

# Bootstrapping Action-Grounded Visual Dynamics in Unified Vision–Language Models

Anonymous ACL submission

## Abstract

Can unified vision–language models (VLMs) perform *forward dynamics prediction* (FDP), i.e., predicting the future state (in image form) given the previous observation and an action (in language form)? We find that VLMs struggle to generate physically plausible transitions between frames from instructions. Nevertheless, we identify a crucial asymmetry in multi-modal grounding: fine-tuning a VLM to learn *inverse dynamics prediction* (IDP)—effectively captioning the action between frames—is significantly easier than learning FDP. In turn, IDP can be used to bootstrap FDP through two main strategies: 1) weakly supervised learning from synthetic data and 2) inference time verification. Firstly, IDP can annotate actions for unlabelled pairs of video frame observations to expand the training data scale for FDP. Secondly, IDP can assign rewards to multiple samples of FDP to score them, effectively guiding search at inference time. We evaluate the FDP resulting from both strategies through the task of *action-centric image editing* on AURORA-BENCH with two families of VLMs. Despite remaining general-purpose, our best model achieves a performance competitive with state-of-the-art image editing models, improving on them by a margin between 7% and 13% according to GPT4o-as-judge, and achieving the best average human evaluation across all subsets of AURORA-BENCH.<sup>1</sup>

## 1 Introduction

*World models* are instrumental in training embodied agents to endow them with specific abilities, such as planning and simulation (Qin et al., 2024; Brohan et al., 2023; Huang et al., 2022; Li et al., 2025; Reed et al., 2022a; Yang et al., 2023; Hafner et al., 2025, *inter alia*). However, learning a specialised world model is challenging due to the scarcity of

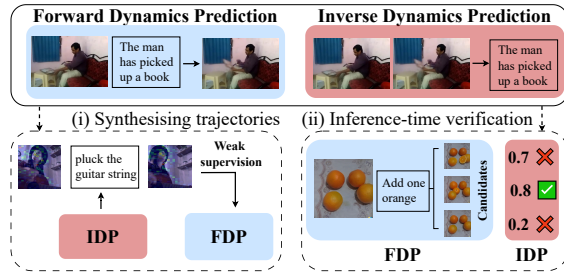


Figure 1: A high-level illustration of our two strategies to bootstrap Forward Dynamics Prediction from Inverse Dynamics Prediction in unified Vision–Language Models: (i) synthesising trajectories for weak supervision (left) and (ii) inference-time verification of candidate future observations (right).

real-world data (Liu et al., 2024; Motamed et al., 2025). Conversely, a promising alternative is endowing pre-existing unified<sup>2</sup> vision-language models (VLMs) with world modelling abilities. In fact, VLMs are already imbued with plenty of real-world knowledge of both action (in language form) and perception (in vision form), because of their large-scale pre-training.

Firstly, we probe whether unified VLMs already contain reliable forward dynamics models (FDP), i.e., the ability to predict the next image frame given the previous image frame and an action expressed as a language instruction. For this assessment, we limit ourselves to single-step trajectories, as a first step towards longer horizons. Based on our evaluation, we empirically demonstrate that existing VLMs do not prefer ground-truth trajectories compared to adversarially generated ones. Hence, we verify that the world model implicit in the original VLMs *per se* is not well grounded on real-world forward dynamics (Gao et al., 2024; Qiu et al., 2024; Abdou et al., 2021).

Surprisingly, we also find that predicting the ac-

<sup>1</sup>The code and models developed in this paper will be made available at [anonymised].

<sup>2</sup>Here ‘unified’ means models capable of interleaving text and images during generation architecturally.

Forward-dynamics Modelling										Inverse-dynamics Modelling									
Rand. Act.	36.7	36.7	31.0	43.5	48.8	56.5	29.8	40.3	44.8	Rand. Act.	58.9	59.7	56.5	60.1	55.6	56.9	59.3	58.1	55.6
Inv. Obs.	53.2	50.4	52.4	55.2	54.4	46.8	46.4	48.4	52.0	Inv. Obs.	53.6	54.8	55.2	50.4	49.6	47.6	48.4	50.8	50.8
Copy. Obs.	54.0	56.9	53.6	81.0	48.8	31.9	74.2	66.9	100.0	Copy. Obs.	54.0	64.5	64.9	67.3	48.4	30.6	55.6	55.2	42.7
Rand. Obs.	42.7	38.3	42.7	36.3	51.6	58.9	35.9	39.9	46.4	Rand. Obs.	47.6	50.4	42.3	43.1	52.0	58.1	54.4	50.4	58.1
	Qwen2-VL-2B	Qwen2-VL-7B	Qwen2.5-VL-3B	Qwen2.5-VL-7B	LLaVA-Next-7B	LLaVA-Interleave	Qwen-Omni-3B	Qwen-Omni-7B	Chameleon-7B		Qwen2-VL-2B	Qwen2-VL-7B	Qwen2.5-VL-3B	Qwen2.5-VL-7B	LLaVA-Next-7B	LLaVA-Interleave	Qwen-Omni-3B	Qwen-Omni-7B	Chameleon-7B

Figure 2: Percentage of times 9 VLMs assign higher probability to observation–action–observation Reference trajectories compared with 4 types of Negative (i.e., adversarially manipulated) trajectories, for both forward dynamics and inverse dynamics prediction. Higher values are better.

tion  $a \in \mathcal{A}$  taking place between observations  $o \in \mathcal{O}$ —also known as *inverse dynamics prediction* (IDP;  $\mathcal{O}_t \times \mathcal{O}_{t+1} \rightarrow \mathcal{A}_t$ )—via supervised fine-tuning is substantially easier than FDP ( $\mathcal{O}_t \times \mathcal{A}_t \rightarrow \mathcal{O}_{t+1}$ ). Inspired by this finding, we propose two strategies to bootstrap the FDP from the IDP in a unified VLM, namely (i) **learning from synthetic trajectories** in videos automatically labelled with actions by the dynamics model; and (ii) **test-time verification** of predicted observations sampled from the FDP through the IDP. Figure 1 illustrates our two strategies.

For the weak supervision strategy, we use a VLM fine-tuned for IDP on the AURORA dataset (Krojer et al., 2024) to annotate with linguistic actions ( $\hat{a} \in \mathcal{A}$ ) motion key-frame pairs extracted from 45 hours of unlabelled real-world videos. These are sourced from movements-in-time (Monfort et al., 2019), Kinetics700 (Kay et al., 2017; Carreira et al., 2019) and UCF-101 (Soomro et al., 2012). Together with the ground-truth trajectories in AURORA, the synthesised trajectory triplets are then used for weakly supervised fine-tuning of the very same VLM for FDP. To encourage training to focus on image regions that change the most, we additionally propose a *loss-weighting method* which weights the loss of each image token according to the visual difference between the ground-truth previous and next observations, as estimated by a recognition model. Instead, in the verification strategy, we use the IDP to assign scores to multiple candidate samples generated by the FDP. Selecting the highest-score prediction can effectively guide search at inference time.

We conduct a thorough evaluation of both strategies to bootstrap FDP from IDP on several datasets from AURORA-BENCH (Krojer et al., 2024): MagicBrush, Action-Genome, Something-Something,

WhatsUp and Kubric. We experiment with two families of unified VLMs for FDP, Chameleon-7B and Liquid-8B. Training on trajectories synthesised by IDP, our FDP can achieve an overall performance superior to state-of-the-art diffusion models specialised for image editing, both in terms of GPT4o-as-a-judge and human evaluation. Inference-time verification can also improve FDP to a comparable degree as trajectory synthesis, providing an effective training-free bootstrapping method.

While limited to single-step action–observation trajectories, this work offers promising early evidence that unified VLMs can be successfully endowed with FDP and hence may be suitably developed into long-horizon world models in the future.

## 2 Unified VLMs Lack a Consistent Preference for Real-World Trajectories

To understand whether off-the-shelf VLMs are already suitable for FDP, the first research question we investigate in this paper is: **To what extent do VLMs exhibit a preference for sequences of actions and observations that align with real-world trajectories?**

To address this question, we evaluate 9 different VLMs on ground-truth trajectories from 5 subsets of AURORA-BENCH (Krojer et al., 2024): MagicBrush, Something-Something, Action-Genome, Whatsup, and Kubric. Each subset contains 50 trajectory triplets of the form  $(o_s, a, o_t)$ , where  $o_s$  is the source observation (image frame),  $a$  the action (text), and  $o_t$  the target observation.<sup>3</sup>

We then manually curate four types of negative trajectories using rule-based manipulations. The first kind is an action-level manipulation. 1) **Ran-**

<sup>3</sup>We choose these 9 VLMs using the following criteria: 1) they are publicly accessible and 2) they have been exposed to interleaved multimodal data during their pre-training.

**dom Action:** for a given pair of observations, we substitute the original action with another randomly sampled within the same subset. We also perform three observation-level manipulations. 2) **Random Observation:** we randomly substitute the target observation with another in the same subset. 3) **Copy Observation:** we copy the source observation as the target observation. 4) **Inverse Observation:** we swap the source and target observations.

In Figure 2, we compare the log-likelihood that VLMs assign to each ground-truth trajectory (Reference) against its corresponding manipulated one (Negatives). We evaluate the VLMs in two tasks: action prediction (i.e., as an inverse-dynamics model) and next-observation prediction (i.e., as a forward-dynamics model). For each kind of negative trajectory, we report the percentage of samples where the model favours the reference trajectory over the negative trajectory. From Figure 2, it emerges that VLMs display no clear preference for the ground-truth trajectories in a zero-shot setting (around 50%).

In the action prediction task (right panel), there is a slightly higher tendency to favour the ground-truth over the group with random actions; however, even in the best case, Qwen2.5-VL-7B prefers the reference in only 60.08% of the samples. The only negative group that seems to be identifiable for VLMs is observation copying, where Qwen2.5-VL-7B has 67.34% of correct preference. In the next-observation prediction task (left panel), the VLM mostly fails in effectively differentiating the ground truth from the negatives. Although the underlying reason remains uncertain, one plausible explanation for this behaviour is that the model’s ability to solve next-observation prediction tasks depends on their alignment with training sequences: for instance, it is plausible that Chameleon’s data rarely features two identical adjacent images. We provide a discussion breaking down Chameleon’s performance in Appendix D.

### 3 Bootstrapping FDP from IDP

Overall, the results from Section 2 reveal that off-the-shelf VLMs are not suitable for FDP and IDP *per se*; however, they also show that IDP is a more feasible task, as VLMs already achieve accuracies above random chance in a zero-shot setting. Possibly, IDP abilities may be even improved with a small amount of fine-tuning. This underlies the key intuition behind our work: **can we leverage**

### this asymmetry to bootstrap forward-dynamics abilities from inverse-dynamics ones within a unified VLM?

We propose two strategies to leverage VLMs fine-tuned for IDP to enhance their own FDP abilities: (i) generating synthetic trajectories by annotating large-scale key-frame pairs from videos with actions predicted by IDP, then using these as weak supervision to train for FDP (Section 3.2); and (ii) using the IDP as a verifier at test time to score candidate next observations sampled from the FDP (Section 3.3).

#### 3.1 Inverse-dynamics Modelling

First, we fine-tune the unified VLM as an inverse-dynamics model (**IDM**)  $p_{\text{IDM}}(a \mid o_s, o_t)$ , which predicts the probability of an action given the previous and next observations. As training data, we rely on high-quality triplets from AURORA (Krojer et al., 2024) and the action recognition track of EPIC-Kitchen (Damen et al., 2018), which is based on videos with an egocentric view. We use the first and last frame in the EPIC-Kitchen video clips as the source and target observation  $o_s$  and  $o_t$  and the annotated action as  $a$ . We provide full details on IDM training data and setting in Appendix H.1.

#### 3.2 Weakly Supervised Learning from Unlabelled Videos

**Synthetic Trajectories.** Taking advantage of the resulting high-quality IDM, we then explore the first of our strategies to bootstrap the FDP in VLMs: we annotate pairs of motion key-frames of unlabelled videos with a textual description of the action with the IDM. To ensure both scale and quality, we collect approximately 45 hours of video from Moments-in-Time (Monfort et al., 2019), Kinetics-700 (Kay et al., 2017; Carreira et al., 2019), and UCF-101 (Soomro et al., 2012), all of which consist of curated clips focused on human actions. To ensure the selected pairs of motion key-frames are meaningful, i.e., they express a valid action, we then calculate the optical flow to quantify the dynamics per frame in the video clips, and select the top- $K_f$  frames while ensuring that the interval between two selected frames is  $I_f$ . Specifically, we set  $I_f = 20$  and  $K_f = 6$  for all three datasets. This results in approximately 20K, 46K, and 21K annotated trajectory triplets from Moments-in-Time, Kinetics-700, and UCF-101, respectively. Finally, we apply a filtering strategy to further guarantee the quality of the resulting triplets. Specifically, we

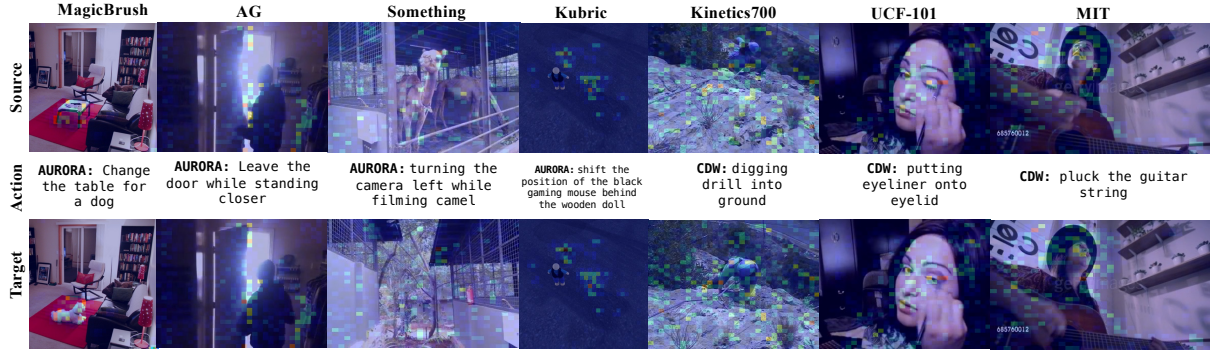


Figure 3: Heatmap visualization of image token weights predicted by the recognition model on examples from AG, Something-Something, MagicBrush, and Kubric, and UCF-101, Kinetics700 and MIT.

apply stratified top- $k$  sampling based on the IDM’s predicted likelihood for each trajectory triplet  $(o_s, \hat{a}_{\text{IDM}}, o_t)$ <sup>4</sup> to select a subset of triplets. We show statistics of the scores and action classes for the selected triplets in Figure 9. We also provide one example for each dataset in Figure 3.

**FDP.** Finally, we fine-tune a VLM as a forward-dynamics model (**FDM**),  $p_{\text{FDM}}(o_t | a, o_s)$  on both AURORA’s supervised triplets  $\mathcal{D}_{\text{sup}}$  and unsupervised triples  $\mathcal{D}_{\text{unsup}}$  with actions sampled from the IDM. Normally, a FDM would be trained with maximum likelihood estimation as an objective:

$$\min_{\theta} \mathbb{E}_{(a, o_s, o_t) \sim \mathcal{D}_{\text{sup}}} [-\log p_{\theta}(o_t | a, o_s)] + \mathbb{E}_{(o_s, o_t) \sim \mathcal{D}_{\text{unsup}}} [\mathbb{E}_{\hat{a} \sim p_{\text{IDM}}(a | o_s, o_t)} [-\log p_{\theta}(o_t | \hat{a}, o_s)]] \quad (1)$$

where  $\theta$  are the parameters for FDM, and  $\hat{a}$  is an action sampled from the IDM.

**Recognition-Weighted Training Loss.** Nevertheless, the objective in Equation 1 is limited by treating all regions (i.e., image patches) of the target observation equally, even if some of them remain identical to the source whereas others change. This may result in degenerate solutions such as always copying the source. As an alternative, we therefore propose a novel training objective for FDM that overcomes this assumption. This objective weights the loss of next-observation image tokens based on their importance. The intuition is that not all image patches in source and target observations contribute equally to modelling real-world transitions; instead, the model should focus on patches most indicative of the action’s consequences. To this end, we leverage a recognition model  $f_{\text{rec}}(w; o_s, o_t)$ , which outputs token-level

weights that represent the similarity between source and target patches. These weights modulate the loss during training, emphasising learning on semantically meaningful regions and down-weighting irrelevant ones. We formulate our objective as:

$$\min_{\theta} \sum_{l=1}^L f_{\text{rec}}(w; o_s, o_t)^{(l)} \cdot \left( -\log p_{\theta}(o_t^{(l)} | o_t^{(<l)}, o_s, a) \right), \quad (2)$$

where  $\theta$  are the parameters of VLM as FDM and a  $L$  is the number of image tokens.  $o_t^{(l)}$  and  $o_t^{(<l)}$  represent the image tokens of  $o_t$  at position  $l$  and the history of previous positions, respectively. For simplicity, we use the pre-trained vector-quantised encoder of the unified VLMs as the recognition model, by computing the squared  $L_2$  norm of pre-quantized features  $\mathbf{z}_{o_s} \in Z_{o_s}$  and  $\mathbf{z}_{o_t} \in Z_{o_t}$  where  $Z_{o_s}$  and  $Z_{o_t}$  are the sets of features of source and target observations, respectively. We visualise the token weights in Figure 3, which capture where acting on the source observation yields the largest effects on the target one.

### 3.3 Test-time Verification

Finally, we introduce an inference-time strategy which harnesses the IDM as a verifier to enhance FDM performance. Inspired by recent work on scaling test-time compute (Muennighoff et al., 2025; Snell et al., 2024), we let the FDM generate  $N$  candidate observations. Each candidate is paired with the source and scored by the IDM, which assigns each a predicted likelihood, interpreted as a reward. The final prediction of the FDM is selected by maximising the IDM’s reward:

$$\hat{o}_t = \operatorname{argmax}_{i \in \{1, \dots, N\}} p_{\text{IDM}}(a | o_s, o_t^{(i)}),$$

where  $o_t^{(i)} \sim p_{\text{FDM}}(o_t | o_s, a)$ .

<sup>4</sup>The details of this algorithm are provided in Appendix E.

where  $\hat{o}_t$  is the selected prediction. While this strategy is training-free, it requires sampling multiple candidates at inference time.

## 4 Experiments and Results

In Section 4.1, we first describe the experimental setup, including benchmarks, baselines, and evaluation metrics. We then report results on the inverse dynamics prediction task in Section 4.2 to verify that fine-tuning on a limited amount of examples is sufficient to develop robust IDP in a VLM. This is followed in Section 4.3 by automatic and human evaluation of both strategies introduced in Section 3 to bootstrap FDP, as well as ablation studies. Finally, we analyse inference-time verification and examine the transfer of forward dynamics prediction to two spatial reasoning benchmark.

### 4.1 Experimental Setting

**Benchmarks.** We select AURORA-BENCH (Krojer et al., 2024) for evaluation of both IDM and FDM. This dataset provides high-quality data for action-centric edits, covering a wide array of phenomena and assessing a model’s alignment with the physical world. We choose 5 subsets: **MagicBrush** for specialised image editing, **Action-Genome (AG)** and **Something-Something (Something)** for real-world actions. **Whatsup** focuses on spatial reasoning, whereas **Kubric** contains samples from a physical engine (Greff et al., 2022).

**Models and Baselines.** We experiment with two unified VLMs, Chameleon-7B (Chameleon Team, 2024) and Liquid-8B (Wu et al., 2024b). The variants of these models fine-tuned on both supervised and weakly supervised data with loss weighting (i.e., with the training-time bootstrapping strategy) are indicated as **C-FDM** and **L-FDM**, respectively. As our first baselines, we use the same two VLMs either zero-shot (**C-ZS** and **L-ZS**) or fine-tuned *only* on AURORA’s supervised data (**C-FT** and **L-FT**). Additionally, we include three state-of-the-art diffusion models specialised for image editing as baselines, such as **PixInstruct** (Brooks et al., 2023), **GoT** (Fang et al., 2025) and **SmartEdit** (Huang et al., 2024). Finally, as a sanity check, we also report the metric scores obtained by simply copying the source observation input as the next-observation prediction (**Copy**).

**Metrics.** For next-observation prediction (FDP) evaluation, following Fang et al. (2025), we rely

	BS	R-1	R-L	BLEU
Chameleon ZS	0.05	0.09	0.08	0.00
Chameleon IDM	<b>0.45</b>	<b>0.45</b>	<b>0.44</b>	0.20
Liquid IDM	0.41	0.41	0.36	<b>0.21</b>

Table 1: Performance of Inverse Dynamics Models on action prediction, measured by text similarity metrics: BERTScore (BS; Zhang et al. 2020), ROUGE (R-1, R-L; Lin 2004) and BLEU (Papineni et al., 2002).

on **GPT4o-as-a-judge** as it is the only metric that reliably penalises Copy. In Appendix B, we show four other metrics, e.g., CLIP, which assign high scores to Copy. GPT4o-as-a-judge scores consider two criteria, one for the editing success rate and one for visual consistency with the original. We take the minimum of the two as the final score. The prompt for the judge is provided in Appendix F.

### 4.2 Inverse Dynamics Prediction

We evaluate the IDM based on the textual similarity of the predicted action with the ground-truth action in AURORA-BENCH. Results are shown in Table 1. Our results demonstrate that a moderate amount of fine-tuning is necessary to let unified VLMs verbalise the dynamics connecting two observations, as evidenced by the large It is also worth noting that both models’ IDP performance (in terms of BLEU) is similar. We leverage these IDP versions of Chameleon and Liquid to bootstrap FDP, whose results are reported in Section 4.3.

### 4.3 Forward Dynamics Prediction

Next, we report FDP results for each of the two bootstrapping strategies: in Section 4.3.1 for trajectory synthesis and in Section 4.3.2 for inference-time verification.

#### 4.3.1 Synthesising Trajectories with IDP

**Automatic evaluation.** To test trajectory synthesis, we evaluate FDM on next-observation prediction for each of the AURORA-BENCH subsets, reporting GPT4o-as-a-judge scores in Table 2. We first notice that the state-of-the-art image editing models (i.e., PixInstruct, GoT, SmartEdit) tend to specialise on the proper image editing subset MagicBrush (5.96 and 6.71 GPT4o scores for GoT and SmartEdit). Nevertheless, in the action-centric subsets, including Action-Genome (AG), Something and Kubric, they are mostly behind bootstrapped models (C/L-FDM) and even the fine-tuned baselines (C/L-FT). Crucially, considering the average performance of C/L-FDM on all subsets of

Dataset	Models								
	PixInst	GoT	SE	C-ZS	C-FT	C-FDM	L-ZS	L-FT	L-FDM
<i>Forward-dynamics Prediction (GPT4o score)</i>									
<b>MagicBrush</b>	3.12	5.96	<b>6.71</b>	0.00	2.52	3.92	0.68	5.56	5.82
<b>AG</b>	1.20	1.61	3.08	0.17	2.48	<b>3.64</b>	0.08	2.59	2.98
<b>Something</b>	0.96	2.62	2.81	0.37	<b>3.11</b>	2.92	0.42	2.72	2.88
<b>WhatsUp</b>	0.00	1.58	0.76	0.15	0.88	0.54	0.82	2.88	<b>3.30</b>
<b>Kubric</b>	1.88	3.92	3.70	0.14	7.30	<b>7.32</b>	2.22	6.28	6.60
<b>AURORA-BENCH Avg.</b>	1.43	3.14	3.41	0.17	3.26	3.67	0.84	4.04	<b>4.32</b>
<i>Spatial Reasoning (Accuracy)</i>									
<b>SpatialMQA</b>	–	–	–	26.1	25.8	27.2	27.7	<b>28.0</b>	27.8
<b>EmbodiedSpatial-Bench</b>	–	–	–	15.1	21.2	17.5	32.6	33.2	<b>33.8</b>

Table 2: Performance on AURORA-BENCH (GPT-4o-as-a-judge score) and spatial reasoning benchmarks (accuracy). We **bold** the best model per dataset. Among our variants, we highlight the **best** and **worst** scores. SE: SmartEdit. Greedy decoding is used for spatial reasoning evaluation and applied only to VLMs.

	Weighted		Standard	
	ES (↑)	ME (↑)	ES (↑)	ME (↑)
<b>MB</b>	3.73	8.17	3.68	8.46
<b>AG</b>	3.18	8.03	2.37	8.13
<b>ST</b>	3.32	7.01	2.78	7.20
<b>WU</b>	0.54	7.25	0.76	7.19
<b>KU</b>	7.75	8.49	7.24	8.70
<b>Avg.</b>	3.71	7.80	3.37	7.94
<b>GPT4o</b>	3.67		3.58	

Table 3: Detailed scores of GPT4o-as-a-judge evaluation for loss-weighting and standard training. We report the scores for **Editing Success (ES)** and **Minimal Editing (ME)**. MB: MagicBrush, AG: Action-Genome, ST: Something-Something, WU: WhatsUp, KU: Kubric. We highlight the **best** and **worst** scores for each category. We report the average of ES and ME as **Avg.** and the final score as **GPT4o**.

AURORA-BENCH reveals the benefit of augmenting the training data with synthetic triplets bootstrapped from the IDM (vs. FT), as it yields a 13% gain for Chameleon and 7% for Liquid, and the benefit of fine-tuning off-the-shelf VLMs for FDP more broadly (vs. ZS).

**Human Evaluation.** Following Krojer et al. (2024), we conduct a blind human evaluation comparing GoT, SmartEdit, C-FT, and C-FDM. We randomly sample 5 examples from each subset within AURORA-BENCH and present the outputs generated by each of the four models. Human annotators are asked to identify the best and worst generated observations based on three criteria: (1) *Realism*: the generated image should exhibit natural textures and lighting while remaining faithful to the input scene; (2) *Instruction-Following Ability*: the edit should clearly reflect the given action;

and (3) *Over-Editing*: the modification should be minimal and focused, altering only what is necessary to fulfil the action. Each model receives +1 point for being selected as the best, -1 for the worst, and 0 otherwise. We compute the average scores over 350 annotated samples, as reported in Table 5. The results are well aligned with automatic evaluations in Table 2: image-editing models excel in the MagicBrush and WhatsUp subsets, but fall short on action-centric datasets such as Action-Genome, Something-Something, and Kubric. In contrast, C-FDM outperforms C-FT (and all other baselines) on all three of these datasets, highlighting its strength in next-observation prediction in real-world, action-centric trajectories despite remaining a general-purpose VLM.

**Ablation Study on Synthetic Trajectories.** To assess the importance of extra supervision from IDM-synthetic trajectories, Table 4 reports GPT-4o’s scores for this ablation. We see performance drops on most datasets—particularly on Something and AG—when the additional training data from unlabelled videos is removed, highlighting the effectiveness of bootstrapping C-FDM with large-scale real-world data via IDM. An exception is the WhatsUp dataset, which focuses on specific actions within a fixed scene; in this case, training in an open-domain setting may not transfer effectively.

**Ablation Study on Loss Weighting.** Based on Table 4, we also observe consistent degradation when loss weighting is removed, demonstrating the benefit of explicitly incorporating the recognition model into visual next-token prediction. To better understand the effect of loss weighting, Table 3 re-

	C-FDM	w/o Synth.	w/o LW
<b>MB</b>	<b>3.48</b>	-0.28	-0.22
<b>AG</b>	<b>3.02</b>	-0.35	-0.08
<b>ST</b>	<b>3.06</b>	-0.18	-0.19
<b>WU</b>	<b>0.46</b>	0.40	0.08
<b>KU</b>	<b>7.14</b>	-0.03	-0.33
<b>All</b>	<b>3.43</b>	-0.09	-0.15

Table 4: Ablation study of synthetic trajectories (Synth.) and loss weighting (LW) in C-FDM. Numbers are GPT-4o-as-judge scores ( $\uparrow$ , average of 3 runs). MB: MagicBrush, AG: Action-Genome, ST: Something-Something, WU: WhatsUp, KU: Kubric.

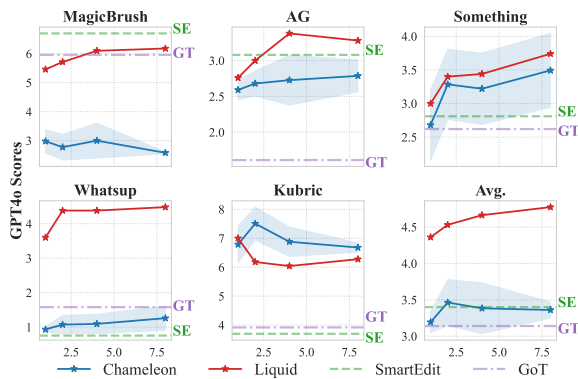


Figure 4: GPT-4o scores for test-time verification with  $K$  samples, where  $K \in \{1, 2, 4, 8\}$ . We use a blue line for C-FT and a red line for L-FT. For C-FT, we plot the standard deviation as the shaded area due to its large variance. We indicate the scores for GoT (GT) and SmartEdit (SE) as horizontal lines.

ports the average scores for two criteria used in the GPT-4o-as-a-judge evaluation separately: Editing Success (ES), which measures how well the model captures the intended action and performs the corresponding edit, and Minimal Editing (ME), which assesses whether the model introduces unnecessary modifications. The full distribution of GPT-4o scores is provided in Appendix G. Our analysis reveals that the primary bottleneck for C-FDM remains its ability to reliably follow the instruction, as reflected by the fact that ES scores are significantly lower than ME scores. Loss weighting partly solves this problem, increasing the editing success and reducing copying behaviour, albeit at the cost of sometimes over-editing the source observation.

**Image Editing as an Auxiliary Task.** Finally, we assess whether enhanced FDP capabilities are beneficial for a broader range of vision-language

	GoT	SE	C-FT	C-FDM
<b>MB</b>	0.06 $\dagger$	<b>0.29<math>\dagger</math></b>	-0.32 $\dagger$	-0.03
<b>AG</b>	-0.23 $\dagger$	-0.46 $\dagger$	0.32	<b>0.37</b>
<b>ST</b>	0.00	-0.37 $\dagger$	0.18	<b>0.20</b>
<b>WU</b>	<b>0.25</b>	-0.38 $\dagger$	0.14	0.00
<b>KU</b>	-0.52 $\dagger$	-0.22 $\dagger$	0.34	<b>0.40</b>
<b>All</b>	-0.09 $\dagger$	-0.23 $\dagger$	0.13	<b>0.19</b>

Table 5: Human evaluation results.  $\dagger$  indicates all results whose gap with respect to C-FDM is significant, based on a Wilcoxon signed-rank test ( $p = 0.05$ ). MB: MagicBrush, AG: Action-Genome, ST: Something-Something, WU: WhatsUp, KU: Kubric.

tasks. Since AURORA and the action-annotated videos contain diverse spatial relations (e.g., left/right orientation), we evaluated whether this supervision helps VLMs generalise beyond editing. We tested our models on two spatial reasoning benchmarks: SpatialMQA (Liu et al., 2025) and EmbodiedSpatial-Bench (Du et al., 2024), and report the corresponding accuracy in Table 2. Both Chameleon and Liquid trained with the FDP objective outperform the zero-shot baseline, demonstrating that the FDP task transfers beyond action-centric editing and highlighting FDP as a signal for enhancing spatial reasoning.

### 4.3.2 Inference-time Verification

We evaluate our test-time strategy using IDM to choose among candidates generated by C/L-FT in Figure 4, using  $K \in \{1, 2, 4, 8\}$ . By increasing exploration on more candidate next observations, Chameleon and Liquid benefit from test-time verification on most datasets with real-world trajectories (e.g., AG, Something, WhatsUp), indicating the reliability of IDM’s trajectory preferences. Increasing  $K$  does not always improve performance (MagicBrush, Kubric), however, suggesting that bootstrapping with IDM that shares the same foundation model backbone may be limiting. In summary, IDM-based verification boosts performance to a similar level as FDM, by leveraging more samples during inference rather than training.

**Qualitative Example.** Figure 5 presents a real-world example demonstrating that C-FDM is also capable of iteratively generating future observations in multiple steps while maintaining consistency with previous frames.

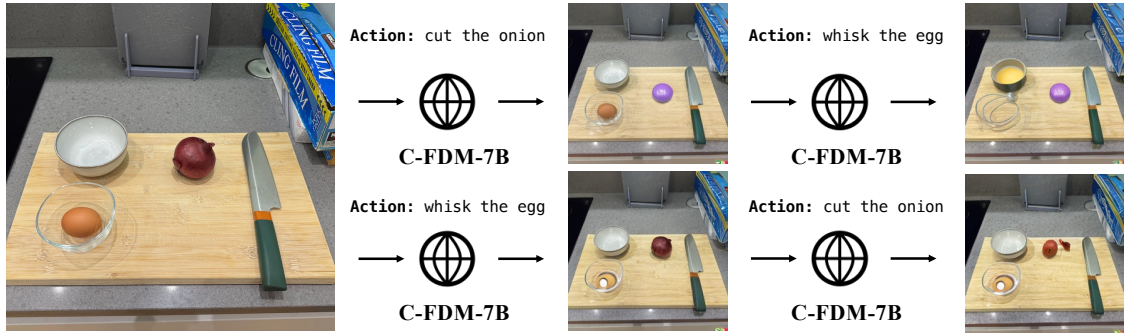


Figure 5: A qualitative case of real-world next-observation prediction, demonstrating C-FDM’s ability to steer predictions using language and perform *sequential* predictions. More cases from AURORA-BENCH are in Appendix C.

## 5 Related Work

Despite the surge in interest for world modelling (Ha and Schmidhuber, 2018; Sutton, 1988; Hafner et al., 2019a), previous works focused mostly on building specialised *ad-hoc* world models. These world models can be explicitly learnt as a visual simulator (Agarwal et al., 2025; Bruce et al., 2024; Brooks et al., 2024), or enable planning with model-based reinforcement learning (Hafner et al., 2019b; Micheli et al., 2022; Robine et al., 2023; Alonso et al., 2024; Hafner et al., 2025). Instead, we focus on leveraging general-purpose, pre-trained vision-language models (Lin et al., 2024; Chen et al., 2025; Wu et al., 2024c) to develop world models, which is more attractive due to the inductive bias they provide from their extensive training.

This is possible thanks to frameworks that integrate observations, actions, and rewards into a unified sequence of tokens in autoregressive Transformers (Wu et al., 2024a), building on pioneering works such as Decision Transformers (Chen et al., 2021) and GATo (Reed et al., 2022b). Related to our work, Chen et al. (2024) initialise the parameters of RL policies with VLMs, thus taking advantage of the abundant and general world knowledge encoded in their representations. 3D-VLA (Zhen et al., 2024) integrates a set of interaction tokens into a Large Language Model to engage with the environment as an embodied agent. Yang et al. (2024) and Soni et al. (2024) explore large-scale self-supervised learning via next token or frame prediction to build a unified model absorbing internet knowledge, learning from interaction via video.

Most similar to our work, Baker et al. (2022) train a dynamics model which aims to uncover the underlying action between video frames in unlabelled video frames from the Minecraft game. Through this model, they synthesise trajectories

to train a policy for sequential decision making. In contrast with Baker et al. (2022), we focus on actin-centric next-observation prediction as a task to evaluate FDP. First, this allows us to port the observation space to real-world frames, rather than simulated ones, hence assessing whether VLMs can eventually be developed into world models. Second, this broadens the space of actions from a few choices to the combinatorially infinite and expressive space of language, capturing a significantly more diverse range of dynamics.

## 6 Conclusions

In this work, we explored whether unified vision-language models (VLMs) can be endowed with the ability to predict forward dynamics, i.e., the next observation in the environment (e.g., an image frame) given the past observation and an action (e.g., a textual instruction). We first show that these models lack a clear preference for ground-truth real-world action-observation trajectories compared with adversarially manipulated ones.

To address this, we leverage an inverse-dynamics model (IDM) fine-tuned from the same VLM, which consists instead of predicting actions taking place between observations and is easier to learn, to bootstrap a better forward-dynamics model (FDM). Specifically, the IDM can be used to 1) automatically annotate pairs of frames from unlabelled videos, which are then used for weakly supervised training of the FDM; or 2) verify the best sample among multiple candidates generated from the FDM at inference time. Experiments confirm the effectiveness of both strategies, with our general-purpose forward-dynamics model achieving state-of-the-art performance compared to existing approaches specialised for image editing.

## 565 Limitations

566 While overall our results demonstrate the effective-  
567 ness of our approaches across AURORA-BENCH,  
568 we would like to highlight a few limitations that  
569 we have discovered:

- 570 • Despite efforts to guide the model via weakly  
571 supervised fine-tuning with loss weighting or  
572 inference-time verification (Table 2), we ob-  
573 serve that the model may still resort to copying  
574 the source observation, especially under low  
575 sampling temperatures or ambiguous instruc-  
576 tions.
- 577 • While we show promising preliminary results  
578 of language-steered next-observation predic-  
579 tion in Figure 5, fine-grained control remains  
580 limited, and understanding subtle instructions  
581 (e.g., spatial or quantitative edits) remains  
582 challenging.
- 583 • We observe high variance across different runs  
584 of experiments for Chameleon, likely due to  
585 the sensitivity of sampling from a weak model.  
586 To address this, we report results averaged  
587 over multiple runs.

## 588 References

589 Mostafa Abdou, Artur Kulmizev, Daniel Hershcovich,  
590 Stella Frank, Ellie Pavlick, and Anders Søgaard.  
591 2021. Can language models encode perceptual struc-  
592 ture without grounding? A case study in color. In  
593 *Proceedings of the 25th Conference on Computa-  
594 tional Natural Language Learning*, pages 109–132,  
595 Online. Association for Computational Linguistics.

596 Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji,  
597 Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay,  
598 Yongxin Chen, Yin Cui, Yifan Ding, and 1 others.  
599 2025. Cosmos world foundation model platform for  
600 physical AI. *arXiv preprint arXiv:2501.03575*.

601 Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kan-  
602 ervisto, Amos J Storkey, Tim Pearce, and François  
603 Fleuret. 2024. Diffusion for world modeling: Visual  
604 details matter in atari. *Advances in Neural Informa-  
605 tion Processing Systems*, 37:58757–58791.

606 Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost  
607 Huizinga, Jie Tang, Adrien Ecoffet, Brandon  
608 Houghton, Raul Sampedro, and Jeff Clune. 2022.  
609 Video PreTraining (VPT): Learning to act by watch-  
610 ing unlabeled online videos. *Advances in Neural  
611 Information Processing Systems*, 35:24639–24654.

Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol  
612 Hausman, Alexander Herzog, Daniel Ho, Julian  
613 Ibarz, Alex Irpan, Eric Jang, Ryan Julian, and 1 oth-  
614 ers. 2023. Do as I can, not as I say: Grounding  
615 language in robotic affordances. In *Conference on  
616 robot learning*, pages 287–318. PMLR. 617

Tim Brooks, Aleksander Holynski, and Alexei A Efros. 618  
2023. InstructPix2Pix: Learning to follow image  
619 editing instructions. In *Proceedings of the IEEE/CVF  
620 conference on computer vision and pattern recogni-  
621 tion*, pages 18392–18402. 622

Tim Brooks, Bill Peebles, Connor Holmes, Will DePue,  
623 Yufei Guo, Li Jing, David Schnurr, Joe Taylor,  
624 Troy Luhman, Eric Luhman, and 1 others. 2024.  
625 Video generation models as world simulators. 2024.  
626 URL [https://openai.com/research/video-generation-  
627 models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators), 3:1. 628

Jake Bruce, Michael D Dennis, Ashley Edwards, Jack  
629 Parker-Holder, Yuge Shi, Edward Hughes, Matthew  
630 Lai, Aditi Mavalankar, Richie Steigerwald, Chris  
631 Apps, and 1 others. 2024. Genie: Generative in-  
632 teractive environments. In *Forty-first International  
633 Conference on Machine Learning*. 634

Joao Carreira, Eric Noland, Chloe Hillier, and An-  
635 drew Zisserman. 2019. A short note on the  
636 kinetics-700 human action dataset. *arXiv preprint  
637 arXiv:1907.06987*. 638

Chameleon Team. 2024. Chameleon: Mixed-modal  
639 early-fusion foundation models. *arXiv preprint  
640 arXiv:2405.09818*. 641

Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee,  
642 Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind  
643 Srinivas, and Igor Mordatch. 2021. Decision trans-  
644 former: Reinforcement learning via sequence mod-  
645 eling. *Advances in neural information processing  
646 systems*, 34:15084–15097. 647

William Chen, Oier Mees, Aviral Kumar, and Sergey  
648 Levine. 2024. Vision-language models provide  
649 promptable representations for reinforcement learn-  
650 ing. In *Automated Reinforcement Learning: Explor-  
651 ing Meta-Learning, AutoML, and LLMs*. 652

Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan,  
653 Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan.  
654 2025. Janus-Pro: Unified multimodal understanding  
655 and generation with data and model scaling. *arXiv  
656 preprint arXiv:2501.17811*. 657

Ethan Chern, Jiadi Su, Yan Ma, and Pengfei Liu. 2024.  
658 ANOLE: An open, autoregressive, native large multi-  
659 modal models for interleaved image-text generation.  
660 *arXiv preprint arXiv:2407.06135*. 661

Dima Damen, Hazel Doughty, Giovanni Maria Farinella,  
662 Sanja Fidler, Antonino Furnari, Evangelos Kazakos,  
663 Davide Moltisanti, Jonathan Munro, Toby Perrett,  
664 Will Price, and 1 others. 2018. Scaling egocentric  
665 vision: The EPIC-KITCHENS dataset. In *Proceed-  
666 ings of the European conference on computer vision  
667 (ECCV)*, pages 720–736. 668



779	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. <i>arXiv preprint arXiv:2501.19393</i> .	832
780		833
781		834
782		835
783		836
784	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	837
785		838
786		839
787		840
788		841
789	Yiran Qin, Zhelun Shi, Jiwen Yu, Xijun Wang, Enshen Zhou, Lijun Li, Zhenfei Yin, Xihui Liu, Lu Sheng, Jing Shao, and 1 others. 2024. WorldSimBench: Towards video generation models as world simulators. <i>arXiv preprint arXiv:2410.18072</i> .	842
790		843
791		844
792		845
793		846
794	Yifu Qiu, Zheng Zhao, Yftah Ziser, Anna Korhonen, Edoardo Ponti, and Shay B Cohen. 2024. Are large language model temporally grounded? In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 7057–7076.	847
795		848
796		849
797		850
798		851
799		852
800		853
801	Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, and 1 others. 2022a. A generalist agent. <i>arXiv preprint arXiv:2205.06175</i> .	854
802		855
803		856
804		857
805		858
806		859
807	Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, and 1 others. 2022b. A generalist agent. <i>arXiv preprint arXiv:2205.06175</i> .	860
808		861
809		862
810		863
811		864
812		865
813	Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. 2023. Transformer-based world models are happy with 100k interactions. <i>arXiv preprint arXiv:2303.07109</i> .	866
814		867
815		868
816		869
817	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling LLM test-time compute optimally can be more effective than scaling model parameters. <i>arXiv preprint arXiv:2408.03314</i> .	870
818		871
819		872
820		873
821	Achint Soni, Sreyas Venkataraman, Abhramil Chandra, Sebastian Fischmeister, Percy Liang, Bo Dai, and Sherry Yang. 2024. VideoAgent: Self-improving video generation. <i>arXiv preprint arXiv:2410.10076</i> .	874
822		875
823		876
824		877
825	Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. 2012. UCF101: A dataset of 101 human actions classes from videos in the wild. <i>arXiv preprint arXiv:1212.0402</i> .	878
826		
827		
828		
829	Richard S Sutton. 1988. Learning to predict by the methods of temporal differences. <i>Machine learning</i> , 3:9–44.	
830		
831		
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. <b>Transformers: State-of-the-art natural language processing</b> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
	Jialong Wu, Shaofeng Yin, Ningya Feng, Xu He, Dong Li, Jianye Hao, and Mingsheng Long. 2024a. <b>iVideoGPT: Interactive VideoGPTs are scalable world models</b> . In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 68082–68119. Curran Associates, Inc.	
	Junfeng Wu, Yi Jiang, Chuofan Ma, Yuliang Liu, Hengshuang Zhao, Zehuan Yuan, Song Bai, and Xiang Bai. 2024b. Liquid: Language models are scalable and unified multi-modal generators. <i>arXiv preprint arXiv:2412.04332</i> .	
	Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang, Dacheng Li, Yunhao Fang, Ligeng Zhu, Enze Xie, Hongxu Yin, Li Yi, and 1 others. 2024c. VILA-U: a unified foundation model integrating visual understanding and generation. <i>arXiv preprint arXiv:2409.04429</i> .	
	Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. 2023. Learning interactive real-world simulators. <i>arXiv preprint arXiv:2310.06114</i> , 1(2):6.	
	Sherry Yang, Jacob Walker, Jack Parker-Holder, Yilun Du, Jake Bruce, Andre Barreto, Pieter Abbeel, and Dale Schuurmans. 2024. Video as the new language for real-world decision making. <i>arXiv preprint arXiv:2402.17139</i> .	
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating text generation with bert. In <i>International Conference on Learning Representations</i> .	
	Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang, Xin Yan, Yilun Du, Yining Hong, and Chuang Gan. 2024. 3D-VLA: A 3D Vision-Language-Action Generative World Model. In <i>International Conference on Machine Learning</i> , pages 61229–61245. PMLR.	

879	<b>A Potential Risks</b>		
880	This work develops models for action-centric image editing for visual world modelling. While our primary aim is to advance fundamental research in world modelling, we acknowledge potential risks, particularly in the generation of realistic future observations.		
881			
882			
883			
884			
885			
886	A core concern is the potential misuse of the models for creating deceptive visual content, including fabricated action sequences or manipulated images that imply false causality. Although the model is not explicitly designed for these tasks, its ability to generate coherent visual predictions from the linguistic action could be adapted for such uses if deployed irresponsibly.		
887			
888			
889			
890			
891			
892			
893			
894	Even in intended use, risks include over-reliance on generated outputs in downstream tasks such as robotic control, or interactive systems. Model failures—e.g., copying artefacts, hallucinations, or broken object continuity—can lead to incorrect inferences or reinforce dataset biases.		
895			
896			
897			
898			
899			
900	To mitigate potential misuse, we limit our model release to research purposes under a non-commercial license and clearly communicate its capabilities and limitations. We urge caution when adapting them for deployment, particularly in settings with high societal or ethical sensitivity.		
901			
902			
903			
904			
905			
906	<b>B Model Performance on</b>		
907	<b>AURORA-BENCH with 5 Metrics</b>		
908	In addition to GPT4o-as-a-judge evaluation, we further employ a diverse set of automatic metrics covering both low-level and semantic fidelity: 1) we compute the <b>L1 distance</b> between the predicted and target observation as a pixel-level metric. 2) We extract visual features and compute the cosine similarity in their respective embedding spaces for several image encoders, including ( <b>CLIP-I</b> and <b>DINO</b> ), to assess semantic similarity. Additionally, to measure alignment between image content and the action semantics, we compute <b>CLIP-T</b> , the similarity between the edited image and its BLIP-generated caption. These metrics are evaluated in addition to GPT4o-as-a-judge metric following previous works in image editing (Huang et al., 2024; Fang et al., 2025; Krojer et al., 2024). We report the detailed results with 5 metrics in Table 6. We notice that copy baseline exhibits the best performance as measured by the distance-based and visual encoder-based approach, as indicated in Table 2. This poses a challenge to the reliability of the traditional met-		
909			
910			
911			
912			
913			
914			
915			
916			
917			
918			
919			
920			
921			
922			
923			
924			
925			
926			
927			
928			
		rics in fairly evaluating the action-centric image editing task. On the other hand, GPT4o-as-a-judge metric robustly assigns 0 score to Copy, indicating its robustness in detecting copying generation while putting GPT-as-a-judge as the most reliable metric to interpret.	929 930 931 932 933 934
	<b>C Qualitative Cases</b>		
	In this section, we present additional qualitative examples from AURORA-BENCH in Figure 6. We observe several common failure modes in image editing models. First, they sometimes fail to preserve the scene from the source observation (e.g., PixInstruct on Action-Genome and MagicBrush). Second, some models generate near-identical copies of the source as the target (e.g., GoT on Something-Something). Third, producing realistic outputs remains difficult, as seen in GoT’s result on Kubric. Finally, maintaining object consistency is also a challenge—SmartEdit alters the object in WhatsUp, and C-FDM does so in Something-Something.	935 936 937 938 939 940 941 942 943 944 945 946 947 948	
	Despite the challenges, we also observe several positive editing behaviours from C-FDM. On Action-Genome, C-FDM correctly predicts spatial changes, such as <i>opening and closing a drawer</i> , which requires a strong understanding of the spatial concepts. In Something-Something, it is the only model to accurately capture the spatial concept of “falling down.” On Kubric, it demonstrates basic counting ability by correctly adding one keyboard. In WhatsUp, C-FDM correctly grounds the action to the laptop, while other models mistakenly edit the monitor.	949 950 951 952 953 954 955 956 957 958 959 960	
	<b>D Detailed Discussion for Chameleon’s Predicted Likelihoods</b>		
	From Figure 7, it emerges that Chameleon-7B displays a very limited preference for the ground-truth trajectories in a zero-shot setting. In the action prediction task (top panel), there is a slightly higher tendency to favour the ground-truth; however, even in the best case (counterfactual action), the model prefers the reference in only 58.1% of the samples. The high correlation in likelihoods indicates that the VLM struggles also on visual manipulations. In the next-observation prediction task (bottom panel), the VLM mostly fails in effectively differentiating the ground truth from the negatives. An exception to this is the copy manipulation, where the model can always tell them apart. Although the underlying reason remains uncertain, one plausible	961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977	

Datasets	Metrics	Models								
		Copy	PixInstruct	GoT	SE	CM	C-FT	+Best-of-3	C-FDM	+Best-of-3
MagicBrush	L1	0.027	0.114	<b>0.063</b>	0.068	0.287	0.075	0.075	0.090	0.078
	CLIP-I	0.959	0.877	0.930	<b>0.937</b>	0.671	0.913	0.914	0.906	0.909
	CLIP-T	0.289	0.275	0.286	0.290	0.227	0.289	0.289	<b>0.291</b>	0.291
	DINO	0.931	0.761	0.881	<b>0.894</b>	0.292	0.883	0.883	0.864	0.864
	GPT-4o	0.000	3.120	5.960	<b>6.710</b>	0.000	2.520	3.270	<b>3.920</b>	<b>3.920</b>
AG	L1	0.069	0.220	0.174	<b>0.137</b>	0.314	0.170	0.168	0.168	0.167
	CLIP-I	0.943	0.757	0.846	0.811	0.609	0.872	0.872	<b>0.881</b>	0.883
	CLIP-T	0.279	0.254	0.280	0.268	0.214	0.280	0.284	<b>0.284</b>	0.284
	DINO	0.929	0.557	0.785	0.774	0.258	0.801	0.817	<b>0.816</b>	0.816
	GPT-4o	0.000	1.200	1.610	3.080	0.170	2.480	2.740	<b>3.640</b>	<b>3.640</b>
Something	L1	0.135	0.232	0.184	<b>0.163</b>	0.293	0.184	0.184	0.196	0.184
	CLIP-I	0.870	0.709	0.807	0.773	0.649	<b>0.820</b>	0.820	0.804	0.804
	CLIP-T	0.275	0.238	0.269	0.265	0.232	<b>0.271</b>	0.269	0.268	0.268
	DINO	0.797	0.453	0.636	0.662	0.297	<b>0.675</b>	0.653	0.666	0.666
	GPT-4o	0.000	0.957	2.620	2.810	0.370	<b>3.110</b>	3.110	2.920	<b>3.310</b>
WhatsUp	L1	0.039	0.138	0.078	0.067	0.251	<b>0.066</b>	0.066	0.070	0.070
	CLIP-I	0.954	0.817	<b>0.923</b>	0.888	0.721	0.877	0.880	0.870	0.883
	CLIP-T	0.326	0.287	<b>0.316</b>	0.312	0.243	0.309	0.310	0.306	0.307
	DINO	0.908	0.615	<b>0.850</b>	0.805	0.424	0.836	0.841	0.831	0.838
	GPT-4o	0.000	0.000	<b>1.580</b>	0.755	0.146	<b>0.880</b>	<b>0.980</b>	0.540	0.540
Kubric	L1	0.011	0.104	<b>0.026</b>	0.064	0.276	0.044	0.044	0.044	0.044
	CLIP-I	0.963	0.796	0.895	0.868	0.660	0.897	0.899	<b>0.897</b>	0.898
	CLIP-T	0.282	0.259	0.281	0.271	0.213	<b>0.287</b>	0.287	0.287	0.288
	DINO	0.955	0.676	0.857	0.798	0.161	<b>0.906</b>	0.906	0.902	0.902
	GPT-4o	0.000	1.880	3.920	3.700	0.140	7.300	7.300	<b>7.320</b>	<b>7.780</b>
All	GPT-4o	0.000	1.430	3.140	3.410	0.165	3.260	3.480	<b>3.670</b>	<b>3.840</b>

Table 6: Model performance at MagicBrush, Action-Genome, Something, WhatsUp and Kubric on AURORA-BENCH. For C-FT and C-FDM We report both the model performance and their performance in the *best-of-N* distribution. We report the average GPT4o scores for each model at the bottom. We highlight the better GPT-4o scores for C-FT and C-FDM. We bold the best performance among all models, except Copy and *best-of-N* performances. SE: SmartEdit.

978 explanation for this behaviour is that the model’s  
979 ability to solve next-observation prediction tasks  
980 depends on their alignment with training sequences:  
981 for instance, it is plausible that Chameleon’s data  
982 rarely features two identical adjacent images. In  
983 Figure 8, we visualize the predicted likelihoods  
984 produced by Chameleon’s fine-tuned inverse dynam-  
985 ics model (IDM). We observe that fine-tuning  
986 on ground-truth trajectories substantially increases  
987 the model’s ability to distinguish ground-truth ac-  
988 tions from negative alternatives. Specifically, the  
989 probability that the ground-truth action is assigned  
990 a higher likelihood increases from 55.6% to 73.2%  
991 under random-action negatives, and from 58.1%  
992 to 72.2% under counterfactual-action negatives.  
993 These results demonstrate the strong potential of  
994 learning effective inverse dynamics models directly  
995 from real-world trajectories. In summary, the zero-

shot Chameleon-7B does not exhibit a preference  
for ground-truth trajectories over negative ones,  
constructed through action- or observation-based  
manipulations. However, it is possible to learn an  
effective IDM from the real-world trajectories.

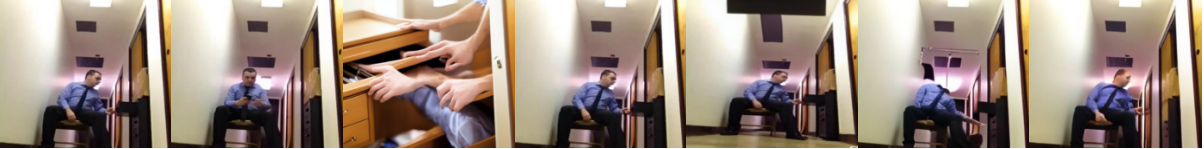
## E Details of Processing IDM Annotations for Unlabelled Videos

We present the raw dataset statistics before sam-  
pling for Movements-in-Time, UCF-101 and Ki-  
netics700 in Table 7. Figure 9 shows the distri-  
bution of IDM’s predicted scores across action  
classes in Movements-in-Time, Kinetics700, and  
UCF-101. The predicted likelihoods are nearly  
uniform within each class, indicating that our sam-  
pling method maintains both class diversity and  
high overall likelihoods. The sampling procedure  
for IDM-annotated trajectories is detailed in Algo-

Action (MagicBrush): Let the laptop screen be blank



Action (Action-Genome): Make him fully close the drawer



Action (Something-Something): Make remote fall down



Action (WhatsUp): Move the book under the table



Action (Kubric): Add 1 white keyboard to the platform



Input

Output

PixInstruct

GoT

SmartEdit-7B

C-FT

CWM

Figure 6: Qualitative examples of the predicted next observation from the state-of-the-art specialised image editing models, and our models including C-FT and C-FDM, on AURORA-BENCH.

rithm 1.

## F Prompt Template for Using GPT4o-as-a-Judge

We provide the prompts used for evaluating image editing performance with GPT-4o in Figure 10. We use GPT-4o-2024-11-20. The final score is the average of the minimum value of the two scores for each sample, as in (Fang et al., 2025).

## G Detailed GPT4o Scores for Editing Success and Minimal Editing

Figure 11 shows the distribution of editing success (ES) and minimal editing (ME) scores for standard training and loss-weighted training. Loss weighting tends to improve editing success, with a modest trade-off in minimal editing quality in most of the datasets.

## H Implementation Details

### H.1 Chameleon Dynamics Model

We fine-tune the Chameleon-7B checkpoint from the Anole-7B version (Chern et al., 2024) to predict the action given a pair of observations, framed as an action-prediction task. The model is trained on a merged dataset from Action-Genome, Kubric, MagicBrush, Something-Something from AURORA’s annotated trajectories, and 15K EPIC-Kitchens processed by us. We downsample Kubric’s trajectories to 10K. Training is performed for 10 epochs with a batch size of 64, using a learning rate of 2e-4

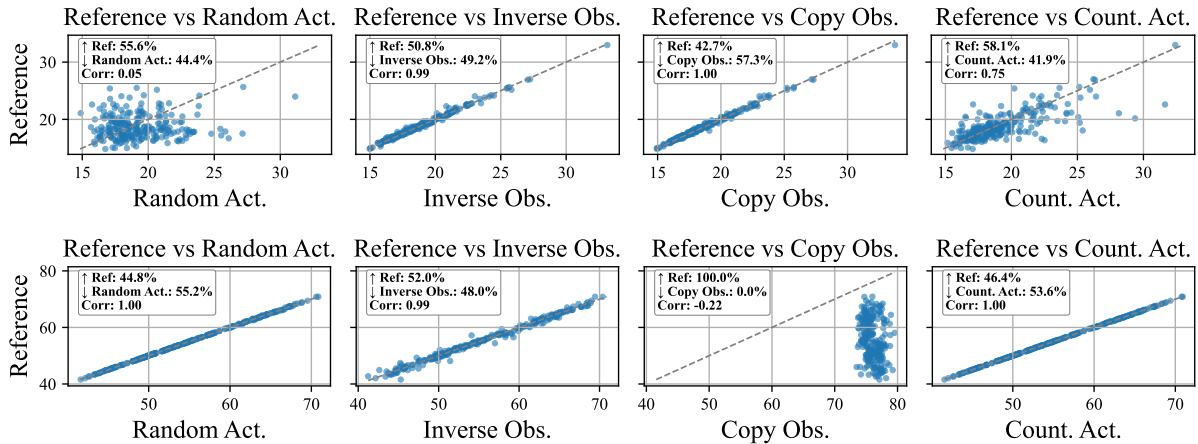


Figure 7: Comparison of predicted negative log-likelihoods (lower values indicate stronger model preference) for ground-truth real-world trajectories versus four types of negative trajectories. **Top:** Action prediction task for the IDP (observation  $\times$  observation  $\rightarrow$  action). **Bottom:** Next observation prediction task for the FDP (observation  $\times$  action  $\rightarrow$  observation). The legend shows the percentage of times the model prefers the ground-truth trajectory ( $\uparrow$ ) over the negatives ( $\downarrow$ ).

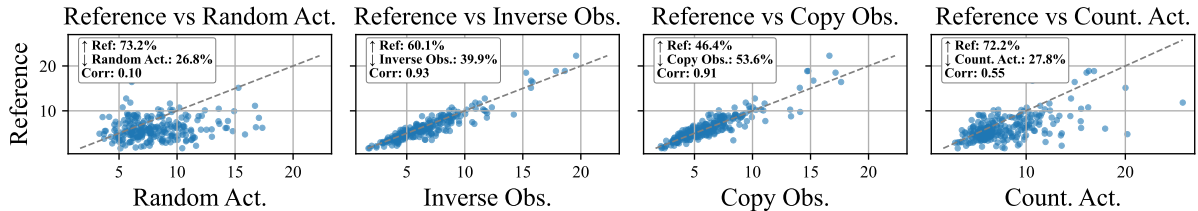


Figure 8: Comparison of negative log-likelihoods (lower values indicate stronger model preference) of the action predicted by IDM for ground-truth trajectories versus four types of negative trajectories.

1041 and cosine scheduling (500 warm-up steps). We  
 1042 use bfloat16 mixed-precision training and apply  
 1043 LoRA (Hu et al., 2022) for parameter-efficient fine-  
 1044 tuning (rank 16,  $\alpha = 32$ , dropout 0.05). Only the  
 1045 completion loss is used to optimise the generation  
 1046 of action. Training is conducted on 4 NVIDIA-  
 1047 H100-80GB-HBM3 GPUs using DeepSpeed for  
 1048 distributed optimisation.

## 1049 H.2 C-FT Baseline

1050 We fine-tune the Chameleon-7B checkpoint from  
 1051 the Anole-7B version (Chern et al., 2024).  
 1052 The model is trained on a combined dataset  
 1053 from Action-Genome, Kubric, MagicBrush, and  
 1054 Something-Something, formatted as the image edit-  
 1055 ing task. We downsample Kubric’s trajectories to  
 1056 10K. Training is conducted for 40 epochs with a  
 1057 batch size of 96 using the AdamW optimiser and  
 1058 a cosine learning rate scheduler (learning rate of  
 1059  $5e-4$ , 400 warm-up steps). We use mixed-precision  
 1060 training with bfloat16 and apply LoRA (Hu et al.,  
 1061 2022) for efficient fine-tuning (rank 16,  $\alpha = 32$ ,  
 1062 dropout 0.05). We only train the model with the

1063 truncated loss from the completion part. We use  
 1064 4 NVIDIA-H100-80GB-HBM3 GPUs with Deep-  
 1065 Speed for distributed training. During inference,  
 1066 we apply a logits processor to mask out non-image  
 1067 tokens, set the temperature to 1, and use top-1 sam-  
 1068 pling. We observe that temperature is critical in  
 1069 controlling model behaviour: lower values often  
 1070 cause the model to copy the source observation  
 1071 instead of generating meaningful edits.

## 1072 H.3 Chameleon FDM

1073 We fine-tune the Chameleon-7B checkpoint from  
 1074 the Anole-7B version (Chern et al., 2024). The  
 1075 model is trained on a combined dataset from  
 1076 Action-Genome, Kubric, MagicBrush, Something-  
 1077 Something from AURORA’s annotated trajectories,  
 1078 together with 7K trajectories from Movements-in-  
 1079 Time, 7K trajectories from UCF-101 and 7K tra-  
 1080 jectories from Kinetics700, formatted as the im-  
 1081 age editing task. Again, we downsample Kubric’s  
 1082 trajectories to 10K. Training is conducted for 40  
 1083 epochs with a batch size of 96 using the AdamW  
 1084 optimiser and a cosine learning rate scheduler

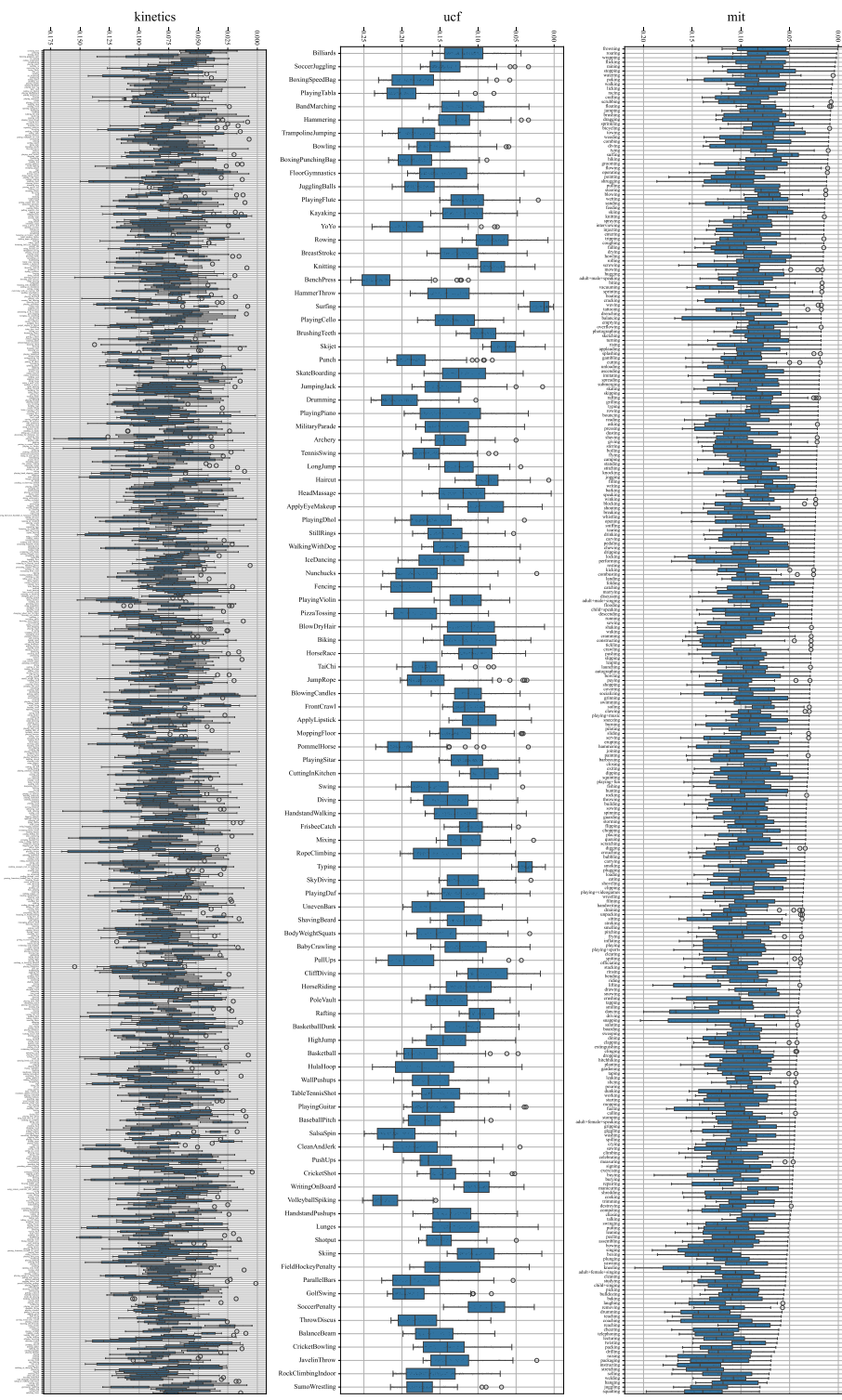


Figure 9: Distributions of triplet log-likelihoods predicted by IDM on Movements-in-Time, UCF-101, and Kinetics-700, based on 7K synthetic triplets per dataset. Triplets are uniformly sampled from each action class while maximising overall predicted likelihoods.

---

**Algorithm 1** Stratified Top-K Sampling with Action Class Uniformity

---

**Require:** Trajectory triplet set  $X = \{(o_s^i, o_t^i, a^i, s^i, c^i)\}_{i=1}^N$ , where  $s_i$  is the predicted likelihood of  $a^i$ ,  $c_i \in \mathcal{C}$  is the class, number of samples  $K$

```
1: Sort  $X$  descending by score  $s_i$ 
2: Initialize  $S \leftarrow \emptyset$ , and  $\text{class\_counts}[c] \leftarrow 0$  for all  $c \in \mathcal{C}$ 
3: while  $|S| < K$  do
4:   for all class  $c \in \mathcal{C}$  in round-robin order do
5:      $X_c \leftarrow$  top unsampled item from class  $c$  in  $X$ 
6:     if  $X_c \neq \emptyset$  then
7:        $S \leftarrow S \cup \{X_c\}$ 
8:       Remove  $X_c$  from  $X$ 
9:        $\text{class\_counts}[c] \leftarrow \text{class\_counts}[c] + 1$ 
10:    end if
11:    if  $|S| = K$  then
12:      break
13:    end if
14:  end for
15: end while
16: return  $S$ 
```

---

Dataset	Video		Triplet			
	Avg. Length	Total Length	#Samples	#Avg. OPV	#Avg. APV	#Avg. WPA
MIT	3.04 seconds	2.57 hours	19,658	2.05	1.05	7.10
UCF-101	7.24 seconds	26 hours	10,965	3.00	2.00	8.96
Kinetics700	9.02 seconds	18 hours	26,959	2.71	1.71	7.39

---

Table 7: Dataset statistics for the video and triplets from the trajectories annotated by IDM. **OPV**: observations (i.e., extracted key-frames) per video, **APV**: actions per video, **WPA**: words per action.

(learning rate of  $5e-4$ , 400 warm-up steps). We use mixed-precision training with bfloat16 and apply LoRA (Hu et al., 2022) for efficient fine-tuning (rank 16,  $\alpha = 32$ , dropout 0.05). We only train the model with the truncated loss from the completion part, but we weight the image tokens using  $L_2$  strategy as introduced in Section 3. We use 4 NVIDIA-H100-80GB-HBM3 GPUs with DeepSpeed for distributed training. We use the same hyperparameters as C-FT during the inference time.

#### H.4 Computing Resources

All training experiments were conducted on a compute node equipped with  $4 \times$  NVIDIA H100 80GB GPUs, 256 CPU cores, and 256GB of memory. The total GPU hours required for training C-FT, C-FDM, and IDM were approximately 200, 400, and 100 hours, respectively.

For inference, we used a single NVIDIA A100 80GB GPU with 8 CPU cores and 128GB memory. Inference for C-FT and C-FDM takes approx-

imately 1 GPU hour per model. When applying verification with  $K = 8$ , inference time increases to around 8 GPU hours. IDM only takes around 0.3 GPU hours for inference.

#### H.5 Assets and Licenses

In this section, we list the public assets we used in this paper and the corresponding links.

**Datasets.** We include the detailed license and URL for the datasets we used in this paper.

- AURORA and AURORA-BENCH (Krojer et al., 2024): MIT license, the reader can find the corresponding version we use in this paper in <https://github.com/McGill-NLP/AURORA>.
- Movements-in-Time (Monfort et al., 2019): BSD-2-Clause license and its own License for Non-Commercial Use, the reader can find the corresponding version we use in this paper in <http://moments.csail.mit.edu/>.

### Prompt Template for GPT4o-as-a-judge Evaluation

You are a professional digital artist. You will have to evaluate the effectiveness of the AI-generated image(s) based on the given rules.

You will have to give your output in a valid way of a Python dictionary format (Keep your reasoning concise and short.):

```
{{"score": [...], "reasoning": "..."} }
```

and don't output anything else. Two images will be provided:

- The first being the original AI-generated image
- The second being an edited version of the first.

The objective is to evaluate how successfully the editing instruction has been executed in the second image.

Note that sometimes the two images might look identical due to a failure in image editing. From a scale of 0 to 10:

- A score from 0 to 10 will be given based on the success of the editing.
- A second score from 0 to 10 will rate the degree of minimal editing.

Editing instruction: {instruction}

Figure 10: Prompt template used for GPT-4o-as-a-judge evaluation.

1124 • UCF-101 (Soomro et al., 2012): unknown  
1125 license, the reader can find the corre-  
1126 sponding version we use in this paper  
1127 in [https://huggingface.co/datasets/  
1128 flwrlabs/ucf101](https://huggingface.co/datasets/flwrlabs/ucf101).

1129 • Kinetics700 (Kay et al., 2017; Carreira et al.,  
1130 2019): Creative Commons Attribution 4.0  
1131 International License, the reader can find  
1132 the corresponding version we use in this  
1133 paper in [https://research.google/pubs/  
1134 the-kinetics-human-action-video-dataset/](https://research.google/pubs/the-kinetics-human-action-video-dataset/).

1135 • EPIC-Kitchens (Damen et al., 2018): Creative  
1136 Commons Attribution-NonCommercial 4.0 In-  
1137 ternational License, the reader can find the  
1138 corresponding version we use in this paper in  
1139 <https://epic-kitchens.github.io/>.

1140 **Implementation.** We use the other following code  
1141 for the implementations:

1142 • Transformers (Wolf et al., 2020): Apache-  
1143 2.0 license. We use the 4.47.0 version, fol-  
1144 lowing the link at [https://github.com/  
1145 huggingface/transformers](https://github.com/huggingface/transformers).

1146 • DeepSpeed: We use the 0.14.4 version, fol-  
1147 lowing the link at [https://github.com/  
1148 deepspeedai/DeepSpeed](https://github.com/deepspeedai/DeepSpeed).

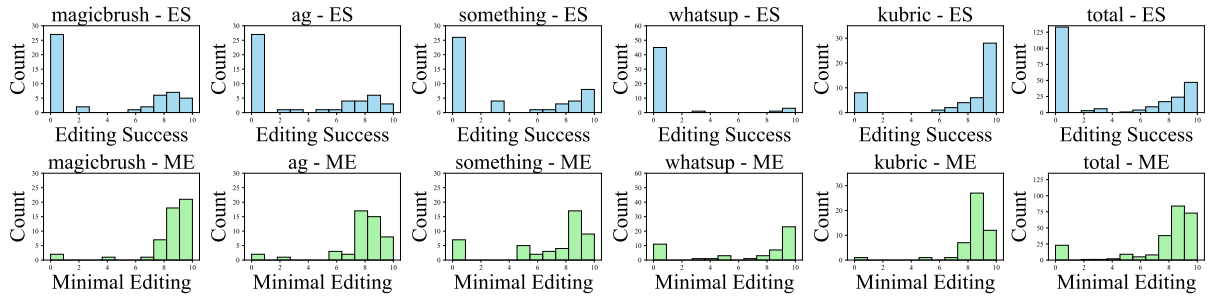
1149 **Model.** We use the following models or check-  
1150 points for the implementations:

1151 • Chameleon (Chameleon Team, 2024):  
1152 Chameleon Research License, the reader  
1153 can find the corresponding version we use  
1154 in this paper in [https://github.com/  
1155 facebookresearch/chameleon](https://github.com/facebookresearch/chameleon).

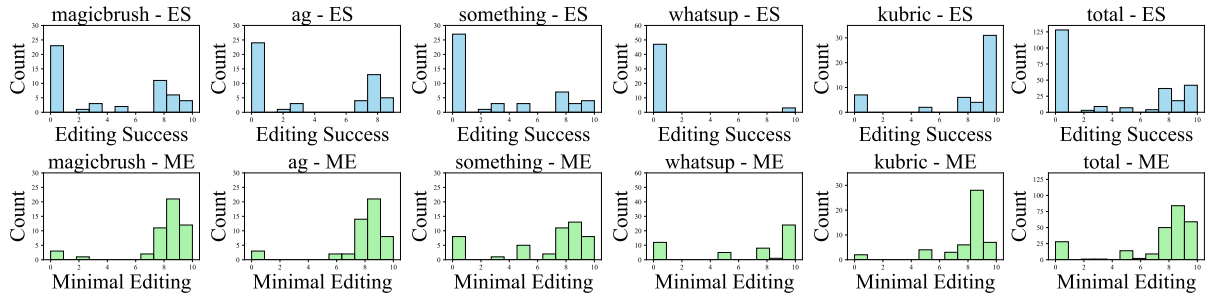
1156 • Anole-7B (Chern et al., 2024): Chameleon  
1157 Research License and MIT License, the  
1158 reader can find the corresponding version we  
1159 use in this paper in [https://github.com/  
1160 GAIR-NLP/anole](https://github.com/GAIR-NLP/anole).

1161 • VILA-U (Chern et al., 2024): MIT License,  
1162 the reader can find the corresponding ver-  
1163 sion we use in this paper in [https://github.  
1164 com/mit-han-lab/vila-u](https://github.com/mit-han-lab/vila-u).

1165 • SmartEdit (Huang et al., 2024): Apache-  
1166 2.0, the reader can find the correspond-



(a) Detailed GPT4o scores for C-FDM trained with the standard loss.



(b) Detailed GPT4o scores for C-FDM trained with the L<sub>2</sub>-weighted loss.

Figure 11: GPT4o scores’ distributions of editing success (ES) and minimal editing (OE) for C-FDM trained with standard loss or our loss-weighting method.

ing version we use in this paper in <https://huggingface.co/TencentARC/SmartEdit-7B>.

- GoT (Fang et al., 2025): MIT License, the reader can find the corresponding version we use in this paper in <https://github.com/rongyaofang/GoT>.
- PixInstruct (Brooks et al., 2023): PixInstruct customised license, the reader can find the corresponding version we use in this paper in <https://github.com/timothybrooks/instruct-pix2pix>.

## I Details of Human Evaluation

We conducted a human evaluation using a custom-built interface, with the full interface and instructions shown in Figure 12. A total of 14 participants were recruited, all of whom are PhD-level graduate students or higher. Participation was voluntary. Each participant was asked to evaluate 25 samples, which typically required 15–20 minutes to complete.

The evaluation process, including recruitment, instructions, and data processing and storage, followed our institution’s ethical guidelines for human subject research. All participants were informed of the purpose of the study and provided consent. No

personally identifiable information was collected, and all data were stored and analysed in accordance with privacy standards.

## J Safeguards

C-FDM performs observation prediction through image generation and, while its outputs are task-specific, we acknowledge that any generative model may carry potential for misuse. To mitigate these risks, we commit to the following safeguards upon release:

The model will be released solely for research purposes under a license that prohibits commercial use or any other harmful applications. The GitHub repository will include clear usage guidelines and terms of use, aligned with responsible AI principles.

We will include a disclaimer that the model is intended only for academic research in controlled environments. The datasets used for training are publicly available, action-centric image editing benchmarks that do not include sensitive or personally identifiable content.

Given the targeted nature of our model and the safeguards in place, we believe the risk of misuse is limited. Nonetheless, we encourage responsible use and welcome feedback from the community regarding potential improvements to safety.

## **K LLMs Usage Declaration**

We declare that the large language model (LLM) was only used to assist in minor tasks, including revising the manuscript for grammatical correctness, improving phrasing, and performing small technical implementations such as debugging code snippets. All scientific ideas, results, analyses, and conclusions presented in this paper are entirely the work of the authors.

## Instruction for Editing: let the chair be red

Input Image



## Anonymous Model Outputs



Model 1

Model 2

Model 3

Model 4

## Select the best and worst model according to these three criteria.

Select the BEST/WORST candidate which satisfies/contradicts with the following criteria as many as possible.

If none of them satisfies the criteria, please prioritise in this order:

Criterion 1: Realism > Criterion 2: Instruction Followed > Criterion 3: Over-editing

### Criterion 1: Realism

- **Good:** The generated image looks like a real photo with natural textures and lighting, mostly follows the scene in the input image.
- **Bad:** Artifacts, distortions, or unnatural results.

### Criterion 2: Instruction Followed

- **Good:** The edit reflects the instruction clearly (e.g., "add a tree" results in a tree in the scene).
- **Bad:** The edit misses the point or wrongly changes something irrelevant.

### Criterion 3: Over-editing

- **Good:** The edit is focused and minimal, changing only what was requested.
- **Bad:** The entire image is edited correctly, but more than what was requested is changed (e.g., adding or altering extra objects).

Select the BEST model:

Model 1  Model 2  Model 3  Model 4

Select the WORST model:

Model 1  Model 2  Model 3  Model 4

Submit Evaluation

Figure 12: The screenshot for the instructions given to participants and the interface developed for conducting the evaluation.