

# CAN WORLD MODELS BENEFIT VLMS FOR WORLD DYNAMICS?

Kevin Zhang<sup>1,\*</sup>, Kuangzhi Ge<sup>1,\*</sup>, Xiaowei Chi<sup>2,\*</sup>, Renrui Zhang<sup>3</sup>, Shaojun Shi<sup>1</sup>,  
Zhen Dong<sup>4</sup>, Sirui Han<sup>2,✉</sup>, Shanghang Zhang<sup>1,✉</sup>

<sup>1</sup> Peking University <sup>2</sup> Hong Kong University of Science and Technology

<sup>3</sup> Chinese University of Hong Kong <sup>4</sup> University of California, Santa Barbara

\* Equal contribution, ✉ Corresponding author

## ABSTRACT

Trained on internet-scale video data, generative world models are increasingly recognized as powerful world simulators that can generate consistent and plausible dynamics over structure, motion, and physics. This raises a natural question: *with the advent of strong video foundational models, might they supplant conventional vision encoder paradigms for general-purpose multimodal understanding?* While recent studies have begun to explore the potential of world models on common vision tasks, these explorations typically lack a systematic investigation of generic, multimodal tasks. In this work, we strive to investigate its current capabilities when these priors are transferred into a Vision-Language Model (VLM): we re-purpose a video diffusion model as a *generative encoder*, queried for a single denoising step, and treat the resulting latents as an additional set of visual embeddings. We empirically investigate this class of models, which we refer to as World-Language Models (WorldLMs), and we find that generative encoders can indeed capture latents useful for downstream understanding, showing distinctions from conventional vision encoders. Naming our best-performing variant **Dynamic Vision Aligner (DyVA)**, we further discover that this method significantly enhances spatial reasoning abilities and enables single-image models to perform multi-frame reasoning. Through the curation of a suite of visual reasoning tasks, we find DyVA to surpass both open-source and proprietary baselines on out-of-domain tasks, achieving state-of-the-art or comparable performance. We attribute these gains to WorldLM’s inherited motion-consistency internalization from video pre-training. Finally, we systematically explore extensive model designs to highlight promising directions for future work. We hope our study can pave the way for a new family of VLMs that leverage priors from world models and are on a promising path towards generalist vision learners.

Project page: <https://dyva-worldlm.github.io/>.

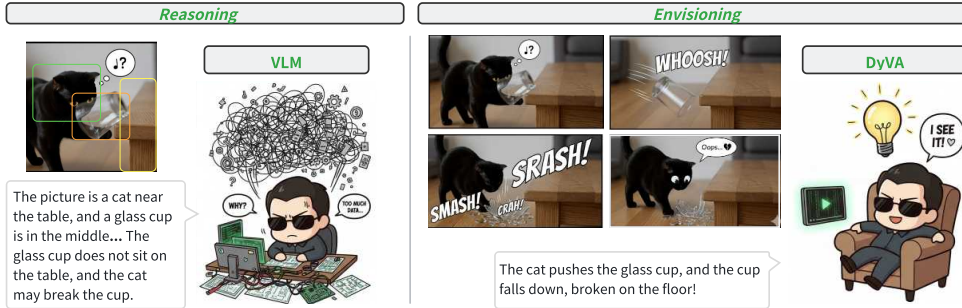


Figure 1: **What will happen?** From reasoning to dynamic intuition — comparing how VLM and WorldLM understand and predict real-world events.

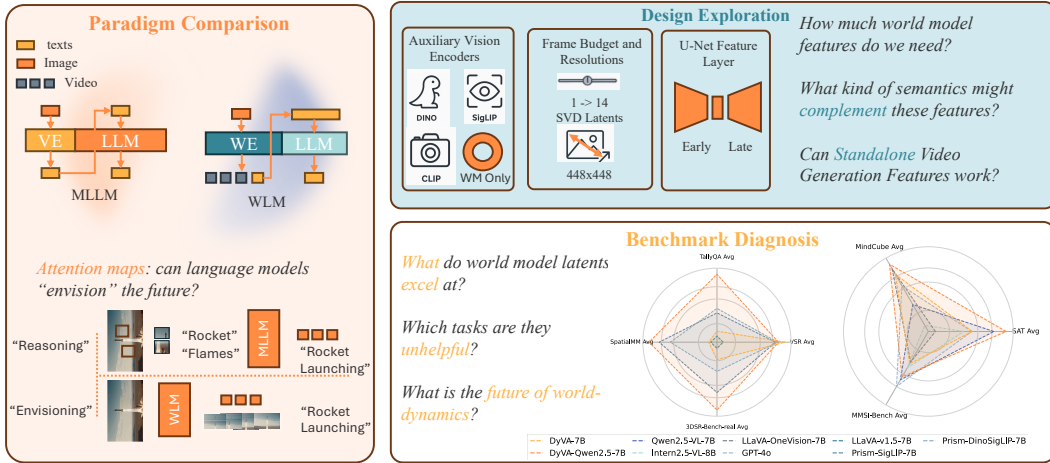


Figure 2: Our analysis is structured around three dimensions: (i) Paradigm comparison between static and generative encoders (e.g., SigLIP vs. SVD); (ii) Benchmark diagnosis, revealing world model latents’ strength (e.g., spatial/multi-frame reasoning) and weaknesses (e.g., language-heavy tasks); and (iii) Design-space exploration, probing different auxiliary encoders, resolutions, and training recipes to understand how world-model features aid visual understanding.

## 1 INTRODUCTION

World models, originally proposed in cognitive science to explain how humans predict and plan in their environments (Tolman, 1948), have recently emerged as powerful tools in machine learning. Generative world models, such as video generation models (VGMs) (Agarwal et al., 2025b; OpenAI, 2024; Wan et al., 2025; Hu et al., 2023; Blattmann et al., 2023; Yang et al., 2025b; Guo\* et al., 2023; 2025; Chen\* et al., 2025) that are trained on internet-scale video data, encode strong priors over objects, spatial layouts, and dynamics. These priors allow them to predict plausible future scenarios that are consistent in 3D structure and physically coherent in motion.

However, a largely overlooked implication of World Models is that the ability to generate coherent futures signals a form of semantic understanding of visual dynamics; this difference between visual generation and understanding has shaped a decade of representation learning. This suggests that world models can be more than generators—they may serve as transferable encoders that enrich downstream tasks with spatial, temporal, and predictive signals. As a result, recent work has attempted to use video generation backbones for visual perception tasks (Acuaviva et al., 2025; Wiedemer et al., 2025).

In this work, we ask a foundational question: *can generative models surpass current vision understanding paradigms for generic, multimodal understanding?*

To empirically investigate the current capabilities of video generation models, we introduce a simple yet effective framework applying them to Vision–Language Models (VLMs). We specifically explore this by evaluating the applicability of predictive world models on a generic multimodal task - Visual Question Answering (VQA) — to assess their broader potential as **generalizable vision encoders**. Currently, mainstream VLMs primarily rely on ViT-based encoders such as CLIP (Radford et al., 2021), SigLIP (Zhai et al., 2023), and DINO (Caron et al., 2021; Oquab et al., 2024), which extract visual semantics from image patches and are then projected as visual tokens into language backbones. While these encoders are semantically aligned, they are limited by temporal reasoning and weaken spatial grounding when multiple views or sequential cues are present. On the other hand, we re-purpose a world model (Stable Video Diffusion, or SVD) as a **Generative Encoder**. Our core mechanism is to extract latent features from a **single denoising step** of its U-Net. This single step, we hypothesize, captures the low-dimensional world-dynamics prior sufficient for downstream understanding. These dynamics-aware latents are then fused with static image features (e.g., SigLIP) and projected into the Large Language Model (LLM). The design is very efficient: all encoders remain frozen, with only lightweight projectors and the LLM being trained.

To this end, we conduct a systematic investigation comparatively evaluating this class of models, which we refer to as World-Language Models (WorldLMs). Our findings are as follows:

- **Shift in Reasoning Paradigm.** The generative prior alters the model’s reasoning process. It moves beyond describing static content to envisioning dynamic possibilities.
- **Zero-shot Multi-Frame Adaptation.** Trained only on single images, the generative encoder enables emergent multi-frame reasoning without multi-image training. On multi-frame visual reasoning, DyVA achieves state-of-the-art or comparable performance with flagship models such as Qwen2.5-VL (Bai et al., 2025) and GPT-4o (OpenAI et al., 2024).
- **We empirically identify the regimes where video priors help.** Our ablations and diagnostics separate the settings in which SVD latents strengthen spatial reasoning from those where they dilute semantic grounding, guiding future designs.

Our best-performing WorldLM variant, **Dynamic Vision Aligner (DyVA)**, exemplifies this paradigm shift. In zero-shot evaluations on challenging multi-frame reasoning benchmarks, DyVA decisively surpasses even proprietary models, for instance, a **28.3%** lead on the **MindCube** benchmark over the GPT-4o model. This provides strong evidence that the ability to predict is a powerful, perhaps essential, foundation for stronger representation learning.

As shown in Figure 2, we systemically organize our investigation revolving around three pillars:

**Paradigm comparison.** World-model encoders versus static encoders reveal distinct strengths: world-model latents benefit spatial and multi-frame reasoning, while static encoders excel on semantics-heavy benchmarks.

**Benchmark diagnostics.** Through curated evaluation sets including MindCube Yin et al. (2025), SAT-Bench Ray et al. (2024), VSR Liu et al. (2023a), we find that DyVA surpasses both open-source and proprietary baselines on out-of-domain tasks, achieving **state-of-the-art performance on MindCube**. Given that SVD is pre-trained on temporally coherent video–text pairs, we show that dynamics-aware latents particularly boost object relations, cross-view understanding, and multi-frame spatial reasoning, while offering less gain on tasks requiring stronger language priors.

**Design-space exploration.** We analyze different encoder setups to identify when predicted latents help or hinder performance and analyzing the co-training of U-Net and VAE layers with text loss, laying the groundwork for a new class of WorldLMs exploiting world-model priors.

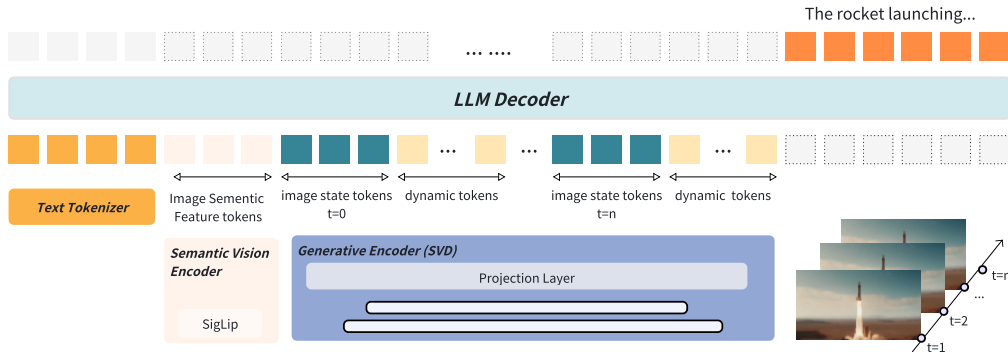


Figure 3: **WorldLM Pipeline.** A SigLIP encoder extracts semantic features from the input image. Concurrently, a Generative Encoder generates dynamic state tokens to capture temporal changes, using evenly spaced keyframe slots. All visual tokens are projected into a shared embedding space, concatenated with text tokens, and then fed into the LLM decoder.

## 2 PRELIMINARY

We first lay a groundwork for our analysis with: 1) a framework to incorporate the dynamic features of a world model into a multimodal language model (which we term **WorldLM**), 2) a training recipe, and 3) an implementation of inference supporting both single and multi-image datasets.

**Framework.** Given an input image  $x_{img} \in \mathbb{R}^{H \times W \times C}$  and a text prompt  $u_{prompt}$ , traditional VLMs such as LLaVA (Liu et al., 2024), QwenVL (Bai et al., 2025), InternVL (Chen et al., 2025), DeepSeekVL (Lu et al., 2024), and Prismatic-VLMs (Karamcheti et al., 2024), process the input with an architecture consisting of three core components:

- **Semantic Vision Encoder.**  $x_{img}$  is processed by a frozen pre-trained ViT-based (Dosovitskiy et al., 2021) encoder  $V_\omega$ , e.g., SigLIP (Zhai et al., 2023), to extract a sequence of feature embeddings  $p_{img} = V_\omega(x_{img})$ , where  $p_{img} \in \mathbb{R}^{L \times d_{vision}}$ , where  $L$  is the token length and  $d_{vision}$  refers to the vision feature dimension.
- **Projector.** The visual features  $p_{img}$  are subsequently mapped into the language model’s embedding space by a projector  $F_\psi$ . This yields a sequence of embeddings  $e_{img} = F_\psi(p_{img})$ , where  $e_{img} \in \mathbb{R}^{L \times d_{text}}$ , where  $d_{text}$  is the text feature dimension. The projector is typically implemented as a simple MLP with GELU activations (Hendrycks & Gimpel, 2023).
- **LLM Backbone.** Finally, the language model  $LM_\theta$  autoregressively generates the textual output  $u_{out}$ . It is conditioned on the concatenated sequence of the projected image features  $e_{img}$  and the text prompt embeddings  $e_{prompt}$ :  $u_{out} = LM_\theta([e_{img}; e_{prompt}])$

On the other hand, in WorldLMs, we employ a Generative Encoder to extract dynamic visual information and motion priors of the input image:

- **Generative Encoder.** We utilize Stable Video Diffusion (SVD) (Blattmann et al., 2023) as our encoder. SVD consists of a VAE (Kingma & Welling, 2022) encoder  $\phi$  and a U-Net (Ronneberger et al., 2015) denoiser  $f_\theta$ . The input image  $x_{img}$  is first embedded by VAE into the latent  $z_0$ , which is then replicated  $T$  times to form the initial video latent  $Z_0$ . A single Euler integration step (Karras et al., 2022) is then applied to yield an updated latent  $Z_1 = Z_0 + \Delta\sigma f_\theta(Z_0, \sigma_0, c)$ . Rather than rendering video frames, the final output  $D_{img} = \text{Hidden}^{\text{mid}}(f_\theta, Z_1)$  is extracted from U-Net’s middle layers.

As shown in Fig. 3, semantic features  $p_{img}$  and dynamic features  $\tilde{H}$  are projected by two separate projectors  $P_{\text{sem}}$  and  $P_{\text{dyn}}$  into the LLM space, yielding  $V_s = P_{\text{sem}}(p_{img}) \in \mathbb{R}^{L_s \times d}$  and  $V_d = P_{\text{dyn}}(\tilde{H}) \in \mathbb{R}^{L_d \times d}$ . The fused sequence is  $V = [V_s; V_d] \in \mathbb{R}^{(L_s + L_d) \times d}$ , which, together with prompt embeddings  $E_{prompt}$ , is fed into the LLM backbone to autoregressively generate answer tokens  $u_{out} = LM_\theta([V; E_{prompt}])$ . By fusing both streams, our WorldLM leverages static semantics (from SigLIP) and dynamics-aware priors (from SVD) for multimodal reasoning.

**Training recipe.** We adopt the training strategy from Prismatic-VLMs (Karamcheti et al., 2024), using single-stage training to align modalities and incorporate dynamic features: We jointly train both the projectors and the LLM on a mixture of multimodal instruction datasets from LLaVA-1.5 (Liu et al., 2023b), together with examples from established vision-language benchmarks (e.g., GQA (Hudson & Manning, 2019), TextCaps (Sidorov et al., 2020)), and language-only samples from ShareGPT (ShareGPT, 2023). This training paradigm not only effectively aligns the representations of the generative encoder with the semantic space of the LLM but also improves its compositional generalization, allowing it to reason over both priors of motion and the static features. Remarkably, the entire training process completes in only 10.3 hours on 16×A800 GPUs ( $\approx 165$  GPU-hours) while achieving competitive performance, underscoring the efficiency of our approach.

**Inference Protocol** During inference, we employ SigLIP-so400m-patch14-224 as the semantic vision encoder and SVD as the generative encoder with an image resolution of  $448 \times 448$ . Shown in Fig. 3, for  $K$  input images, we allocate key frames using evenly spaced indices within the  $T$ -frame latent tensor, replacing the corresponding slots with encoded keyframes before the Euler step, and reuse the resulting latents as visual tokens. For the semantic vision encoder, only the first input image

is encoded and concatenated with the input of the generative encoder. Unless otherwise specified, the number of frames ( $T$ ) is set to 8 for both single-image and multi-image inputs.

Following the proposed framework, training setup, and inference principles, we train a family of WorldLM models and designate the optimal ones in **Dynamic Vision Alignment** as **DyVA**.

### 3 PARADIGM COMPARISON

#### Do WorldLM Encoders Entail Visual Semantics Understanding?

In this section, we explore how world model latents can benefit visual understanding by contrasting two differentiating encoder paradigms: (i) conventional static encoders such as CLIP and SigLIP that prioritize multimodal semantic alignment, and (ii) WorldLM encoders based on video generation models that generate dynamics-aware latents. We begin by comparing the most intuitive design to test if WorldLMs can work, by directly replacing the CLIP vision encoder of LLaVA 1.5 (Liu et al., 2024) with a Generative Encoder (e.g., SVD) following the WorldLM pipeline settings in Fig. 3.

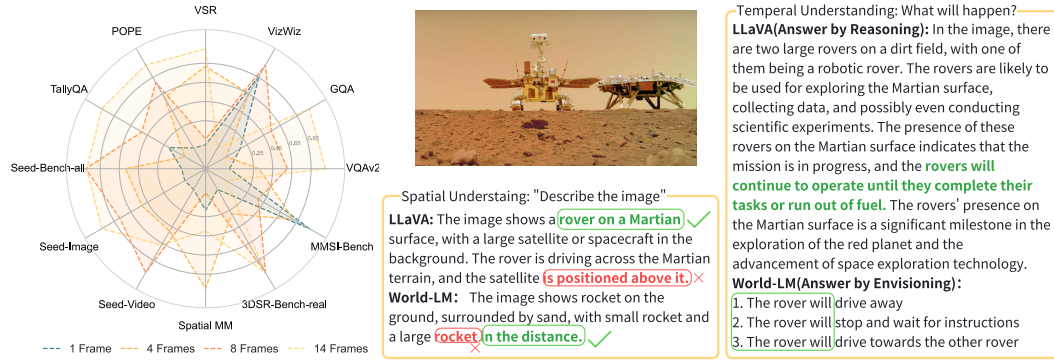


Figure 4: **Paradigm Comparison.** We evaluate predicting 1, 4, 8, and 14 frames with a straight-forward WorldLM setup. The radar chart (left) demonstrates that more frames boosts performance across various tasks, especially in visual reasoning. The qualitative example (right) illustrates that our WorldLM exhibits a distinct reasoning paradigm by envisioning, offering concise descriptions, stronger spatial grounding, and more structured temporal foresight compared to LLaVA.

#### Generative encoders exhibit fundamentally different performance.

We begin with a motivating case study, illustrated in Fig. 4. Models leveraging static encoders, such as LLaVA, adopt a *reasoning* paradigm. The output of LLaVA tends to be more descriptive, describing in depth the details of the given image input. WorldLM, on the other hand, employs an *envisioning* paradigm, which not only encodes the current state of the image, but it also performs a prediction of plausible future conditions (e.g., “will drive away”, and “drive to the other rover”). This case reveals an intrinsic difference between the two paradigms: VLM reasons by the given image’s embeddings, whereas WorldLM attends to depict the embeddings of generated predictions.

**Multi-frame capture more useful semantic features than Single Frame.** The quantitative 4 comparison between using different numbers of generated dynamic latents shows its effect on downstream tasks. When the generated frames of the video prediction model increase from 1 to 14, we see a general rising trend of performance on all tasks.

**Meanwhile, the vanilla WorldLM performs great on spatial-reasoning tasks.** Notably, the gains are most pronounced on benchmarks requiring visual reasoning through space and time, such as *SeedBench*, *VSR*, and *TallyQA*. This demonstrates the potential of using world models as dynamics-aware encoders to allow VLMs gain a deeper and more grounded level of spatial understanding.

**Limitations of WorldLMs.** Despite the clear advantages in spatial reasoning, our empirical study reveals a critical trade-off in Fig. 4, where its performance relative to LLaVA is much lower across all tasks. The case study in Fig. 4 offers a qualitative explanation for this phenomenon. While our world model correctly grounds the spatial structure of the scene (e.g., “rocket on the ground... large

rocket in the distance”), it hallucinates the semantic identity of the objects, misidentifying the Mars lander and rover as “rockets”. Therefore, we believe that using a world model as an encoder has the potential to enhance predictive and spatial reasoning tasks, but requires further improvement to ensure basic semantic capabilities.

## 4 BENCHMARK ANALYSIS: INVESTIGATION

### 4.1 EXPERIMENTAL SETUP

We document the configurations, datasets, and training protocols underlying our study. Unless otherwise noted, all settings use a 7B-parameter LLaMA-2 LLM backbone, with both SigLIP and SVD encoders frozen during a single-stage instruction tuning. Training updates are restricted to lightweight projection layers and the language backbone.

### 4.2 DATASETS AND EVALUATION TARGETS

Benchmarks vary widely in their emphasis on *spatial grounding*, *temporal coherence*, and *semantic understanding*. To assess these dimensions, we curate a suite of open-source **out-of-domain (OOD)** datasets on which our models have not been trained. This allows us to isolate the transferability of world-model priors.

**Single-image spatial reasoning.** We evaluate on benchmarks that probe relational and spatial understanding without temporal context, including VSR (Liu et al., 2023a), TallyQA (Acharya et al., 2018), SpatialMM-Obj (Shiri et al., 2024), and 3DSR-Bench-real (Ma et al., 2025). Baselines include LLaVA-1.5 (Liu et al., 2024), Prism-SigLIP-7B (Karamcheti et al., 2024), and Prism-DinoSigLIP-7B (Karamcheti et al., 2024).

**Multi-image and temporal reasoning.** To assess robustness to sequential inputs and temporal structure, we use MMSI-Bench (Yang et al., 2025a), SAT-Synthetic (Ray et al., 2024), and MindCube (Yin et al., 2025). These benchmarks require models to integrate cues across frames or view-points, testing whether world-model latents can enable multi-frame reasoning. We compare against both open-source and proprietary large-scale VLMs, including Qwen-2.5-VL-7B (Bai et al., 2025), InternVL-2.5-7B (Chen et al., 2025), LLaVA-OneVision-7B (Li et al., 2024), and GPT-4o (OpenAI et al., 2024). Note that all of the compared benchmarks are trained with multi-frame or video data, whereas we train on single images only.

### 4.3 EXPERIMENTAL ANALYSIS AND INSIGHTS

Table 1: Performance comparison between DyVA and state-of-the-art methods on multi-image benchmarks SAT Synthetic, MMSI-Bench, and MindCube. DyVA outperforms baselines in these OOD tasks without training on multi-image datasets. The highest average values are in bold.

Model	SAT Synthetic						MindCube				
	Obj Move.	Act. Seq.	Act. Cons.	Goal Aim	Persp.	Avg.	Rot.	Among	Around	Avg.	
Qwen2.5-VL-7B	79.29	84.70	47.83	25.84	35.17	53.16	38.76	29.50	21.35	29.26	
Intern2.5-VL-8B	77.74	55.49	53.74	15.03	32.61	48.06	18.68	36.45	18.20	18.68	
LLaVA-OneVision-7B	71.10	21.64	49.85	31.76	35.43	43.24	36.45	48.42	44.09	47.43	
GPT-4o	61.50	33.20	47.60	67.50	37.50	49.40	40.17	29.16	38.81	38.81	
DyVA-7B	49.15	57.81	49.25	53.38	40.44	49.51	37.70	43.10	49.00	44.62	
DyVA-Qwen2.5-7B	78.83	62.13	49.85	51.86	41.72	55.24	37.20	39.10	51.70	49.80	

MMSI-Bench												
Model	Positional Relationship						Attribute		Motion		MSR	Avg.
	Cam-Cam	Obj-Obj	Reg-Reg	Cam-Obj	Obj-Reg	Cam-Reg	Means	Appr	Cam	Obj		
Qwen2.5-VL-7B	32.3	27.7	29.6	32.6	24.7	32.5	26.6	27.3	16.2	31.6	30.3	28.70
Intern2.5-VL-8B	24.7	24.5	24.7	25.6	29.4	26.5	25.0	18.2	20.3	39.5	25.8	25.90
LLaVA-OneVision-7B	20.4	33.0	29.6	29.1	25.9	30.1	29.7	25.8	18.9	34.2	11.6	24.50
GPT-4o	34.4	24.5	23.5	19.8	37.6	27.7	32.8	31.8	35.1	36.8	30.8	30.30
DyVA-7B	21.5	30.9	25.9	31.4	27.1	20.5	35.9	24.2	13.5	19.7	24.2	24.90
DyVA-Qwen2.5-7B	15.1	33.0	25.9	33.7	35.3	30.1	32.8	25.8	17.6	27.6	29.3	28.00



Table 2: Performance comparison of DyVA variants against baselines on various single-image spatial reasoning benchmarks, including VSR, TallyQA, SpatialMM-Obj, and 3DSR-Bench-real. These are Out-of-Domain tasks where models are not trained and perform zero-shot inference. Our results surpass all baseline models. Highest values are highlighted in bold.

Models	Data	VSR							
		Topo.	Prox.	Proj.	Direc.	Adj.	Orien.	Unall.	Avg.
LLaVA-v1.5-7B	558k+665k	52.24	50.00	54.77	50.00	50.86	48.98	57.50	52.94
Prism-SigLIP-7B	665k	67.48	62.50	65.63	66.67	55.17	55.10	67.50	64.97
Prism-DinoSigLIP-7B	665k	71.34	59.38	65.63	64.29	53.45	48.98	52.50	65.46
<b>DyVA-7B</b>	665k	68.90	68.75	66.74	66.67	66.38	61.22	57.50	<b>67.10</b>
<b>DyVA-Qwen2.5-7B</b>	665k	66.67	71.88	68.74	61.90	62.93	40.82	55.00	65.63

Models	TallyQA	SpatialMM-Obj			3DSR-Bench-real				
	Avg.	1-obj	2-obj	Avg.	H.	L.	O.	M.	Avg.
LLaVA-v1.5-7B	58.74	57.37	44.87	48.91	55.42	57.82	26.09	39.42	45.02
Prism-SigLIP-7B	62.25	62.54	46.77	51.86	52.28	60.22	27.23	42.17	46.55
Prism-DinoSigLIP-7B	62.93	58.56	47.72	51.22	56.85	59.42	27.23	38.97	45.82
<b>DyVA-7B</b>	59.47	54.78	46.29	49.03	53.71	57.60	27.23	40.80	45.41
<b>DyVA-Qwen2.5-7B</b>	<b>68.11</b>	62.74	47.53	<b>52.44</b>	52.57	54.51	27.23	49.60	<b>47.16</b>

Tab. 1 and 2 present representative results under both single and multi-image settings. This framing allows us to disentangle how world-model features contribute across different reasoning regimes.

As presented in Tab. 1 and 2, we evaluate the OOD performance of DyVA-7B and DyVA-Qwen2.5-7B. We examine DyVA’s performance relative to existing vision-language models across various spatial reasoning tasks. The key differences lie in DyVA’s use of Generative Encoders versus baselines that use only standard visual embeddings. Below, we discuss the strengths and weaknesses of DyVA in each benchmark category, drawing on the reported results of these tasks and models.

**DyVA can enable single-image trained WorldLMs to perform multi-image tasks exceptionally well.** As in Tab. 1, our best variant can perform strongly in multi-frame spatial understanding tasks.

Specifically, on the MindCube benchmark (Tab. 1), DyVA-Qwen2.5-7B achieves a new state-of-the-art performance with the highest overall score (49.8% vs. 47.4% for the runner-up baseline). It particularly excels in “Around” (rotating viewpoint) tasks (51.7% vs. 44.1%) and matches or slightly exceeds baselines on “Rot” tasks (37% vs. 36%). These results suggest that DyVA latents significantly aid in tasks requiring mental rotation and perspective-taking, likely because they encode cross-view consistency. Specifically, these margins are consistent with the motion-consistency priors inherited from SVD’s pre-training on LVD-F video-text pairs Blattmann et al. (2023) that include how an object may appear from different angles.

This achievement is especially noteworthy considering the training efficiency. Compared to baselines where LLaVA-One-Vision is trained on 4M multi-frame images, Intern 2.5-VL is pretrained with 16.3M samples, including multi-image and video data, and Qwen-2.5-VL is also pre-trained with a variety of data comprising videos and multi-images. These baselines also have several complex methods for image preprocessing, such as patchifying (Li et al., 2024), processing at different fps (Bai et al., 2025), and high-res processing (Chen et al., 2025). In contrast, we trained our DyVA model using only the most basic processing methods with minimal amount of data.

Our modest training budget and intuitive multi-image inference method suggest that world model latents strongly enhance the spatial understanding on multi-image benchmarks. We also believe that the fusion of SVD with SigLIP is a key factor that directly improves multi-image reasoning abilities.

**DyVA excels in handling spatial relations, counting and object queries, and 3D Scenes.** In Single-Image Spatial Reasoning, DyVA’s world-model features boost performance on tasks emphasizing geometric and relational spatial reasoning (orientation, adjacency, multi-object spatial layouts), reflecting improved 3D awareness.

1. **Visual Spatial Relations (VSR):** DyVA (SigLIP+SVD) achieves the highest average score (67.1%) across VSR subtasks (topology, proximity, projection, direction, adjacency, orientation, unaligned), outperforming the SigLIP-only baselines (64.9–65.5%) in Tab. 2. In particular, DyVA significantly improves orientation reasoning (61.2% vs 55–49% for baselines) and proximity/topology, suggesting it can better encode spatial layouts and object alignment.

2. **Counting and Object Queries (TallyQA, SpatialMM-Obj):** On TallyQA (visual counting), DyVA-Qwen2.5 excels (68.1% average), well above Prismatic baselines (62–63%) and LLaVA (58.7%) Tab. 2. For the SpatialMM-Obj task (single- vs multi-object queries), DyVA-Qwen2.5 again slightly outperforms others (52.4% vs 51.8% baseline) on the combined 1- and 2-object questions.

3. **3D Scene Reasoning (3DSR-Bench-real):** This benchmark measures 3D spatial and depth understanding in real images. Notably, DyVA greatly improves the “Multiple objects” (M) subset (49.6% vs 40% for baselines). This aligns with the conception that SVD latents capture implicit depth and occlusion cues learned from video modeling.

**Limitations and Areas for Improvement.** Despite its strengths in spatial reasoning, DyVA exhibits certain limitations, particularly on tasks that rely heavily on semantic language priors, non-canonical object arrangements, or temporal sequence understanding.

1. **Weakened Performance on Language-Intensive Tasks:** The fusion of world-model tokens can dilute the semantic precision required for certain tasks. As shown in Tab. 3, on benchmarks such as VQAv2 and TextVQA, which demand strong language priors and OCR capabilities, DyVA underperforms compared to SigLIP-only baselines. This suggests that while SVD latents enhance spatial awareness, they can interfere with fine-grained semantic grounding and text recognition, where the original visual features are more direct and precise.

2. **Bias Towards Canonical Scene Structures:** As previously noted in the VSR analysis, DyVA’s performance drops significantly on the “Unaligned” subtask (57.5% vs. 67.5%). This indicates that embedding world-model context can be detrimental when objects lack canonical alignments. The model’s latent prior appears biased toward common or expected scene structures, hindering its ability to reason about novel or unusual spatial configurations.

3. **Less Reliable Sequential and Temporal Reasoning:** The current SVD latents are less effective for understanding dynamic sequences. This is evidenced by a large performance drop in SAT Action Sequence and mixed results on MMSI. These outcomes suggest that the latents, while powerful for static scenes, are less reliable for predicting discrete action orders or interpreting rapid changes over time, marking a clear area for future improvement.

## 5 DESIGN-SPACE EXPLORATION: WHY DYVA WORKS

### **Generative Encoders rely on both dynamic frames and text-aligned semantics as support.**

Building on the strong spatial performance demonstrated in both single-image and multi-image tasks in our experiments, we further analyze two key design axes to investigate the sources of WorldLM’s benefits: (i) the choice of different semantic vision encoders, and (ii) the potential of leveraging text-loss to supervise the joint-training of VAE and U-Net in SVD.

#### 5.1 WHY DO VAE, DINO, SVD-ONLY NOT WORK, BUT SIGLIP+SVD DOES?

To investigate the respective roles of the generative encoder and the semantic vision encoder within WorldLM, we conduct a two-stage ablation study. **First**, in a setting without the semantic vision encoder, we decouple the generative encoder into its component VAE and the complete generative encoder architecture. We then train and comparatively evaluate the performance of two distinct encoding approaches: one employing only the VAE for encoding and the other utilizing the entire generative encoder (SVD). **Second**, while keeping the generative encoder fixed, we systematically substitute the backbone of the semantic vision encoder with various alternative architectures to analyze its impact on the model’s overall performance. Our quantitative experimental results are presented in Tab. 3.

**Prediction Matters.** The inference protocol for the SVD encoder is detailed in Sec. 2. A similar inference process is employed when using VAE as the generative encoder. In contrast to extracting



Table 3: **Performance Comparison of SVD-based Vision Models.** Benchmark scores across a set of VQA, reasoning, and spatio-temporal tasks. All experiments use the LLaMA-2 7B backbone. The highest score in each column is marked in **bold**, and the second-highest is underlined. Align: one-time alignment on LAION-558k Schuhmann et al. (2022). F1: one-time finetuning. Fused: 3-layer MLP projector.

Model	Align	VQAv2	GQA	VizWiz	VSR	POPE	TallyQA	SeedBench	SpatialMM	3DSR
VAE-Only	×	46.98	40.53	38.90	52.04	66.42	39.55	38.18	38.81	44.15
	✓	50.70	43.26	48.67	52.29	60.80	42.48	41.53	37.3	43.43
SVD-Only	×	63.51	55.18	44.95	57.93	82.38	49.75	50.15	42.03	42.93
	✓	61.82	50.20	50.60	53.60	75.61	53.27	52.55	40.60	43.50
U-Net Trainable	✓	63.36	54.49	50.24	57.93	79.88	51.51	52.76	40.80	43.43
U-Net & VAE Trainable	✓	60.99	49.80	50.17	52.53	77.08	53.75	52.33	39.50	44.00
Dino + SVD	×	68.77	58.50	50.73	62.52	85.25	52.78	55.19	44.79	44.26
	✓	68.44	55.57	51.13	59.41	85.54	54.15	56.49	43.40	45.07
SigLIP + SVD	×	<b>75.36</b>	<b>61.52</b>	<b>55.95</b>	<b>67.10</b>	85.97	<b>59.47</b>	<b>66.61</b>	<b>49.03</b>	45.40
	✓	73.63	58.89	54.63	61.62	84.37	56.98	62.09	45.40	<u>45.49</u>
U-Net Trainable	✓	74.02	59.86	54.60	62.27	85.61	57.42	63.39	45.95	44.11
CLIP + SVD	×	73.51	59.67	53.14	64.89	85.80	58.25	<u>65.45</u>	46.07	<b>46.13</b>
	✓	72.99	<u>60.74</u>	<u>55.89</u>	<u>65.38</u>	85.80	55.37	65.33	46.70	44.42
DinoSigLIP + SVD	×	74.28	60.16	54.13	64.81	<b>87.27</b>	57.42	64.54	<u>48.65</u>	44.15
	✓	72.42	59.28	54.47	61.29	<u>86.75</u>	54.98	61.54	47.00	45.14

features from the layer before the middle block of U-Net, we directly use the features encoded by the VAE. To align the feature dimensionality with that of the SVD, we prepend several convolutional layers to the projector. As evidenced by our experimental results in Tab. 3, the model employing only VAE for encoding exhibits a performance degradation across nearly all benchmarks when compared to models using SVD. This finding underscores the significance of the predicted dynamics for the WorldLM.

**WorldLMs need a text-aligned encoder.** Although SigLIP (Zhai et al., 2023) has recently shown dominant performance as an emerging vision encoder in current state-of-the-art VLMs, such as LLaVA-One-Vision (Li et al., 2024) and Prismatic-VLM (Karamcheti et al., 2024), in this study, we investigate the respective roles of SigLIP, CLIP (Radford et al., 2021), DINOv2 (Oquab et al., 2024), and a combined DINO-SigLIP architecture as the semantic vision encoder. To ensure a fair comparison, we selected the ViT-L version for each model, all configured for a  $224 \times 224$  input resolution. Furthermore, we adopted a consistent image processing strategy, which involves scaling and then cropping all images to uniform resolutions.

As demonstrated in Tab. 3, models that utilize SigLIP (including the DINO-SigLIP combination) or CLIP as the semantic vision encoder significantly outperform the model using DINOv2. Furthermore, when considering the aforementioned investigation of the generative encoder, the model with DINOv2 as the semantic vision encoder shows better performance than the generative-encoder-only architecture.

This leads to a key insight: for our WorldLM framework that is trained with text-loss supervision, in addition to predicted dynamic features, DyVA requires supplementary visual-semantic information from a model pre-trained on language-vision tasks (i.e., a text-aligned model). This insight also paves the way for future explorations: Can the generative encoder alone suffice to replace the semantic vision encoder? And is text-loss supervision the answer to WorldLM training?

## 5.2 CAN DYVA BENEFIT FROM U-NET & VAE TRAINING ON TEXT-LOSS?

We investigated the efficacy of fine-tuning the SVD’s core components (U-Net and VAE) using only a text-loss signal. Our experimental results indicate this strategy is largely ineffective.

**Text supervision failed to help VQA tasks.** As shown in Tab. 3, making only the U-Net trainable yields inconsistent and marginal performance changes, while allowing both the U-Net and VAE to be trainable leads to a distinct and widespread degradation in performance across the benchmarks.

This suggests the high-level semantic supervision from the text-loss is ill-suited for adapting the low-level generative priors of these components. This constitutes one of the limitations of our current work. An alternative approach, inspired by methods like RAPE-E (Leng et al., 2025), involves

aligning the features from the VAE and U-Net with the visual features from a semantic encoder such as DINOv2. Exploring such an alignment strategy is a promising direction for future research.

## 6 DISCUSSIONS AND OUTLOOKS

Over the recent past, video foundation models have demonstrated remarkable performance in key areas such as consistency and content generation. Through our empirical investigation on multimodal general tasks through a VLM framework, we find that:

- (1) Paradigm comparisons reveal that WorldLM latents are powerful: these latents enable effective spatial and multi-view reasoning.
- (2) Design-space explorations clarify which architectural choices benefit WorldLMs, while benchmark diagnostics explain where DyVA excels.
- (3) WorldLM encoders unlock visual reasoning through leveraging SVD’s predictive pre-training supplies transferable camera-motion and interaction priors, yet semantic gaps persist until the generative encoders are co-trained or better aligned with language signals.

Overall, we observe that WorldLM encoders offer a reliable pathway to stronger spatial and multi-view reasoning, and scaling trends in video generation (Wiedemer et al., 2025; Chi et al., 2025) suggest the semantic deficit may narrow as these models progress. Closing that gap for VLMs will likely require tighter alignment between dynamics-rich latents and language-grounded objectives.

**Outlooks.** Promising next steps include: (i) exploring text-to-video generators as encoders to test whether text-aligned priors further boost visual understanding; and (ii) designing specialized training that aligns generative latents with semantics without eroding their physical fidelity.

## REFERENCES

- Manoj Acharya, Kushal Kafle, and Christopher Kanan. Tallyqa: Answering complex counting questions, 2018. URL <https://arxiv.org/abs/1810.12440>.
- Pablo Acuaiviva, Aram Davtyan, Mariam Hassan, Sebastian Stapf, Ahmad Rahimi, Alexandre Alahi, and Paolo Favaro. From generation to generalization: Emergent few-shot learning in video diffusion models, 2025. URL <https://arxiv.org/abs/2506.07280>.
- Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chatopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, Daniel Dworakowski, Jiaojiao Fan, Michele Fenzi, Francesco Ferroni, Sanja Fidler, Dieter Fox, Songwei Ge, Yunhao Ge, Jinwei Gu, Siddharth Gururani, Ethan He, Jiahui Huang, Jacob Huffman, Pooya Jannaty, Jingyi Jin, Seung Wook Kim, Gergely Klár, Grace Lam, Shiyi Lan, Laura Leal-Taixé, Anqi Li, Zhaoshuo Li, Chen-Hsuan Lin, Tsung-Yi Lin, Huan Ling, Ming-Yu Liu, Xian Liu, Alice Luo, Qianli Ma, Hanzi Mao, Kaichun Mo, Arsalan Mousavian, Seungjun Nah, Sriharsha Niverty, David Page, Despoina Paschalidou, Zeeshan Patel, Lindsey Pavao, Morteza Ramezanali, Fitsum Reda, Xiaowei Ren, Vasanth Rao Naik Sabavat, Ed Schmerling, Stella Shi, Bartosz Stefaniak, Shitao Tang, Lyne Tchampi, Przemysław Tredak, Wei-Cheng Tseng, Jibin Varghese, Hao Wang, Haoxiang Wang, Heng Wang, Ting-Chun Wang, Fangyin Wei, Xinyue Wei, Jay Zhangjie Wu, Jiashu Xu, Wei Yang, Lin Yen-Chen, Xiaohui Zeng, Yu Zeng, Jing Zhang, Qinsheng Zhang, Yuxuan Zhang, Qingqing Zhao, and Artur Zabkowski. Cosmos world foundation model platform for physical ai. *CoRR*, abs/2501.03575, 2025a. doi: 10.48550/arXiv.2501.03575.
- Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chatopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025b.
- Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Mojtaba Komeili, Matthew J. Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zhoulus, Sergio Arnaud, Abha Gejji, Ada Martin, François R. Hogan, Daniel Dugas, Piotr Bojanowski, Vasil Khalidov, Patrick Labatut, Francisco Massa, Marc Szafraniec, Kapil Krishnakumar, Yong Li, Xiaodong

- Ma, Sarath Chandar, Franziska Meier, Yann LeCun, Michael Rabbat, and Nicolas Ballas. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *CoRR*, abs/2506.09985, 2025. doi: 10.48550/arXiv.2506.09985.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL <https://arxiv.org/abs/2502.13923>.
- Yutong Bai, Xinyang Geng, Karttikeya Mangalam, Amir Bar, Alan L. Yuille, Trevor Darrell, Jitendra Malik, and Alexei A. Efros. Sequential modeling enables scalable learning for large vision models. *arXiv preprint arXiv:2312.00785*, 2024. doi: 10.48550/arXiv.2312.00785.
- Dmitry Baranchuk, Ivan Rubachev, Andrey Voynov, Valentin Khruikov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models, 2022. URL <https://arxiv.org/abs/2112.03126>.
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large datasets. *CoRR*, abs/2311.15127, 2023. doi: 10.48550/arXiv.2311.15127.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers, 2021. URL <https://arxiv.org/abs/2104.14294>.
- Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. Diffusiondet: Diffusion model for object detection. *arXiv preprint arXiv:2211.09788*, 2022.
- Xinyan Chen\*, Renrui Zhang\*, Dongzhi Jiang, Aojun Zhou, Shilin Yan, Weifeng Lin, and Hongsheng Li. Mint-cot: Enabling interleaved visual tokens in mathematical chain-of-thought reasoning. *arXiv preprint arXiv:2506.05331*, 2025.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang, Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo, Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Botian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong Wu, Huipeng Deng, Jiaye Ge, Kai Chen, Kaipeng Zhang, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhao Wang. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling, 2025. URL <https://arxiv.org/abs/2412.05271>.
- Xiaowei Chi, Peidong Jia, Chun-Kai Fan, Xiaozhu Ju, Weishi Mi, Kevin Zhang, Zhiyuan Qin, Wanxin Tian, Kuangzhi Ge, Hao Li, Zezhong Qian, Anthony Chen, Qiang Zhou, Yueru Jia, Jiaming Liu, Yong Dai, Qingpo Wu, Chengyu Bai, Yu-Kai Wang, Ying Li, Lizhang Chen, Yong Bao, Zhiyuan Jiang, Jiacheng Zhu, Kai Tang, Ruichuan An, Yulin Luo, Qiuxuan Feng, Siyuan Zhou, Chi min Chan, Chengkai Hou, Wei Xue, Sirui Han, Yike Guo, Shanghang Zhang, and Jian Tang. Wow: Towards a world omniscient world model through embodied interaction, 2025. URL <https://arxiv.org/abs/2509.22642>.
- Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified multimodal pretraining, 2025. URL <https://arxiv.org/abs/2505.14683>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL <https://arxiv.org/abs/2010.11929>.
- Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng Zhang, Han Hu, Dong Chen, and Baining Guo. Instructdiffusion: A generalist modeling interface for vision tasks. *CoRR*, abs/2309.03895, 2023. doi: 10.48550/arXiv.2309.03895.

- Zhangxuan Gu, Haoxing Chen, Zhuoer Xu, Jun Lan, Changhua Meng, and Weiqiang Wang. Diffusioninst: Diffusion model for instance segmentation, 2022. URL <https://arxiv.org/abs/2212.02773>.
- Ziyu Guo\*, Renrui Zhang\*, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023.
- Ziyu Guo\*, Renrui Zhang\*, Chengzhuo Tong\*, Zhizheng Zhao\*, Peng Gao, Hongsheng Li, and Pheng-Ann Heng. Can we generate images with cot? let’s verify and reinforce image generation step by step. *CVPR 2025*, 2025.
- David Ha and Jürgen Schmidhuber. World models. *CoRR*, abs/1803.10122, 2018. doi: 10.48550/arXiv.1803.10122.
- Danijar Hafner, Timothy P. Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. *CoRR*, abs/1811.04551, 2018. doi: 10.48550/arXiv.1811.04551.
- Danijar Hafner, Timothy P. Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *CoRR*, abs/1912.01603, 2019. doi: 10.48550/arXiv.1912.01603.
- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023. URL <https://arxiv.org/abs/1606.08415>.
- Geoffrey E. Hinton. To recognize shapes, first learn to generate images. In Paul Cisek, Trevor Drew, and John F. Kalaska (eds.), *Computational Neuroscience: Theoretical Insights into Brain Function*, volume 165 of *Progress in Brain Research*, pp. 535–547. Elsevier, 2007. doi: [https://doi.org/10.1016/S0079-6123\(06\)65034-6](https://doi.org/10.1016/S0079-6123(06)65034-6). URL <https://www.sciencedirect.com/science/article/pii/S0079612306650346>.
- Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving, 2023. URL <https://arxiv.org/abs/2309.17080>.
- Drew A. Hudson and Christopher D. Manning. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models, 2024. URL <https://arxiv.org/abs/2402.07865>.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models, 2022. URL <https://arxiv.org/abs/2206.00364>.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022. URL <https://arxiv.org/abs/1312.6114>.
- Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng. Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers, 2025. URL <https://arxiv.org/abs/2504.10483>.
- Alexander C. Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion model is secretly a zero-shot classifier, 2023. URL <https://arxiv.org/abs/2303.16203>.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024. URL <https://arxiv.org/abs/2408.03326>.

- Yijing Lin, Mengqi Huang, Shuhan Zhuang, and Zhendong Mao. Realgeneral: Unifying visual generation via temporal in-context learning with video models. *arXiv preprint arXiv:2503.10406*, 2025. URL <https://arxiv.org/abs/2503.10406>.
- Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. *Transactions of the Association for Computational Linguistics*, 11:635–651, 2023a. doi: 10.1162/tacl.a.00566. URL <https://aclanthology.org/2023.tacl-1.37/>.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b. URL <https://arxiv.org/abs/2304.08485>.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2024. URL <https://arxiv.org/abs/2310.03744>.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. Deepseek-vl: Towards real-world vision-language understanding, 2024.
- Wufei Ma, Haoyu Chen, Guofeng Zhang, Yu-Cheng Chou, Jieneng Chen, Celso M de Melo, and Alan Yuille. 3dsrbench: A comprehensive 3d spatial reasoning benchmark, 2025. URL <https://arxiv.org/abs/2412.07825>.
- OpenAI. Video generation models as world simulators, 2024. URL <https://openai.com/index/video-generation-models-as-world-simulators/>. Accessed: 2025-09-10.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codisposi, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Gierler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khosravan, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood,

- Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljube, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyei Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. URL <https://arxiv.org/abs/2304.07193>.
- Mihir Prabhudesai, Tsung-Wei Ke, Alexander C. Li, Deepak Pathak, and Katerina Fragkiadaki. Diffusion-tta: Test-time adaptation of discriminative models via generative feedback, 2023. URL <https://arxiv.org/abs/2311.16102>.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. doi: 10.48550/arXiv.2103.00020.
- Arijit Ray, Jiafei Duan, Reuben Tan, Dina Bashkirova, Rose Hendrix, Kiana Ehsani, Aniruddha Kembhavi, Bryan A. Plummer, Ranjay Krishna, Kuo-Hao Zeng, and Kate Saenko. Sat: Spatial aptitude training for multimodal language models. *arXiv preprint arXiv:2412.07755*, 2024.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation, 2015. URL <https://arxiv.org/abs/1505.04597>.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev.



- Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. URL <https://arxiv.org/abs/2210.08402>.
- ShareGPT. Sharegpt, 2023.
- Fatemeh Shiri, Xiao-Yu Guo, Mona Golestan Far, Xin Yu, Reza Haf, and Yuan-Fang Li. An empirical analysis on spatial reasoning capabilities of large multimodal models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 21440–21455, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1195. URL <https://aclanthology.org/2024.emnlp-main.1195/>.
- Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioning with reading comprehension, 2020. URL <https://arxiv.org/abs/2003.12462>.
- Edward C Tolman. Cognitive maps in rats and men. *Psychological Review*, 55(4):189–208, 1948.
- Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Fei Wu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- Zhendong Wang, Yifan Jiang, Yadong Lu, Yelong Shen, Pengcheng He, Weizhu Chen, Zhangyang Wang, and Mingyuan Zhou. In-context learning unlocked for diffusion models. *arXiv preprint arXiv:2305.01115*, 2023. URL <https://arxiv.org/abs/2305.01115>.
- Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane Gu, Nick Matarese, Kevin Swersky, Been Kim, Priyank Jaini, and Robert Geirhos. Video models are zero-shot learners and reasoners, 2025. URL <https://arxiv.org/abs/2509.20328>.
- Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models, 2023. URL <https://arxiv.org/abs/2303.04803>.
- Sihan Yang, Runsen Xu, Yiman Xie, Sizhe Yang, Mo Li, Jingli Lin, Chenming Zhu, Xiaochen Chen, Haodong Duan, Xiangyu Yue, Dahua Lin, Tai Wang, and Jiangmiao Pang. Mmsi-bench: A benchmark for multi-image spatial intelligence. *arXiv preprint arXiv:2505.23764*, 2025a.
- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Yuxuan Zhang, Weihang Wang, Yean Cheng, Bin Xu, Xiaotao Gu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video diffusion models with an expert transformer. In *The Thirteenth International Conference on Learning Representations*, 2025b. URL <https://openreview.net/forum?id=LQzN6TRFg9>.
- Baiqiao Yin, Qineng Wang, Pingyue Zhang, Jianshu Zhang, Kangrui Wang, Zihan Wang, Jieyu Zhang, Keshigeyan Chandrasegaran, Han Liu, Ranjay Krishna, Saining Xie, Manling Li, Jijun Wu, and Li Fei-Fei. Spatial mental modeling from limited views. *arXiv preprint arXiv:2506.21458*, 2025.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. URL <https://arxiv.org/abs/2303.15343>.
- Gaoyue Zhou, Hengkai Pan, Yann LeCun, and Lerrel Pinto. Dino-wm: World models on pre-trained visual features enable zero-shot planning. *CoRR*, abs/2411.04983, 2024. doi: 10.48550/arXiv.2411.04983.

## OUTLINE

- **Related Work (Section A):** Reviews prior work in three key areas: (1) predictive World Models, (2) diffusion-based Generalist Models for in-context learning, and (3) the application of diffusion models to discriminative vision tasks.
- **Model Formalization:** Details our architecture, including:
  - Static visual features from a SigLIP encoder (Eq. 1).
  - Dynamic features from SVD U-Net hidden states (Eq. 2, 3).
  - The fusion mechanism for static and dynamic tokens (Eq. 4).
- **Training Hyperparameters:** Specifies all training configurations, which are listed in Table 4.
- **Design Space Explorations:** Presents key ablation studies, demonstrating:
  - The model’s sensitivity to temporal frames over spatial resolution (Table 5).
  - The rationale for our SVD feature fusion strategy, with comparative results in Table 6.

## A RELATED WORK

### A.1 WORLD MODELS

Various methods have been developed to learn predictive models of visual dynamics. Ha & Schmidhuber (2018) proposed the original World Models framework, which learns a compressed latent representation of an environment’s dynamics using generative RNNs (Ha & Schmidhuber, 2018). Hafner introduced PlaNet (Hafner et al., 2018) and later Dreamer (Hafner et al., 2019), which use latent space dynamics models trained on pixel observations for planning and control. More recently, large-scale self-supervised video models have emerged. For example, Stability AI’s Stable Video Diffusion trains a high-capacity latent video diffusion model on vast video datasets for high-quality text-to-video and image-to-video generation (Blattmann et al., 2023). Zhou (2024) introduced DINO-WM, a world model that leverages pretrained DINOv2 patch features to enable zero-shot goal-reaching via planning in feature space (Zhou et al., 2024). Meta’s V-JEPA 2 (Assran et al., 2025) and NVIDIA’s Cosmos platform (Agarwal et al., 2025a) provide video foundation models that enable understanding, prediction, and planning from raw visual data.

### A.2 GENERALIST MODELS

Recent work has explored using diffusion-based generative models for flexible multi-task and in-context learning. Wang et al. (2023) presented Prompt Diffusion, a method that enables in-context learning in diffusion models by conditioning on example input-output image pairs and a text prompt. Geng et al. (2023) proposed InstructDiffusion, a unified framework that casts diverse vision tasks as a pixel-space image manipulation guided by human instructions, learned via a diffusion process. Bai et al. (2024) introduced a sequential modeling approach that represents images and annotations as “visual sentences,” enabling training a single large vision model across many tasks without using any language data. Lin et al. (2025) presented RealGeneral, which reformulates image generation as conditional frame prediction analogous to LLM in-context learning: using video diffusion models with novel modules, they unify multiple image-generation tasks (e.g., custom generation, canny-to-image) within one framework. Recently, Bagel (Deng et al., 2025) further extends these ideas by introducing novel techniques for improving the generalization and efficiency of multi-task learning in diffusion models.

### A.3 DIFFUSION MODELS ON VISION TASKS

Recently, diffusion models, having established state-of-the-art performance in image generation, are increasingly being explored for their potential in discriminative vision tasks. This trend continues the historical trajectory to leverage generative models for discriminative tasks (Hinton, 2007). Current research has primarily followed three strategies for repurposing these models. The first utilizes them as potent feature extractors, leveraging the rich internal representations from frozen, large-scale text-to-image models for tasks like open-vocabulary panoptic segmentation (Baranchuk et al.,

2022; Xu et al., 2023). The second employs them at inference time as probabilistic world models, providing generative feedback to adapt discriminative models (Prabhudesai et al., 2023). A third strategy directly leverages the model’s likelihood estimation capabilities, reframing classification as an ”analysis-by-synthesis” problem to perform zero-shot classification without additional training (Li et al., 2023).

More recently, a fundamental paradigm shift has emerged, reformulating core discriminative tasks as conditional denoising problems. This moves beyond using diffusion models as auxiliary tools, making the generative process itself the core mechanism for prediction. Seminal works in this area include DiffusionDet (Chen et al., 2022), which frames object detection as a ”noise-to-box” process of refining random boxes into precise detections, and DiffusionInst (Gu et al., 2022), which formulates instance segmentation as a ”noise-to-filter” denoising process. This unified ”denoising-as-prediction” framework replaces task-specific architectures (e.g., RPNs, query-based heads) with a single generative principle, marking a significant convergence and evolution in the modeling of discriminative vision tasks.

## B MODEL FORMALIZATION

**VLM basics.** A frozen SigLIP image encoder  $E_{\text{siglip}}$  maps an image  $x \in \mathbb{R}^{H \times W \times 3}$  to a grid of patch embeddings  $S \in \mathbb{R}^{N \times C_s}$ , where  $N$  is the number of patches and  $C_s$  the channel width. A lightweight projector  $P_{\text{siglip}} : \mathbb{R}^{C_s} \rightarrow \mathbb{R}^d$  aligns these to the LLM token space:

$$V_s = P_{\text{siglip}}(S) = \text{MLP}_s(S) \in \mathbb{R}^{N \times d}, \quad (1)$$

where  $\text{MLP}_s$  is a 3-layer MLP with GELU activations.

**SVD for single-image  $\rightarrow$  video.** Stable Video Diffusion (SVD) consists of a VAE encoder  $\phi$  and a U-Net denoiser  $f_\theta$  operating over a continuous noise scale  $\sigma$  (Karras et al.). Given a conditioning image  $x$ , we compute a latent  $z_0 = \phi(x)$ . To form a video latent tensor, we replicate  $z_0$  across  $T$  frames:

$$Z_0 = [z_0, \dots, z_0] \in \mathbb{R}^{T \times C \times H' \times W'}.$$

Let  $\sigma_0$  denote the initial noise level from the SVD schedule. We perform one explicit Euler integration step over the ODE at  $\sigma_0$  (classifier-free guidance disabled):

$$Z_1 = Z_0 + \Delta\sigma f_\theta(Z_0, \sigma_0, c), \quad (2)$$

where  $c$  denotes SVD conditioning (e.g., time/frame embeddings, text/image prompts), and  $\Delta\sigma$  is the step size.

We do not render frames; instead, we extract a U-Net hidden activation at the lowest spatial resolution on the downsampling path before the mid-block:

$$H \in \mathbb{R}^{T \times H_d \times W_d \times C_h} = \text{Hidden}^{\text{pre-mid}}(f_\theta, Z_1). \quad (3)$$

**Multi-image extension.** For multiple images  $\{x_k\}_{k=1}^K$ , we first compute their latents  $\{z_0^{(k)}\}$ . These are inserted as keyframes within  $T$  frames at indices  $i_k = \text{round}(\text{linspace}(0, T-1, K))$ . We initialize  $Z_0$  with copies of  $z_0^{(1)}$  and set  $(Z_0)_{i_k} \leftarrow z_0^{(k)}$  before the Euler step, yielding multi-image-aware  $H$ .

**Static+dynamics token fusion.** We convert  $H$  into a token sequence by flattening spatial locations:  $L = H_d W_d$ ,  $\tilde{H} \in \mathbb{R}^{(T \cdot L) \times C_h}$ . A projector  $P_{\text{svd}} : \mathbb{R}^{C_h} \rightarrow \mathbb{R}^d$  maps these to the LLM token space:

$$V_d = P_{\text{svd}}(\tilde{H}) = \text{MLP}_d(\tilde{H}) \in \mathbb{R}^{M \times d}, \quad (4)$$

where  $M = T \cdot L$ .

The SigLIP tokens  $\hat{V}_s$  (Eq. 1) are concatenated with  $\hat{V}_d$  to form the visual sequence:

$$V = [\hat{V}_s; \hat{V}_d].$$

## C TRAINING HYPERPARAMETERS

We adopt the hyperparameters in Table 4 for all our models (for both DyVA-7B and DyVA-Qwen2.5-7B).

Table 4: Training Hyperparameters

Hyperparameter	Value
Batch Size	128
Max Gradient Norm	1.0
Weight Decay	0.1
Learning Rate	2e-5
Optimizer	AdamW
Scheduler	Warmup & Cosine Decay
Warmup Ratio	0.03

## D MORE DESIGN SPACE EXPLORATIONS

**WorldLM is sensitive to temporal information but demonstrates robustness to spatial resolution.** This dual characteristic is evident from our ablation studies. First, as detailed in Table 5, increasing the number of input frames from 1 to 14 yields a consistent and significant improvement across most benchmarks, such as VQAv2 (59.38 to 61.73). This highlights the model’s proficiency in leveraging richer temporal context. Conversely, the impact of spatial resolution appears marginal. By comparing the results in Table 5 (at 576×1024 resolution) with those in Table 3, we find that variations in resolution do not lead to substantial performance changes. These combined findings suggest that our model architecture prioritizes temporal patterns over high-frequency spatial details for the evaluated tasks.

Table 5: Model Performance Across Different Frame Numbers. These are DyVA with SVD only encoders using a image resolution of  $576 \times 1024$

Frames	Pretrain	Tuning	VQAv2	GQA	VizWiz	VSR	POPE	TallyQA	SeedBench	SpatialMM-Obj	3DSR-Bench-real
1	558k	665k	59.38	47.75	48.74	52.12	75.74	50.97	51.12	38.81	45.40
4	558k	665k	60.10	47.36	46.24	53.19	77.60	50.68	52.24	42.48	45.67
8	558k	665k	60.80	48.63	50.25	52.20	78.15	51.46	52.81	37.98	46.32
14	558k	665k	61.73	49.71	38.68	53.43	78.80	52.19	53.28	39.78	46.32

**Fusing SVD latents after the U-Net’s middle block substantially improves performance over the baseline fusion strategy.** We further explore the optimal strategy for integrating SVD-derived temporal latents into the model architecture. Specifically, we compare our baseline DyVA-SVD model with a variant, DyVA-SVD-Post-MiddleBlock, which injects the latents after the U-Net’s middle block. The results, presented in Table 6, indicate that the Post-MiddleBlock fusion strategy yields significant performance gains across most benchmarks. Notably, we observe substantial improvements on GQA (+4.1), VSR (+4.09), and POPE (+4.56), strongly advocating for this modified fusion approach and highlighting the critical impact of architectural choices in temporal feature integration. However, given the absence of across-the-board performance gains, and in consideration of inference efficiency, we ultimately adopted the ”pre-mid” implementation.

Table 6: **SVD vs. SVD-MiddleBlock.** Comparison of different fusion strategies using SVD latents.

Model	VQAv2	GQA	VizWiz	VSR	POPE	TallyQA	SeedBench	Spatial	3DSR
SVD-Only	61.82	50.20	50.60	53.60	75.61	53.27	52.55	40.60	43.50
SVD-Only-Post-MiddleBlock	62.86	54.30	51.41	57.69	80.17	51.36	52.50	41.13	43.84