Appendix

A. Details of ELVA Model

A.1. Training Details

Tab 1 summarizes the hyperparameters used across different training stages. During the spatial pretraining stage, we adopt a low number of frames, increasing to 32 frames for both the spatial-temporal pretraining and supervised finetuning (SFT) stages.

Table 1. **Hyper-parameter Settings for Training Details.** PE denotes Patch Embedding, TC represents the Temporal-Capture Block, TH refers to the Task Head, and LM indicates the language model.

Hyperparameter	Stage 1	Stage 2	Stage 3
Data Scale	4M	3M	843K / 3M
Batch Size	256	256	128 / 256
Video Frame	1	16	32
Hierarchical Merge	×	×	\checkmark
Learning Rate (lr)	4e-5	4e-5	2e-5
LR Schedule	cosine decay	cosine decay	cosine decay
LR Warmup Ratio	0.03	0.01	0.01
Epoch	1	2	1
Weight Decay		0	
Optimizer		AdamW	
DeepSpeed stage		2	

We utilize a total of 4M image samples, comprising 1M from Densefusion and 3M from re-annotated CC3M and COCO in stage 1. For stage 2, we employ 3M re-annotated samples, including 2M from WebVid [2] and 1M from VALOR [4]. See Table 2 for a detailed breakdown of data sources.

A.2. Prompt Engineering

We utilize the following prompt in table 3 to generate detailed captions for the provided images and videos using Qwen2-VL (7B) [14]. For image data, we limit the maximum pixel count to $1280 \times 28 \times 28$ to ensure computational efficiency. Using 16 Nvidia A100 GPUs, generating 3 million high-quality image descriptions takes approximately two days. For video data, we process frames at a rate of 1 fps, with the maximum pixel count per frame set to 360×420 . Generating 3 million video captions under these settings requires three days with 16 Nvidia A100 GPUs.

B. Evaluations on Image-Language Benchmarks.

Evaluation on Image-Language Benchmarks. We evaluate ELVA on a series of general visual understanding benchmarks including GQA [9], SEED-Bench [10], MME [8], MMBench [12]. Part of the image results of Chameleon

Table 2. Data used in pre-training and multimodal supervised fine-tuning stages. * indicates the data is used only in ELVA-7B (HD).

Stage	Dataset	Scale	Source	
Stage 1	ELVA-Image	3M	CC-3M, COCO	
	ELVA-Image	1M	DenseFusion	
Stage 2	ELVA-Video	3M	Webvid-2.5M, VALOR-1M	
	LLaVA-Video	178K	NeXT-QA, ActivityNetQA,	
			PerceptionTest, LLaVA-Hound	
Stage 3	LLaVA-665K /	665K /	COCO, VG, OCR-VQA,	
			GQA, TextVQA	
	LLaVA-OneVision*	3M*	High-Quality Single-Image	

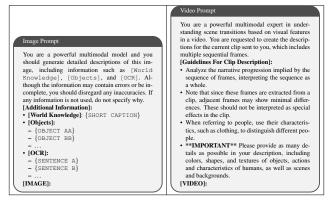


Table 3. Prompt for Caption Engine.

Table 4. **Evaluation on image-language benchmarks.** Our evaluation involves GQA, MME, MM Bench (MMB), and SEED. The results in **bold** and <u>underline</u> are the best and second-best results among encoder-free models, respectively.

Model	Data	GQA	SEEDI	MME	MMB
Encoder-based Models					
Qwen-VL [1]	7.2B	57.5	58.2	1487.5	60.6
LLaVA-v1.5 [11]	0.4B +	62.0	58.6	1510.7	64.3
LLaVA-1.6 (HD) [11]	0.4B +	64.2	64.7	1519.3	67.4
Encoder-free Models					
Fuyu-8B [3]	-	-	-	-	10.7
Chameleon [13]	1.4B+	-	30.6	170	31.1
EVE [6]	33M	60.8	61.3	1217.3	49.5
SOLO [5]	43.7M	-	64.4	1260	-
EVE (HD) [6]	33M	62.6	64.6	1305.7	52.3
ELVA	7M	60.5	<u>66.5</u>	1262.1	<u>53.8</u>
ELVA (HD)	7M	<u>61.1</u>	67.2	<u>1291.5</u>	58.2

and EVE are evaluated with VLMEvalKit [7] or from the OpenCompass. Following the EVE approach, we introduce a high-resolution training strategy exclusively during the SFT stage to develop the ELVA-7B (HD) model, where the longest image edge is set to 1,344 pixels.

B.1. Visualization on Data

We provide additional examples from the ELVA datasets in Table 5. For enhanced visualization, different colors are used to highlight distinct types of information within the descriptions.

References

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv* preprint arXiv:2308.12966, 2023. 1
- [2] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *ICCV*, pages 1728–1738, 2021.
- [3] Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. Introducing our multimodal models, 2023. 1
- [4] Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weining Wang, Jinhui Tang, and Jing Liu. Valor: Vision-audio-language omni-perception pretraining model and dataset. arXiv preprint arXiv:2304.08345, 2023. 1
- [5] Yangyi Chen, Xingyao Wang, Hao Peng, and Heng Ji. Solo: A single transformer for scalable vision-language modeling, 2024.
- [6] Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. Unveiling encoder-free vision-language models. arXiv preprint arXiv:2406.11832, 2024. 1
- [7] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. In *Proceedings* of the 32nd ACM International Conference on Multimedia, pages 11198–11201, 2024. 1
- [8] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. arXiv preprint arXiv:2306.13394, 2023. 1
- [9] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In CVPR, 2019.
- [10] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal Ilms with generative comprehension. arXiv preprint arXiv:2307.16125, 2023.
- [11] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 1
- [12] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? arXiv preprint arXiv:2307.06281, 2023.
- [13] Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- [14] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024. 1

Table 5. Visualizations of the descriptions in ELVA-Image and ELVA-Video. For enhanced clarity, information related to *objects/attributes*, *spatial positions/scene changes*, and *text information* is highlighted using distinct colors.

Visualization



Detailed Caption

The image depicts a red and black train, stationed at a platform. The train is numbered 96104 and is positioned on the tracks, ready to depart or having just arrived. The platform is equipped with a yellow safety line and a metal fence on the right side. In the background, there is a pedestrian bridge and some buildings, indicating an urban setting. The weather appears to be clear, and the overall scene is well-lit.



The image depicts a cozy living room with a modern and somewhat industrial aesthetic. The room features two black leather sofas, one of which is positioned against the wall and the other facing the television. The sofas are arranged around a wooden coffee table, which has a few items on it, including a small plant and a cup. ... On the right side of the room, there is a flat-screen television mounted on the wall. Below the television, there is a small shelf with a few items on it. In the background, there is a poster on the wall with the text "ROXY PRO FIJI" visible, suggesting a connection to surfing or a surfing event. The room also contains a guitar leaning against the wall, adding a personal touch to the space.



The video begins with a title screen that reads "A LIFESTYLE." The scene then shifts to a man wearing a white t-shirt, black cap, and sunglasses, standing in front of a car at night. The man appears to be talking or speaking, and the background shows a brightly lit gas station with an American flag and a sign that reads "2 for \$2.22." The scene suggests a casual, relaxed atmosphere, possibly indicating a lifestyle that involves leisure activities and spending time outdoors.



The video opens with a close-up shot of a white fireplace with a warm fire burning inside. The flames flicker and dance, casting a cozy glow on the surrounding area. As the camera pans out, we see that the fireplace is situated in a well-decorated room, with a potted plant and a few decorative items placed on the mantel. The walls are painted a warm beige color, and the floor is covered with a soft, plush carpet. The overall atmosphere is one of warmth and comfort, with the fire providing a focal point for the room.