict the next observation  $o_t$  given the past observations  $o_{<t}$  and actions  $a_{<t}$ :  $W : (o_{<t}, a_{<t}) \rightarrow o_t$ , where  $o_t \in \mathcal{O}$  represents the sensory observation at timestep  $t$ , typically an RGB image. Rollouts consist of temporally grounded sequences that reflect the causal effects of control inputs. These outputs are semantically rich and visually complex, requiring timestamp-level assessment of correctness.

To enable automatic evaluation, we propose UNIVERSE, an adapted Vision-Language Model (VLM) that serves as a structured evaluator for world model rollouts. Formally, it operates as a function:  $E : (V, Q) \rightarrow \hat{A}$ , where  $V = (o_{t_1}, \dots, o_{t_k}) \in \mathcal{O}^k$  is a sequence of frames from a rollout,  $Q \in \mathcal{L}$  is a natural language question, and  $\hat{A} \in \mathcal{L}$  is the predicted answer. Evaluation quality is measured by comparing  $\hat{A}$  to the reference answer  $A$  using semantic similarity metrics.

**Evaluation Protocol.** We define two structured recognition tasks: (i) *Action Recognition (AR)*: Assesses whether generated sequences accurately reflect the effects of agent actions given the segment; (ii) *Character Recognition (CR)*: Evaluates whether entities maintain consistent identity and appearance across time. Each task is framed as a visual QA problem: the evaluator receives a sequence of frames and a natural language prompt (binary, multiple-choice, or open-ended), and generates a textual response. Outputs are scored using Exact Match (EM) and ROUGE-F<sub>1</sub> (ROUGE), capturing both literal and semantic alignment with the reference answer. Metric details are in Appendix E.2.

**Dataset Construction.** Effective VLM adaptation for rollout evaluation requires a dataset that (i) captures realistic human behavior in interactive environments, and (ii) aligns with prior work in simulated settings to support comparability and reproducibility. To satisfy these constraints, we partnered with Ninja Theory and curated a dataset from both internal and public *Bleeding Edge* gameplay recordings, focusing on the *Skygarden* environment (Kanervisto et al., 2025). This dataset provides high visual and behavioral diversity (Pearce et al., 2025), includes a publicly available evaluation split, and is closely aligned with prior work in the domain (Kanervisto et al., 2025; Pearce et al., 2025; Tot et al., 2025; Sharma et al., 2024; Devlin et al., 2021), enabling fair comparison.

Data preparation proceeds in three stages: (i) *Preprocessing*: Segment gameplay into 14-frame clips with synchronized video, control logs, and metadata; (ii) *Description Generation*: Convert structured annotations (e.g., actions, agent states) into natural language summaries; (iii) *Question-Answer Pair Construction*: Generate six complimentary QA pairs per clip (binary, multiple-choice, and open-ended) spanning both AR and CR tasks. The final dataset contains 32.453 training clips and 8.113 validation clips, yielding 194.718 and 48.678 QA pairs, respectively. See Appendix D for details.

**Model Architecture.** We adapt a model from the PaliGemma family (Beyer et al., 2024; Steiner et al., 2024), consisting of a vision encoder  $\mathcal{M}_V$ , a projection head  $\mathcal{M}_P$ , and a language decoder  $\mathcal{M}_L$ . Based on initial zero-shot evaluations (Appendix G.2), we use a single configuration for all experiments—PaliGemma 2 3b, which includes a 2B-parameter Gemma 2 decoder pretrained on 2T tokens. Input frames are resized to  $224 \times 224$  and tokenized into 256 patches each. Model architecture details are in Appendix E.1.

Each model input sequence  $S = \{S_{\mathcal{I}}, S_{\mathcal{T}}^{\text{PREF}}, S_{\mathcal{T}}^{\text{SUFF}}\}$  consists of: visual tokens  $S_{\mathcal{I}}$  from  $k$  frames, a textual prefix  $S_{\mathcal{T}}^{\text{PREF}}$  containing the task-language cue and question, and a suffix  $S_{\mathcal{T}}^{\text{SUFF}}$  with the expected answer (used only during training). This format allows the decoder to attend jointly over visual and textual context. Full prompt details are provided in Appendix E.1.

**Training Objective.** We optimize a causal language modeling loss on the answer suffix:

$$\mathcal{L}(S) = - \sum_{t=1}^{T_{\text{SUFF}}} \log P(s_t^{\text{SUFF}} | S_{<t'}) \quad (1)$$

where  $s_t^{\text{SUFF}}$  is the  $t$ -th token in the suffix, and  $t' = T_{\mathcal{I}} + T_{\text{PREF}} + t$  is the token position in the flattened sequence.

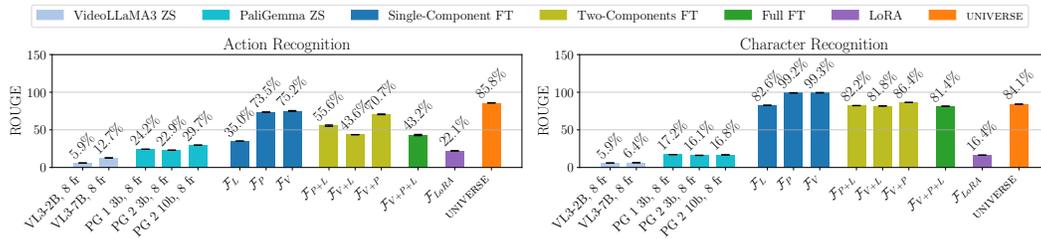


Figure 2: Comparison of UNIVERSE and baseline models on Action and Character Recognition, all models trained for 1 epoch. **Left:** UNIVERSE outperforms all baselines on AR. **Right:** On CR, it ranks third, behind models with either full vision encoder tuning or task-specific training with greater supervision. Trained under a unified protocol with minimal parameter updates (0.07%) and reduced per-task data, UNIVERSE delivers strong performance across both tasks.

**Adaptation Strategies.** We explore a broad design space for adapting pretrained VLMs to temporally grounded rollout evaluation. Our study spans three core dimensions: *fine-tuning configurations*, *frame sampling policy*, and *supervision composition*.

*Fine-Tuning Configurations.* We compare five adaptation strategies varying in parameter count and modularity: (i) *Zero-shot prompting*: No tuning; model is prompted directly. (ii) *Full fine-tuning*: All parameters  $\theta = \theta_V \cup \theta_P \cup \theta_L$  are updated end-to-end. (iii) *Dual-component fine-tuning*: Two of three modules are trained (e.g.,  $\theta_P \cup \theta_L$ ). (iv) *Single-component fine-tuning*: Only one module—vision, projection, or language—is updated. (v) *Parameter-efficient fine-tuning*: We apply LoRA (Hu et al., 2022) adapters to attention and MLP layers in vision and language components:  $\mathbf{W} \leftarrow \mathbf{W} + \frac{\alpha}{r} \mathbf{A}\mathbf{B}$ ,  $\alpha = 8$ ,  $r \in \{8, 16, 32, 48, 64\}$ .

*Frame Sampling Policy.* We vary both the number of input frames and their sampling strategy. Specifically, we sweep over  $k \in [1, 8]$ , and evaluate two selection methods: (i) *First- $n$* : selecting the first  $k$  frames from each rollout; (ii) *Uniform- $n$* : sampling  $k$  frames uniformly across the full clip.

*Supervision Composition.* To support generalization across QA formats and tasks, we construct a multi-task dataset covering binary, multiple-choice, and open-ended prompts across both AR and CR. We perform a three-stage grid search to optimize the data mixture: (i) Varying AR/CR task ratios ( $\alpha_{AR}, \alpha_{CR}$ ) while fixing QA type proportions ( $\beta_{Binary}, \beta_{MC}, \beta_{OE}$ ); (ii) Tuning the proportion of open-ended supervision ( $\beta_{OE}$ ) for best performance; (iii) Adjusting  $\beta_{Binary}$  and  $\beta_{MC}$ .

**UNIVERSE: UNIFIED Vision-language Evaluator for Rollouts in Simulated Environments.** We distill our empirical findings into UNIVERSE, a lightweight and scalable adaptation method for temporally grounded evaluation of world model rollouts using VLMs. Designed for constrained compute and limited supervision, UNIVERSE delivers strong generalization across our evaluation protocol using a single, partially tuned model. The method combines three main components: I *Partial fine-tuning*: We update only the projection head ( $\theta_P$ ), training just 0.07% of model parameters. Despite this minimal footprint, it achieves the second-best performance among all strategies—trailing only vision encoder tuning, which requires  $\sim 11\%$  of parameters and incurs significantly higher compute cost. II *Efficient frame sampling*: Each input sequence includes  $k = 8$  frames sampled uniformly from a 14-frame rollout. This sparsity-aware strategy maintains long-range temporal structure while reducing token count and enabling efficient batching. III *Mixed supervision*: We train on a hierarchical mixture of tasks and QA formats. The task distribution favors Action Recognition ( $\alpha_{AR} = 0.8$ ) as it converges slower. Within each task, we emphasize open-ended questions ( $\beta_{OE} = 0.8$ ), while maintaining smaller proportions of binary ( $\beta_{binary} = 0.15$ ) and multiple-choice ( $\beta_{MC} = 0.05$ ) examples.

#### 4 EXPERIMENTS

We evaluate UNIVERSE on ground truth video data, focusing on AR and CR across binary, multiple-choice, and open-ended formats. Our goals are twofold: (i) to benchmark performance against zero-shot and fine-tuned baselines, and (ii) to assess the trade-offs between adaptation strategies under constrained supervision and compute.

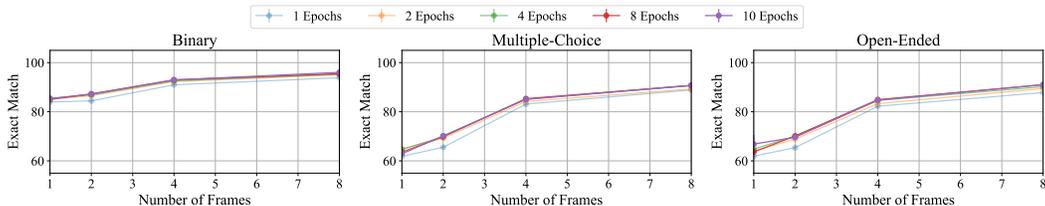


Figure 3: Action Recognition performance as a function of training supervision (epochs) and temporal context (number of frames), evaluated across all formats. Performance improves along both axes, with highest accuracy achieved when both dimensions are scaled.

**Baselines.** We compare UNIVERSE against two classes of baselines: (i) *Zero-shot VLMs*: Seven off-the-shelf models, including VideoLLaMA3 (2B, 7B) (Boqiang Zhang, 2025) and three PaliGemma models: version 1 (3B) and version 2 (3B and 10B) (Beyer et al., 2024; Steiner et al., 2024), evaluated without domain adaptation using an 8-frame visual context window.<sup>2</sup> (ii) *Fine-tuned PaliGemma 2*: Variants adapted via full, partial, and parameter-efficient tuning. This backbone is selected based on a sweep over PaliGemma variants, using zero-shot performance as a guide (Appendix G.2). The adaptation space includes 8 primary baselines: (i) *Single-component fine-tuning*: tuning only the vision encoder ( $\mathcal{F}_V$ ), the multimodal projector ( $\mathcal{F}_P$ ), or the language head ( $\mathcal{F}_L$ ); (ii) *Two-component fine-tuning*: jointly tuning pairs of components— $\mathcal{F}_{V+P}$ ,  $\mathcal{F}_{V+L}$ , and  $\mathcal{F}_{P+L}$ ; (iii) *Full-model fine-tuning*: tuning all components simultaneously ( $\mathcal{F}_{V+P+L}$ ); (iv) *LoRA-based tuning*: Parameter-efficient adaptation with rank  $r = 8$ , selected after observing minimal performance variation across  $r \in \{8, 16, 32, 48, 64\}$  (see Appendix G.4 for details). All models are trained using 8-frame clips and a single epoch.

**Results.** Figure 2 (left, center) summarizes performance across Action Recognition (AR) and Character Recognition (CR). Zero-shot VLMs perform poorly: VideoLLaMA3 variants stay below 12.7% (AR) and 6.4% (CR), while PaliGemma reaches 29.7% (AR) and 17.2% (CR), confirming that general-purpose models lack the temporal grounding and domain-specific semantics needed for rollout evaluation. In contrast, UNIVERSE outperforms all models on AR and ranks third on CR—despite tuning only the 2.66M-parameter projector (0.07% of the model) under a unified protocol spanning both tasks, all prompt formats, and reduced supervision. The two stronger CR baselines fine-tune either the full 400M-parameter vision encoder or use  $5\times$  more task-specific CR data. Its performance under these constraints underscores the efficiency and generality of our adaptation strategy for temporally grounded evaluation.

## 5 ANALYSIS

We find a consistent performance gap between AR and CR, highlighting the greater temporal and causal complexity of action understanding. This motivates our focus on AR as the more challenging and diagnostic task. Below, we analyze how adaptation choices shape performance on AR.

**Supervision and Temporal Context.** We begin by analyzing how supervision (training budget) and temporal input (number of frames) influence UNIVERSE performance. By independently and jointly varying the number of training epochs and input frames, we disentangle the contributions of model capacity and temporal context to task success.

**Results.** CR converges rapidly, achieving over 97% exact match after 12.5% of an epoch ( $\sim 4K$  samples; Figure 5, bottom), and shows minor improvement with further training. In contrast, AR improves only modestly under extended training when limited to a single frame (Figure 5, top), suggesting that supervision alone is insufficient in the absence of temporal information. Motivated by this, we jointly scale both supervision and input length, varying the number of frames and epochs. As shown in Figure 3, performance on AR improves consistently across all formats, with the best results achieved under combined scaling.

<sup>2</sup>We also experimented with CLIPScore-based evaluation (Appendix G.3); results underperformed relative to selected baselines and were constrained to predefined candidate sets, further underscoring the need for model adaptation.

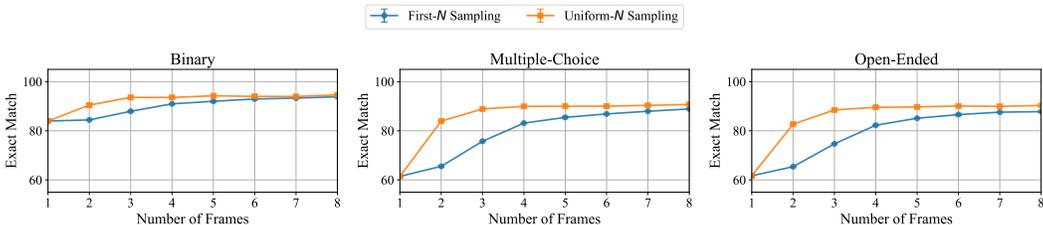


Figure 4: Effect of frame sampling strategy on Action Recognition performance across all formats. Uniform- $n$  sampling (orange) consistently outperforms first- $n$  (blue), with especially large gains at low frame counts, and maintains an advantage as temporal context increases.

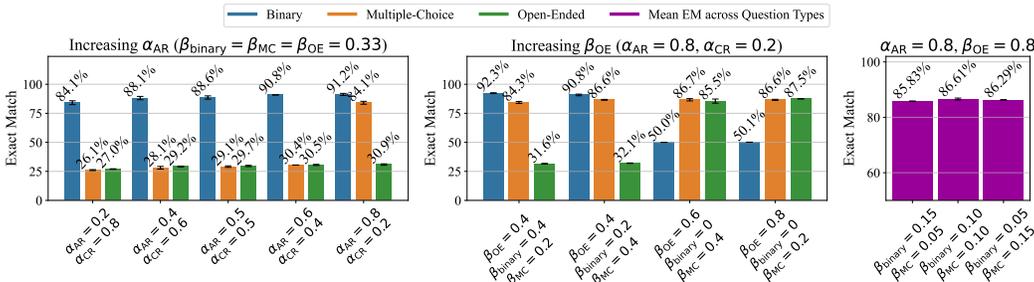


Figure 6: Hierarchical ablation of training data composition for UNIVERSE. **Left:** Varying task-level ratio  $\alpha$  (AR vs. CR) with uniform format distribution ( $\beta = 1/3$ ) shows that increasing  $\alpha_{AR}$  improves AR performance, especially for multiple-choice, while open-ended remains flat. **Center:** Sweeping format-level ratio  $\beta_{OE}$  with fixed  $\alpha_{AR} = 0.8$  reveals that oversampling open-ended data ( $\beta_{OE} = 0.8$ ) improves AR-OE performance. **Right:** Fine-tuning binary and MC proportions under  $\beta_{OE} = 0.8$  shows performance is stable across mixes, with slight gains from  $\beta_{binary} = 0.15, \beta_{MC} = 0.05$ .

**Temporal Sampling Strategies.** Following the observation that AR requires both extended supervision and temporally rich input, we examine how frame selection impacts performance. We compare first- $n$  sampling, which selects the first  $n$  consecutive frames from each rollout, to uniform- $n$  sampling, which draws  $n$  evenly spaced frames across the entire sequence. We conduct experiments at varying context lengths, using  $n \in \{1, 2, \dots, 8\}$  frames. to evaluate the impact of both sampling method and input horizon.

**Results.** As shown in Figure 4, uniform- $n$  consistently outperforms first- $n$  across all evaluation formats. The effect is most pronounced at low frame counts. With only 2 input frames, uniform sampling improves exact match accuracy from 84.42% to 90.47% in Binary, from 65.53% to 83.93% in Multiple-Choice, and from 65.38% to 82.68% in Open-Ended formats. Gains persist even at 8 frames, where uniform sampling maintains an advantage across formats.

**Optimizing Data Mix for Unified Multi-Task Evaluation.** We analyze how training data composition affects multi-task performance in UNIVERSE, with the goal of enabling a single model to generalize across AR and CR. Specifically, we study how the task-level ratio  $\alpha$  and format-level ratio  $\beta$  influence performance across evaluation settings. We first conduct a hierarchical ablation to identify an optimized data mixture, then assess its impact by comparing against a default task mix with uniform sampling.

**Data Mix Optimization.** To determine an effective training mixture for UNIVERSE, we perform a hierarchical ablation over task-level and format-level data ratios. We begin by varying the task-level proportion  $\alpha$  (AR vs. CR), holding the format distribution fixed at

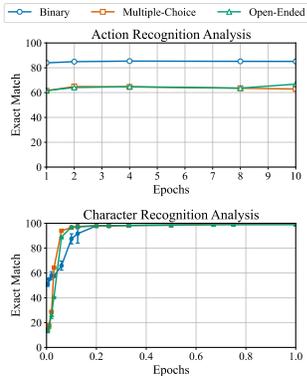


Figure 5: Exact Match accuracy for Action Recognition and Character Recognition.

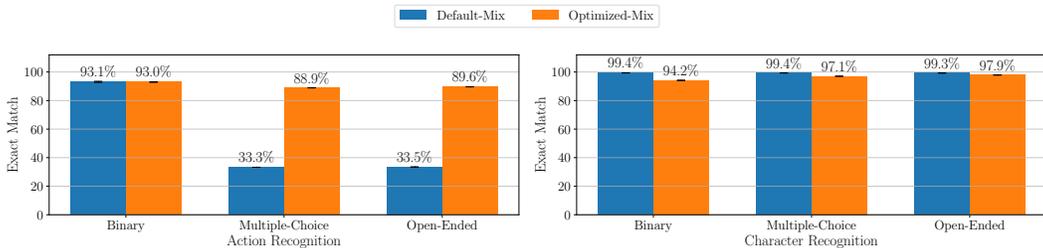


Figure 7: Comparison of two training regimes: a default data mix (equal task and format proportions) and the optimized mix derived from hierarchical tuning. The optimized configuration yields substantial gains on AR, while maintaining strong CR performance.

$\beta_{\text{binary}} = \beta_{\text{MC}} = \beta_{\text{OE}} = 1/3$ . As shown in Figure 6 (left), increasing  $\alpha_{\text{AR}}$  improves AR performance—especially for multiple-choice—while CR remains stable, with a favorable tradeoff reached at  $\alpha_{\text{AR}} = 0.8$ . However, open-ended accuracy shows little change, motivating format-specific rebalancing. Fixing  $\alpha_{\text{AR}} = 0.8$ , we sweep the format ratio  $\beta_{\text{OE}}$ , and observe in Figure 6 (center) that AR-OE accuracy improves substantially with increased open-ended coverage, peaking at  $\beta_{\text{OE}} = 0.8$ , albeit at the cost of binary performance. To restore balance, we fix  $\beta_{\text{OE}} = 0.8$  and allocate the remaining budget across binary and multiple-choice formats. As shown in Figure 6 (right), performance remains robust across configurations, with a slight preference for  $\beta_{\text{binary}} = 0.15$  and  $\beta_{\text{MC}} = 0.05$ . Based on these findings, we adopt the following optimized data composition:  $\alpha_{\text{AR}} = 0.8$ ,  $\alpha_{\text{CR}} = 0.2$ ;  $\beta_{\text{binary}} = 0.15$ ,  $\beta_{\text{MC}} = 0.05$ , and  $\beta_{\text{OE}} = 0.8$ .

*Effectiveness of the Optimized Mix.* Having identified an optimized training mixture through hierarchical ablation, we now evaluate its impact in practice. We compare the final UNIVERSE model—trained with this optimized mix—to a baseline trained with a default task and format distribution. We train both models on 4 epochs. As shown in Figure 7, the optimized configuration yields substantial gains on AR, particularly for multiple-choice and open-ended formats, while maintaining competitive performance on CR. These results underscore the importance of data composition in enabling robust multi-task learning within a unified evaluator.

## 6 EVALUATING WORLD MODEL ROLLOUTS WITH UNIVERSE

We evaluate the reliability of UNIVERSE through a human study spanning eight distinct settings that vary in model scale, training data diversity, and output resolution. Our analysis considers two axes: (i) *in-domain accuracy*, measured on Skygarden, and (ii) *generalization* to six previously unseen environments.

Concretely, we study rollouts generated by: (i) a large-scale model trained across multiple environments with higher-resolution rollouts ( $300 \times 180$ ), and (ii) a smaller model trained on a single environment with lower-resolution rollouts ( $128 \times 128$ ). These two model families expose complementary challenges: resolution mismatch, domain coverage, and rollout fidelity. We construct eight evaluation settings: settings 1–7 draw from the large-scale model across diverse environments, while setting 8 uses the smaller model in the fine-tuning environment. Each setting contains 30 rollouts, paired with six natural-language questions from our evaluation protocol, yielding 240 rollouts in total. UNIVERSE answers each question via majority vote over five greedy decoding samples. Human annotators rate responses on a four-point ordinal scale (*Correct*, *Partially Correct*, *Incorrect*, *Unclear*), with double annotation and adjudication on disagreements. Inter-rater reliability is measured using Cohen’s  $\kappa$ . Full details of the annotation protocol are provided in Appendix F.

*Results.* Figure 9 reports graded accuracy across all settings. Rollouts from the smaller, single-environment model yield lower evaluation accuracy, likely due to resolution mismatch, while the larger, multi-environment model provides higher-quality inputs. These results demonstrate that UNIVERSE generalizes across model scales, rollout fidelities, and environments, while remaining closely aligned with human judgments.

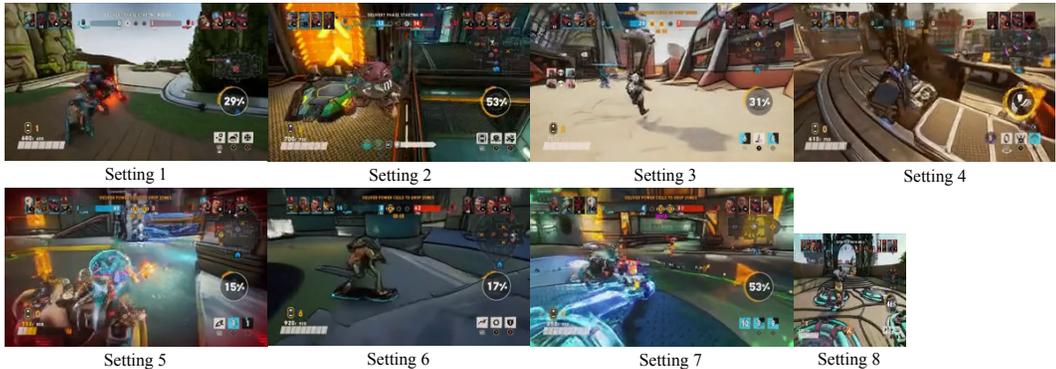


Figure 8: Example frames from the eight evaluation settings, spanning different model scales, rollout fidelities, and environments. Note the resolution difference: Settings 1-7 feature  $300 \times 180$  rollouts, whereas Setting 8 uses  $128 \times 128$ .

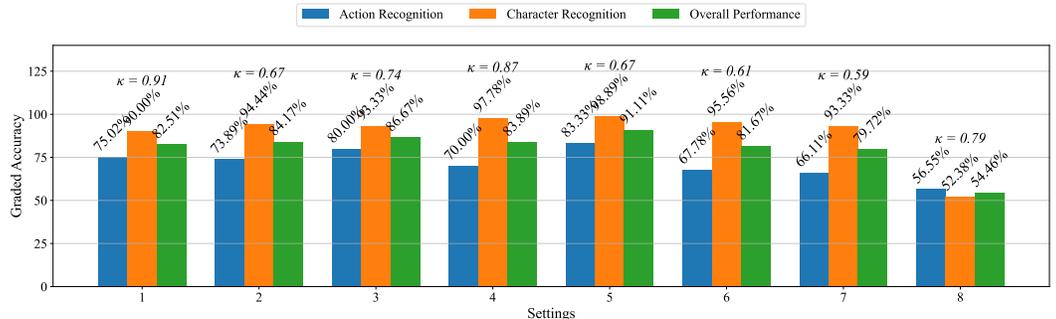


Figure 9: Graded accuracy of UNIVERSE across the eight evaluation settings. Performance improves with higher-fidelity rollouts and remains stable across unseen environments (2–7). Cohen’s  $\kappa$  indicates substantial inter-rater agreement.

## 7 CONCLUSION

World model evaluation remains a fundamental challenge, requiring fine-grained assessment of semantic consistency and action alignment – capabilities poorly addressed by existing metrics. In this work, we investigate whether lightweight VLM adaptation can provide reliable semantic evaluation of world model rollouts through a comprehensive case study. We introduce a structured evaluation protocol centered on action and character recognition tasks across binary, multiple-choice, and open-ended formats. To support this, we propose UNIVERSE, a unified method for adapting VLMs to this setting through mixed supervision, efficient frame sampling, and lightweight fine-tuning. Our large-scale study demonstrates that UNIVERSE matches the performance of task-specific checkpoints using a single unified model and aligns closely with human judgments, establishing it as a lightweight, semantics-aware evaluator for evaluating world models.

**Limitations.** While our experiments focus on simulated environments, chosen both for their ground-truth availability and their direct relevance to large video-game and interactive-entertainment industries (Kanervisto et al., 2025), evaluating UNIVERSE in real-world settings remains an important next step. Our protocol focuses on foundational semantic tasks, and extending it to cover higher-level reasoning represents an important direction for future work. Scaling UNIVERSE to long-horizon rollouts is also challenging, as fine-grained reasoning becomes harder with larger visual context. Our frame-sampling analysis suggests that more intelligent sampling or hierarchical summarization could address this, and we plan to explore such strategies. Finally, as with all pretrained VLMs, UNIVERSE may inherit biases and exhibit reduced reliability on ambiguous cases (Bleeker et al., 2024). While we examined them partially during evaluation, a deeper investigation into bias propagation is an exciting research direction for future work.

## REFERENCES

- 486  
487  
488 Marah I. Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany  
489 Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha  
490 Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu  
491 Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon,  
492 Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider,  
493 Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos  
494 Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee,  
495 Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik  
496 Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid  
497 Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli  
498 Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma,  
499 Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael  
500 Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong  
501 Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and  
502 Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- 503 Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chat-  
504 topadhyay, Yongxin Chen, Yin Cui, Yifan Ding, Daniel Dworakowski, Jiaojiao Fan, Michele  
505 Fenuz, Francesco Ferroni, Sanja Fidler, Dieter Fox, Songwei Ge, Yunhao Ge, Jinwei Gu, Siddharth  
506 Gururani, Ethan He, Jiahui Huang, Jacob Samuel Huffman, Pooya Jannaty, Jingyi Jin, Seung Wook  
507 Kim, Gergely Klár, Grace Lam, Shiyi Lan, Laura Leal-Taixé, Anqi Li, Zhaoshuo Li, Chen-Hsuan  
508 Lin, Tsung-Yi Lin, Huan Ling, Ming-Yu Liu, Xian Liu, Alice Luo, Qianli Ma, Hanzi Mao, Kaichun  
509 Mo, Arsalan Mousavian, Seungjun Nah, Sriharsha Niverty, David Page, Despoina Paschalidou,  
510 Zeeshan Patel, Lindsey Pavao, Morteza Ramezani, Fitsum Reda, Xiaowei Ren, Vasanth Rao Naik  
511 Sabavat, Ed Schmerling, Stella Shi, Bartosz Stefaniak, Shitao Tang, Lyne Tchapmi, Przemek  
512 Tredak, Wei-Cheng Tseng, Jibin Varghese, Hao Wang, Haoxiang Wang, Heng Wang, Ting-Chun  
513 Wang, Fangyin Wei, Xinyue Wei, Jay Zhangjie Wu, Jiashu Xu, Wei Yang, Lin Yen-Chen, Xiaohui  
514 Zeng, Yu Zeng, Jing Zhang, Qinsheng Zhang, Yuxuan Zhang, Qingqing Zhao, and Artur Zólkowski.  
515 Cosmos world foundation model platform for physical AI. *arXiv preprint arXiv:2501.03575*, 2025.
- 516 Ibrahim M Alabdulmohsin, Xiaohua Zhai, Alexander Kolesnikov, and Lucas Beyer. Getting vit  
517 in shape: Scaling laws for compute-optimal model design. *Advances in Neural Information  
518 Processing Systems*, 36:16406–16425, 2023.
- 519 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel  
520 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford,  
521 Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick,  
522 Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski,  
523 Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a visual  
524 language model for few-shot learning. *Advances in neural information processing systems*, 35:  
525 23716–23736, 2022.
- 526 Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos J Storkey, Tim Pearce, and  
527 François Fleuret. Diffusion for world modeling: Visual details matter in atari. *Advances in Neural  
528 Information Processing Systems*, 37:58757–58791, 2024.
- 529 Artificial Analysis. Video generation arena leaderboard, 2024. URL <https://huggingface.co/spaces/ArtificialAnalysis/Video-Generation-Arena-Leaderboard>.  
530 Accessed: 24 March 2025.
- 531  
532  
533 Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon  
534 Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (VPT): Learning to act by watching  
535 unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654,  
536 2022.
- 537 Charles Beattie, Joel Z. Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler,  
538 Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, Julian Schrittwieser, Keith Anderson,  
539 Sarah York, Max Cant, Adam Cain, Adrian Bolton, Stephen Gaffney, Helen King, Demis Hassabis,  
Shane Legg, and Stig Petersen. Deepmind lab. *arXiv preprint arXiv:1612.03801*, 2016.

- 540 Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environ-  
541 ment: An evaluation platform for general agents. *Journal of artificial intelligence research*, 47:  
542 253–279, 2013.
- 543 Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz,  
544 Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas  
545 Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey A. Gritsenko,  
546 Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer,  
547 Matko Bosnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan  
548 Puigcerver, Pinelopi Papalampidi, Olivier J. Hénaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen,  
549 and Xiaohua Zhai. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*,  
550 2024.
- 551 Lukas Biewald. Experiment tracking with weights and biases, 2020. URL <https://www.wandb.com/>.  
552 Software available from wandb.com.
- 553 Mikolaj Binkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD  
554 gans. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC,*  
555 *Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- 556 Steven Bird and Edward Loper. NLTK: The natural language toolkit. In *Proceedings of the*  
557 *ACL Interactive Poster and Demonstration Sessions*, pp. 214–217, Barcelona, Spain, July 2004.  
558 Association for Computational Linguistics.
- 559 Maurits Bleeker, Mariya Hendriksen, Andrew Yates, and Maarten de Rijke. Demonstrating and  
560 reducing shortcuts in vision-language representation learning. *Transactions on Machine Learning*  
561 *Research*, 2024.
- 562 Zesen Cheng Zhiqiang Hu Yuqian Yuan Guanzheng Chen Sicong Leng Yuming Jiang Hang Zhang  
563 Xin Li Peng Jin Wenqi Zhang Fan Wang Lidong Bing Deli Zhao Boqiang Zhang, Kehan Li.  
564 Videollama 3: Frontier multimodal foundation models for image and video understanding. *arXiv*  
565 *preprint arXiv:2501.13106*, 2025. URL <https://arxiv.org/abs/2501.13106>.
- 566 G. Bradski. The OpenCV Library. *Dr. Dobb’s Journal of Software Tools*, 2000.
- 567 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
568 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
569 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 570 Jake Bruce, Michael D. Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,  
571 Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, Yusuf Aytar, Sarah Bechtle, Feryal  
572 M. P. Behbahani, Stephanie C. Y. Chan, Nicolas Heess, Lucy Gonzalez, Simon Osindero, Sherjil  
573 Ozair, Scott E. Reed, Jingwei Zhang, Konrad Zolna, Jeff Clune, Nando de Freitas, Satinder Singh,  
574 and Tim Rocktäschel. Genie: Generative interactive environments. In *Forty-first International*  
575 *Conference on Machine Learning*, 2024.
- 576 João Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics  
577 dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017,*  
578 *Honolulu, HI, USA, July 21-26, 2017*, pp. 4724–4733. IEEE Computer Society, 2017.
- 579 Soravit Changpinyo, Doron Kukliansky, Idan Szpektor, Xi Chen, Nan Ding, and Radu Soricut. All  
580 you may need for VQA are image captions. *arXiv preprint arXiv:2205.01883*, 2022.
- 581 Dongping Chen, Ruoxi Chen, Shilin Zhang, Yaochen Wang, Yinuo Liu, Huichi Zhou, Qihui Zhang,  
582 Yao Wan, Pan Zhou, and Lichao Sun. Mllm-as-a-judge: Assessing multimodal llm-as-a-judge with  
583 vision-language benchmark. In *Forty-first International Conference on Machine Learning*, 2024.
- 584 Jingye Chen, Yuzhong Zhao, Yupan Huang, Lei Cui, Li Dong, Tengchao Lv, Qifeng Chen, and Furu  
585 Wei. Model as a game: On numerical and spatial consistency for generative games. *arXiv preprint*  
586 *arXiv:2503.21172*, 2025.
- 587 Wenhu Chen, Hexiang Hu, Xi Chen, Pat Verga, and William Cohen. Murag: Multimodal retrieval-  
588 augmented generator for open question answering over images and text. In *Proceedings of the*  
589 *2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5558–5570, 2022a.

- 594 Xi Chen, Xiao Wang, Soravit Changpinyo, A. J. Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian  
595 Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver,  
596 Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish V. Thapliyal, James  
597 Bradbury, and Weicheng Kuo. Pali: A jointly-scaled multilingual language-image model. *arXiv  
598 preprint arXiv:2209.06794*, 2022b.
- 599  
600 Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Car-  
601 los Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, Siamak Shakeri, Mostafa Dehghani,  
602 Daniel Salz, Mario Lucic, Michael Tschannen, Arsha Nagrani, Hexiang Hu, Mandar Joshi, Bo Pang,  
603 Ceslee Montgomery, Paulina Pietrzyk, Marvin Ritter, A. J. Piergiovanni, Matthias Minderer, Filip  
604 Pavetic, Austin Waters, Gang Li, Ibrahim Alabdulmohsin, Lucas Beyer, Julien Amelot, Kenton  
605 Lee, Andreas Peter Steiner, Yang Li, Daniel Keysers, Anurag Arnab, Yuanzhong Xu, Keran Rong,  
606 Alexander Kolesnikov, Mojtaba Seyedhosseini, Anelia Angelova, Xiaohua Zhai, Neil Houlsby,  
607 and Radu Soricut. Pali-x: On scaling up a multilingual vision and language model. *arXiv preprint  
608 arXiv:2305.18565*, 2023.
- 609 Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. Lon-  
610 glora: Efficient fine-tuning of long-context large language models. In *The Twelfth International  
611 Conference on Learning Representations*.
- 612 Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Moham-  
613 madreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin  
614 Bransom, Kiana Ehsani, Huong Ngo, Yen-Sung Chen, Ajay Patel, Mark Yatskar, Chris Callison-  
615 Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne  
616 Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron  
617 Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt,  
618 Dirk Groeneveld, Jen Dumas, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick,  
619 Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross B. Gir-  
620 shick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for  
621 state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*, 2024.
- 622 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning  
623 of quantized llms. *Advances in neural information processing systems*, 36:10088–10115, 2023.
- 624  
625 Sam Devlin, Raluca Georgescu, Ida Momennejad, Jaroslaw Rzepecki, Evelyn Zuniga, Gavin Costello,  
626 Guy Leroy, Ali Shaw, and Katja Hofmann. Navigation turing test (ntt): Learning to evaluate  
627 human-like navigation. In *International Conference on Machine Learning*, pp. 2644–2653. PMLR,  
628 2021.
- 629  
630 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
631 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,  
632 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.  
633 In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria,  
634 May 3-7, 2021*, 2021.
- 635 Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,  
636 Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar,  
637 Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc  
638 Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. PaLM-E: An embodied  
639 multimodal language model. In *International Conference on Machine Learning*, pp. 8469–8488.  
640 PMLR, 2023.
- 641  
642 Yilun Du, Sherry Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet, Tianhe  
643 Yu, Pieter Abbeel, Joshua B. Tenenbaum, Leslie Pack Kaelbling, Andy Zeng, and Jonathan  
644 Tompson. Video language planning. In *The Twelfth International Conference on Learning  
645 Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.
- 646  
647 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image  
synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
pp. 12873–12883, 2021.

- 648 Shen yuan Gao, Siyuan Zhou, Yilun Du, Jun Zhang, and Chuang Gan. Adaworld: Learning adaptable  
649 world models with latent actions. *arXiv preprint arXiv:2503.18938*, 2025.  
650
- 651 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
652 Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In Zoubin Ghahramani,  
653 Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (eds.), *Advances in  
654 Neural Information Processing Systems 27: Annual Conference on Neural Information Processing  
655 Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pp. 2672–2680, 2014.
- 656 Junliang Guo, Yang Ye, Tianyu He, Haoyu Wu, Yushu Jiang, Tim Pearce, and Jiang Bian. Mineworld:  
657 a real-time and open-source interactive world model on minecraft. *arXiv preprint arXiv:2504.08388*,  
658 2025.  
659
- 660 William H Guss, Mario Ynocente Castro, Sam Devlin, Brandon Houghton, Noboru Sean Kuno,  
661 Crissman Loomis, Stephanie Milani, Sharada Mohanty, Keisuke Nakata, Ruslan Salakhutdinov,  
662 et al. The minerl 2020 competition on sample efficient reinforcement learning using human priors.  
663 *arXiv preprint arXiv:2101.11071*, 2021.
- 664 David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in  
665 Neural Information Processing Systems*, 31:2451–2463, 2018.  
666
- 667 Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with  
668 discrete world models. In *International Conference on Learning Representations*.  
669
- 670 Danijar Hafner, Timothy P. Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learn-  
671 ing behaviors by latent imagination. In *8th International Conference on Learning Representations,  
672 ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- 673 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains  
674 through world models. *arXiv preprint arXiv:2301.04104*, 2023.  
675
- 676 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse control tasks  
677 through world models. *Nature*, pp. 1–7, 2025.
- 678 Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning  
679 for large models: A comprehensive survey. *Transactions on Machine Learning Research*, 2024.  
680
- 681 Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David  
682 Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti  
683 Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández  
684 del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy,  
685 Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming  
686 with NumPy. *Nature*, 585(7825):357–362, September 2020.
- 687 Mariya Hendriksen, Maurits Bleeker, Svitlana Vakulenko, Nanne Van Noord, Ernst Kuiper, and  
688 Maarten De Rijke. Extending clip for category-to-image retrieval in e-commerce. In *European  
689 Conference on Information Retrieval*, pp. 289–303. Springer, 2022.  
690
- 691 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint  
692 arXiv:1606.08415*, 2016.  
693
- 694 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-  
695 free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical  
696 Methods in Natural Language Processing*, pp. 7514–7528, 2021.
- 697 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans  
698 trained by a two time-scale update rule converge to a local nash equilibrium. In Isabelle Guyon,  
699 Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and  
700 Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference  
701 on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp.  
6626–6637, 2017.

- 702 Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton,  
703 and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving. *arXiv preprint*  
704 *arXiv:2309.17080*, 2023a.
- 705 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
706 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth International*  
707 *Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*, 2022.
- 708 Ziniu Hu, Ahmet Iscen, Chen Sun, Zirui Wang, Kai-Wei Chang, Yizhou Sun, Cordelia Schmid,  
709 David A Ross, and Alireza Fathi. Reveal: Retrieval-augmented visual-language pre-training with  
710 multi-source multimodal knowledge memory. In *Proceedings of the IEEE/CVF conference on*  
711 *computer vision and pattern recognition*, pp. 23369–23379, 2023b.
- 712 Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing  
713 Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video  
714 generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
715 *Recognition*, pp. 21807–21818, 2024.
- 716 John D Hunter. Matplotlib: A 2d graphics environment. *Computing in science & engineering*, 9(03):  
717 90–95, 2007.
- 718 Sadeep Jayasumana, Srikumar Ramalingam, Andreas Veit, Daniel Glasner, Ayan Chakrabarti, and  
719 Sanjiv Kumar. Rethinking FID: Towards a better evaluation metric for image generation. In  
720 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA,*  
721 *USA, June 16-22, 2024*, pp. 9307–9315. IEEE, 2024.
- 722 Woojeong Jin, Yu Cheng, Yelong Shen, Weizhu Chen, and Xiang Ren. A good prompt is worth  
723 millions of parameters: Low-resource prompt-based learning for vision-language models. In  
724 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume*  
725 *1: Long Papers)*, pp. 2763–2775, 2022.
- 726 Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H. Campbell, Konrad  
727 Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin,  
728 Ryan Sepassi, George Tucker, and Henryk Michalewski. Model based reinforcement learning for  
729 atari. *ICLR*, 1:2, 2020.
- 730 Anssi Kanervisto, Christian Scheller, and Ville Hautamäki. Action space shaping in deep reinforce-  
731 ment learning. In *2020 IEEE conference on games (CoG)*, pp. 479–486. IEEE, 2020.
- 732 Anssi Kanervisto, David Bignell, Linda Yilin Wen, Martin Grayson, Raluca Georgescu, Sergio Val-  
733 carcel Macua, Shan Zheng Tan, Tabish Rashid, Tim Pearce, Yuhan Cao, Abdelhak Lemkhenter,  
734 Chentian Jiang, Gavin Costello, Gunshi Gupta, Marko Tot, Shu Ishida, Tarun Gupta, Udit Arora,  
735 Ryan W. White, Sam Devlin, Cecily Morrison, and Katja Hofmann. World and human action  
736 models towards gameplay ideation. *Nature*, 638(8051):656–663, 2025.
- 737 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
738 adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*  
739 *2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 4401–4410. Computer Vision Foundation /  
740 IEEE, 2019.
- 741 Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword  
742 tokenizer and detokenizer for neural text processing. In Eduardo Blanco and Wei Lu (eds.),  
743 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing,*  
744 *EMNLP 2018: System Demonstrations, Brussels, Belgium, October 31 - November 4, 2018*, pp.  
745 66–71. Association for Computational Linguistics, 2018.
- 746 Seongyun Lee, Seungone Kim, Sue Park, Geewook Kim, and Minjoon Seo. Prometheus-vision:  
747 Vision-language model as a judge for fine-grained evaluation. In *Findings of the Association for*  
748 *Computational Linguistics ACL 2024*, pp. 11286–11315, 2024.
- 749 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
750 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-  
751 tion for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:  
752 9459–9474, 2020.

- 756 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image  
757 pre-training with frozen image encoders and large language models. In *International conference*  
758 *on machine learning*, pp. 19730–19742. PMLR, 2023.
- 759 Mingxiang Liao, Qixiang Ye, Wangmeng Zuo, Fang Wan, Tianyu Wang, Yuzhong Zhao, Jingdong  
760 Wang, Xinyu Zhang, et al. Evaluation of text-to-video generation models: A dynamics perspective.  
761 *Advances in Neural Information Processing Systems*, 37:109790–109816, 2024.
- 762 Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and  
763 Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *European*  
764 *Conference on Computer Vision*, pp. 366–384. Springer, 2024.
- 765 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction  
766 tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
767 pp. 26296–26306, 2024a.
- 768 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-  
769 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. In *Forty-first*  
770 *International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*.
- 771 Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu,  
772 Tiejiong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large  
773 video generation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
774 *Pattern Recognition*, pp. 22139–22149, 2024b.
- 775 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International*  
776 *Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*,  
777 2019.
- 778 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *Interna-*  
779 *tional Conference on Learning Representations, 2022*.
- 780 Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. Cheap and  
781 quick: Efficient vision-language instruction tuning for large language models. *Advances in neural*  
782 *information processing systems (NeurIPS)*, 2023.
- 783 Oscar Mañas, Benno Krojer, and Aishwarya Agrawal. Improving automatic vqa evaluation using large  
784 language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp.  
785 4171–4179, 2024.
- 786 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin  
787 Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods, 2022. URL <https://github.com/huggingface/peft>.
- 788 Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter,  
789 Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, Anton Belyi, Haotian Zhang, Karanjeet  
790 Singh, Doug Kang, Hongyu He, Max Schwarzer, Tom Gunter, Xiang Kong, Aonan Zhang, Jianyu  
791 Wang, Chong Wang, Nan Du, Tao Lei, Sam Wiseman, Mark Lee, Zirui Wang, Ruoming Pang,  
792 Peter Grasch, Alexander Toshev, and Yinfei Yang. MM1: Methods, analysis and insights from  
793 multimodal LLM pre-training. In *European Conference on Computer Vision*, pp. 304–323. Springer,  
794 2024.
- 795 Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre,  
796 Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha  
797 Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie  
798 Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Char-  
799 line Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David  
800 Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Cristian Mu-  
801 raru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin,  
802 James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy  
803 Chen, Johan Ferret, Justin Chiu, and et al. Gemma: Open models based on gemini research and  
804 technology. *arXiv preprint arXiv:2403.08295*, 2024.

- 810 Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world models.  
811 In *The Eleventh International Conference on Learning Representations, 2023*.  
812
- 813 Atsuyuki Miyai, Qing Yu, Go Irie, and Kiyoharu Aizawa. Locoop: Few-shot out-of-distribution  
814 detection via prompt learning. *Advances in Neural Information Processing Systems*, 36:76298–  
815 76310, 2023.
- 816 Ivona Najdenkoska, Xiantong Zhen, and Marcel Worring. Meta learning to bridge vision and language  
817 models for multimodal few-shot learning. In *The Eleventh International Conference on Learning  
818 Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023, 2023*.
- 819
- 820 Jingcheng Ni, Yuxin Guo, Yichen Liu, Rui Chen, Lewei Lu, and Zehuan Wu. Maskgwm: A gener-  
821 alizable driving world model with video mask reconstruction. *arXiv preprint arXiv:2502.11663*,  
822 2025.
- 823 Jack Parker-Holder, Philip Ball, Jake Bruce, Vibhavari Dasagi, Kristian Holsheimer, Chris-  
824 tos Kaplanis, Alexandre Moufarek, Guy Scully, Jeremy Shar, Jimmy Shi, Stephen Spencer,  
825 Jessica Yung, Michael Dennis, Sultan Kenjeyev, Shangbang Long, Vlad Mnih, Harris  
826 Chan, Maxime Gazeau, Bonnie Li, Fabio Pardo, Luyu Wang, Lei Zhang, Frederic Besse,  
827 Tim Harley, Anna Mitenkova, Jane Wang, Jeff Clune, Demis Hassabis, Raia Hadsell,  
828 Adrian Bolton, Satinder Singh, and Tim Rocktäschel. Genie 2: A large-scale founda-  
829 tion world model, 2024. URL [https://deepmind.google/discover/blog/  
830 genie-2-a-large-scale-foundation-world-model/](https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-world-model/). Accessed: 20 March 2025.
- 831
- 832 Tim Pearce, Tabish Rashid, Dave Bignell, Raluca Georgescu, Sam Devlin, and Katja Hofmann.  
833 Scaling laws for pre-training agents and world models. In *Proceedings of the 42nd International  
834 Conference on Machine Learning (ICML), 2025*.
- 835
- 836 AJ Piergiovanni, Weicheng Kuo, and Anelia Angelova. Pre-training image-language transformers for  
837 open-vocabulary tasks. *arXiv preprint arXiv:2209.04372*, 2022.
- 838
- 839 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
840 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 841
- 842 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
843 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.  
844 Learning transferable visual models from natural language supervision. *Proceedings of the  
845 International Conference on Machine Learning (ICML), 2021*.
- 846
- 847 Morgane Rivière, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard  
848 Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya  
849 Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy  
850 Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt  
851 Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna  
852 Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda  
853 Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian,  
854 Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty,  
855 Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar,  
856 Dominika Rogozinska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira,  
857 Evan Senter, Evgenii Eltyshv, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus  
858 Martins, Hadi Hashemi, Hanna Klimczak-Plucinska, Harleen Batra, Harsh Dhand, Ivan Nardini,  
859 Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana  
860 Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon,  
861 Josh Lipschultz, Josh Newlan, Ju-yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie  
862 Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjöstrand,  
863 Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, and Lilly McNealus. Gemma 2:  
Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- 862
- 863 Lloyd Russell, Anthony Hu, Lorenzo Bertoni, George Fedoseev, Jamie Shotton, Elahe Arani, and  
Gianluca Corrado. Gaia-2: A controllable multi-view generative world model for autonomous  
driving. *arXiv preprint arXiv:2503.20523*, 2025.

- 864 Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.  
865 Improved techniques for training gans. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg,  
866 Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems*  
867 *29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016,*  
868 *Barcelona, Spain*, pp. 2226–2234, 2016.
- 869 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon  
870 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy P. Lillicrap,  
871 and David Silver. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*,  
872 588(7839):604–609, 2020.
- 873 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,  
874 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th*  
875 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.  
876 2556–2565, 2018.
- 877 Sugandha Sharma, Guy Davidson, Khimya Khetarpal, Anssi Kanervisto, Udit Arora, Katja Hofmann,  
878 and Ida Momennejad. Toward human-ai alignment in large-scale multi-player games. In *ACL 2024*  
879 *Worldplay Workshop*, 2024.
- 880 Haoyu Song, Li Dong, Weinan Zhang, Ting Liu, and Furu Wei. Clip models are few-shot learners:  
881 Empirical studies on vqa and visual entailment. In *Proceedings of the 60th Annual Meeting of the*  
882 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6088–6100, 2022.
- 883 Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. Wit:  
884 Wikipedia-based image text dataset for multimodal multilingual machine learning. In *Proceedings*  
885 *of the 44th international ACM SIGIR conference on research and development in information*  
886 *retrieval*, pp. 2443–2449, 2021.
- 887 Andreas Steiner, André Susano Pinto, Michael Tschannen, Daniel Keysers, Xiao Wang, Yonatan  
888 Bitton, Alexey A. Gritsenko, Matthias Minderer, Anthony Sherbondy, Shangbang Long, Siyang  
889 Qin, R. Reeve Ingle, Emanuele Bugliarello, Sahar Kazemzadeh, Thomas Mesnard, Ibrahim  
890 Alabdulmohsin, Lucas Beyer, and Xiaohua Zhai. Paligemma 2: A family of versatile vlms for  
891 transfer. *arXiv preprint arXiv:2412.03555*, 2024.
- 892 Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart*  
893 *Bulletin*, 2(4):160–163, 1991.
- 894 Suramya Tomar. Converting video formats with ffmpeg. *Linux Journal*, 2006(146):10, 2006.
- 895 Marko Tot, Shu Ishida, Abdelhak Lemkhenter, David Bignell, Pallavi Choudhury, Chris Lovett, Luis  
896 França, Matheus Ribeiro Furtado de Mendonça, Tarun Gupta, Darren Gehring, et al. Adapting  
897 a world model for trajectory following in a 3d game. In *ICLR 2025 Workshop on World Models*,  
898 2025.
- 899 Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill.  
900 Multimodal few-shot learning with frozen language models. *Advances in Neural Information*  
901 *Processing Systems*, 34:200–212, 2021.
- 902 P Umesh. Image processing in python. *CSI Communications*, 23, 2012.
- 903 Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and  
904 Sylvain Gelly. Towards accurate generative models of video: A new metric and challenges. *arXiv*  
905 *preprint arXiv:1812.01717*, 2018.
- 906 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
907 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*  
908 *systems*, 30, 2017.
- 909 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,  
910 Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng  
911 Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-VL: Enhancing vision-language model’s  
912 perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

- 918 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from  
919 error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612,  
920 2004.
- 921 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
922 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. Huggingface’s  
923 transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*,  
924 2019.
- 925  
926 Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and  
927 Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint*  
928 *arXiv:2104.14806*, 2021.
- 929  
930 Mingrui Wu, Xinyue Cai, Jiayi Ji, Jiale Li, Oucheng Huang, Gen Luo, Hao Fei, Guannan Jiang,  
931 Xiaoshuai Sun, and Rongrong Ji. Controlmlm: Training-free visual prompt learning for multimodal  
932 large language models. *Advances in Neural Information Processing Systems*, 37:45206–45234,  
933 2024a.
- 934  
935 Xun Wu, Shaohan Huang, and Furu Wei. Mixture of lora experts. In *ICLR*, 2024b.
- 936  
937 Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen,  
938 Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language  
939 models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna,*  
940 *Austria, May 7-11, 2024*, 2024.
- 941  
942 Sherry Yang, Yilun Du, Seyed Kamyar Seyed Ghasemipour, Jonathan Tompson, Leslie Pack Kael-  
943 bling, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. In *The*  
944 *Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May*  
945 *7-11, 2024*, 2024.
- 946  
947 Peng Ye, Yongqi Huang, Chongjun Tu, Minglei Li, Tao Chen, Tong He, and Wanli Ouyang.  
948 Partial fine-tuning: A successor to full fine-tuning for vision transformers. *arXiv preprint*  
949 *arXiv:2312.15681*, 2023.
- 950  
951 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language  
952 image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*,  
953 pp. 11975–11986, 2023.
- 954  
955 Bingchen Zhao, Haoqin Tu, Chen Wei, Jieru Mei, and Cihang Xie. Tuning layernorm in attention:  
956 Towards efficient multi-modal llm finetuning. In *ICLR*, 2024.
- 957  
958 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
959 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
960 chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.
- 961  
962 Sharon Zhou, Mitchell L. Gordon, Ranjay Krishna, Austin Narcomey, Li Fei-Fei, and Michael S.  
963 Bernstein. HYPE: A benchmark for human eye perceptual evaluation of generative models. In  
964 Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox,  
965 and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual*  
966 *Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019,*  
967 *Vancouver, BC, Canada*, pp. 3444–3456, 2019.
- 968  
969 Yiwei Zhou, Xiaobo Xia, Zhiwei Lin, Bo Han, and Tongliang Liu. Few-shot adversarial prompt  
970 learning on vision-language models. *Advances in Neural Information Processing Systems*, 37:  
971 3122–3156, 2024.

---

# ADAPTING VISION-LANGUAGE MODELS FOR EVALUATING WORLD MODELS

## APPENDIX

---

### TABLE OF CONTENTS

---

972		
973		
974		
975		
976		
977		
978		
979		
980		
981		
982		
983		
984	<b>A Broader Impact</b>	<b>20</b>
985		
986	<b>B Reproducibility Statement</b>	<b>20</b>
987		
988		
989	<b>C UNIVERSE: Additional Details</b>	<b>21</b>
990	C.1 Implementation Overview . . . . .	21
991	C.2 Adaptation to New Domains . . . . .	22
992	C.3 Inference . . . . .	23
993		
994		
995	<b>D Dataset</b>	<b>23</b>
996		
997	D.1 Construction Process . . . . .	23
998	D.2 Release Details . . . . .	25
999		
1000	<b>E Experimental Details</b>	<b>26</b>
1001		
1002	E.1 Model . . . . .	26
1003	E.2 Evaluation Metrics . . . . .	29
1004	E.3 Results . . . . .	29
1005		
1006		
1007	<b>F Human Annotation Study</b>	<b>31</b>
1008		
1009	F.1 Study Design . . . . .	31
1010	F.2 Evaluation Metrics . . . . .	37
1011	F.3 Results . . . . .	38
1012		
1013	<b>G Supplementary Experimental Results</b>	<b>39</b>
1014		
1015	G.1 GPT-5 Performance . . . . .	39
1016	G.2 Zero-Shot Performance of PaliGemma Models . . . . .	39
1017	G.3 CLIPScore Comparisons . . . . .	39
1018	G.4 Low-Rank Adaptation Comparisons . . . . .	42
1019		
1020		
1021	<b>H Comparison with Existing Evaluation Baselines</b>	<b>43</b>
1022		
1023	H.1 Comparison with FVD . . . . .	43
1024	H.2 Correlation with VBench Metrics . . . . .	44
1025		

---

## 1026 A BROADER IMPACT

1027  
1028 As world models become integral to simulation, planning, and decision-making in interactive en-  
1029 vironments, evaluation remains a key bottleneck for both research progress and safe deployment.  
1030 We address this challenge by introducing a unified, sample-efficient framework for evaluating world  
1031 model rollouts using adapted VLMs, designed for fine-grained, temporally grounded, and semanti-  
1032 cally coherent assessment.

1033 This capability has direct implications for high-impact domains such as neural game engines (Kan-  
1034 ervisto et al., 2025; Guo et al., 2025; Gao et al., 2025; Chen et al., 2025), embodied AI (Du et al.; Yang  
1035 et al., 2024), and autonomous driving (Russell et al., 2025; Hu et al., 2023a; Ni et al., 2025), where  
1036 world models simulate environment dynamics and support downstream control and generalization. In  
1037 such contexts, precise and interpretable evaluation is critical not only for benchmarking, but also for  
1038 diagnosing failure modes and ensuring alignment with intended behaviors.

1039 By reducing dependence on human annotation and task-specific fine-tuning, UNIVERSE offers a  
1040 scalable alternative that lowers the computational and environmental costs of rollout evaluation.  
1041 However, reliance on automated evaluators introduces risks: adapted VLMs may inherit biases  
1042 from pretraining, struggle under distributional shift, or yield unreliable judgments in edge cases.  
1043 These risks are amplified in safety-critical settings, where miscalibrated evaluations can propagate  
1044 downstream errors.

1045 We therefore advocate for cautious deployment, accompanied by human oversight, rigorous validation,  
1046 and transparent reporting. While UNIVERSE advances the automation of world model evaluation,  
1047 it must be situated within evaluation pipelines that foreground robustness, interpretability, and  
1048 accountability.

## 1050 B REPRODUCIBILITY STATEMENT

1051  
1052 To support reproducibility and facilitate future research, we provide detailed instructions for repro-  
1053 ducing all main experiments. Detailed descriptions of model architectures, training procedures, and  
1054 dataset construction are provided in Section 4 and Appendix E. A high-level overview of the overall  
1055 implementation framework is included in Appendix C.1. All experiments have been repeated for  
1056 three runs. Plots and tables with quantitative results show the standard deviation across these runs.

1057 *Use of Existing Assets.* We experiment with a range of open-weight VLMs, including three  
1058 PaliGemma variants (version 1 (3B) (Beyer et al., 2024) and version 2 (3B and 10B) (Steiner  
1059 et al., 2024)), VideoLLaMA3 (2B, 7B) (Boqiang Zhang, 2025), and CLIP (Radford et al., 2021)  
1060 with the following vision encoder configurations: ViT-B/32, ViT-B/16, ViT-L/14, and ViT-L/14  
1061 with  $336 \times 336$  resolution. UNIVERSE is built on top of PaliGemma v2 (3B), using publicly re-  
1062 leased checkpoints for initialization. Further architectural and implementation details are provided  
1063 in Appendix E.1. For our software stack, we use Matplotlib (Hunter, 2007) for plotting, NumPy  
1064 (Harris et al., 2020) for data handling, openCV (Bradski, 2000), FFmpeg (Tomar, 2006) and PIL  
1065 (Umesh, 2012) for video and image processing, and NLTK (Bird & Loper, 2004) for text processing.  
1066 Parameter-efficient fine-tuning is implemented using the PEFT library (Mangrulkar et al., 2022). We  
1067 log our experiments using Weights and Biases (Biewald, 2020).

1068 *Use of Large Language Models.* Portions of the manuscript were polished with the assistance of large  
1069 language models (LLMs). The use of LLMs was limited to only improving readability and style; all  
1070 ideas, and experimental designs were developed by the authors.

1071 *Compute Resources.* All experiments were conducted using NVIDIA A100 GPUs (40GB memory)  
1072 on an internal compute cluster. Each model was trained and/or evaluated using 8 GPUs. The  
1073 compute breakdown is as follows: zero-shot evaluation experiments consumed approximately 136  
1074 GPU-days; baseline fine-tuning experiments required around 864 GPU-days; analysis experiments  
1075 contributed the bulk of usage, totaling 2,554 GPU-days. Human evaluation experiments—including  
1076 rollout generation and response annotation using UNIVERSE—incur an additional 1.125 GPU-days.  
1077 Additional compute was required for preliminary experiments, and failed runs not included in the  
1078 final paper. These development activities accounted for an estimated 1,599 GPU-days. In total, all  
1079 experiments amounted to approximately 5,153 GPU-days, equivalent to 14.12 GPU-years.

## C UNIVERSE: ADDITIONAL DETAILS

### C.1 IMPLEMENTATION OVERVIEW

This section outlines the implementation of UNIVERSE in Python, presented as high-level pseudocode. The system is structured around two main stages: (i) *Adaptation*: fine-tuning a VLM on task-specific question-answer (QA) supervision derived from ground truth; (ii) *Evaluation*: using the adapted model to assess new rollouts via structured, prompt-based recognition tasks.

#### Adaptation Pipeline.

The adaptation stage can be implemented as two modules: `AdaptationDatasetBuilder` and `VLMAdapter`.

**AdaptationDatasetBuilder.** This class constructs an adaptation dataset from raw ground truth data, initialized via `load_ground_truth_data` (see Section 3 and Appendix D). The core method, `build`, takes four arguments: `alpha_task`, which specifies the task mixture ratio; `beta_format`, which controls the distribution over QA prompt formats; `context_length`, which determines the number of frames per QA instance; and `sampling_strategy`, which defines how frames are sampled from rollouts. The builder first applies `stratified_sample` to select a subset of annotated samples that match the specified configuration. For each sample, it invokes `_sample_visual_context` to extract the relevant frames, and constructs a triplet consisting of frames, question, and answer.

**VLMAdapter.** This class applies an adaptation strategy to a base VLM, passed via the `base_vlm` argument. Given an adaptation dataset `adaptation_data`, a tuning strategy specified by the `strategy` parameter, and a fixed number of training steps `num_steps`, the adapter trains the model by iteratively sampling a batch, computing the loss via `compute_loss`, and applying updates with `update_model`.

```

1108 class AdaptationDatasetBuilder:
1109     def __init__(self, raw_data_path):
1110         self.samples = load_ground_truth_data(raw_data_path)
1111
1112     def build(self, alpha_task, beta_format, context_length,
1113             ↪ sampling_strategy):
1114         formatted = stratified_sample(
1115             samples=self.samples,
1116             task_proportions=alpha_task,
1117             format_proportions=beta_format
1118         )
1119         dataset = []
1120         for sample in formatted:
1121             visual_ctx = self._sample_visual_context(
1122                 sample["frames"], context_length, sampling_strategy
1123             )
1124             dataset.append({
1125                 "frames": visual_ctx,
1126                 "question": sample["question"],
1127                 "answer": sample["answer"]
1128             })
1129         return dataset
1130
1131 class VLMAdapter:
1132     def __init__(self, base_vlm):
1133         self.base_vlm = base_vlm
1134
1135     def adapt(self, adaptation_data, strategy, num_steps):
1136         configure_adaptation(self.base_vlm, strategy)
1137         for step in range(num_steps):

```

```

1134         batch = sample_from(adaptation_data)
1135         loss = compute_loss(self.base_vlm, batch)
1136         update_model(self.base_vlm, loss)
1137     return self.base_vlm

```

### Evaluation Pipeline.

The evaluation stage can be implemented via two additional modules: `RolloutsGenerator` and `Universe`.

**RolloutsGenerator.** This component autoregressively samples rollout trajectories from a world model (`textttworld_model`). Given an initial observation `o_initial` and an action sequence `a_seq`, the `rollout` method generates a sequence of predicted observations by maintaining lists of past observations (`o_lt`) and actions (`a_lt`). At each timestep, it calls `predict_next_observation` to obtain the next predicted frame, appends it to the rollout sequence `o_seq`, and continues until `timesteps` is reached. This process produces a full trajectory simulating environment dynamics.

**Universe.** This module serves as the inference engine of our framework. It wraps an adapted VLM passed via `adapted_vlm`. Given a generated rollout and an evaluation specification, the method `evaluate_rollout` constructs a prompt using `generate_question`, parameterized by a recognition target and complexity level. It then calls `evaluate`, which queries the VLM with the resulting rollout and question, returning the model’s answer.

```

1157 class RolloutsGenerator:
1158     def __init__(self):
1159         self.world_model = WorldModel(...)
1160
1161     def predict_next_observation(self, o_lt, a_lt):
1162         return self.world_model(o_lt, a_lt)
1163
1164     def rollout(self, o_initial, a_seq, timesteps):
1165         o_seq = [o_initial]
1166         o_lt, a_lt = [o_initial], []
1167         for t in range(timesteps):
1168             a_lt.append(a_seq[t])
1169             o_t = self.predict_next_observation(o_lt, a_lt)
1170             o_seq.append(o_t)
1171             o_lt.append(o_t)
1172         return o_seq
1173
1174 class Universe:
1175     def __init__(self, adapted_vlm):
1176         self.vlm = adapted_vlm
1177
1178     def evaluate(self, rollout, question):
1179         return self.vlm(rollout, question)
1180
1181     def evaluate_rollout(self, rollout, target, complexity):
1182         question = generate_question(rollout, target, complexity)
1183         return self.evaluate(rollout, question)

```

## C.2 ADAPTATION TO NEW DOMAINS

To adapt UNIVERSE to a new environment, one could collect a small set of reference trajectories that provide ground-truth observations and actions for the evaluation dimensions of interest. Question and answer templates can then be instantiated from these ground-truth signals to construct a mixed-supervision adaptation dataset.

In practice, three complementary routes exist for obtaining such reference data: (i) *Environment-side metadata*: partnering with environment developers (as in our work) provides high-fidelity ground truth through controller logs and environment state. This route yields the highest-quality labels but requires coordination with studios or simulation platforms. (ii) *Manual annotation*: when metadata is unavailable, a one-time human annotation effort can label a small set of trajectories. While this incurs upfront cost, it is less expensive than continuous human evaluation of every generated rollout. (iii) *Synthetic data generation*: for domains with publicly available gameplay videos or recordings, generative models can provide initial noisy labels. These synthetic labels can bootstrap the adaptation dataset, either directly or after human verification of a subset.

Once reference data is obtained through either route, question-answer templates are instantiated from ground-truth signals, and the resulting dataset can be used for model adaptation. Once adapted, the evaluator could provide dimension-wise assessments that can be aggregated or used as feedback in downstream world-model training.

Although our experiments use 14-frame inputs, this is a design choice rather than a fundamental constraint. The same procedure extends to longer rollouts by increasing the visual context window or applying a sliding-window scheme over successive 14-frame segments.

Overall, deploying UNIVERSE as a standalone evaluator in novel domain would require: (i) a small adaptation dataset of reference trajectories with ground-truth labels, (ii) the task specification, and (iii) the lightweight fine-tuning recipe.

### C.3 INFERENCE

At inference time, the evaluator processes rollouts by emitting natural-language responses along specified dimensions. These outputs can be used either as (i) structured feedback for world model improvement or (ii) mapped to numerical scores and aggregated into trajectory-level quality metrics.

Specifically, our method supports two complementary evaluation setups:

- (i) *Ground-truth-aligned evaluation*: given a generated rollout  $r$  and ground truth information  $GT$  (e.g., actions conditioning the world model that we want to evaluate), we use  $GT$  to instantiate  $QA$  pairs with known correct answers  $(A, \hat{A})$ . UNIVERSE generates predictions  $\hat{A}$  for each question  $Q$  and we compute EM and ROUGE for  $(A, \hat{A})$  to quantify alignment.
- (ii) *Open-ended evaluation*: given a generated rollout  $r$  and a dimension of interest (e.g., action execution, character presence), we instantiate questions from the world model’s conditioning information without reference answers. UNIVERSE’s responses provide structured semantic feedback, e.g., binary questions can identify which frames contain target actions or characters, enabling downstream analysis. While this setting lacks ground-truth alignment scores, the evaluator’s responses can be aggregated into meaningful statistics.

## D DATASET

This section details the construction and release of the dataset used to adapt VLMs for fine-grained evaluation of world model rollouts. We curate a realistic, human-centered dataset derived from actual gameplay in a complex multi-agent environment. Designed to provide temporally grounded and semantically structured supervision, the dataset aligns with the downstream evaluation setting and supports adaptation to both action and character recognition tasks across all QA formats. We describe the data construction pipeline, QA generation process, and release format below.

### D.1 CONSTRUCTION PROCESS

The ground truth dataset for adapting the evaluator (see Section 3) was developed in collaboration with *Ninja Theory* using human gameplay recordings from *Bleeding Edge*, a 4v4 multiplayer combat game. Data use was governed by a formal agreement with the studio, and collection adhered to the game’s End User License Agreement (EULA). All protocols were approved by our Institutional Review Board (IRB), and personally identifiable information (PII) was removed prior to analysis.

Each gameplay session is represented as a tuple  $s = (v, c, m)$ , where  $v$  is a high-resolution MP4 video (60 FPS),  $c$  is the synchronized controller action log, and  $m$  contains structured metadata (e.g., player roles, agent identities, action categories, and map context). The full set of gameplay sessions is denoted by  $\mathcal{S} = \{(v_i, c_i, m_i)\}_{i=1}^{|\mathcal{S}|}$ .

The dataset construction pipeline proceeds in three stages:

- (i) *Preprocessing*. We begin by filtering out corrupted applying or inactive sessions and synchronizes the video, controller logs, and metadata streams using internal game timestamps:  $\mathcal{S}_{\text{valid}} = \text{Preprocessing}(\mathcal{S})$ . Each valid session is segmented into non-overlapping clips of fixed length  $L = 14$  frames, each paired with controller input and shared metadata; formally, for a session  $s = (v, c, m) \in \mathcal{S}_{\text{valid}}$ , the segmentation produces  $\text{Segment}(v, c, m, L) = \{(f^{(1:L)}, c^{(1:L)}, m)\}$ , where  $f^{(1:L)}$  denotes the sequence of frames,  $c^{(1:L)}$  the aligned controller inputs, and  $m$  the associated metadata. The complete set of extracted clips across all valid sessions is defined as  $\mathcal{V} = \bigcup_{s \in \mathcal{S}_{\text{valid}}} \text{Segment}(s, L)$ , where each element  $v \in \mathcal{V}$  is a triplet  $(f^{(1:L)}, c^{(1:L)}, m)$  consisting of video frames, corresponding controller inputs, and metadata.
- (ii) *Description Generation*. Next, for each sequence of frames  $f^{(1:L)} \in \mathcal{V}$ , we use the associated control log  $c^{(1:L)}$  to extract action information and the metadata  $m$  to obtain character-related attributes. These are combined to generate a structured natural language description via  $d = \text{Describe}(c^{(1:L)}, m)$ , where `Describe` is a rule-based procedure that transforms the logged actions and metadata into textual descriptions used for constructing the QA supervision. This yields a set of paired video–text examples:  $\mathcal{Z} = \{(f^{(1:L)}, d) \mid f^{(1:L)} \in \mathcal{V}\}$ .
- (iii) *Question-Answer Pair Construction*. Finally, we generate six QA pairs per clip, spanning two predefined tasks (AR and CR), each instantiated in three question formats: binary, multiple-choice, and open-ended. To enable this, we define task-specific answer spaces using `GetAnswerSpace( $\mathcal{Z}$ )`, which returns  $\mathcal{Y}_{\text{AR}}$  for action categories and  $\mathcal{Y}_{\text{CR}}$  for character identities, based on all video–text pairs in  $\mathcal{Z}$ . For each clip, we extract the task-specific ground-truth answer from the corresponding description as  $y = \text{ExtractLabel}(d, t)$ , where  $t \in \{\text{AR}, \text{CR}\}$ . Each QA format is constructed as follows: (i) *Binary*: Two binary question-answer pairs are generated per instance using `FormatBinaryPrompt`. The positive question  $Q^{\text{pos}}$  is constructed using the correct label  $y \in \mathcal{Y}^{(t)}$  and paired with the positive answer  $A^{\text{pos}}$ . The negative question  $Q^{\text{neg}}$  is constructed using an incorrect label  $\tilde{y} \sim \text{SampleDistractor}(\mathcal{Y}^{(t)} \setminus \{y\})$  and paired with the negative answer  $A^{\text{neg}}$ . (ii) *Multiple-Choice*: A question  $Q$  is generated using the full set of candidate options, formatted via `FormatOptions( $\mathcal{Y}_t$ )`. The question is constructed with `FormatMCPrompt( $t, O$ )` and paired with the correct answer  $y \in \mathcal{Y}_t$ . (iii) *Open-Ended*: A free-form question  $Q$  is generated using `FormatOEPrompt( $t$ )`, prompting the model to produce the correct label  $y \in \mathcal{Y}_t$  without access to predefined answer choices.

The final dataset is represented as  $\mathcal{D} = \{(f_i^{(1:L)}, QA_i)\}_{i=1}^{|\mathcal{D}|}$ , where each  $f^{(1:L)}$  is a video clip and  $QA = \{(Q_j, A_j)\}_{j=1}^6$  is the associated set of question–answer pairs, covering all combinations of three question formats (binary, multiple-choice, open-ended) and two tasks (Action Recognition and Character Recognition). A detailed data pipeline is provided in Algorithm 1.

---

```

1296 Algorithm 1 Dataset Construction Process
1297
1298 Procedure DatasetCreation ( $\mathcal{S}, L$ ):
1299    $\mathcal{S}_{\text{valid}} \leftarrow \text{Preprocessing}(\mathcal{S})$ 
1300    $\mathcal{V} \leftarrow \emptyset$ 
1301   for  $(v, c, m) \in \mathcal{S}_{\text{valid}}$  do
1302      $\mathcal{V}_s \leftarrow \text{Segment}(v, c, m, L)$ 
1303      $\mathcal{V} \leftarrow \mathcal{V} \cup \mathcal{V}_s$ 
1304    $\mathcal{Z} \leftarrow \emptyset$ 
1305   for  $(f^{(1:L)}, c^{(1:L)}, m) \in \mathcal{V}$  do
1306      $d \leftarrow \text{Describe}(m, c^{(1:L)})$ 
1307      $\mathcal{Z} \leftarrow \mathcal{Z} \cup \{(f^{(1:L)}, d)\}$ 
1308    $\mathcal{D} \leftarrow \emptyset$ 
1309    $\mathcal{Y}_{\text{AR}}, \mathcal{Y}_{\text{CR}} \leftarrow \text{GetAnswerSpace}(\mathcal{Z})$ 
1310   for  $(f^{(1:L)}, d) \in \mathcal{V}$  do
1311      $\mathcal{QA} \leftarrow \text{GenerateQAPairs}(d, \mathcal{Y}_{\text{AR}}, \mathcal{Y}_{\text{CR}})$ 
1312     for  $(Q, A) \in \mathcal{QA}$  do
1313        $\mathcal{D} \leftarrow \mathcal{D} \cup \{(f^{(1:L)}, Q, A)\}$ 
1314   return  $\mathcal{D}$ 
1315
1316 Procedure GenerateQAPairs ( $d, \mathcal{Y}_{\text{AR}}, \mathcal{Y}_{\text{CR}}$ ):
1317    $\mathcal{QA} \leftarrow \emptyset$ 
1318   for  $t \in \{\text{AR}, \text{CR}\}$  do
1319      $y \leftarrow \text{ExtractLabel}(d, t)$ 
1320      $QA_{\text{bin}}^{\text{pos}}, QA_{\text{bin}}^{\text{neg}} \leftarrow \text{CreateBinaryQA}(t, y)$ 
1321      $\mathcal{QA} \leftarrow \mathcal{QA} \cup \{QA_{\text{bin}}^{\text{pos}}, QA_{\text{bin}}^{\text{neg}}\}$ 
1322      $QA_{\text{mc}} \leftarrow \text{CreateMCQA}(t, y, \mathcal{Y}_t)$ 
1323      $\mathcal{QA} \leftarrow \mathcal{QA} \cup QA_{\text{mc}}$ 
1324      $QA_{\text{oe}} \leftarrow \text{CreateOpenEndedQA}(t, y)$ 
1325      $\mathcal{QA} \leftarrow \mathcal{QA} \cup QA_{\text{oe}}$ 
1326   return  $\mathcal{QA}$ 
1327
1328 Procedure CreateBinaryQA ( $t, y$ ):
1329    $\tilde{y} \leftarrow \text{SampleDistractor}(\mathcal{Y}_t \setminus \{y\})$ 
1330    $Q^{\text{pos}} \leftarrow \text{FormatBinaryPrompt}(t, y)$ 
1331    $Q^{\text{neg}} \leftarrow \text{FormatBinaryPrompt}(t, \tilde{y})$ 
1332   return  $\{(Q^{\text{pos}}, A^{\text{pos}}), (Q^{\text{neg}}, A^{\text{neg}})\}$ 
1333
1334 Procedure CreateMCQA ( $t, y, \mathcal{Y}_t$ ):
1335    $O \leftarrow \text{FormatOptions}(\mathcal{Y}_t)$ 
1336    $Q \leftarrow \text{FormatMCPrompt}(t, O)$ 
1337   return  $Q, y$ 
1338
1339 Procedure CreateOpenEndedQA ( $t, y$ ):
1340    $Q \leftarrow \text{FormatOEPrompt}(t)$ 
1341   return  $Q, y$ 

```

---

## 1339 D.2 RELEASE DETAILS

1341 To support reproducibility and further research, we aim to release a subset of our evaluation data.  
1342 This includes sampled human gameplay segments, aligned action vectors and environment states,  
1343 natural language descriptions, and QA annotations spanning binary, multiple-choice, and open-ended  
1344 formats. The dataset is included in the supplementary ZIP file and will be publicly released following  
1345 the publication of the paper.

1346 **File Layout.** The data is organized as follows:

- 1347 • human-gameplay-segments/: directory of .npz files, each containing image frames  
1348 along with frame-aligned actions and states;

- `annotations.jsonl`: line-delimited JSON file containing natural language descriptions, QA prompts, and ground truth answers.

**Structure.** Each dataset instance corresponds to a short human gameplay segment stored as a NumPy archive (`.npz`), containing:

- (i) `images`  $\in \mathbb{R}^{14 \times 3 \times 180 \times 300}$ : a sequence of 14 RGB frames in channel-first (CHW) format;
- (ii) `actions`  $\in \mathbb{R}^{14 \times 16}$ : frame-aligned control vectors;
- (iii) `states`  $\in \mathbb{R}^{14 \times 56}$ : frame-aligned environment states.

**Annotation Format.** Annotations are provided in `annotations.jsonl`, a line-delimited JSON file where each entry corresponds to a single gameplay segment. Each entry includes structured prompts and ground truth answers spanning all tasks and formats.

Specifically, each annotation entry includes:

- `filename`: Unique identifier of the associated `.npz` file containing visual observations (frames), action vectors, and states.
- `description`: Natural language summary of the video segment.
- `ar_binary_pos_q`, `ar_binary_pos_a`: Affirmative binary question and corresponding answer, evaluating recognition of the correct action.
- `ar_binary_neg_q`, `ar_binary_neg_a`: Negative binary question and corresponding answer, targeting rejection of an incorrect action.
- `ar_mc`: Multiple-choice question prompting the model to select the correct action from a list of candidate classes.
- `ar_oe`: Open-ended question prompting free-form generation of the observed action.
- `ar_answer`: Ground truth action label corresponding to both `ar_mc` and `ar_oe`.
- `cr_binary_pos_q`, `cr_binary_pos_a`: Affirmative binary question and corresponding answer for identifying the correct character.
- `cr_binary_neg_q`, `cr_binary_neg_a`: Negative binary question and corresponding answer targeting an incorrect character identity.
- `cr_mc`: Multiple-choice question prompting identification of the correct character from a candidate set.
- `cr_oe`: Open-ended question prompting free-form naming of the character.
- `cr_answer`: Ground truth character label shared across both `cr_mc` and `cr_oe`.

## E EXPERIMENTAL DETAILS

In this section, we provide a detailed description of the dataset preparation process, model architecture, prompt templates, training procedure. Additionally, we provide an overview of all results presented in the main paper in numerical table form, an report additional experimental results leveraging alternate fine-tuning solutions.

### E.1 MODEL

This section provides extended details on the architecture, pretraining configuration, and input formatting of the vision-language models used in our experiments. Our primary backbone is PaliGemma (Beyer et al., 2024; Steiner et al., 2024).

#### E.1.1 OVERVIEW

PaliGemma is a VLM that processes both images and text as input and autoregressively generates natural language output. It follows the training paradigm of PaLI-3 (Chen et al., 2023), combining a ViT-based vision encoder (Dosovitskiy et al., 2021) with a decoder-only Transformer language model.

1404 Table 1: Detailed architecture of the PaliGemma model, comprising a SigLIP-So400m vision tower,  
 1405 a multimodal projection head, and a Gemma-based language decoder. All transformer layers follow  
 1406 standard design and include residual connections around attention and MLP blocks.  
 1407

1408	Component	Configuration
1409	<i>Vision Tower: SigLIP-So400m</i>	
1410	Patch Embedding	Conv2d(in=3, out=1152, kernel=14, stride=14)
1411	Position Embedding	Embedding(num_embeddings=256, emb_dim=1152)
1412	Encoder	27 × Transformer Encoder Layers
1413	Self-Attention	—
1414	Query / Key / Value projection	Linear(1152 → 1152, bias=True)
1415	Layer Normalization	LayerNorm((1152,), eps=1e-6)
1416	MLP Block	—
1417	Activation Function	GELU-Tanh
1418	Feedforward layer (up)	Linear(1152 → 4304, bias=True)
1419	Feedforward layer (down)	Linear(4304 → 1152, bias=True)
1420	Layer Normalization	LayerNorm((1152,), eps=1e-6)
1421	Post-Encoder Layer Norm	LayerNorm((1152,), eps=1e-6)
1422	<i>Multimodal Projection Head</i>	
1423	Linear Projection	Linear(1152 → 2304, bias=True)
1424	<i>Language Model: Gemma</i>	
1425	Token Embedding	Embedding(vocab=257216, dim=2304)
1426	Decoder Stack	26 × Transformer Decoder Layers
1427	Self-Attention	—
1428	Query projection	Linear(2304 → 2048, bias=False)
1429	Key projection	Linear(2304 → 1024, bias=False)
1430	Value projection	Linear(2304 → 1024, bias=False)
1431	Output projection	Linear(2048 → 2304, bias=False)
1432	MLP Block	—
1433	Gating projection	Linear(2304 → 9216, bias=True)
1434	Down projection	Linear(2304 → 9216, bias=True)
1435	Up projection	Linear(9216 → 2304, bias=True)
1436	Activation Function	GELU-Tanh
1437	Normalization Layers	—
1438	Input Norm	RMSNorm(2304, eps=1e-6)
1439	Post-Attn Norm	RMSNorm(2304, eps=1e-6)
1440	Pre-FFN Norm	RMSNorm(2304, eps=1e-6)
1441	Post-FFN Norm	RMSNorm(2304, eps=1e-6)
1442	Rotary Embeddings	GemmaRotaryEmbedding
1443	LM Head	Linear(2304 → 257216, bias=False)

1444 The model is publicly available (Wolf et al., 2019). The architecture is fully modular, comprising  
 1445 three parameterized components: (i) *Vision encoder* ( $\mathcal{M}_V$ ): based on SigLIP (Zhai et al., 2023),  
 1446 specifically the “shape optimized” So400m (Alabdulmohsin et al., 2023). (ii) *Multimodal projection*  
 1447 *head* ( $\mathcal{M}_P$ ): a single linear layer for projecting visual features into the language decoder’s embedding  
 1448 space. (iii) *Language decoder* ( $\mathcal{M}_L$ ): a Transformer-based autoregressive model from the Gemma  
 1449 family (Mesnard et al., 2024; Rivière et al., 2024). Below, we discuss the architecture in more details,  
 1450 the general layer-level overview is also provided in Table 1.  
 1451

1452 **Vision Encoder: SigLIP-So400m.** The visual backbone  $\mathcal{M}_V$  is a ViT-style encoder pretrained  
 1453 using a Sigmoid contrastive loss (SigLIP). It processes input images by dividing them into non-  
 1454 overlapping  $14 \times 14$  patches. Each patch is linearly projected into a 1152-dimensional embedding  
 1455 via a convolutional stem. To encode spatial structure, learned positional embeddings are added before  
 1456

Table 2: Component-wise parameter overview of the PaliGemma model.

Component	Model / Variant	Details	# Params
Vision Encoder	SigLIP-So400m	Input resolutions: 224px <sup>2</sup> , 448px <sup>2</sup> , 896px <sup>2</sup>	400M
Multimodal Projection	—	Connects vision and language components	2.66M
Language Model	PG 1	Gemma 1 2B, pre-trained on 6T tokens	3B
	PG 2	Gemma 2 2B, pre-trained on 2T tokens	3B
	PG 3	Gemma 2 9B, pre-trained on 8T tokens	9.7B

the representation is passed through a stack of 27 SigLIP encoder layers. Each encoder layer contains multi-head self-attention with projection layers for queries, keys, and values, followed by an MLP block with GELU-Tanh activations. All transformer blocks use LayerNorm and residual connections. The vision tower supports multiple input resolutions (224, 448, 896), though our experiments fix resolution at 224px<sup>2</sup> for consistency and efficiency.

**Multimodal Projection Head.** The projection head  $\mathcal{M}_P$  is a lightweight linear mapping from the vision encoder’s output dimension (1152) to the language decoder’s input dimension (2304). It contains approximately 2.66M parameters and is initialized with zero-mean weights. This head enables alignment between visual and linguistic modalities and is important for bridging the representation gap between the vision and language components.

**Language Decoder: Gemma.** The language module  $\mathcal{M}_L$  is a decoder-only Transformer with 26 layers and 2304-dimensional hidden states. Token embeddings are learned over a vocabulary of 257,216 tokens, encoded using the SentencePiece tokenizer (Kudo & Richardson, 2018). Each Transformer block contains a self-attention mechanism with separate linear projections for queries, keys, and values. The MLP block follows a gated architecture, where the input is processed through parallel down projection and gating projection layers, modulated by a GELU-Tanh activation (Hendrycks & Gimpel, 2016), combined via elementwise multiplication, and then passed through an up projection to return to the model’s hidden dimension. RMSNorm is applied before and after both attention and MLP sublayers to stabilize training. Rotary positional embeddings are added to enable relative position encoding. Output tokens are produced via a tied language modeling head that projects back to the vocabulary space.

### E.1.2 CONFIGURATIONS

Table 2 summarizes the architecture components and parameter counts of the PaliGemma configurations available for experimentation. While we focus on the PaliGemma 2 3b variant in our study, we include all publicly released configurations for completeness and to clarify how our selected model compares to other available options. All three variants share the same vision encoder and multimodal integration strategy, differing only in the language decoder. The first configuration, PaliGemma 1 3b, pairs the visual encoder with Gemma 1 (2B), pretrained on 6 trillion tokens, resulting in a total model size of approximately 3 billion parameters. The second configuration, PaliGemma 2 3b, replaces the decoder with Gemma 2 (2B), pretrained on 2 trillion tokens, and maintains a comparable total parameter count. The third and largest variant, PaliGemma 2 10b, uses Gemma 2 (9B) as the decoder, pretrained on 8 trillion tokens, yielding a total model size of approximately 9.7 billion parameters.

### E.1.3 PROMPT FORMAT

To generate textual responses, we adopt a unified prompt format for the decoder. Each input sequence consists of image tokens  $S_I$ , a textual prefix  $S_T^{\text{PREFIX}}$  containing the question, and a suffix  $S_T^{\text{SUFFIX}}$  containing the expected answer. The model autoregressively generates the answer tokens, and training loss is applied only to the suffix.

Let  $n$  denote the number of input frames and  $p$  the number of visual tokens (patch embeddings) per frame. In our setting, each frame is encoded as  $p = 256$  visual tokens. The overall input schema is as follows:

$$\begin{aligned}
S = & \underbrace{\langle \text{image} \rangle_1^{(1)}, \dots, \langle \text{image} \rangle_p^{(1)}, \dots, \langle \text{image} \rangle_1^{(n)}, \dots, \langle \text{image} \rangle_p^{(n)}}_{S_{\mathcal{I}}: \text{Visual tokens from } n \text{ frames, each represented as } p \text{ patches}} \\
& \underbrace{\langle \text{BOS} \rangle, \text{answer en}, \langle \text{QUESTION} \rangle, \langle \text{SEP} \rangle}_{S_{\mathcal{T}}^{\text{PREFIX}}: \text{Prefix (cue + question)}} \\
& \underbrace{\langle \text{ANSWER} \rangle, \langle \text{EOS} \rangle, \langle \text{PAD} \rangle, \dots, \langle \text{PAD} \rangle}_{S_{\mathcal{T}}^{\text{SUFFIX}}: \text{Suffix (answer)}}
\end{aligned}$$

Here,  $S_{\mathcal{I}}$  contains visual tokens produced by the vision encoder  $\mathcal{M}_V$ , and projected into  $\mathcal{M}_L$  space using  $\mathcal{M}_P$ . The prefix  $S_{\mathcal{T}}^{\text{PREFIX}}$  starts with a special  $\langle \text{BOS} \rangle$  token and includes a task-language cue (e.g., “answer en”), the question, and a separator  $\langle \text{SEP} \rangle$ . The suffix  $S_{\mathcal{T}}^{\text{SUFFIX}}$  contains the target answer, terminated with  $\langle \text{EOS} \rangle$  and padded with  $\langle \text{PAD} \rangle$  tokens for batching.

#### E.1.4 PRETRAINING DATA AND FILTERING

PaliGemma is pretrained on a mixture of large-scale vision-language datasets, including WebLI (Chen et al., 2022b), CC3M-35L (Sharma et al., 2018), VQ<sup>2</sup>A-CC3M-35L (Changpinyo et al., 2022), OpenImages (Piergiovanni et al., 2022), and WIT (Srinivasan et al., 2021). Data quality and safety are maintained through pornographic content filtering, text safety and toxicity filtering, and privacy-preserving measures.

## E.2 EVALUATION METRICS

In this section, we provide additional details on metrics used for quantitative evaluation. We employ two complementary metrics: *Exact Match (EM)* and *ROUGE-F<sub>1</sub> (ROUGE)*, which together capture both syntactic precision and semantic alignment.

**Exact Match Accuracy (EM)** measures whether the generated answer is identical to the expected answer, providing a high-precision signal for correctness. Formally, it is defined as:

$$EM = \mathbb{1}(\hat{A} = A) \quad (2)$$

where  $\hat{A}$  is the model’s prediction and  $A$  is the corresponding ground-truth answer. This metric is especially informative for binary and multiple-choice formats where the output space is well-defined.

**ROUGE F<sub>1</sub> (ROUGE)** captures token-level semantic overlap between generated and reference responses by computing the harmonic mean of precision and recall. This allows us to account for partially correct or paraphrased answers. For binary questions, we compute the metric on the bigram level, while for multiple-choice and open-ended formats, we use trigram-level evaluation.

Formally, let  $G$  and  $GT$  denote the sets of  $n$ -grams in the generated and ground truth answers, respectively. Precision and recall are defined as:

$$P = \frac{|G \cap GT|}{|G|}, \quad R = \frac{|G \cap GT|}{|GT|} \quad (3)$$

where  $|G \cap GT|$  counts overlapping  $n$ -grams. The ROUGE score is then computed as:

$$\text{ROUGE} = 2 \times \frac{P \times R}{P + R} \quad (4)$$

Together, these metrics provide a robust view of model performance: EM reflects exact correctness, while ROUGE provides a softer measure of semantic fidelity, particularly useful for evaluating open-ended generations.

## E.3 RESULTS

**Hyperparameters.** Table 3 summarizes the core training hyperparameters used across all adaptation experiments. We train all models on 8 NVIDIA A100 GPUs with a batch size of 1 per device and

Table 3: Summary of hyperparameters used in our experiments.

Hyperparameter	Value
Input resolution	$224 \times 224$
Image frames per input	1–8
Number of epochs	1–10
Batch size (per device)	1
Gradient accumulation steps	4
Optimizer	AdamW Loshchilov & Hutter (2019)
Learning rate	$5 \times 10^{-5}$ , cosine annealing
Learning rate warmup	10%
Weight decay	$1 \times 10^{-6}$
Gradient clipping	Global norm, threshold 1.0
VLM backbone	PaliGemma 2 (3B) Beyer et al. (2024)

accumulate gradients over 4 steps, yielding an effective batch size of 32. Each epoch corresponds to a full pass over the adaptation dataset, and no early stopping is applied. Models were trained for 1–10 epochs depending on task and setting. Optimization is performed using AdamW (Loshchilov & Hutter, 2019) with parameters  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , a base learning rate of  $5 \times 10^{-5}$ , and weight decay of  $1 \times 10^{-6}$ . We use cosine learning rate annealing (Loshchilov & Hutter, 2022) with a linear warmup over the first 10% of training steps. To stabilize training, we apply gradient clipping with a global norm threshold of 1.0. All models use PaliGemma 2 (3B) (Beyer et al., 2024) as the vision-language backbone unless otherwise noted. We vary the number of input frames between 1 and 8 depending on task, and all images are resized to a fixed resolution of  $224 \times 224$ . Training is conducted in `bfloat16` precision using data parallelism. Model selection is based on final validation accuracy.

**Tabular Results Summary.** The following tables summarize primary experimental findings across our study. Each entry corresponds to a core evaluation or analysis in the paper, organized by experimental section and aligned with the corresponding table description.

- *Zero-Shot Evaluation* (Section 4): Table 4 reports ROUGE-F<sub>1</sub> zero-shot performance of pretrained PaliGemma and VideoLLaMA3 models on Action and Character Recognition tasks. Models are evaluated in a zero-shot setting with 1 or 8 input frames, across binary, multiple-choice, and open-ended formats.
- *Fine-Tuned Baselines* (Section 4): Table 5 reports ROUGE-F<sub>1</sub> and Exact Match performance of PaliGemma 2 variants fine-tuned using full, partial, and parameter-efficient strategies. All models are trained on a single frame for one epoch, and evaluated across binary, multiple-choice, and open-ended formats.
- *Analysis: Supervision and Temporal Context* (Section 5): Table 6 examines early-stage learning dynamics on Character Recognition (CR), with evaluation at sub-epoch intervals. Table 7 reports AR performance as a function of training budget, scaling the number of epochs with a single input frame. Table 8 extends this analysis to jointly vary training epochs and the number of input frames, disentangling the effects of temporal context and supervision on AR.
- *Analysis: Temporal Sampling Strategies* (Section 5): Table 9 compares first- $n$  and uniform- $n$  frame sampling strategies for Action Recognition, evaluating model performance across varying temporal context lengths ( $n \in \{1, \dots, 8\}$ ).
- *Analysis: Optimizing Data Mix for Unified Multi-Task Evaluation* (Section 5): This analysis spans three tables. Table 10 explores task-level trade-offs when jointly training on Action and Character Recognition by varying  $\alpha_{AR}$  vs.  $\alpha_{CR}$ , with format distribution held uniform. Table 11 fixes  $\alpha_{AR} = 0.8$  and searches over format-level ratios ( $\beta$ ), revealing the impact of increased open-ended (OE) supervision. Table 12 further investigates this high-OE regime, balancing the remaining budget between binary and multiple-choice for optimal performance.

Table 4: Zero-shot ROUGE-F<sub>1</sub>-based evaluation of PaliGemma (PG) and VideoLLaMA3 (VL3) models on Action and Character Recognition tasks using 1 and 8 input frames. “MC” denotes multiple-choice and “OE” open-ended formats.

Fr	Model	Action Recognition			Character Recognition		
		Binary	MC	OE	Binary	MC	OE
1	PG 1 3B	50.43 ± 0.13	8.12 ± 0.02	10.83 ± 0.01	50.73 ± 0.38	0.46 ± 0.06	0.00 ± 0.00
	PG 2 3B	44.69 ± 0.03	9.30 ± 0.17	12.64 ± 0.01	48.58 ± 0.07	0.28 ± 0.06	0.01 ± 0.00
	PG 2 10B	50.04 ± 0.03	26.98 ± 0.00	12.35 ± 0.21	50.08 ± 0.07	8.33 ± 0.50	0.00 ± 0.00
	VL3-2B	3.24 ± 0.00	18.52 ± 0.06	6.27 ± 0.05	8.76 ± 0.04	3.44 ± 0.08	0.50 ± 0.01
	VL3-7B	45.02 ± 0.28	15.53 ± 0.05	6.54 ± 0.04	39.09 ± 0.73	6.21 ± 0.05	0.51 ± 0.02
8	PG 1 3B	51.67 ± 0.02	10.68 ± 0.00	10.32 ± 0.00	51.39 ± 0.07	0.25 ± 0.00	0.00 ± 0.00
	PG 2 3B	47.61 ± 0.19	6.73 ± 0.04	14.52 ± 0.00	48.37 ± 0.18	0.03 ± 0.00	0.01 ± 0.00
	PG 2 10B	50.02 ± 0.06	26.93 ± 0.01	12.12 ± 0.00	50.09 ± 0.06	0.22 ± 0.00	0.00 ± 0.00
	VL3-2B	13.92 ± 0.13	3.47 ± 0.02	0.32 ± 0.01	13.92 ± 0.13	3.46 ± 0.04	0.32 ± 0.01
	VL3-7B	15.05 ± 0.21	16.67 ± 0.35	6.35 ± 0.06	12.76 ± 0.52	5.88 ± 0.01	0.54 ± 0.01

Table 5: Performance of fine-tuned PaliGemma 2 variants on Action and Character Recognition tasks. We compare full, partial, and parameter-efficient tuning strategies. “MC” denotes multiple-choice and “OE” open-ended formats.

Model	Binary		Multiple-choice		Open-ended	
	EM	ROUGE	EM	ROUGE	EM	ROUGE
<b>Action Recognition</b>						
$\mathcal{F}_L$	50.00 ± 0.00	50.00 ± 0.00	13.13 ± 0.00	27.57 ± 0.00	13.13 ± 0.00	27.57 ± 0.00
$\mathcal{F}_P$	83.97 ± 0.02	83.97 ± 0.02	61.43 ± 0.58	68.05 ± 0.70	61.68 ± 0.35	68.46 ± 0.19
$\mathcal{F}_V$	83.70 ± 0.97	83.70 ± 0.97	63.40 ± 0.45	69.87 ± 0.44	66.03 ± 0.10	71.92 ± 0.08
$\mathcal{F}_{P+L}$	74.47 ± 1.64	74.47 ± 1.64	13.13 ± 0.00	27.57 ± 0.00	55.74 ± 0.70	64.83 ± 0.29
$\mathcal{F}_{V+L}$	75.80 ± 0.16	75.80 ± 0.16	13.13 ± 0.00	27.57 ± 0.00	13.13 ± 0.00	27.57 ± 0.00
$\mathcal{F}_{V+P}$	73.46 ± 0.85	73.46 ± 0.85	61.21 ± 0.23	67.57 ± 0.21	64.70 ± 0.02	70.93 ± 0.01
$\mathcal{F}_{all}$	74.35 ± 1.37	74.35 ± 1.37	13.13 ± 0.00	27.57 ± 0.00	13.13 ± 0.00	27.57 ± 0.00
$\mathcal{F}_{LoRA}$	44.66 ± 0.21	44.66 ± 0.21	0.02 ± 0.01	9.21 ± 0.01	0.00 ± 0.00	12.49 ± 0.00
<b>Character Recognition</b>						
$\mathcal{F}_L$	50.00 ± 0.00	50.00 ± 0.00	98.92 ± 0.00	98.92 ± 0.00	98.98 ± 0.00	98.99 ± 0.01
$\mathcal{F}_P$	99.09 ± 0.11	99.09 ± 0.11	99.22 ± 0.33	99.22 ± 0.33	99.15 ± 0.07	99.15 ± 0.07
$\mathcal{F}_V$	99.31 ± 0.01	99.31 ± 0.01	99.14 ± 0.42	99.14 ± 0.42	99.61 ± 0.12	99.61 ± 0.12
$\mathcal{F}_{P+L}$	50.00 ± 0.00	50.00 ± 0.00	98.28 ± 0.00	98.30 ± 0.02	98.39 ± 0.00	98.39 ± 0.00
$\mathcal{F}_{V+L}$	50.00 ± 0.00	50.00 ± 0.00	96.88 ± 0.00	96.88 ± 0.00	98.45 ± 0.00	98.45 ± 0.00
$\mathcal{F}_{V+P}$	60.32 ± 0.02	60.32 ± 0.02	99.22 ± 0.00	99.22 ± 0.00	99.79 ± 0.00	99.79 ± 0.00
$\mathcal{F}_{all}$	50.00 ± 0.00	50.00 ± 0.00	97.67 ± 0.06	97.67 ± 0.06	96.55 ± 0.01	96.55 ± 0.01
$\mathcal{F}_{LoRA}$	48.76 ± 0.00	48.76 ± 0.00	0.00 ± 0.00	0.32 ± 0.00	0.00 ± 0.00	0.01 ± 0.01

## F HUMAN ANNOTATION STUDY

This section provides full details of our human annotation study, including rollout generation, annotation procedures, inter-annotator agreement, and evaluation metrics. The goal is to validate the adapted VLM’s fine-grained predictions on generated video rollouts.

### F.1 STUDY DESIGN

**Task Overview.** Human annotators were presented with short video clips generated by a world model, each paired with a natural language question and an answer generated by the VLM. They were asked to judge whether the model’s answer accurately described what was shown in the video. Each QA pair was rated using one of four categories: *Correct* (score = 1), *Partially Correct* (0.5), *Incorrect* (0), or *Unclear / Cannot Tell* (excluded from accuracy computation).

Table 6: *Supervision and Temporal Context*: Training budget analysis for Character Recognition. Models are fine-tuned for sub-epoch durations and evaluated across binary, multiple-choice (MC), and open-ended (OE) formats.

Ep	Binary		Multiple-choice		Open-ended	
	EM	ROUGE	EM	ROUGE	EM	ROUGE
0.005	50.84 ± 1.70	50.84 ± 1.70	14.27 ± 0.00	14.27 ± 0.00	13.16 ± 0.00	13.27 ± 0.00
0.01	54.92 ± 0.84	54.92 ± 0.84	17.85 ± 0.00	17.85 ± 0.00	16.78 ± 0.14	16.88 ± 0.01
0.02	57.91 ± 3.19	57.91 ± 3.19	28.64 ± 0.00	28.64 ± 0.00	26.29 ± 2.96	28.38 ± 0.01
0.03	57.45 ± 0.58	57.45 ± 0.58	64.19 ± 0.00	64.19 ± 0.00	40.76 ± 0.01	40.98 ± 0.00
0.06	65.88 ± 3.71	65.88 ± 3.71	93.96 ± 0.00	93.96 ± 0.00	88.89 ± 0.00	88.95 ± 0.01
0.10	87.47 ± 4.11	87.47 ± 4.11	96.54 ± 0.00	96.54 ± 0.00	97.01 ± 0.38	97.02 ± 0.40
0.125	91.63 ± 7.71	91.63 ± 7.71	97.08 ± 0.00	97.08 ± 0.00	97.28 ± 0.00	97.30 ± 0.00
0.20	97.96 ± 0.45	97.96 ± 0.45	97.89 ± 0.28	97.89 ± 0.28	98.12 ± 0.00	98.14 ± 0.02
0.25	97.75 ± 0.74	97.75 ± 0.74	98.08 ± 0.00	98.08 ± 0.00	98.12 ± 0.00	98.15 ± 0.00
0.33	98.42 ± 0.20	98.42 ± 0.20	98.19 ± 0.00	98.19 ± 0.00	98.30 ± 0.00	98.35 ± 0.00
0.50	98.74 ± 0.06	98.74 ± 0.06	98.45 ± 0.00	98.45 ± 0.00	98.51 ± 0.00	98.54 ± 0.00
0.67	99.11 ± 0.10	99.11 ± 0.10	98.99 ± 0.08	98.99 ± 0.08	99.09 ± 0.00	99.09 ± 0.00
0.75	99.03 ± 0.04	99.03 ± 0.04	99.15 ± 0.00	99.15 ± 0.00	99.21 ± 0.00	99.22 ± 0.01
1	99.09 ± 0.11	99.09 ± 0.11	99.22 ± 0.33	99.22 ± 0.33	99.15 ± 0.07	99.15 ± 0.07

Table 7: *Supervision and Temporal Context*: Training budget analysis for Action Recognition, with models fine-tuned for up to 10 epochs. Evaluated using across binary, multiple-choice (MC), and open-ended (OE) formats.

Ep	Binary		Multiple-Choice		Open-Ended	
	EM	ROUGE	EM	ROUGE	EM	ROUGE
1	83.97 ± 0.02	83.97 ± 0.02	61.43 ± 0.58	68.05 ± 0.70	61.68 ± 0.35	68.46 ± 0.19
2	84.92 ± 0.23	84.92 ± 0.23	64.90 ± 0.15	71.17 ± 0.01	64.05 ± 0.01	70.36 ± 0.01
4	85.37 ± 0.30	85.37 ± 0.30	64.58 ± 0.69	70.89 ± 0.55	64.79 ± 0.31	70.88 ± 0.27
8	85.18 ± 0.20	85.18 ± 0.20	63.53 ± 0.35	69.95 ± 0.37	63.43 ± 0.28	69.91 ± 0.18
10	85.11 ± 0.41	85.11 ± 0.41	62.88 ± 1.35	69.34 ± 1.20	66.82 ± 3.75	72.34 ± 3.04



Figure 10: Reference slides shown to annotators during the human annotation study, illustrating the two recognition targets: *actions* (left) and *characters* (center and right). The slides include 20 exemplar videos (7 actions, 13 characters) to support consistent evaluation of VLM-generated responses.

**Annotation Setup and Interface.** Annotations were collected using a custom PowerPoint-based interface (see Figure 11). Each slide presented a short video, a question, and a generated answer. Annotators selected a rating from a predefined rubric. The full annotation guidelines – including action and character definitions and rating instructions – were embedded in the annotation deck for reference. For completeness, we also provide them in Table 13 and Figure 10. The annotation study was carried out by a subset of the authors with prior experience in the environment. Judging correctness required non-trivial familiarity with the visual dynamics and task ontology, making expert annotation necessary. All annotators were compensated above local minimum wage rates. Each QA pair was independently rated by two primary annotators. In cases of disagreement or if either

Table 8: *Supervision and Temporal Context*: Training budget and temporal context analysis for Action Recognition. Models are fine-tuned for up to 10 epochs and evaluated with up to 8 input frames.

Ep	Fr	Binary		Multiple-choice		Open-ended	
		EM	ROUGE	EM	ROUGE	EM	ROUGE
1	1	83.97 ± 0.02	83.97 ± 0.02	61.43 ± 0.58	68.05 ± 0.70	61.68 ± 0.35	68.46 ± 0.19
	2	84.42 ± 0.06	84.42 ± 0.06	65.53 ± 0.27	72.03 ± 0.06	65.38 ± 0.06	71.74 ± 0.23
	4	90.97 ± 0.10	90.97 ± 0.10	83.11 ± 0.08	87.13 ± 0.04	82.26 ± 0.14	87.02 ± 0.06
	8	93.85 ± 0.28	93.85 ± 0.28	88.89 ± 0.14	93.40 ± 0.76	87.80 ± 0.20	92.23 ± 0.16
2	1	85.10 ± 0.02	85.10 ± 0.02	64.93 ± 0.06	71.09 ± 0.11	64.05 ± 0.01	70.36 ± 0.01
	2	86.53 ± 0.45	86.40 ± 0.26	69.20 ± 0.88	75.02 ± 0.67	68.83 ± 0.47	72.45 ± 2.76
	4	92.26 ± 0.34	92.26 ± 0.34	84.19 ± 0.10	88.46 ± 0.11	83.34 ± 0.15	87.84 ± 0.06
	8	95.05 ± 0.15	95.05 ± 0.15	89.27 ± 0.14	93.30 ± 0.22	89.42 ± 0.22	93.25 ± 0.15
4	1	85.37 ± 0.30	85.37 ± 0.30	64.58 ± 0.69	70.89 ± 0.55	64.79 ± 0.31	70.88 ± 0.27
	2	86.89 ± 0.06	86.89 ± 0.06	69.89 ± 0.83	75.49 ± 0.69	70.04 ± 0.08	75.74 ± 0.06
	4	92.58 ± 0.18	92.58 ± 0.18	85.13 ± 0.19	89.28 ± 0.17	84.61 ± 0.07	88.81 ± 0.04
	8	95.29 ± 0.07	95.29 ± 0.07	90.64 ± 0.00	94.09 ± 0.04	90.18 ± 0.16	93.81 ± 0.11
8	1	85.04 ± 0.00	85.04 ± 0.00	63.53 ± 0.35	69.95 ± 0.37	63.62 ± 0.00	70.03 ± 0.00
	2	87.27 ± 0.40	87.27 ± 0.40	69.84 ± 0.66	75.44 ± 0.45	70.11 ± 0.65	75.75 ± 0.52
	4	92.97 ± 0.49	92.97 ± 0.49	85.32 ± 0.08	89.27 ± 0.08	84.93 ± 0.21	89.05 ± 0.11
	8	95.48 ± 0.21	95.48 ± 0.21	90.71 ± 0.14	94.15 ± 0.13	91.02 ± 0.28	93.96 ± 0.71
10	1	85.40 ± 0.00	85.40 ± 0.00	62.88 ± 1.35	69.34 ± 1.20	66.82 ± 3.75	72.34 ± 3.04
	2	87.17 ± 0.22	87.17 ± 0.22	70.18 ± 0.00	75.64 ± 0.00	69.59 ± 0.20	75.37 ± 0.07
	4	92.96 ± 0.37	92.96 ± 0.37	85.02 ± 0.58	89.05 ± 0.49	84.71 ± 0.08	88.94 ± 0.06
	8	96.03 ± 0.05	96.03 ± 0.05	90.75 ± 0.04	94.22 ± 0.05	91.00 ± 0.11	94.33 ± 0.09

Table 9: Comparison of frame sampling strategies for Action Recognition. We evaluate first- $n$  vs. uniform- $n$  sampling across varying temporal context lengths ( $n \in \{1, \dots, 8\}$ ).

Fr	Binary		Multiple-choice		Open-ended		
	EM	ROUGE	EM	ROUGE	EM	ROUGE	
First-N	1	83.97 ± 0.02	83.97 ± 0.02	61.43 ± 0.58	68.05 ± 0.70	61.68 ± 0.35	68.46 ± 0.19
	2	84.42 ± 0.06	84.42 ± 0.06	65.53 ± 0.27	72.03 ± 0.06	65.38 ± 0.06	71.74 ± 0.23
	3	87.93 ± 0.28	87.93 ± 0.28	75.73 ± 0.18	81.07 ± 0.06	74.68 ± 0.16	80.21 ± 0.08
	4	90.97 ± 0.10	90.97 ± 0.10	83.11 ± 0.08	87.13 ± 0.04	82.26 ± 0.14	87.02 ± 0.06
	5	92.00 ± 0.30	92.00 ± 0.30	85.46 ± 0.34	89.84 ± 0.18	85.10 ± 0.16	89.47 ± 0.10
	6	92.95 ± 0.30	92.95 ± 0.30	86.86 ± 0.08	91.13 ± 0.06	86.59 ± 0.39	90.82 ± 0.30
	7	93.31 ± 0.03	93.31 ± 0.03	87.95 ± 0.08	92.06 ± 0.06	87.58 ± 0.17	91.82 ± 0.08
	8	93.85 ± 0.28	93.85 ± 0.28	88.89 ± 0.14	93.40 ± 0.76	87.80 ± 0.20	92.23 ± 0.16
Uniform-N	1	83.97 ± 0.02	83.97 ± 0.02	61.43 ± 0.58	68.05 ± 0.70	61.68 ± 0.35	68.46 ± 0.19
	2	90.47 ± 0.62	90.47 ± 0.62	83.93 ± 0.04	88.36 ± 0.08	82.68 ± 0.19	87.33 ± 0.01
	3	93.59 ± 0.07	93.59 ± 0.07	88.90 ± 0.11	92.85 ± 0.10	88.49 ± 0.24	92.57 ± 0.04
	4	93.57 ± 0.39	93.57 ± 0.39	89.94 ± 0.04	93.65 ± 0.01	89.56 ± 0.42	93.49 ± 0.28
	5	94.25 ± 0.04	94.25 ± 0.04	89.99 ± 0.18	93.70 ± 0.13	89.72 ± 0.10	93.63 ± 0.23
	6	94.01 ± 0.57	94.01 ± 0.57	90.03 ± 0.04	93.73 ± 0.06	90.09 ± 0.18	93.88 ± 0.16
	7	93.96 ± 0.16	93.96 ± 0.16	90.34 ± 0.23	94.00 ± 0.10	89.94 ± 0.11	93.73 ± 0.10
	8	94.62 ± 0.48	94.62 ± 0.48	90.72 ± 0.12	94.30 ± 0.10	90.30 ± 0.04	94.01 ± 0.02

annotator marked the example as *Unclear*, a third, more experienced adjudicator reviewed the pair and assigned a final rating.

**Selected World Models.** For our study, we aim to evaluate rollouts generated across different model scales, training diversities, and output resolutions, while keeping the underlying architecture general enough to apply broadly. To this end, we select two autoregressive world models (Kanervisto et al., 2025). The autoregressive formulation offers a flexible and widely adopted framework, and is the basis for many state-of-the-art private world models.

Table 10: *Optimizing Data Mix for Unified Multi-Task Evaluation*: Performance tradeoffs under varying task-level allocation ratios for Action ( $\alpha_{AR}$ ) vs. Character Recognition ( $\alpha_{CR}$ ), with a fixed format distribution ( $\beta = 1/3$  per format). Evaluated across binary, multiple-choice (MC), and open-ended (OE) formats.

$\alpha_{AR}$	$\alpha_{CR}$	Binary		Multiple-choice		Open-ended		
		EM	ROUGE	EM	ROUGE	EM	ROUGE	
<b>Action Recognition</b>								
0.20	0.80	84.13 $\pm$ 1.66	84.13 $\pm$ 1.66	26.10 $\pm$ 0.43	39.04 $\pm$ 0.10	27.03 $\pm$ 0.21	39.60 $\pm$ 0.13	
0.40	0.60	88.11 $\pm$ 1.44	88.11 $\pm$ 1.44	28.13 $\pm$ 1.15	40.93 $\pm$ 0.57	29.17 $\pm$ 0.40	41.66 $\pm$ 0.51	
0.50	0.50	88.59 $\pm$ 1.41	88.59 $\pm$ 1.41	29.10 $\pm$ 0.65	41.20 $\pm$ 0.16	29.66 $\pm$ 0.54	41.17 $\pm$ 0.22	
0.60	0.40	90.80 $\pm$ 0.04	90.80 $\pm$ 0.04	30.44 $\pm$ 0.03	42.54 $\pm$ 0.01	30.55 $\pm$ 0.49	42.32 $\pm$ 0.60	
0.80	0.20	91.23 $\pm$ 0.91	91.23 $\pm$ 0.91	84.06 $\pm$ 1.26	89.42 $\pm$ 1.00	30.88 $\pm$ 0.63	42.85 $\pm$ 0.42	
<b>Character Recognition</b>								
0.20	0.80	98.57 $\pm$ 0.47	98.57 $\pm$ 0.47	98.95 $\pm$ 0.16	98.95 $\pm$ 0.16	98.94 $\pm$ 0.22	98.97 $\pm$ 0.21	
0.40	0.60	98.51 $\pm$ 0.53	98.51 $\pm$ 0.53	98.77 $\pm$ 0.16	98.77 $\pm$ 0.16	98.98 $\pm$ 0.06	98.98 $\pm$ 0.06	
0.50	0.50	96.33 $\pm$ 1.81	96.33 $\pm$ 1.81	98.03 $\pm$ 0.06	98.03 $\pm$ 0.06	98.23 $\pm$ 0.06	98.23 $\pm$ 0.06	
0.60	0.40	93.22 $\pm$ 3.38	93.22 $\pm$ 3.38	96.91 $\pm$ 1.81	96.94 $\pm$ 1.77	97.93 $\pm$ 0.02	97.93 $\pm$ 0.02	
0.80	0.20	80.53 $\pm$ 0.49	80.53 $\pm$ 0.49	89.08 $\pm$ 0.49	89.08 $\pm$ 0.49	89.02 $\pm$ 0.25	89.18 $\pm$ 0.39	

Table 11: *Optimizing Data Mix for Unified Multi-Task Evaluation*: Performance on Action and Character Recognition under varying format-level sampling ratios ( $\beta$ ) for Binary, Multiple-choice (MC), and Open-ended (OE) questions. We fix  $\alpha_{AR} = 0.8$  and train all models on the first 8 frames.

Ep	$\beta_{Binary}$	$\beta_{MC}$	$\beta_{OE}$	Binary		Multiple-Choice		Open-Ended	
				EM	ROUGE	EM	ROUGE	EM	ROUGE
<b>Action Recognition</b>									
1	0.4	0.2	0.4	92.32 $\pm$ 0.37	92.32 $\pm$ 0.37	84.32 $\pm$ 0.89	89.54 $\pm$ 0.64	31.61 $\pm$ 0.23	43.51 $\pm$ 0.09
	0.2	0.4	0.4	90.80 $\pm$ 0.71	90.80 $\pm$ 0.71	86.60 $\pm$ 0.38	91.27 $\pm$ 0.42	32.06 $\pm$ 0.18	43.58 $\pm$ 0.54
	0.0	0.4	0.6	49.98 $\pm$ 0.04	49.98 $\pm$ 0.04	86.65 $\pm$ 0.99	91.26 $\pm$ 0.87	85.51 $\pm$ 1.78	90.13 $\pm$ 1.53
	0.0	0.2	0.8	50.11 $\pm$ 0.06	50.11 $\pm$ 0.06	86.58 $\pm$ 0.47	91.42 $\pm$ 0.21	87.45 $\pm$ 0.09	91.78 $\pm$ 0.09
2	0.4	0.2	0.4	93.14 $\pm$ 0.48	93.14 $\pm$ 0.48	86.77 $\pm$ 0.00	91.28 $\pm$ 0.00	32.96 $\pm$ 0.00	44.00 $\pm$ 0.00
	0.2	0.4	0.4	92.89 $\pm$ 0.16	92.89 $\pm$ 0.16	87.83 $\pm$ 0.00	92.13 $\pm$ 0.00	33.37 $\pm$ 0.00	44.13 $\pm$ 0.00
	0.0	0.4	0.6	41.22 $\pm$ 0.00	41.30 $\pm$ 0.01	89.17 $\pm$ 0.07	93.12 $\pm$ 0.02	88.55 $\pm$ 0.00	93.65 $\pm$ 0.00
	0.0	0.2	0.8	49.98 $\pm$ 0.03	49.99 $\pm$ 0.02	88.68 $\pm$ 0.00	92.71 $\pm$ 0.00	88.59 $\pm$ 0.00	92.56 $\pm$ 0.00
4	0.4	0.2	0.4	94.33 $\pm$ 0.34	94.33 $\pm$ 0.34	92.67 $\pm$ 0.07	93.27 $\pm$ 0.71	33.94 $\pm$ 0.01	43.96 $\pm$ 0.00
	0.2	0.4	0.4	94.19 $\pm$ 0.13	94.19 $\pm$ 0.13	93.04 $\pm$ 0.01	93.57 $\pm$ 0.74	33.78 $\pm$ 0.00	44.73 $\pm$ 0.00
	0.0	0.4	0.6	50.19 $\pm$ 0.27	50.19 $\pm$ 0.27	89.78 $\pm$ 0.11	93.52 $\pm$ 0.03	88.95 $\pm$ 0.05	92.75 $\pm$ 0.06
	0.0	0.2	0.8	49.98 $\pm$ 0.02	49.98 $\pm$ 0.02	89.25 $\pm$ 0.00	93.13 $\pm$ 0.00	89.41 $\pm$ 0.00	93.16 $\pm$ 0.00
<b>Character Recognition</b>									
1	0.4	0.2	0.4	86.42 $\pm$ 0.25	86.42 $\pm$ 0.25	94.77 $\pm$ 0.14	94.77 $\pm$ 0.14	94.76 $\pm$ 0.08	94.63 $\pm$ 0.21
	0.2	0.4	0.4	77.57 $\pm$ 0.01	77.57 $\pm$ 0.01	94.93 $\pm$ 0.03	94.93 $\pm$ 0.03	94.51 $\pm$ 0.57	94.15 $\pm$ 0.00
	0.0	0.4	0.6	50.37 $\pm$ 0.52	50.37 $\pm$ 0.52	96.56 $\pm$ 0.28	96.56 $\pm$ 0.28	96.81 $\pm$ 0.39	96.82 $\pm$ 0.37
	0.0	0.2	0.8	50.51 $\pm$ 0.03	50.51 $\pm$ 0.03	96.88 $\pm$ 0.27	96.88 $\pm$ 0.27	97.39 $\pm$ 0.18	97.39 $\pm$ 0.18
2	0.4	0.2	0.4	89.95 $\pm$ 0.46	89.95 $\pm$ 0.46	87.35 $\pm$ 0.00	87.36 $\pm$ 0.00	88.64 $\pm$ 0.00	88.64 $\pm$ 0.00
	0.2	0.4	0.4	91.28 $\pm$ 0.20	91.28 $\pm$ 0.20	93.90 $\pm$ 0.00	93.93 $\pm$ 0.04	93.06 $\pm$ 0.00	93.09 $\pm$ 0.00
	0.0	0.4	0.6	47.37 $\pm$ 0.02	47.48 $\pm$ 0.04	97.69 $\pm$ 0.00	97.70 $\pm$ 0.00	98.07 $\pm$ 0.00	98.07 $\pm$ 0.00
	0.0	0.2	0.8	50.07 $\pm$ 0.04	50.07 $\pm$ 0.04	97.75 $\pm$ 0.00	97.79 $\pm$ 0.00	98.07 $\pm$ 0.00	98.07 $\pm$ 0.00
4	0.4	0.2	0.4	97.71 $\pm$ 0.05	97.71 $\pm$ 0.05	97.70 $\pm$ 0.00	97.70 $\pm$ 0.00	97.88 $\pm$ 0.01	97.91 $\pm$ 0.03
	0.2	0.4	0.4	96.55 $\pm$ 0.06	96.55 $\pm$ 0.06	98.51 $\pm$ 0.00	98.51 $\pm$ 0.00	98.46 $\pm$ 0.00	98.46 $\pm$ 0.00
	0.0	0.4	0.6	51.25 $\pm$ 1.56	51.25 $\pm$ 1.56	98.21 $\pm$ 0.13	98.21 $\pm$ 0.13	98.48 $\pm$ 0.06	98.49 $\pm$ 0.06
	0.0	0.2	0.8	50.93 $\pm$ 0.10	50.93 $\pm$ 0.10	98.79 $\pm$ 0.00	98.79 $\pm$ 0.00	99.08 $\pm$ 0.00	99.09 $\pm$ 0.00

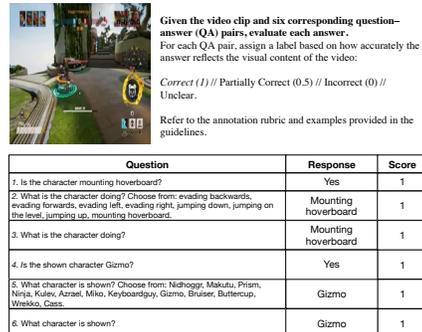
The selected models generate sequences of visual frames and controller actions without textual supervision, using a decoder-only transformer (Radford et al., 2019; Vaswani et al., 2017) trained autoregressively on discrete tokens. Visual frames are encoded with a VQGAN (Esser et al., 2021), while joystick actions are tokenized using a learned discretization scheme based on action bucketization (Kanervisto et al., 2020).

Table 12: *Optimizing Data Mix for Unified Multi-Task Evaluation: Performance on Action and Character Recognition under high open-ended (OE) supervision, with  $\beta_{OE} = 0.8$  and remaining budget split between Binary and Multiple-choice (MC).*

Ep	$\beta_{BN}$	$\beta_{MC}$	$\beta_{OE}$	Binary		Multiple-Choice		Open-Ended	
				EM	ROUGE	EM	ROUGE	EM	ROUGE
<b>Action Recognition</b>									
1	0.15	0.05	0.80	88.85 $\pm$ 0.04	88.85 $\pm$ 0.04	81.93 $\pm$ 0.00	89.08 $\pm$ 0.00	86.72 $\pm$ 0.00	91.33 $\pm$ 0.00
	0.10	0.10	0.80	87.38 $\pm$ 0.59	87.38 $\pm$ 0.59	85.58 $\pm$ 0.00	90.50 $\pm$ 0.00	86.88 $\pm$ 0.00	91.25 $\pm$ 0.00
	0.05	0.15	0.80	85.67 $\pm$ 0.19	85.67 $\pm$ 0.19	86.38 $\pm$ 0.09	91.21 $\pm$ 0.06	86.84 $\pm$ 0.02	91.34 $\pm$ 0.03
2	0.15	0.05	0.80	92.45 $\pm$ 0.11	92.45 $\pm$ 0.11	87.52 $\pm$ 0.00	91.90 $\pm$ 0.00	87.97 $\pm$ 0.63	92.19 $\pm$ 0.41
	0.10	0.10	0.80	92.11 $\pm$ 0.18	92.11 $\pm$ 0.18	88.42 $\pm$ 0.00	92.54 $\pm$ 0.00	88.50 $\pm$ 0.00	92.54 $\pm$ 0.00
	0.05	0.15	0.80	91.98 $\pm$ 0.24	91.98 $\pm$ 0.24	88.72 $\pm$ 0.00	92.78 $\pm$ 0.00	88.56 $\pm$ 0.00	92.66 $\pm$ 0.00
4	0.15	0.05	0.80	92.98 $\pm$ 0.21	92.98 $\pm$ 0.21	88.93 $\pm$ 0.00	93.02 $\pm$ 0.00	89.64 $\pm$ 0.00	93.34 $\pm$ 0.00
	0.10	0.10	0.80	92.81 $\pm$ 0.11	92.81 $\pm$ 0.11	91.40 $\pm$ 2.81	93.88 $\pm$ 0.70	89.43 $\pm$ 0.00	93.20 $\pm$ 0.00
	0.05	0.15	0.80	91.52 $\pm$ 0.37	91.52 $\pm$ 0.37	89.80 $\pm$ 0.00	93.54 $\pm$ 0.01	89.81 $\pm$ 0.01	93.49 $\pm$ 0.06
<b>Character Recognition</b>									
1	0.15	0.05	0.80	59.75 $\pm$ 0.04	59.75 $\pm$ 0.04	95.45 $\pm$ 0.00	95.45 $\pm$ 0.00	97.16 $\pm$ 0.00	97.16 $\pm$ 0.00
	0.10	0.10	0.80	56.31 $\pm$ 0.53	56.31 $\pm$ 0.53	94.55 $\pm$ 0.00	94.55 $\pm$ 0.00	95.96 $\pm$ 0.00	95.96 $\pm$ 0.00
	0.05	0.15	0.80	50.86 $\pm$ 0.08	50.86 $\pm$ 0.08	96.87 $\pm$ 0.02	96.87 $\pm$ 0.02	97.12 $\pm$ 0.01	97.12 $\pm$ 0.01
2	0.15	0.05	0.80	80.18 $\pm$ 0.09	80.18 $\pm$ 0.09	95.41 $\pm$ 0.00	95.41 $\pm$ 0.00	96.91 $\pm$ 0.00	96.91 $\pm$ 0.00
	0.10	0.10	0.80	70.20 $\pm$ 0.21	70.20 $\pm$ 0.21	98.02 $\pm$ 0.00	98.02 $\pm$ 0.00	98.15 $\pm$ 0.00	98.15 $\pm$ 0.00
	0.05	0.15	0.80	69.67 $\pm$ 0.36	69.67 $\pm$ 0.36	97.37 $\pm$ 0.00	97.37 $\pm$ 0.00	97.79 $\pm$ 0.00	97.79 $\pm$ 0.00
4	0.15	0.05	0.80	94.16 $\pm$ 0.12	94.16 $\pm$ 0.12	97.06 $\pm$ 0.00	97.07 $\pm$ 0.00	97.91 $\pm$ 0.00	97.91 $\pm$ 0.00
	0.10	0.10	0.80	86.67 $\pm$ 0.13	86.67 $\pm$ 0.13	98.50 $\pm$ 0.00	98.50 $\pm$ 0.00	98.77 $\pm$ 0.00	98.77 $\pm$ 0.00
	0.05	0.15	0.80	71.79 $\pm$ 0.01	71.79 $\pm$ 0.01	98.22 $\pm$ 0.00	98.22 $\pm$ 0.00	98.57 $\pm$ 0.00	98.57 $\pm$ 0.00

Each world model is implemented as a decoder-only transformer (Radford et al., 2019; Vaswani et al., 2017), trained to predict discrete tokens representing visual observations and actions. Visual frames are first compressed with a VQGAN (Esser et al., 2021), while joystick actions are tokenized using a learned discretization scheme based on action bucketization (Kanervisto et al., 2020). The objective is next-token prediction conditioned on prior visual and action tokens. We focus on the following model variants that differ in capacity, training diversity, and output resolution: (i) a large-scale model: 1.6B parameters, trained for 200K steps on gameplay from seven diverse environments (including Skygarden) at  $300 \times 180$ , and (ii) a smaller model: 140M parameters, trained for 100K steps on gameplay from a single environment (Skygarden) at  $128 \times 128$  resolution.

**Rollouts Generation.** Rollout generation follows a consistent protocol for both world models: at inference time, the model is conditioned on 1 second of ground-truth gameplay (visual and action tokens), after which it generates 10 seconds of future gameplay conditioned only on a sequence of held-out controller actions. For rollout generation, we use reference trajectories from the game dataset, which include time-aligned controller inputs and image frames. For each generated rollout, we randomly sample a 1-second context window (video frames and corresponding action tokens) from these trajectories. The world model is then conditioned on this context, which consists of an interleaved sequence of tokenized images and actions, and proceeds autoregressively to generate latent image tokens and discretized action tokens. The generated rollout is then split into 14-frame chunks. This setup enables a comprehensive analysis of the UNIVERSE’s evaluation capabilities across two axes: (i) *in-domain performance*: evaluating on Skygarden, the environment used for fine-tuning; (ii) *generalization*: assessing performance on six unseen environments. It also allows comparison across generation quality and model capacity. We generate 82 rollouts for each model-environment setting, resulting in 656 rollouts in total.



Given the video clip and six corresponding question-answer (QA) pairs, evaluate each answer. For each QA pair, assign a label based on how accurately the answer reflects the visual content of the video:

Correct (1) // Partially Correct (0.5) // Incorrect (0) // Unclear.

Refer to the annotation rubric and examples provided in the guidelines.

Question	Response	Score
1. Is the character mounting hoverboard?	Yes	1
2. What is the character doing? Choose from: evading backwards, evading forwards, evading left, evading right, jumping down, jumping on the level, jumping up, mounting hoverboard.	Mounting hoverboard	1
3. What is the character doing?	Mounting hoverboard	1
4. Is the shown character Gizmo?	Yes	1
5. What character is shown? Choose from: Nidhoggr, Makulu, Prism, Niqa, Kulu, Ansel, Nika, Keybladejay, Gizmo, Brusar, Buttercup, Wreiko, Cass.	Gizmo	1
6. What character is shown?	Gizmo	1

Figure 11: Annotation interface example. Each instance includes a video clip, task instructions, and a table with: *Question* (generated via evaluation protocol), *Response* (VLM output), and *Score* (human-assigned label).

1890 Table 13: Annotation instructions provided to human raters as part of the study. The interface  
 1891 outlines task context, scoring criteria, general guidelines, and reference definitions for supported  
 1892 action categories.

1893	
1894	<b>1. Task Overview</b>
1895	You will be presented with:
1896	• A short video clip;
1897	• A natural language question about the video;
1898	• An answer generated by a vision-language model.
1899	Your task is to evaluate whether the model’s answer accurately describes the events depicted in the
1900	video.
1901	<b>2. How to Rate Each Answer</b>
1902	Assign one of the following categories:
1903	• <i>Correct (1.0)</i> : Fully matches the event in the video;
1904	• <i>Partially Correct (0.5)</i> : Captures the general idea but contains a minor error;
1905	• <i>Incorrect (0.0)</i> : Wrong, hallucinated, or mismatched with the visual evidence;
1906	• <i>Unclear / Cannot Tell</i> : Not enough evidence to confidently decide.
1907	<b>3. General Guidelines</b>
1908	• Watch the full video before rating;
1909	• Base your decision solely on visible content;
1910	• Use provided action and character references;
1911	• If multiple plausible interpretations exist and the answer matches one, mark as <i>Correct</i> ;
1912	• If unsure even after review, mark <i>Unclear / Cannot Tell</i> ;
1913	• Optionally leave comments for ambiguous or interesting cases.
1914	<b>5. Action Label Definitions</b>
1915	• <i>Evading Backwards</i> : Moves backwards to avoid threat or reposition.
1916	• <i>Evading Forwards</i> : Moves forwards.
1917	• <i>Evading Left / Right</i> : Lateral movement left or right.
1918	• <i>Jumping Down</i> : Jumps from a higher to a lower platform or level.
1919	• <i>Jumping on the Level</i> : Jumps without elevation change.
1920	• <i>Jumping Up</i> : Jumps upward to reach a higher platform.
1921	• <i>Mounting Hoverboard</i> : Begins riding or is seen riding a hoverboard.



1922 (i) Early, Uninformative Segments



1928 (ii) Rollouts Dominated by Occlusion



1934 (iii) Sequence with No Visible Agents

1940 Figure 12: Randomly sampled examples of rollouts excluded from the human evaluation study.

1942  
1943

**Rollout Filtering.** To keep evaluation focused on challenging sequences, we remove only those

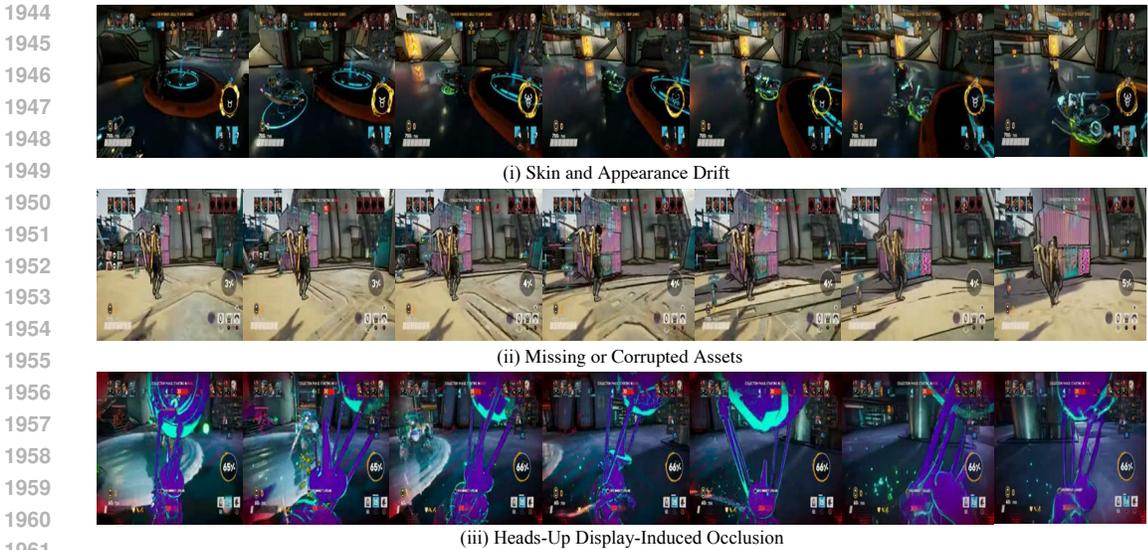


Figure 13: Randomly sampled examples of rollouts retained for human evaluation.

rollouts that are either uninterpretable or uninformative. Specifically, annotators discard: (i) early segments with no meaningful interactions; (ii) rollouts dominated by occlusion; and (iii) sequences with no visible agents. These criteria target cases where neither humans nor models can extract actionable dynamics, while retaining all visually imperfect, ambiguous, and failure-prone generations. Thus, filtering eliminates only sequences that provide no evaluable signal. Figure 12 shows randomly sampled examples from each excluded category.

**Complexity of Retained Rollouts.** Despite filtering out uninterpretable/uninformative sequences, many retained rollouts exhibit visual imperfections that make evaluation challenging. Annotators highlighted recurring generation failures including: (i) skin and appearance drift, where model-generated agent skins diverge from canonical silhouettes; (ii) missing or corrupted assets, such as hoverboards implied only through motion cues; and (iii) HUD-induced occlusion, where oversized or misplaced UI elements block important scene content. These imperfections arise from the underlying world models and introduce ambiguity while still permitting meaningful assessment. Figure 13 illustrates representative failure-prone cases, emphasizing that UNIVERSE is evaluated on realistic, imperfect rollouts rather than sanitized inputs.

**UNIVERSE’s Response Generation.** To obtain responses from UNIVERSE, we provide it with a video segment (resized to match the evaluator’s input resolution) along with its corresponding question. We then sample five responses using greedy decoding and select the most frequent response as the final answer. In cases where all five responses are unique (i.e., no majority), one response is selected at random. The resulting dataset comprises rollouts from 8 settings: rollouts generated by a smaller model on Skygarden, and rollouts generated by a larger model across seven distinct environments. For each model–environment pair, we sample 30 rollouts. Each rollout is annotated with 6 question–answer (QA) pairs, along with a corresponding response from the adapted evaluator. Each of the resulting 1,440 QA instances was rated by 3 annotators, yielding 4,320 total human judgments.

## F.2 EVALUATION METRICS

We report two accuracy-based metrics using the adjudicated labels:

*Strict Accuracy.*: The proportion of QA pairs labeled as *Correct*:

$$Acc_{\text{Strict}} = \frac{N_{\text{Correct}}}{N_{\text{Answerable}}}, \tag{5}$$

Table 14: Inter-annotator agreement and valid QA coverage across environments. We report Cohen’s  $\kappa$  between the two primary annotators for each world model–map pair. The total number of valid examples excludes QA pairs marked as *Unclear* by at least one annotator.

Setting	Valid QA Pairs	Cohen’s $\kappa$
1	29	0.91
2	28	0.67
3	28	0.74
4	29	0.87
5	30	0.67
6	29	0.61
7	30	0.59
8	24	0.79

Table 15: Graded/strict accuracy of UNIVERSE on Action and Character Recognition tasks, evaluated by human annotators across different environments and question formats. We report results for Binary, Multiple-Choice (MC), and Open-Ended (OE) prompts, disaggregated by task and world model. All metrics are based on final adjudicated ratings.

Setting	Binary	Action Recognition			Character Recognition	
		MC	OE	Binary	MC	OE
1	98.3 / 96.7	51.7 / 46.7	75.0 / 73.3	93.3 / 93.3	83.3 / 83.3	93.3 / 93.3
2	96.7 / 96.7	60.0 / 60.0	65.0 / 60.0	99.9 / 99.9	90.0 / 90.0	93.3 / 93.3
3	96.7 / 96.7	63.3 / 63.3	80.0 / 80.0	99.9 / 99.9	86.7 / 86.7	93.3 / 93.3
4	93.3 / 93.3	43.3 / 43.3	73.3 / 73.3	96.7 / 96.7	96.7 / 96.7	99.9 / 99.9
5	80.0 / 76.7	76.7 / 73.3	93.3 / 93.3	96.7 / 96.7	99.9 / 99.9	99.8 / 99.8
6	71.7 / 70.0	56.7 / 56.7	75.0 / 70.0	96.7 / 96.7	93.3 / 93.3	96.7 / 96.7
7	68.3 / 66.7	50.0 / 46.7	80.0 / 76.7	93.3 / 93.3	90.0 / 90.0	96.7 / 96.7
8	92.9 / 89.3	35.7 / 32.1	41.1 / 39.3	85.7 / 85.7	10.7 / 10.7	60.7 / 60.7

*Graded Accuracy.*: Partial credit given to *Partially Correct* responses:

$$\text{Acc}_{\text{Graded}} = \frac{N_{\text{Correct}} + 0.5 \times N_{\text{Partial}}}{N_{\text{Answerable}}}. \quad (6)$$

Only examples not marked *Unclear* by adjudication are included in  $N_{\text{Answerable}}$ .

*Inter-Annotator Agreement.* To quantify rating consistency, we compute Cohen’s  $\kappa$  between the two primary annotators. The adjudicator’s label is used only when disagreement occurs and is excluded from agreement computation. Results are shown in Table 14.

*Sample Size Justification.* We annotate 30 rollouts per model–environment pair. Assuming a standard deviation of  $\sigma \approx 0.2$  and a 95% confidence level, the confidence interval (CI) width is given by  $\text{CI Width} = z_{1-\frac{C}{2}} \cdot \frac{\sigma}{\sqrt{n}}$ . This yields an estimated CI of  $\sim 7.1\%$  for individual model–environment pairs ( $n = 30$ ), and  $\sim 2.5\%$  when aggregating across all eight pairs ( $n = 240$ ), offering sufficient precision for comparative evaluation.

### F.3 RESULTS

Table 15 reports graded and strict accuracy across environments, recognition targets (Action and Character Recognition), and question formats (Binary, Multiple-Choice, Open-Ended). We observe a clear gap in performance between rollouts generated by the two world models. UNIVERSE struggles with outputs from WHAM 140M, achieving substantially lower accuracy compared to WHAM 1.6B. This is likely due to a mismatch in image resolution: WHAM 140M generates frames at  $128 \times 128$  resolution, which must be upsampled to the UNIVERSE’s expected input of  $224 \times 224$ . Despite resizing, the resulting frames often lack sharpness, making actions and characters harder to recognize. In

contrast, UNIVERSE performs well on rollouts from WHAM 1.6B, even across diverse environments. On the in-domain setting (Environment A), the model achieves strong results—averaging 75.02% graded accuracy for AR and 90.00% for CR. When evaluating on the six unseen environments (Environments B–G), performance for AR drops slightly (from 75.02% to 73.52%), while CR remains stable or improves, suggesting strong generalization in character grounding and visual consistency tracking.

**Qualitative Examples.** Figure 14 illustrates the diversity of generated rollouts across environments. WHAM 1.6B captures greater visual variation and scene composition compared to WHAM 140M.

## G SUPPLEMENTARY EXPERIMENTAL RESULTS

This section presents additional experimental results that support the main findings but are omitted from the main paper for clarity and space. These include: (i) a zero-shot analysis of PaliGemma variants to motivate backbone selection, (ii) CLIPScore-based baselines to contextualize performance without adaptation, and (iii) a study of low-rank adaptation (LoRA) across different rank values. While these results are not central to the unified evaluation framework proposed in the main text, they provide valuable insight into model selection, adaptation efficiency, and the limitations of standard evaluation proxies in our setting.

### G.1 GPT-5 PERFORMANCE

To demonstrate the complexity of the tasks that comprise our protocol, we conducted an evaluation of GPT-5 on randomly selected examples. We deliberately chose the simplest evaluation regime (binary action recognition) to test whether the model can succeed without adaptation.

Figure 15 demonstrates that out of six random samples, GPT-5 produced incorrect answers in five cases. This consistent failure highlights the difficulty of the task.

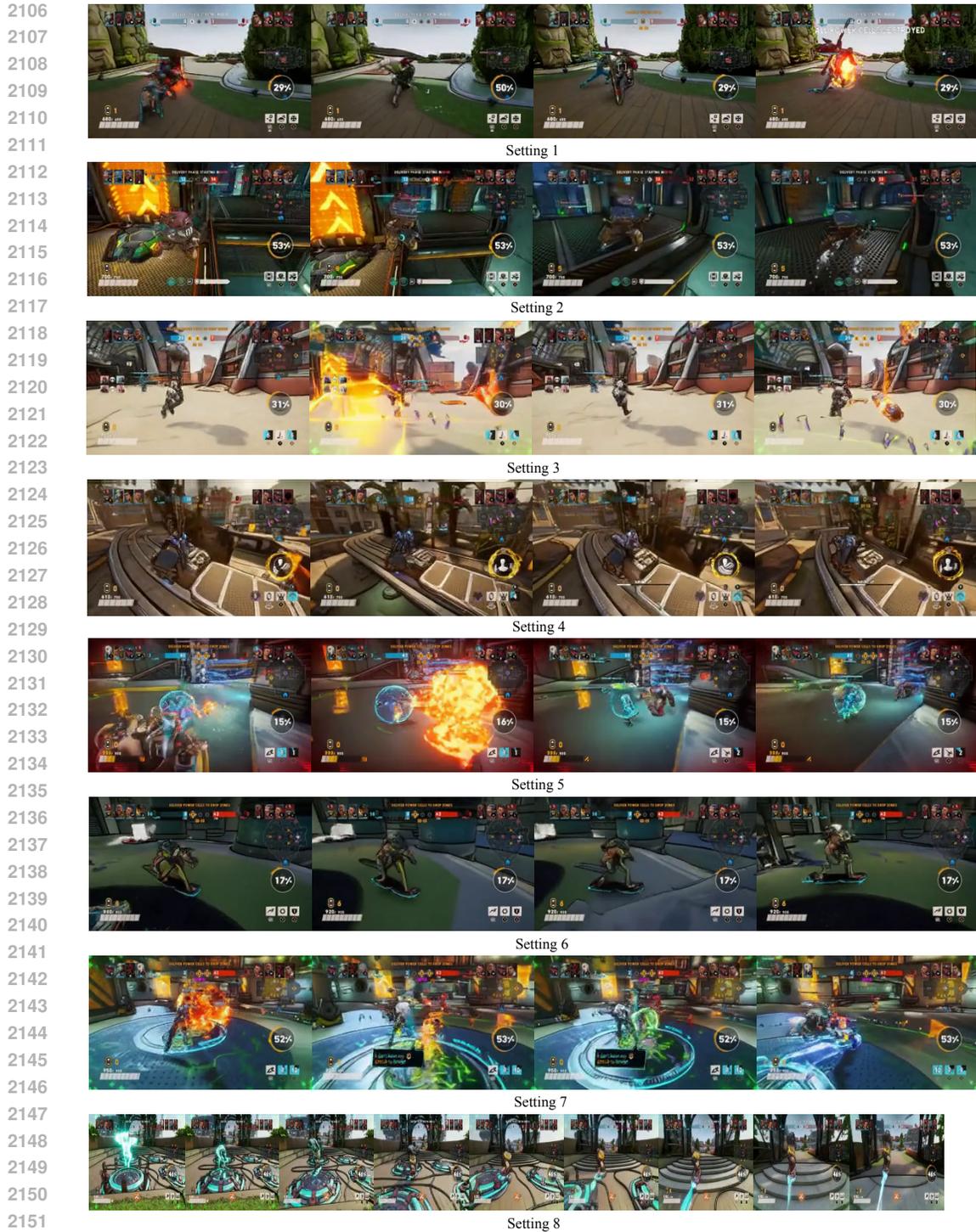
### G.2 ZERO-SHOT PERFORMANCE OF PALIGEMMA MODELS

In this section, we benchmark three pretrained configurations—PaliGemma 1 3b, PaliGemma 2 3b, and PaliGemma 2 10b—under our proposed protocol and motivate our choice of PaliGemma 2 3b as the default backbone for subsequent experiments. Each model receives a natural language prompt along with either 1 or 8 image frames as input and produces a textual response. This experiment probes both model capacity and the role of temporal visual context in zero-shot settings.

**Results.** Figure 16 reports ROUGE scores across task types, question formats, and visual context lengths. While zero-shot performance reveals some capacity for structured reasoning—particularly in the multiple-choice setting—it remains limited overall. Binary accuracy hovers near chance, and open-ended responses frequently lack specificity. Performance is strongest on action recognition (AR), likely reflecting pretrained models’ familiarity with generic visual dynamics. In contrast, character recognition (CR) lags behind, underscoring a lack of grounding in domain-specific entities. Increasing the number of input frames modestly improves AR, but yields diminishing returns for CR. Among the evaluated configurations, PaliGemma 2 10b performs best in absolute terms. However, the margin over PaliGemma 2 3b is narrow, and PaliGemma 2 3b offers a substantially smaller footprint while using a newer Gemma 2 decoder architecture. We therefore adopt PaliGemma 2 3b as the default model for all subsequent adaptation experiments, balancing performance, compute efficiency, and architectural recency.

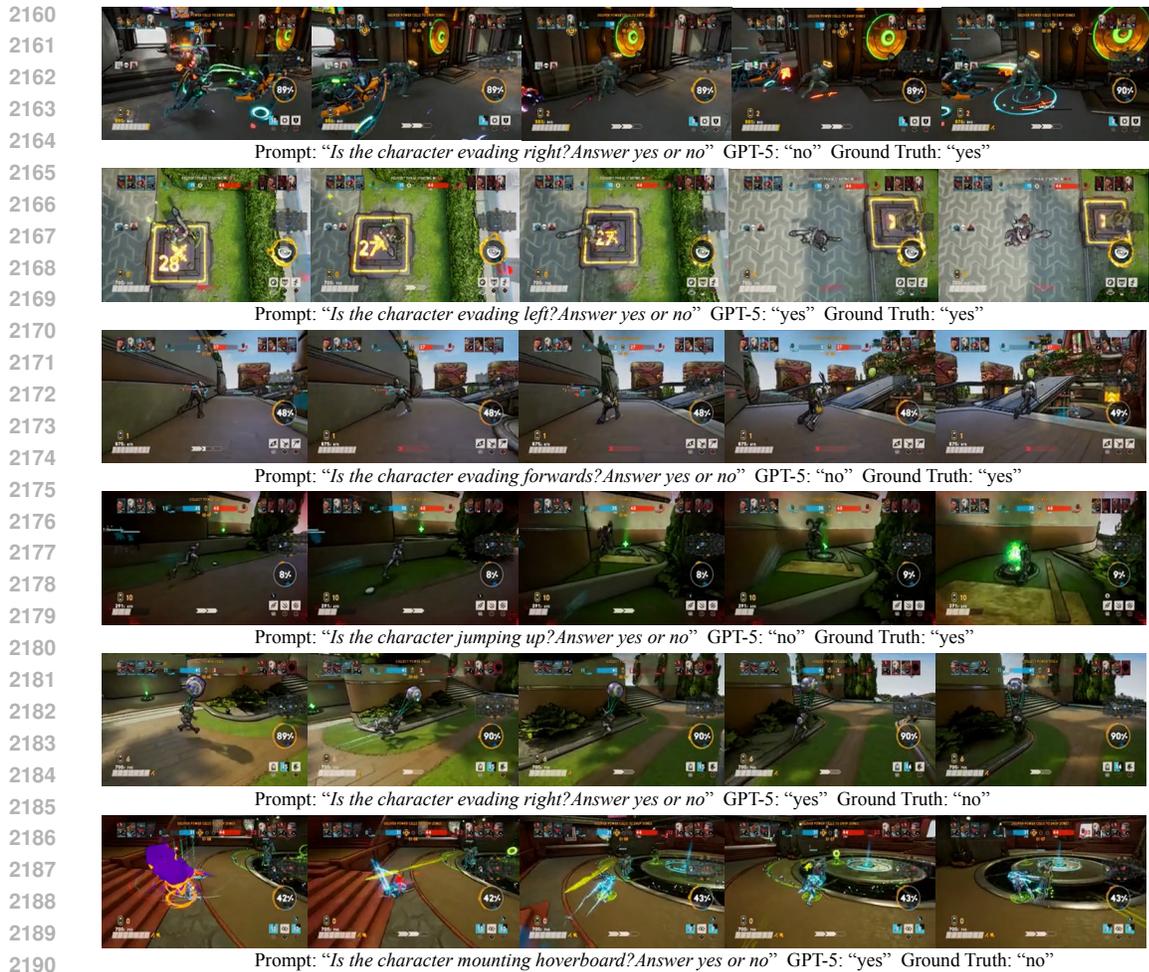
### G.3 CLIPSCORE COMPARISONS

To further evaluate zero-shot recognition capabilities without adaptation, we apply CLIPScore to our rollout evaluation protocol. Specifically, we assess four pretrained CLIP variants – ViT-B/32, ViT-B/16, ViT-L/14, and ViT-L/14-336 – across both Action Recognition (AR) and Character Recognition (CR) tasks using 1-frame and 8-frame visual inputs. For each evaluation instance, we extract either 1 or 8 frames from the video segment and compute the cosine similarity between each image and a predefined set of textual labels (i.e., action verbs for AR, character names for CR). For single-frame settings, we select the label with the highest similarity score as the predicted class. In the multi-frame

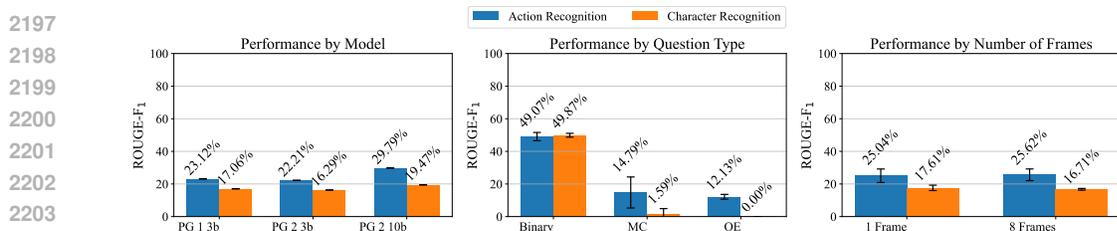


2153 Figure 14: Representative frames from rollouts across the eight evaluation settings, spanning different  
2154 environments, scales, and resolutions.

2155  
2156  
2157 setting, we compute predictions for each frame independently and use a majority vote to produce  
2158 the final prediction. We also report two reference baselines for context: a random classifier, which  
2159 achieves 12.5% on AR and 7.7% on CR, and a majority-class predictor, which yields 35.5% and  
17.6% respectively. These are included only for calibration.



2192 Figure 15: Performance of GPT-5 on six randomly sampled binary action recognition questions. Each  
2193 panel shows the trial prompt, the model’s response, and the ground-truth label. Despite the apparent  
2194 simplicity of the setup, GPT-5 fails on 5/6 trials, underscoring the difficulty of the task and the need  
2195 for task-specific adaptation.



2205 Figure 16: Zero-shot evaluation results for PaliGemma variants across tasks, prompt formats, and  
2206 visual context sizes. Overall performance remains limited, indicating the need for task-specific  
2207 adaptation.

2210 **Results.** Table 16 demonstrates the results. While CLIP ViT-B/16 performs relatively well on AR in  
2211 both input settings, performance remains inconsistent across model scales and tasks. In particular, CR  
2212 accuracy remains low, reflecting CLIP’s limited grounding in domain-specific visual semantics and  
2213 fine-grained identity resolution. Larger CLIP models such as ViT-L/14 do not consistently outperform  
smaller variants, and 8-frame inputs provide only marginal gains over single-frame inputs.

Table 16: Zero-shot accuracy-based evaluation of CLIP models and baseline methods on Action and Character Recognition tasks using 1 and 8 input frames.

Fr	Model	Action Recognition	Character Recognition
1	CLIP ViT-B/32	24.04 ± 0.00	13.32 ± 0.00
	CLIP ViT-B/16	52.67 ± 0.00	16.47 ± 0.00
	CLIP ViT-L/14	24.60 ± 0.00	9.95 ± 0.00
	CLIP ViT-L/14-336	12.17 ± 0.00	8.85 ± 0.05
8	CLIP ViT-B/32	36.22 ± 0.00	14.41 ± 0.00
	CLIP ViT-B/16	57.36 ± 0.00	17.24 ± 0.00
	CLIP ViT-L/14	17.57 ± 0.00	10.10 ± 0.00
	CLIP ViT-L/14-336	23.12 ± 0.00	8.64 ± 0.00

Table 17: Performance on Action and Character Recognition tasks after LoRA-based adaptation with varying ranks ( $r \in \{8, 16, 32, 48, 64\}$ ). Adapters are applied to attention and MLP layers in both vision and language components.

Rank	Binary		Multiple-choice		Open-ended	
	EM	ROUGE	EM	ROUGE	EM	ROUGE
<b>Action Recognition</b>						
8	44.66 ± 0.21	44.66 ± 0.21	0.02 ± 0.00	9.21 ± 0.00	0.00 ± 0.00	12.49 ± 0.00
16	44.47 ± 0.43	44.47 ± 0.43	0.02 ± 0.00	9.21 ± 0.00	0.00 ± 0.00	12.49 ± 0.00
32	44.59 ± 0.03	44.59 ± 0.03	0.02 ± 0.00	9.21 ± 0.00	0.00 ± 0.00	12.49 ± 0.00
48	46.71 ± 3.20	46.71 ± 3.20	0.02 ± 0.00	9.21 ± 0.00	0.00 ± 0.00	12.49 ± 0.00
64	48.67 ± 0.13	48.67 ± 0.13	0.02 ± 0.00	9.21 ± 0.00	0.00 ± 0.00	12.49 ± 0.00
<b>Character Recognition</b>						
8	48.76 ± 0.00	48.76 ± 0.00	0.00 ± 0.00	0.32 ± 0.00	0.00 ± 0.00	0.01 ± 0.01
16	48.62 ± 0.23	48.62 ± 0.23	0.00 ± 0.00	0.14 ± 0.00	0.00 ± 0.00	0.05 ± 0.00
32	48.98 ± 0.08	48.98 ± 0.08	0.00 ± 0.00	0.14 ± 0.00	0.00 ± 0.00	0.05 ± 0.00
48	48.91 ± 0.09	48.91 ± 0.09	0.00 ± 0.00	0.14 ± 0.00	0.00 ± 0.00	0.05 ± 0.00
64	48.72 ± 0.06	48.72 ± 0.06	0.00 ± 0.00	0.14 ± 0.00	0.00 ± 0.00	0.05 ± 0.00

Overall, these results suggest that while CLIPScore offers a lightweight and scalable evaluation proxy, it lacks the temporal grounding and semantic specificity required for structured rollout evaluation. Performance falls short relative to our selected baselines, and the method is inherently constrained to predefined candidate sets—limiting its applicability to open-ended or compositional tasks. As such, we exclude CLIP-based scores from our primary comparisons and instead focus on adapted, generative VLM-based evaluators.

#### G.4 LOW-RANK ADAPTATION COMPARISONS

This section presents an extended analysis of low-rank adaptation (LoRA) as a parameter-efficient strategy for adapting vision-language models to our protocol. We systematically vary the rank parameter  $r$  and measure its impact on Action and Character Recognition performance across all prompt formats. All experiments in this section are conducted using PaliGemma 2 (3B) as the backbone model, consistent with the main fine-tuning results. These experiments assess whether increasing rank provides meaningful gains, and inform our decision to report only the rank-8 setting in the main paper.

**Results.** Table 17 presents the performance of LoRA-based adaptation across a range of rank values ( $r \in \{8, 16, 32, 48, 64\}$ ) for both Action Recognition (AR) and Character Recognition (CR) tasks, across all prompt formats. We report exact match (EM) and ROUGE-F<sub>1</sub> averaged over three runs. Increasing the rank beyond  $r = 8$  yields no consistent improvements across tasks or formats. Performance on binary prompts remains close to random, while performance on multiple-choice and

Table 18: Pearson correlation between FVD and UNIVERSE across experimental settings. None of the correlations are statistically significant (all  $p > 0.05$ ).

Setting	Action Recognition	Character Recognition
1	0.09	0.03
2	-0.07	-0.08
3	-0.17	0.06
4	-0.03	-0.10
5	0.07	-0.25
6	-0.01	0.21
7	0.32	0.03
8	-0.10	0.33

open-ended formats stays near zero across all ranks. These results suggest that LoRA, even with increased capacity, is insufficient for capturing the fine-grained temporal and semantic dependencies required by our evaluation protocol. Given the lack of benefit from increasing rank—and the added parameter cost—it is inefficient to scale LoRA rank beyond  $r = 8$ . Accordingly, all results reported in the main paper use  $r = 8$ , while extended comparisons with higher ranks are presented here for completeness.

## H COMPARISON WITH EXISTING EVALUATION BASELINES

To validate that UNIVERSE captures evaluation dimensions beyond existing automated metrics, we analyze its correlation with two representative baselines: Fréchet Video Distance (FVD), a standard distributional metric for video generation quality, and VBench, a comprehensive multi-dimensional benchmark for text-to-video evaluation. These comparisons assess whether UNIVERSE’s focus on action-conditioned semantic alignment provides complementary signal to metrics designed primarily for perceptual quality and temporal coherence.

### H.1 COMPARISON WITH FVD

FVD Unterthiner et al. (2018) is a widely used reference-free metric for evaluating generative video models, defined as a distance between feature distributions of real and generated videos. Although commonly treated as a default quantitative proxy for rollout quality in world-model evaluation, FVD is agnostic to task semantics and action alignment—the dimensions UNIVERSE is specifically designed to capture. Across the eight environments, raw FVD values are tightly clustered (means between 3.23 and 3.27 with standard deviations below 0.05), reflecting limited dynamic range and strong dataset-dependent scale effects. As a result, raw scores are not directly comparable across settings, making correlations a more informative basis for analysis. We therefore assess the relationship between the two metrics by computing Pearson correlations between FVD and UNIVERSE’s human-aligned scores across all environments for both Action Recognition and Character Recognition tasks.

**Results.** For each environment (Settings 1–8) and each task, we obtain FVD scores on the same set of rollouts evaluated by UNIVERSE, report the scores and compute the correlation between the two. The resulting coefficients are reported in Table 18. Correlations are small in magnitude and unstable in sign, ranging from  $-0.25$  to  $0.33$ , with most values close to zero. The largest positive correlations are modest (Action Recognition, Setting 7:  $r = 0.32$ ; Character Recognition, Setting 8:  $r = 0.33$ ), and several settings exhibit weak negative correlations (e.g., Character Recognition, Setting 5:  $r = -0.25$ ). None of these correlations are statistically significant ( $p > 0.05$  for all settings), indicating UNIVERSE captures complementary semantic cues. These results demonstrate that FVD and UNIVERSE provide complementary evaluation signals: FVD captures distribution-level perceptual quality and temporal coherence, while UNIVERSE quantifies task-specific semantic alignment through human-grounded assessment.

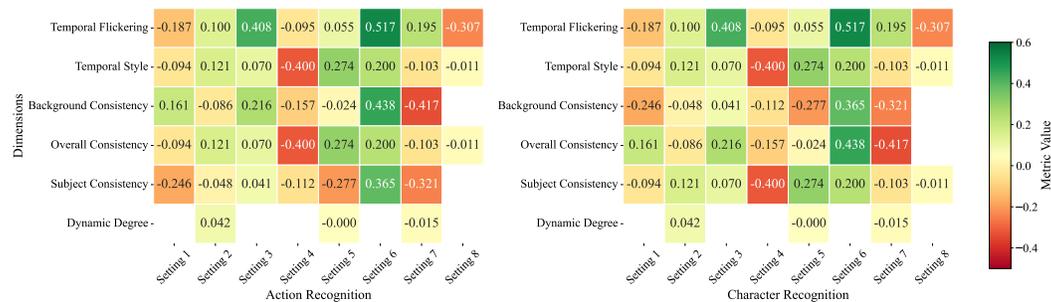


Figure 17: Correlation between UNIVERSE and VBench evaluation suite across eight environments and six dimensions. The *human action* dimension is omitted as it was constant and non-informative. Correlations remain consistently low to moderate ( $|r| < 0.4$ ), with none of the correlations are statistically significant (all  $p > 0.05$ ).

## H.2 CORRELATION WITH VBENCH METRICS

To validate that UNIVERSE captures evaluative dimensions beyond existing automated benchmarks, we compare its human-aligned evaluation scores against VBench Huang et al. (2024), a comprehensive state-of-the-art benchmark for generative video evaluation. We computed correlations on identical rollout sets across eight experimental environments using six evaluation configurations, parameterized by task (Action Recognition or Character Recognition) and prompt type (binary, multiple-choice, open-ended). The analysis spans seven VBench dimensions: Temporal Flickering, Temporal Style, Background Consistency, Overall Consistency, Subject Consistency, Dynamic Degree, and Human Action. The Human Action dimension is omitted from visualization as it exhibited no variance across samples.

**Results.** Figure 17 shows that correlations between UNIVERSE and VBench metrics remain consistently low to moderate ( $|r| < 0.4$ ) across all valid comparisons. Average correlations range from 0.037 for Subject Consistency to 0.06 for Temporal Flickering, with most values falling within 0.1–0.4. Notably, the sign variability across settings (including both positive and negative correlations) indicates that UNIVERSE captures complementary semantic cues not reflected in VBench’s quality dimensions. This divergence is expected: while VBench focuses on low-level perceptual quality and temporal coherence, UNIVERSE explicitly targets semantic alignment.

**Summary.** The consistent orthogonality observed across both analyses demonstrates that UNIVERSE provides a complementary evaluation dimension to existing benchmarks. While standard distributional metrics (FVD) capture visual fidelity and text-to-video benchmarks (VBench) assess perceptual quality, UNIVERSE quantifies action-conditioned semantic alignment through human-grounded evaluation. Comprehensive world model assessment therefore benefits from combining these approaches: existing metrics for perceptual coherence and UNIVERSE for task-relevant semantic fidelity.