

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 QUICKSVIEWER: AN LMM FOR EFFICIENT VIDEO UNDERSTANDING VIA REINFORCED COMPRESSION OF VIDEO CUBES

Anonymous authors

Paper under double-blind review

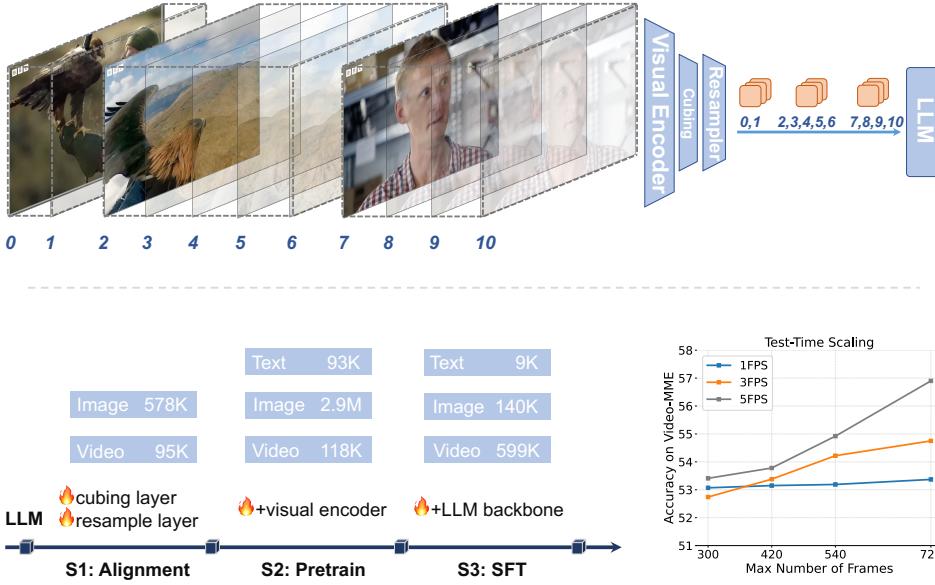


Figure 1: Quicksviewer involves a cubing network that partitions a video into nonuniform cubes, followed by a 3D resampler to gather a fixed number of visual tokens per cube. This efficiency enables **Large Receptive Field** (420 frames) with **High Compression Rate** (64 $\times$ ) during all training stages, and scaling laws on extended frames in inference.

## ABSTRACT

Large Multimodal Models (LMMs) uniformly perceive video frames, creating computational inefficiency for videos with inherently varying temporal information density. This paper present **Quicksviewer**, an LMM with new perceiving paradigm that partitions a video of nonuniform density into varying cubes using Gumbel Softmax, followed by a unified resampling for each cube to achieve efficient video understanding. This simple and intuitive approach dynamically compress video online based on its temporal density, significantly reducing spatiotemporal redundancy (overall 45 $\times$  compression rate), while enabling efficient training with large receptive field. We train the model from a language backbone through three progressive stages, each incorporating lengthy videos on average of 420s/1fps thanks to the perceiving efficiency. With only 0.8M total video-text samples for training, our model outperforms the direct baseline employing a fixed partitioning strategy by a maximum of 8.72 in accuracy, demonstrating the effectiveness in performance. On Video-MME, Quicksviewer achieves competitive performance compared to models of similar size while utilizing just up to 5% of tokens per frame required by baselines. With this paradigm, scaling up the number of input frames reveals a clear power law of the model capabilities. It is also empirically verified that the segments generated by the cubing network can help for analyzing continuous events in videos.

054  
055  
056  
057  
1 INTRODUCTION058  
059  
060  
061  
062  
063  
064  
Large Multimodal Models (LMMs) (Deepmind, 2024; OpenAI, 2024; Bai et al., 2025) have shown  
promising progress in video understanding, paving the way for general intelligence in physical world.  
These models build on Large Language Models (LLMs) and are trained in stages with large-scale  
image and video data, encoding video frames in the same manner as images before feeding them  
into the LLM for inference. At the core of these models is the efficient perception of input videos,  
which is crucial in tackling the persistent contradiction between the **temporal redundancy** of video  
streams (Buckler et al., 2018; Wenger, 1997) and the **computational efficiency** of LMMs with long  
context (Fu, 2024).065  
066  
067  
068  
069  
070  
071  
Extensive studies have been striving to develop LMMs for solving this fundamental issue. Building  
on devise of frame sampling methods, trailblazing efforts typically involve dedicated token merging  
strategies (Bai et al., 2025; Wang et al., 2025; Shen et al., 2024; Zohar et al., 2024; Li et al., 2024d;  
Zhang et al., 2025) and adapted parallel training infrastructures (Zhang et al., 2024b; Chen et al.,  
2024c; Shen et al., 2025). However, the arbitrary frames sampling and tokens merging introduces  
inevitable information loss, while marginal compression limits the number of frames in large-scale  
pre-training.072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
The velocity of content change in videos is inherently nonuniform, suggesting that the density varies  
across different temporal cubes. For example in Figure 1, the initial short period features rapidly  
changing scenes of a researcher attaching a camera to an eagle’s back, followed by an extended  
sequence of stable footage from the camera, and a largely static interview. Inspired by the way that  
humans adjust their perception speed based on content changes, this paper explores how LMMs  
can perform video understanding on the nonuniform cubes to achieve dynamic compression, and  
significantly reduce the spatiotemporal complexity and enhance the efficiency. For practical scenarios  
where videos typically originate from lengthy offline recordings or online video streams, we thereby  
aim for the model to (1) **learn from unlabeled data**, perform (2) **online cube partitioning**, and  
establish a (3) **unified perception paradigm** for images and videos.082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
We present Quicksviewer, an LMM that perceives nonuniform video cubes for efficient video  
understanding. Given a video passed from a visual encoder, a small cubing network first partitions it  
into nonuniform cubes based on the momentum of semantic feature differences between frames, a  
process that can be conducted online in streaming scenarios. Next, a unified resampling is employed  
to the cubes to gather a fixed number of tokens for adaptive compression. Finally, these visual tokens,  
along with absolute timestamps, are fed into the LLM for inference. We integrates the learning  
of the cubing network into the end-to-end training of the LMM using the Gumbel Softmax (Jang  
et al., 2016; Herrmann et al., 2020) method with an improved noise annealing mechanism. This  
reinforced approach not only enables efficient learning on videos without boundary labels but also  
insures sufficient sampling over the cubes distribution with continuous gradient during training. The  
nonuniform perception paradigm, which is solely driven by the properties of input video, together  
with the subsequent resampling enables an efficient video encoding with  $45\times$  compression reate,  
large temporal receptive filed of 420 frames for pre-training, and a consistent representation for both  
images and videos.095  
096  
097  
098  
099  
100  
101  
102  
103  
We train our models starting from the LLM backbones through three progressive stages, each incorpo-  
rating lengthy videos averaging 420s/1fps by benefiting from the efficient perception mechanism. The  
resulting model, which we coined as Quicksviewer, is an efficient LMM capable of understanding  
single/multi-images and long videos. We also find that our network is efficient in learning. With  
only 0.8M video-text samples in total for training, our model achieves competitive performance on  
Video-MME (Fu et al., 2024a) using just up to 5% of the tokens per frame required by baselines. In  
addition, to facilitate training on ultra-long videos (e.g., over 1hour), we developed a training infra-  
structure supporting dynamic changes of sequences lengths based on an existing effort (Chen et al., 2024c), to  
further facilitate potential explorations in future.104  
105  
106  
107  
We evaluate Quicksviewer on various video understanding benchmarks, ranging the duration from 16  
seconds to 1 hour. Results show that our model outperforms the direct baseline employing a fixed  
partitioning strategy by a maximum of 8.72 in accuracy, suggesting the utility of the nonuniform  
perception. We further analyze the cubes partitioning in videos, which demonstrates the emergence  
of "Visual Lag" phenomenon when the model perceives videos phase-by-phase. We also conduct

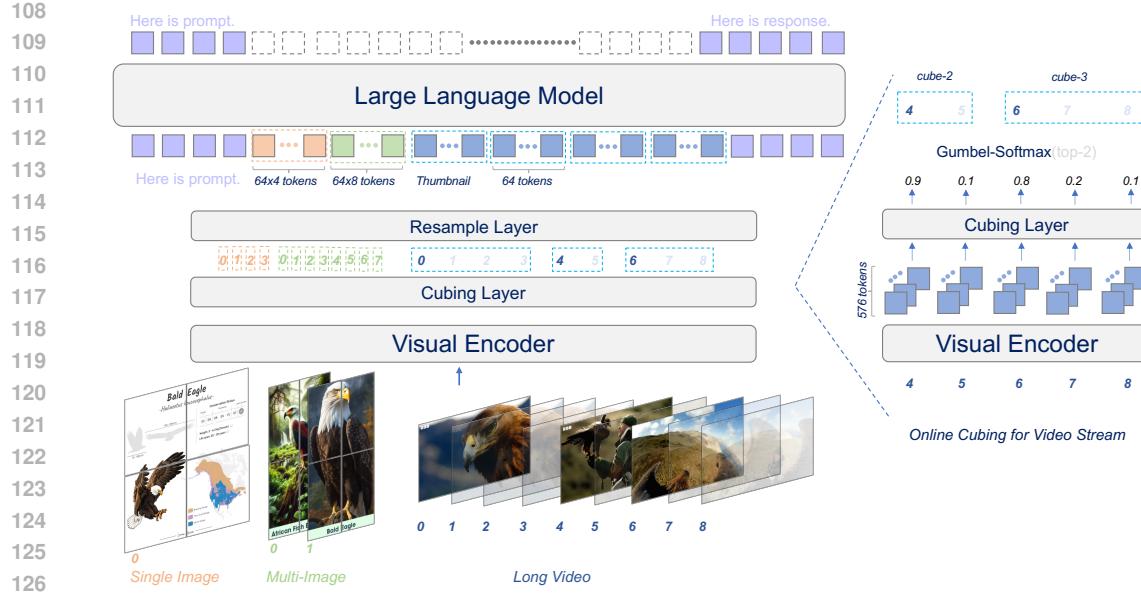


Figure 2: Left: The network architecture of Quicksviewer, that performs unified understanding of videos and images through visual tokens from cascaded modules. Right: The cubing network, that partitions an online video into nonuniform cubes based on Gumbel Softmax.

extensive ablation studies to show the effectiveness of presented components, including the cubing approach, 3D positional encoding, loss penalty, and the annealing strategy.

## 2 APPROACH

The overall architecture of the presented model is shown in Figure 2. We build an efficient LMM that can receive both images and videos as input, where the video is perceived based on nonuniform cubes partitioned through a small cubing network. The model consists of four basic components: (1) a visual encoder  $f_v(\cdot)$  that encode image slices or video frames into visual tokens, (2) a cubing network  $f_q(\cdot)$  that partitions the video frames into  $N_Q$  cubes, (3) a resampler  $f_r(\cdot)$  which compress each slice or cube into a fixed number of tokens, and (4) an LLM  $f_l(\cdot)$  which accept concatenation of visual tokens and user prompt for response generation. Note that we introduce **FPQ**, the average number of frames per cube, which regulates the perception granularity and enables adaptive number of cubes  $N_Q$  according to video duration.

### 2.1 CUBING USING GUMBEL SOFTMAX

**Visual Encoding** Given an input video to our model, we first uniformly sample  $N_F$  frames with a fixed FPS to form  $[F_i]_{i=1,\dots,N_F}$ . And then, each frame is firstly encoded using the visual encoder  $f_v$  to obtain  $N_1$  visual tokens  $\mathbf{F}_j = [\mathbf{v}_i]_{i=1,\dots,N_1}$ , where  $\mathbf{v}_i \in \mathbb{R}^d$  is the token representation. The large number of visual tokens provides fine-grained semantic representations for each frame, while the concatenation of all tokens (*i.e.*,  $N_F \times N_1$ ) across frames results in an unacceptable sequence length.

**Cubing** It is natural to leverage the semantic difference of features to find keyframes. For compatible with online streaming scenarios and also considering the event-level long-term changes, we utilize the momentum accumulated from the representations of previous frames to track the semantic changes. Given the visual tokens of frames, the cubing network calculate the  $i^{th}$  momentum representation as:

$$\Delta_i = \alpha(\mathbf{F}_i - \mathbf{F}_{i-1}) + (1 - \alpha)\Delta_{i-1} \quad (1)$$

162 where  $\Delta_i = [\delta_j]_{j=1, \dots, N_1} \in \mathbb{R}^{N_1 \times d}$  captures the accumulated semantic changes discounted by factor  
 163  $\alpha \in [0, 1]$ . A 2-layers MLP with LayerNorm (Ba et al., 2016) is then applied to the mean of visual  
 164 tokens momentums to quantify the significance of the frame:

$$167 \quad \mathbf{z}_i = \text{MLP}_{\times 2}(\text{LayerNorm}\left(\frac{1}{N_1} \sum_{j=1}^{N_1} \delta_j\right)), \quad \mathbf{z}_i \in \mathbb{R}^2 \quad (2)$$

170 where the 2-dimensional vector  $\mathbf{z} = [w_0, w_1]$  forms a gate, and the sigmoid of their difference directly  
 171 reflects the keyframe probability  $p = \frac{1}{1 + \exp(w_0 - w_1)}$  (Herrmann et al., 2020). We next apply the  
 172 *top-k* operation on the first dimension of  $[\mathbf{z}_i]_{i=2, \dots, N_F}$  vectors to obtain the indices of  $N_Q - 1$  largest  
 173 values  $[l_i]_{i=1, \dots, N_Q - 1}$ , which are selected as keyframes. The keyframe  $F_{l_i}$ , along with its subsequent  
 174 consecutive non-keyframes, forms the cube  $Q_{i+1} = [F_{l_i}, \dots, F_{l_{i+1}}]$ . Note that the first sampled frame  
 175 consistently serves as the keyframe for deriving the first cube  $Q_1$ .

177 **Sampling for Training** During training, we expect the model to perform sufficient sampling  
 178 exploration while ensuring gradient continuity. We achieve this using Gumbel Softmax with the  
 179 Straight-Through trick (Gumbel, 1954; Jang et al., 2016):

$$182 \quad \mathbf{z}_i = \text{softmax}(\mathbf{z}_i - \log(-\log(\epsilon)/\tau)), \quad \epsilon \sim U(0, 1) \quad (3)$$

184 where the log term approximates the sampling process and regulates the degree of exploration.

185 In experiments, we found that persistent exploration prevents the model from establishing a stable  
 186 cubing paradigm for subsequent reasoning, leading to sustained loss oscillations. We propose to add  
 187 an learning rate  $\eta$  before the Gumbel noise  $\eta \log(-\log(\epsilon)/\tau)$ , which is **annealed** from  $\eta_0 = 1.0$  to  
 188  $\eta_T = 0.001$  during training using a cosine scheduler.

## 190 2.2 RESAMPLING WITH 3D POSITIONAL ENCODING

192 Based on the partitions from cubing network, a unified 3D resampler is adopted to compress each  
 193 cube of arbitrary length into a fixed number of  $N_2$  tokens.

195 **Resampling Video Cubes** We employ the same resampler architecture as (Yao et al., 2024) to  
 196 compress each cube into a fixed number of dense tokens. We extend the original 2D positional  
 197 encoding by incorporating a temporal dimension to form 3D position encoding. As a result, each  
 198 video token is assigned three positional coordinates  $(x, y, z)$ , representing time, width, and height. We  
 199 then unfold each cube into tokens sequence along the frame dimension. After adding 3D positional  
 200 embeddings, an unified resampling is performed to obtain  $N_2 = 64$  visual tokens for each cube. For  
 201 images, we first adopt the AnyRes (Liu et al., 2024a) to divide high resolution images into slices,  
 202 and apply the same resampling to each slice to obtain visual tokens and finally concatenate the tokens.

203 **Resampling for Video Thumbnail** Using cubes quantized from the cubing network for response  
 204 tokens generation introduces a fundamental problem: ***How does the NTP training objective optimize***  
 205 ***the boundary prediction of the cubing network?*** We introduce video thumbnail to resolve this  
 206 problem meanwhile provide effective global representation. Specifically, we first (1) multiply the 0-1  
 207 discretized first dimension of vectors  $[\mathbf{z}_i]_{i=1, \dots, N_F}$  with their corresponding frame representations  
 208  $[\mathbf{F}_i]_{i=1, \dots, N_F}$ , then (2) average across the frame dimension to obtain  $N_1$  visual tokens. A further  
 209 resampling is performed to get final thumbnail representation containing  $N_2$  tokens. This simple  
 210 approach allows gradients to be directly propagated back to the cubing boundaries. The final  
 211 representation of a video is a concatenation of the representations of the thumbnail and cubes.

## 213 2.3 LLM INFERENCE WITH AUXILIARY LOSS

215 Following resampling of nonuniform video cubes, the tokens of each cube span varying temporal  
 windows. We prepend each cube with an absolute timestamp as a float number in 0.01-second units,

216 enabling explicit temporal awareness. We also enclose the video, thumbnail, and image tokens with  
 217 their corresponding special tokens to enable explicit content differentiation.  
 218

219 During training, we observed that excessively large values of  $z$  cause overly high gradients, impairing  
 220 convergence. To address this, we introduce an auxiliary  $L_2$  norm loss with  $\beta = 0.001$  penalty weight  
 221 on them to constrain its values within a reasonable range.  
 222

### 223 3 TRAINING PROCESS

225 We train our models with three progressive stages starting from LLM backbones, each stage incorpo-  
 226 rating lengthy videos on average of 420s by benefiting from the efficient perception approach.  
 227

228 **Stage-1: Multimodal Alignment** We utilize both interleaved and captioning image-text corpuses,  
 229 and video-text captioning corpus to train our models, establishing fundamental alignment between  
 230 visual encoder and LLM backbones with in-context learning capabilities. We sample a subset  
 231 of 20K sequences from OBELICS (Laurençon et al., 2023), with each containing more than two  
 232 interleaved pairs. We utilize LCS (Li et al., 2024a), a re-captioned dataset consisting of 558K detailed  
 233 descriptions from the CC3M (Sharma et al., 2018). The video-text training data incorporates a  
 234 sampled subset of 87K captioning pairs from FineVideo (Farré et al., 2024) and 8K captioning pairs  
 235 from ANetCaptions (Krishna et al., 2017). We train parameters of the cubing network and resampler  
 236 while keeping all other parameters frozen to establish a stable projection. The models are trained for  
 237 1 epoch with a  $lr$  that warms up to  $1e^{-4}$  over the first 2% of steps, then gradually decays to 0.  
 238

239 **Stage-2: Pre-training** We employ large-scale pretraining data, primarily consisting of image-text  
 240 multi-task data, to pre-train models establishing general multimodal capabilities across broad visual  
 241 scenarios. We utilize a subset of 2.99M samples from LLaVA-OneVision-SingleImage (Li et al.,  
 242 2024a) as training corpus, which incorporates 2.9M image-text pairs and 93K textual instruction-  
 243 tuning samples from Evo-Instruct (Chen et al., 2024a). For video-text corpus, we utilize a sampled  
 244 subset of 75K video QAs from FineVideo (Farré et al., 2024) and 38K captioning pairs from  
 245 ShareGPT4Video (Chen et al., 2024b). To mitigate catastrophic forgetting, we retain 5% of the  
 246 previous image and video data in our training corpus. Alongside the cubing network and resampler,  
 247 we also unfreeze the visual encoder to improve the visual representation. We train models for 1 epoch  
 248 with a  $1e^{-5}$  initial learning rate, with the same warmup and decay schedule as stage-1.  
 249

250 **Stage-3: Supervised Fine-tuning** We primarily leverage extensive video-text paired corpus to  
 251 train our models in this stage, enabling robust video understanding capabilities. We primarily utilize  
 252 a subset of 476K video-text samples sourced from VideoChat2-IT (Li et al., 2024c), and a subset of  
 253 79K samples from ShareGPTVideo (Zhang et al., 2024c) as the video corpus. To enhance adaptation  
 254 to long video scenarios, we further integrate 5K samples from MovieChat (Song et al., 2024) and  
 255 39K samples derived by (Chen et al., 2024c) from the Shot2Story dataset (Han et al., 2023). The  
 256 image-text corpus incorporates a sampled subset of 100K multi-image, multi-task understanding  
 257 samples from LLaVA-OneVision-MultiImages (Li et al., 2024a). We also preserve a subset of training  
 258 data from the previous stage, consisting of 40K text-image pairs and 9K textual instruction-tuning  
 259

	Stage-1	Stage-2	Stage-3
<i>Vision</i>			
<b>Resolution</b>	$384 \times 384$	$384 \times 384$	$384 \times 384$
FPS, #Frames	1, Max 420	1, Max 420	1, Max 420
<i>Data</i>			
<b>Image-Text</b>	LCS, OBELICS	LLaVAOV-SingleImage	LLaVAOV-MultiImages
#samples	558K, 20K	2.99M	100K
<b>Video-Text</b>	FineVideo, ANetCaptions	FineVideo, ShareGPT4Video, ANetCaptions	Sec. 3. 3
#samples	87K, 8K	118K	599K
<i>Model</i>			
<b>Trainable</b>	Cubing, Resampler	Cubing, Resampler, ViT	Full Model
#Parameters	75M	500M	8B
<i>Training</i>			
<b>Anneal:</b> $\eta_0, \eta_T$ , ratio	1.0, 0.01, 0.8	1.0, 0.01, 0.6	1.0, 0.01, 0.6
<b>LR:</b> $\theta_c, \theta_r, \theta_v, \theta_l$	$1e^{-4}, 1e^{-4}, -, -$	$2e^{-5}, 2e^{-5}, 2e^{-5}, -$	$1e^{-5}$
<b>Epoch</b>	1	1	1

260  
 261 Table 1: Detailed configuration for each training stage.  
 262  
 263  
 264  
 265  
 266  
 267  
 268  
 269

270 samples. We train all parameters for 1 epoch with a learning rate that warms up to  $1e^{-5}$  over the 0.02  
 271 epoch, followed by gradual decay to 0 for the remaining duration.  
 272

273 During training, all videos sampled at 1FPS to extract full frames. For videos exceeding 420s, we  
 274 uniformly extract 420 frames to maintain computational tractability. Images are processed using  
 275 AnyRes with a resolution of  $384 \times 384$ . For all stages, the Gumbel noise learning rate  $\eta$  (initialized  
 276 at 1.0) undergoes cosine annealing to 0.01 within: 0.8 epoch (Stage 1) or 0.6 epochs (Stages 2-3).  
 277

## 278 4 EXPERIMENT

280 Models	281 Size	282 #Tokens	283 #Train	284 MMBench-Video	285 MVBench	286 MLVU	287 Video-MME
288 Duration	289 /Frame	290 V-T	291 3 min	292 16 sec	293 3~120 min	294 1~60 min	295
<i>Proprietary Models</i>							
GPT4-V (OpenAI, 2023)	-	-	-	1.53	43.7	-	60.7
GPT4-o (OpenAI, 2024)	-	-	-	1.63	64.6	66.2	77.2
<i>Open-Source Video LMMs</i>							
LLaMA-VID (Li et al., 2024e)	7B	2	0.4M	1.08	41.5	33.2	-
LongLLaVA (Wang et al., 2024c)	9B	144	0.5M	-	49.1	-	43.7
Chat-UniVI (Jin et al., 2024)	7B	112	100K	1.06	42.9	-	45.9
ShareGPT4Video (Chen et al., 2024b)	8B	144	4.8M	1.05	51.2	46.4	43.6
LLaVA-NeXT-Video (Zhang et al., 2024d)	7B	144	100K	1.14	33.7	-	46.5
VideoLLaMA2 (Cheng et al., 2024)	7B	32	10.7M	1.08	54.6	48.5	46.6
LongVA (Zhang et al. (2024b))	7B	144	-	-	-	56.3	54.3
VideoChat2 (Li et al., 2024c)	7B	64	2.8M	1.22	<b>60.4</b>	47.9	54.6
mPLUG-Owl3 (Ye et al., 2024)	8B	729	134K	<b>1.35</b>	54.5	-	53.5
Fixed-LLama3.1	8B	12.8	0.8M	0.71	45.2	50.2	45.0
Quicksviewer-LLama3.1	8B	12.8	0.8M	0.87	53.9	58.6	47.6
Quicksviewer	8B	12.8	0.8M	<u>1.24</u>	<u>55.6</u>	<b>61.5</b>	<b>56.9</b>

296 Table 2: Video benchmarking results between Quicksviewer and baselines under comparable total  
 297 sequence length. Quicksviewer achieves multiple SOTA performance while using fewer tokens per  
 298 frame (up to 5% of baseline) and substantially less video-text training samples.  
 299

### 300 4.1 IMPLEMENTATION DETAILS

301 We use SigLIP (Zhai et al., 2023) (soo400m-path14-384) as our visual encoder inconsistent with  
 302 previous works. We adopt Qwen2.5 (Yang et al., 2024) as the language backbone for our standard  
 303 implementation (i.e., Quicksviewer), while utilizing Llama3.1 (Touvron et al., 2023a) as the alter-  
 304 native LLM for another version (i.e., Quicksviewer-Llama3.1) for comprehensive exploration. We  
 305 use AdamW (Loshchilov & Hutter, 2017) optimizer with a cosine scheduler for all training stages.  
 306 The number of tokens generated from visual encoder and resampler are  $N_1 = 576$ , and  $N_2 = 64$ ,  
 307 respectively. The discounting factor of momentum is set to  $\alpha = 0.9$ . The penalty weight to the  
 308 auxiliary loss is set to  $\beta = 0.001$ . We use FPQ=5 for all models. Our models is trained on 48  
 309 NVIDIA A100 GPUs.  
 310

### 311 4.2 EXPERIMENTS ON VIDEO UNDERSTANDING

312 We train a direct baseline, Fixed-Llama3.1, which utilizes uniform temporal partitioning with the  
 313 same FPQ of input videos. For an unbiased comparison, we evaluate with baselines configured with  
 314 comparable total sequence lengths, maintaining equivalent computational budgets.  
 315

316 **Benchmarks and Metrics** We evaluate the our models on widely used video understanding  
 317 benchmarks Video-MME (Fu et al., 2024a), MVBench (Li et al., 2024c), and MLVU (Zhou et al.,  
 318 2024) to investigate the effectiveness. VideoMME is a general video understanding benchmark  
 319 that collect videos (1min~1hour) from Youtube with manual annotations. MVBench covers 20  
 320 challenging tasks ranging from perception to cognition. MLVU (3mins~2hours) refers to an long  
 321 video understanding benchmark for long-term inference.  
 322

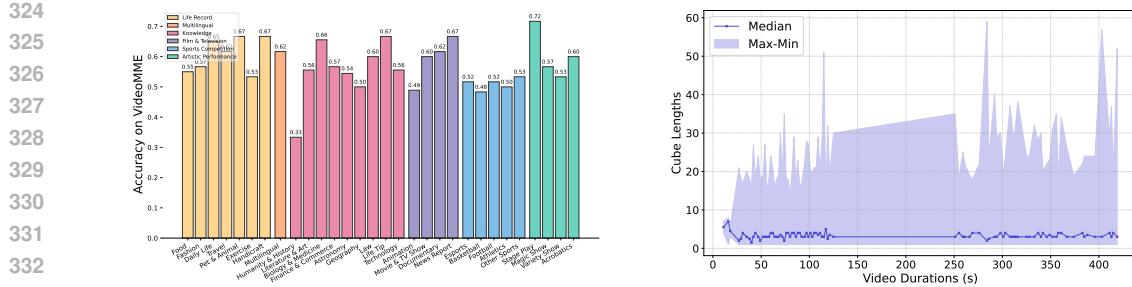


Figure 3: (a) Left: Performance of Quicksviewer on particular domains and categories of Video-MME. (b) Right: Distribution of cube lengths across Video-MME videos.

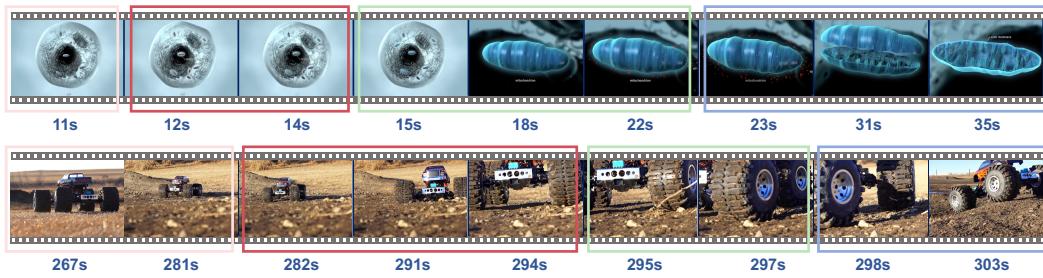


Figure 4: The "Visual Lag" phenomenon occurring during the model's cube-based segmental comprehension, where current cubes incorporate terminal frames from preceding event scenes to enable retrospective understanding.

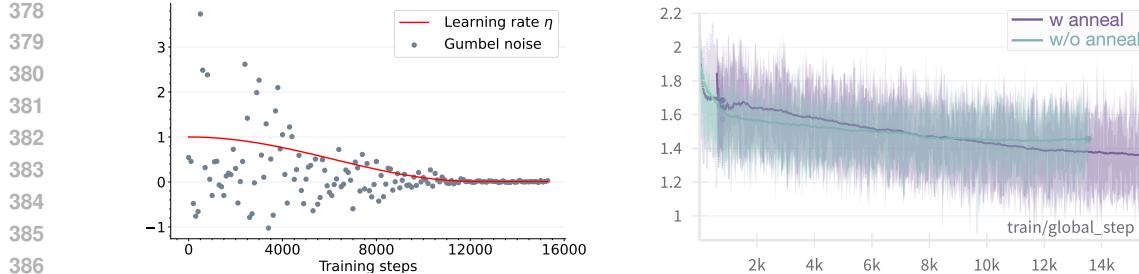
**Quantitative Results** We adopt empirically optimized configuration: 5 FPS with a maximum of 720 frames for all benchmarks. As the main result shown in Table 2, our standard model achieves competitive performance on Video-MME and MLVU, and the competitive performance on MVbench with significantly fewer tokens and training volumes. Specifically, our model achieves competitive performance on Video-MME, albeit with slightly inferior results on long videos. This demonstrates that the encoding paradigm harness the scaling benefits by high frame rate. In comparison with the direct baseline with fixed cubing strategy, our model obtain large improvements, suggesting the effectiveness of the cuing strategy. Our approach achieves competitive performance on MLVU while demonstrating competitive performance on MVBench, despite utilizing substantially less training data (only 28% of VideoChat2’s and 7.5% of VideoLLaMA2’s requirements). This evidences our network’s exceptional learning efficiency.

**Analysis** We further analyze the model performance across distinct domain categories in Video-MME, systematically examine both capability advantages and limitations. As illustrated in Figure 3 (a), bars sharing identical colors belong to the same domain. Primarily, we observe consistent model performance across all domains, with mean scores of 0.61, 0.62, 0.55, 0.60, 0.51, and 0.61 respectively, suggesting limited domain-specific variation in question difficulty. Secondly, the model demonstrates suboptimal performance (below 50%) in three categories: Humanity & History, Animation, and Basketball. This may indicate persistent challenges in fine-grained character recognition that require further improvement.

We further analyze the distribution of cube lengths on Video-MME, with results shown in Figure 3 (b). Based on the predefined FPQ, we found the median cube length approximates 5 frames. Notably, the model demonstrates a tendency to partition diverse length of cubes for longer videos, which aligns with the variable viewing speeds in human perception of lengthy videos.

#### 4.3 ANALYSIS OF THE CUBES PARTITIONING

To investigate how the trained model partitions cubes for understanding, we analyzed two representative video cases by examining cubes relative to content transitions. As shown in Figure 3, each box



378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

Figure 5: (a) Left: Gumbel noise progressively anneals to 0.001 following the decaying learning rate with cosine scheduler. (b) Right: Compared to non-annealed training (cyan curve), adding Gumbel noise annealing (purple curve) yields more stable and superior loss convergence.

refers a cube spanning from its start to end timestamps (with similar intermediate frames omitted). We reveal a "**Visual Lag**" phenomenon during cube-based video perception of the model: terminal frames from preceding event scenes are incorporated into cubes containing subsequent event scenes. For example in the first video, the initial frames of cubes 2-4 respectively contain content from three event scenes: 1) cellular details, 2) mitochondrial positioning, and 3) ATP synthesis exhibiting in cubes 1-3 respectively. We posit this mechanism enables the model to retain partial memory of preceding scenes to facilitate current scene understanding.

#### 4.4 ANALYSIS OF THE ANNEALING STRATEGY

Traditional Gumbel-Softmax training controls sampling randomness exclusively through temperature adjustment, making it unsuitable for training the cubing network as a component of an LMM. To resolve this issue, we propose annealing the Gumbel noise, which substantially improves both training stability and effectiveness. To further evaluate the performance of the proposed annealing mechanism, we examine the evolution of Gumbel noise values throughout a training epoch in with the annealed learning rate. For clear visualization, we uniformly sample training steps at 100-step intervals, as illustrated in Figure 5 (a). From the figure, we observe that in the early training stages, larger Gumbel noise effectively facilitates exploration for the cubing network. As training progresses, the Gumbel noise gradually converges to the predefined value of 0.01. This allows the model to leverage its learned segmentation mechanism for video understanding in later stages, stabilizing the training process and achieving optimal performance.

Figure 5 (b) compares the loss trajectories of models trained with and without the annealing mechanism. Initialized from the same checkpoint from Stage-2, we train parallel models using both approaches and monitor loss variations throughout one epoch to assess learning efficiency and stability. Our analysis reveals that the model benefits from the progressive annealing of Gumbel noise in the later stages. During this phase, the model effectively utilizes its learned cubing mechanism to accelerate loss minimization, achieving superior convergence efficiency, demonstrating the effectiveness and training stability.

#### 4.5 ABLATION STUDIES

We conduct comprehensive ablation studies to evaluate the efficacy of the components leveraged in Quicksviewer. To establish a simple baseline, we first train a Llama3.1 (Touvron et al., 2023a) model through Stages 1-2 using only image-text data introduce in Sec. 3, deliberately excluding video inputs. This image-only pretrained checkpoint then serves as the initialization point for systematically investigating various Stage-3 configurations with video-text data.

**Cubing network with ViT** To accelerate cube processing, we investigate the feasibility of using only the initial  $n$  layers of ViT for the cubing network. Our ablation study employs the first 2 ViT layers for cube feature generation while maintaining all other model components unchanged. As demonstrated in Table 7, this configuration results in significant performance degradation, indicating that shallow visual features are insufficient for effective cubes partitioning.

432	433					434	
	435					436	437
438	439	440	441	442	443	444	445
PE	Cubing Network	Penalty	Annealing	Trainable	Overall	Long	Video-MME
					1~60 min	30~60 min	
2D	ViT first 2-layers	0.1	-	$\theta_c, \theta_r, \theta_v$	33.92	35.56	
2D	ViT full	0.1	-	$\theta_c, \theta_r, \theta_v$	41.22	38.67	
3D	ViT full	0.1	-	$\theta_c, \theta_r, \theta_v$	44.37	40.67	
3D	ViT full	0.001	-	$\theta_c, \theta_r, \theta_v$	44.66	40.44	
3D	ViT full	0.001	annealing	$\theta_c, \theta_r, \theta_v$	45.44	43.44	
3D	ViT full	0.001	annealing	All	<b>45.96</b>	<b>38.44</b>	

Table 3: Ablation results of Stage-3 training initialized from a checkpoint pretrained only with image data (Stage 1-2). The optimal configuration: 3D positional encoding, Gumbel noise annealing with 0.001 penalty weight, and full trainable parameters, demonstrating superior performance.

**3D positional encoding** We systematically evaluate the impact of 3D positional encoding compared to the original 2D formulation. Implementing this modification while keeping all other parameters fixed in Stage-3 training, our experiments demonstrate a consistent accuracy improvement of +3.15% (Table 7), confirming the benefits of spatiotemporal position awareness for video understanding.

**Penalty weight to the auxiliary loss** The penalty weight applied to the  $L_2$ -norm of  $\mathbf{z}$  values regulates the scale of logistic outputs, consequently influencing partition behavior. Through systematic ablation while maintaining other parameters fixed, we identify  $\beta = 0.001$  as the optimal weighting that simultaneously: (1) maintains  $\mathbf{z}$  values in an effective operational range, and (2) achieves superior model performance (Table 7).

**The annealing strategy** Having established the optimal penalty weight, we proceed to evaluate the efficacy of our proposed Gumbel noise annealing strategy. This approach systematically reduces exploration randomness during training, transitioning from aggressive parameter space exploration to fine-tuned optimization. Comparative results in Table 7 demonstrate consistent performance improvements over the fixed-noise baseline, validating the benefits of noise scheduling.

**The tuning parameters** We examine the impact of trainable parameters on video understanding by comparing two training regimes: (1) our baseline approach that only fine-tune the ViT and resampler parameters during Stage-3, versus (2) a full-parameter optimization strategy that additionally fine-tunes the LLM backbone. As evidenced in Table 7, comprehensive parameter training yields superior benchmark performance, suggesting that joint visual-linguistic optimization enhances multimodal alignment for video understanding tasks.

## 5 CONCLUSION

In this paper, we introduced Quicksviewer, an LMM designed for efficient video understanding through a nonuniform perception paradigm. By dynamically partitioning videos into nonuniform cubes and applying adaptive resampling, our approach achieves a  $45\times$  compression rate while maintaining a consistent representation for both images and videos. We demonstrated that integrating the cubing network into end-to-end training via Gumbel Softmax with an improved noise annealing mechanism, enables efficient learning without boundary labels. Furthermore, our model, trained on just 0.8M videos, achieves competitive performance on VideoMME with significantly fewer tokens per frame than the baseline methods. To support training on ultra-long videos, we also developed an infra that allows dynamic sequence lengths. These contributions pave the way for efficient and scalable LMMs, facilitating future research in long video understanding.

Our approach demonstrates strong potential for future applications in video segmentation and long-sequence processing. In video tasks, cubing-based networks enable fine-grained segmentation of videos into multiple events, supporting effective video preprocessing. Moreover, our method is highly generalizable and can be adapted for various sequence segmentation and compression scenarios, providing a foundation for efficient long-sequence analysis.

486 REFERENCES  
487

488 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel  
489 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language  
490 model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736,  
491 2022.

492 Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrahman, Essam Sleiman, Deyao Zhu, Jian Ding, and  
493 Mohamed Elhoseiny. Minigpt4-video: Advancing multimodal llms for video understanding with  
494 interleaved visual-textual tokens. *arXiv preprint arXiv:2404.03413*, 2024.

495 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint  
496 arXiv:1607.06450*, 2016.

497 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
498 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

500 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,  
501 Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*,  
502 2025.

503 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
504 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
505 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

506 Mark Buckler, Philip Bedoukian, Suren Jayasuriya, and Adrian Sampson. Eva<sup>2</sup>: Exploiting temporal  
507 redundancy in live computer vision. In *2018 ACM/IEEE 45th Annual International Symposium on  
508 Computer Architecture (ISCA)*, pp. 533–546. IEEE, 2018.

509 Tyler A Chang, Catherine Arnett, Zhuowen Tu, and Benjamin K Bergen. Goldfish: Monolingual  
510 language models for 350 languages. *arXiv preprint arXiv:2408.10441*, 2024.

511 Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang,  
512 Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized  
513 data for lite vision-language models. *arXiv preprint arXiv:2402.11684*, 2024a.

514 Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan,  
515 Zhenyu Tang, Li Yuan, et al. Sharegpt4video: Improving video understanding and generation with  
516 better captions. *Advances in Neural Information Processing Systems*, 37:19472–19495, 2024b.

517 Yukang Chen, Fuzhao Xue, Dacheng Li, Qinghao Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian  
518 Tang, Shang Yang, Zhijian Liu, et al. Longvila: Scaling long-context visual language models for  
519 long videos. *arXiv preprint arXiv:2408.10188*, 2024c.

520 Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong  
521 Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal  
522 models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024d.

523 Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong  
524 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning  
525 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision  
526 and pattern recognition*, pp. 24185–24198, 2024e.

527 Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi  
528 Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and  
529 audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024.

530 Google Deepmind. Gemini 2.0. <https://deepmind.google/technologies/gemini/>,  
531 2024.

532 Miquel Farré, Andi Marafioti, Lewis Tunstall, Leandro Von Werra, and Thomas Wolf. Finevideo.  
533 <https://huggingface.co/datasets/HuggingFaceFV/finevideo>, 2024.

540 Chaoyou Fu, Yuhang Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu  
 541 Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation  
 542 benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024a.  
 543

544 Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Shaoqi Dong, Xiong  
 545 Wang, Di Yin, Long Ma, et al. Vita: Towards open-source interactive omni multimodal llm. *arXiv  
 546 preprint arXiv:2408.05211*, 2024b.

547 Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Haoyu Cao,  
 548 Zuwei Long, Heting Gao, Ke Li, et al. Vita-1.5: Towards gpt-4o level real-time vision and speech  
 549 interaction. *arXiv preprint arXiv:2501.01957*, 2025.  
 550

551 Yao Fu. Challenges in deploying long-context transformers: A theoretical peak performance analysis.  
 552 *arXiv preprint arXiv:2405.08944*, 2024.

553 Emil Julius Gumbel. *Statistical theory of extreme values and some practical applications: a series of  
 554 lectures*, volume 33. US Government Printing Office, 1954.

555 Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang,  
 556 Jianyu Jiang, Jiawei Wang, et al. Seed1. 5-vl technical report. *arXiv preprint arXiv:2505.07062*,  
 557 2025.

558 Mingfei Han, Linjie Yang, Xiaojun Chang, and Heng Wang. Shot2story20k: A new benchmark for  
 559 comprehensive understanding of multi-shot videos. *arXiv preprint arXiv:2312.10300*, 2023.

560 Charles Herrmann, Richard Strong Bowen, and Ramin Zabih. Channel selection using gumbel  
 561 softmax. In *European conference on computer vision*, pp. 241–257. Springer, 2020.

562 Wenyi Hong, Weihan Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng,  
 563 Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvilm2: Visual language models for image and video  
 564 understanding. *arXiv preprint arXiv:2408.16500*, 2024a.

565 Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan  
 566 Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents.  
 567 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
 568 14281–14290, 2024b.

569 Md Mohaiminul Islam, Tushar Nagarajan, Huiyu Wang, Gedas Bertasius, and Lorenzo Torresani.  
 570 Bimba: Selective-scan compression for long-range video question answering. *arXiv preprint  
 571 arXiv:2503.09590*, 2025.

572 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv  
 573 preprint arXiv:1611.01144*, 2016.

574 Peng Jin, Ryuichi Takanobu, Wancai Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified  
 575 visual representation empowers large language models with image and video understanding. In  
 576 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
 577 13700–13710, 2024.

578 Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning  
 579 events in videos. In *Proceedings of the IEEE international conference on computer vision*, pp.  
 580 706–715, 2017.

581 Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov,  
 582 Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open  
 583 web-scale filtered dataset of interleaved image-text documents. *Advances in Neural Information  
 584 Processing Systems*, 36:71683–71702, 2023.

585 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan  
 586 Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint  
 587 arXiv:2408.03326*, 2024a.

594 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.  
 595 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv*  
 596 *preprint arXiv:2407.07895*, 2024b.

597 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image  
 598 pre-training with frozen image encoders and large language models. In *International conference*  
 599 *on machine learning*, pp. 19730–19742. PMLR, 2023a.

600 KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhui Wang, Ping Luo, Yali Wang, Limin Wang, and  
 601 Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023b.

602 Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,  
 603 Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In  
 604 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
 605 22195–22206, 2024c.

606 Xinhao Li, Yi Wang, Jiashuo Yu, Xiangyu Zeng, Yuhua Zhu, Haian Huang, Jianfei Gao, Kunchang  
 607 Li, Yinan He, Chenting Wang, et al. Videochat-flash: Hierarchical compression for long-context  
 608 video modeling. *arXiv preprint arXiv:2501.00574*, 2024d.

609 Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language  
 610 models. In *European Conference on Computer Vision*, pp. 323–340. Springer, 2024e.

611 Zhiqi Li, Guo Chen, Shilong Liu, Shihao Wang, Vibashan VS, Yishen Ji, Shiyi Lan, Hao Zhang,  
 612 Yilin Zhao, Subhashree Radhakrishnan, et al. Eagle 2: Building post-training data strategies from  
 613 scratch for frontier vision-language models. *arXiv preprint arXiv:2501.14818*, 2025.

614 Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On  
 615 pre-training for visual language models. In *Proceedings of the IEEE/CVF conference on computer*  
 616 *vision and pattern recognition*, pp. 26689–26699, 2024.

617 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in*  
 618 *neural information processing systems*, 36:34892–34916, 2023.

619 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.  
 620 Llavanext: Improved reasoning, ocr, and world knowledge, 2024a.

621 Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi,  
 622 Shiyi Cao, Yuxian Gu, Dacheng Li, et al. Nvila: Efficient frontier visual language models. *arXiv*  
 623 *preprint arXiv:2412.04468*, 2024b.

624 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
 625 *arXiv:1711.05101*, 2017.

626 Muhammad Maaz, Hanooma Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt:  
 627 Towards detailed video understanding via large vision and language models. *arXiv preprint*  
 628 *arXiv:2306.05424*, 2023.

629 OpenAI. Gpt-4v(ision) system card. <https://openai.com/research/gpt-4v-system-card>, 2023.

630 OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2024.

631 Ji Qi, Ming Ding, Weihan Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li,  
 632 Yuxiao Dong, et al. Cogcom: Train large vision-language models diving into details through chain  
 633 of manipulations. *arXiv preprint arXiv:2402.04236*, 2024.

634 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
 635 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
 636 models from natural language supervision. In *International conference on machine learning*, pp.  
 637 8748–8763. PMLR, 2021.

638 Shuhuai Ren, Sishuo Chen, Shicheng Li, Xu Sun, and Lu Hou. Testa: Temporal-spatial token  
 639 aggregation for long-form video-language understanding. *arXiv preprint arXiv:2310.19060*, 2023.

648 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,  
 649 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th*  
 650 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.  
 651 2556–2565, 2018.

652 Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng Wu, Jun Chen, Chenchen Zhu, Zechun  
 653 Liu, Fanyi Xiao, Balakrishnan Varadarajan, Florian Bordes, et al. Longvu: Spatiotemporal adaptive  
 654 compression for long video-language understanding. *arXiv preprint arXiv:2410.17434*, 2024.

655 Yunhang Shen, Chaoyou Fu, Shaoqi Dong, Xiong Wang, Peixian Chen, Mengdan Zhang, Haoyu Cao,  
 656 Ke Li, Xiawu Zheng, Yan Zhang, et al. Long-vita: Scaling large multi-modal models to 1 million  
 657 tokens with leading short-context accuray. *arXiv preprint arXiv:2502.05177*, 2025.

658 Min Shi, Fuxiao Liu, Shihao Wang, Shijia Liao, Subhashree Radhakrishnan, Yilin Zhao, De-An  
 659 Huang, Hongxu Yin, Karan Sapra, Yaser Yacoob, et al. Eagle: Exploring the design space for  
 660 multimodal llms with mixture of encoders. *arXiv preprint arXiv:2408.15998*, 2024.

661 Yan Shu, Zheng Liu, Peitian Zhang, Minghao Qin, Junjie Zhou, Zhengyang Liang, Tiejun Huang,  
 662 and Bo Zhao. Video-xl: Extra-long vision language model for hour-scale video understanding.  
 663 *arXiv preprint arXiv:2409.14485*, 2024.

664 Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe  
 665 Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for  
 666 long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
 667 *Pattern Recognition*, pp. 18221–18232, 2024.

668 Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, Yuxuan  
 669 Wang, and Chao Zhang. video-salmonn: Speech-enhanced audio-visual large language models.  
 670 *arXiv preprint arXiv:2406.15704*, 2024.

671 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
 672 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
 673 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

674 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
 675 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
 676 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

677 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,  
 678 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the  
 679 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024a.

680 Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,  
 681 Lei Zhao, Song XiXuan, et al. Cogvilm: Visual expert for pretrained language models. *Advances*  
 682 *in Neural Information Processing Systems*, 37:121475–121499, 2024b.

683 Xidong Wang, Dingjie Song, Shunian Chen, Chen Zhang, and Benyou Wang. Longllava: Scal-  
 684 ing multi-modal llms to 1000 images efficiently via a hybrid architecture. *arXiv preprint*  
 685 *arXiv:2409.02889*, 2024c.

686 Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan  
 687 Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative and  
 688 discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022.

689 Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng,  
 690 Zun Wang, Yansong Shi, et al. Internvideo2: Scaling foundation models for multimodal video  
 691 understanding. In *European Conference on Computer Vision*, pp. 396–416. Springer, 2024d.

692 Yi Wang, Xinhao Li, Ziang Yan, Yinan He, Jiashuo Yu, Xiangyu Zeng, Chenting Wang, Changlian  
 693 Ma, Haian Huang, Jianfei Gao, et al. Internvideo2. 5: Empowering video mllms with long and rich  
 694 context modeling. *arXiv preprint arXiv:2501.12386*, 2025.

695 Stephan Wenger. Video redundancy coding in h. 263+. In *1997 International Workshop on Audio-  
 696 Visual Services over Packet Networks*. Citeseer, 1997.

702 Lin Xu, Yilin Zhao, Daquan Zhou, Zhiping Lin, See Kiong Ng, and Jiashi Feng. Pllava: Parameter-free  
 703 llava extension from images to videos for video dense captioning. *arXiv preprint arXiv:2404.16994*,  
 704 2024a.

705 Mingze Xu, Mingfei Gao, Zhe Gan, Hong-You Chen, Zhengfeng Lai, Haiming Gang, Kai Kang, and  
 706 Afshin Dehghan. Slowfast-llava: A strong training-free baseline for video large language models.  
 707 *arXiv preprint arXiv:2407.15841*, 2024b.

708 Fuzhao Xue, Yukang Chen, Dacheng Li, Qinghao Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian  
 709 Tang, Shang Yang, Zhiping Liu, et al. Longvila: Scaling long-context visual language models for  
 710 long videos. *arXiv preprint arXiv:2408.10188*, 2024.

711 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,  
 712 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint  
 713 arXiv:2412.15115*, 2024.

714 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,  
 715 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint  
 716 arXiv:2408.01800*, 2024.

717 Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren  
 718 Zhou. mplug-owl3: Towards long image-sequence understanding in multi-modal large language  
 719 models. *arXiv preprint arXiv:2408.04840*, 2024.

720 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language  
 721 image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*,  
 722 pp. 11975–11986, 2023.

723 Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu, Yuqian Yuan, Guanzheng Chen, Sicong Leng,  
 724 Yuming Jiang, Hang Zhang, Xin Li, et al. Videollama 3: Frontier multimodal foundation models  
 725 for image and video understanding. *arXiv preprint arXiv:2501.13106*, 2025.

726 Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong  
 727 Duan, Bin Wang, Linke Ouyang, et al. Internlm-xcomposer-2.5: A versatile large vision language  
 728 model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*, 2024a.

729 Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue  
 730 Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision.  
 731 *arXiv preprint arXiv:2406.16852*, 2024b.

732 Ruohong Zhang, Liangke Gui, Zhiqing Sun, Yihao Feng, Keyang Xu, Yuanhan Zhang, Di Fu,  
 733 Chunyuan Li, Alexander Hauptmann, Yonatan Bisk, et al. Direct preference optimization of video  
 734 large multimodal models from language model reward. *arXiv preprint arXiv:2404.01258*, 2024c.

735 Y Zhang, B Li, H Liu, Y Lee, L Gui, D Fu, J Feng, Z Liu, and C Li. Llava-next: A strong zero-shot  
 736 video understanding model. 2024d.

737 Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video  
 738 instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024e.

739 Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang,  
 740 Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video  
 741 understanding. *arXiv preprint arXiv:2406.04264*, 2024.

742 Orr Zohar, Xiaohan Wang, Yann Dubois, Nikhil Mehta, Tong Xiao, Philippe Hansen-Estruch,  
 743 Licheng Yu, Xiaofang Wang, Felix Juefei-Xu, Ning Zhang, et al. Apollo: An exploration of video  
 744 understanding in large multimodal models. *arXiv preprint arXiv:2412.10360*, 2024.

745

746

747

748

749

750

751

752

753

754

755

756  
757  
758  
759  

## Appendix

760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  

### A RELATED WORKS

776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  

**LMMs for General Vision Understanding** Building upon pretrained large language models (Brown et al., 2020; Touvron et al., 2023a;b) and vision foundation models (Radford et al., 2021; Zhai et al., 2023) with exceptional capabilities, Large Multimodal Models (LMMs) (Alayrac et al., 2022; Li et al., 2023a; Liu et al., 2023; Bai et al., 2023; Wang et al., 2024b; Yao et al., 2024) were initially proposed to learn mappings from visual inputs to linguistic representations through training on vast amounts of image-text paired data. These models have demonstrated remarkable visual-language understanding abilities. Subsequent studies has witnessed significant advancements in this field, with notable contributions including superior visual grounding Qi et al. (2024); Wang et al. (2024a), for web content recognition and visual agents (Hong et al., 2024b), enhanced multimodal reasoning capabilities (Chen et al., 2024e;d; Fu et al., 2024b; 2025; Shi et al., 2024; Zhang et al., 2024a), and the efforts contributed substantially to the open-source community (Lin et al., 2024; Liu et al., 2024b; Li et al., 2024a;b). These developments collectively represent substantial progress in expanding the applicability and performance of multimodal systems across diverse scenarios. Recent research has witnessed a surge of work developing LMMs with video understanding capabilities by incorporating video frames as training data (Chen et al., 2024d; Hong et al., 2024a; Ye et al., 2024; Yang et al., 2024; Wang et al., 2022; 2024d; Guo et al., 2025; Li et al., 2025; 2023b; 2024c).

787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  

**LMMs for Efficient Video Understanding** To handle longer video inputs, numerous efficient LMMs have been proposed. Early optimizations drew inspiration from long-context techniques in language models, including the adoption of sequence parallelism frameworks (Zhang et al., 2024b; Xue et al., 2024; Shen et al., 2025). Subsequently, a series of token merging techniques were developed to reduce spatiotemporal redundancy by consolidating the expanded token sequences resulting from longer video input (Li et al., 2024e; Maaz et al., 2023; Xu et al., 2024a; Shen et al., 2024; Sun et al., 2024; Shu et al., 2024; Ren et al., 2023). Additionally, significant research efforts have focused on video frame sampling strategies to minimize inter-frame redundancy (Ataallah et al., 2024; Zhang et al., 2024e; Xu et al., 2024b; Zohar et al., 2024; Chen et al., 2024b; Chang et al., 2024). More recently, novel architectures with enhanced long-sequence memory capabilities have emerged to facilitate the understanding of extended video sequences (Wang et al., 2024c; Islam et al., 2025).

805  
806  
807  
808  
809  

### B LIMITATIONS

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
10010  
10011  
10012  
10013  
10014  
10015  
10016  
10017  
10018  
10019  
10020  
10021  
10022  
10023  
10024  
10025  
10026  
10027  
10028  
10029  
10030  
10031  
10032  
10033  
10034  
10035  
10036  
10037  
10038  
10039  
10040  
10041  
10042  
10043  
10044  
10045  
10046  
10047  
10048  
10049  
10050  
10051  
10052  
10053  
10054  
10055  
10056  
10057  
10058  
10059  
10060  
10061  
10062  
10063  
10064  
10065  
10066  
10067  
10068  
10069  
10070  
10071  
10072  
10073  
10074  
10075  
10076  
10077  
10078  
10079  
10080  
10081  
10082  
10083  
10084  
10085  
10086  
10087  
10088  
10089  
10090  
10091  
10092  
10093  
10094  
10095  
10096  
10097  
10098  
10099  
100100  
100101  
100102  
100103  
100104  
100105  
100106  
100107  
100108  
100109  
100110  
100111  
100112  
100113  
100114  
100115  
100116  
100117  
100118  
100119  
100120  
100121  
100122  
100123  
100124  
100125  
100126  
100127  
100128  
100129  
100130  
100131  
100132  
100133  
100134  
100135  
100136  
100137  
100138  
100139  
100140  
100141  
100142  
100143  
100144  
100145  
100146  
100147  
100148  
100149  
100150  
100151  
100152  
100153  
100154  
100155  
100156  
100157  
100158  
100159  
100160  
100161  
100162  
100163  
100164  
100165  
100166  
100167  
100168  
100169  
100170  
100171  
100172  
100173  
100174  
100175  
100176  
100177  
100178  
100179  
100180  
100181  
100182  
100183  
100184  
100185  
100186  
100187  
100188  
100189  
100190  
100191  
100192  
100193  
100194  
100195  
100196  
100197  
100198  
100199  
100200  
100201  
100202  
100203  
100204  
100205  
100206  
100207  
100208  
100209  
100210  
100211  
100212  
100213  
100214  
100215  
100216  
100217  
100218  
100219  
100220  
100221  
100222  
100223  
100224  
100225  
100226  
100227  
100228  
100229  
100230  
100231  
100232  
100233  
100234  
100235  
100236  
100237  
100238  
100239  
100240  
100241  
100242  
100243  
100244  
100245  
100246  
100247  
100248  
100249  
100250  
100251  
100252  
100253  
100254  
100255  
100256  
100257  
100258  
100259  
100260  
100261  
100262  
100263  
100264  
100265  
100266  
100267  
100268  
100269  
100270  
100271  
100272  
100273  
100274  
100275  
100276  
100277  
100278  
100279  
100280  
100281  
100282  
100283  
100284  
100285  
100286  
100287  
100288  
100289  
100290  
100291  
100292  
100293  
100294  
100295  
100296  
100297  
100298  
100299  
100300  
100301  
100302  
100303  
100304  
100305  
100306  
100307  
100308  
100309  
100310  
100311  
100312  
100313  
100314  
100315  
100316  
100317  
100318  
100319  
100320  
100321  
100322  
100323  
100324  
100325  
100326  
100327  
100328  
100329  
100330  
100331  
100332  
100333  
100334  
100335  
100336  
100337  
100338  
100339  
100340  
100341  
100342  
100343  
100344  
100345  
100346  
100347  
100348  
100349  
100350  
100351  
100352  
100353  
100354  
100355  
100356  
100357  
100358  
100359  
100360  
100361  
100362  
100363  
100364  
100365  
100366  
100367  
100368  
100369  
100370  
100371  
100372  
100373  
100374  
100375  
100376  
100377  
100378  
100379  
100380  
100381  
100382  
100383  
100384  
100385  
100386  
100387  
100388  
100389  
100390  
100391  
100392  
100393  
100394  
100395  
100396  
100397  
100398  
100399  
100400  
100401  
100402  
100403  
100404  
100405  
100406  
100407  
100408  
100409  
100410  
100411  
100412  
100413  
100414  
100415  
100416  
100417  
100418  
100419  
100420  
100421  
100422  
100423  
100424  
100425  
100426  
100427  
100428  
100429  
100430  
100431  
100432  
100433  
100434  
100435  
100436  
100437  
100438  
100439  
100440  
100441  
100442  
100443  
100444  
100445  
100446  
100447  
100448  
100449  
100450  
100451  
100452  
100453  
100454  
100455  
100456  
100457  
100458  
100459  
100460  
100461  
100462  
100463  
100464  
100465  
100466  
100467  
100468  
100469  
100470  
100471  
100472  
100473  
100474  
100475  
100476  
100477  
100478  
100479  
100480  
100481  
100482  
100483  
100484  
100485  
100486  
100487  
100488  
100489  
100490  
100491  
100492  
100493  
100494  
100495  
100496  
100497  
100498  
100499  
100500  
100501  
100502  
100503  
100504  
100505  
100506  
100507  
100508  
100509  
100510  
100511  
100512  
100513  
100514  
100515  
100516  
100517  
100518  
100519  
100520  
100521  
100522  
100523  
100524  
100525  
100526  
100527  
100528  
100529  
100530  
100531  
100532  
100533  
100534  
100535  
100536  
100537  
100538  
100539  
100540  
100541  
100542  
100543  
100544  
100545  
100546  
100547  
100548  
100549  
100550  
100551  
100552  
100553  
100554  
100555  
100556  
100557  
100558  
100559  
100560  
100561  
100562  
100563  
100564  
100565  
100566  
100567  
100568  
100569  
100570  
100571  
100572  
100573  
100574  
100575  
100576  
100577  
100578  
100579  
100580  
100581  
100582  
100583  
100584  
100585  
100586  
100587  
100588  
100589  
100590  
100591  
100592  
100593  
100594  
100595  
100596  
100597  
100598  
100599  
100600  
100601  
100602  
100603  
100604  
100605  
100606  
100607  
100608  
100609  
100610  
100611  
100612  
100613  
100614  
100615  
100616  
100617  
100618  
100619  
100620  
100621  
100622  
100623  
100624  
100625  
100626  
100627  
100628  
100629  
100630  
100631  
100632  
100633  
100634  
100635  
100636  
100637  
100638  
100639  
100640  
100641  
100642  
100643  
100644  
100645  
100646  
100647  
100648  
100649  
100650  
100651  
100652  
100653  
100654  
100655  
100656  
100657  
100658  
100659  
100660  
100661  
100662  
100663  
100664  
100665  
100666  
100667  
100668  
100669  
100670  
100671  
100672  
100673  
100674  
100675  
100676  
100677  
100678  
100679  
100680  
100681  
100682  
100683  
100684  
100685  
100686  
100687  
100688  
100689  
100690  
100691  
100692  
100693  
100694  
100695  
100696  
100697  
100698  
100699  
100700  
100701  
100702  
100703  
100704  
100705  
100706  
100707  
100708  
100709  
100710  
100711  
100712  
100713  
100714  
100715  
100716  
100717  
100718  
100719  
100720  
100721  
100722  
100723  
100724  
100725  
100726  
100727  
100728  
100729  
100730  
100731  
100732  
100733  
100734  
100735  
100736  
100737  
100738  
100739  
100740  
100741  
100742  
100743  
100744  
100745  
100746  
100747  
100748  
100749  
100750  
100751  
100752  
100753  
100754  
100755  
100756  
100757  
100758  
100759  
100760  
100761  
100762  
100763  
100764  
100765  
100766  
100767  
100768  
100769  
100770  
100771  
100772  
100773  
100774  
100775  
100776  
100777  
100778  
100779  
100780  
100781  
100782  
100783  
100784  
100785  
100786  
100787  
100788  
100789  
100790  
100791  
100792  
100793  
100794  
100795  
100796  
100797  
100798  
100799  
100800  
100801  
100802  
100803  
100804  
100805  
100806  
100807  
100808  
100809  
100810  
100811  
100812  
100813  
100814  
100815  
100816  
100817  
100818  
100819  
100820  
100821  
100822  
100823  
100824  
100825  
100826  
100827  
100828  
100829  
100830  
100831  
100832  
100833  
100834  
100835  
100836  
100837  
100838  
100839  
100840  
100841  
100842  
100843  
100844  
100845  
100846  
100847  
100848  
100849  
100850  
100851  
100852  
100853  
100854  
100855  
100856  
100857  
100858  
100859  
100860  
100861  
100862  
100863  
100864  
100865  
100866  
100867  
100868  
100869  
100870  
100871  
100872  
100873  
100874  
100875  
100876  
100877  
100878  
100879  
100880  
100881  
100882  
100883  
100884  
100885  
100886  
100887  
100888  
100889  
100890  
100891  
100892  
100893  
100894  
100895  
100896  
100897  
100898  
100899  
100900  
100901  
100902  
100903  
100904  
100905  
100906  
100907  
100908  
100909  
100910  
100911  
100912  
100913  
100914  
100915  
100916  
100917  
100918  
100919  
100920  
100921  
100922  
100923  
100924  
100925  
100926  
100927  
100928  
100929  
100930  
100931  
100932  
100933  
100934  
100935  
100936  
100937  
100938  
100939  
100940  
100941  
100942  
100943  
100944  
100945  
100946  
100947  
100948  
100949  
100950  
100951  
100952  
100953  
100954  
100955  
100956  
100957  
100958  
100959  
100960  
100961  
100962  
100963  
100964  
100965  
100966  
100967  
100968  
100969  
100970  
100971  
100972  
100973  
100974  
100975  
100976  
100977  
100978  
100979  
100980  
100981  
100982  
100983  
100984  
100985  
100986  
100987  
100988  
100989  
100990  
100991  
100992  
100993  
100994  
100995  
100996  
100997  
100998  
100999  
1001000  
100101  
100102  
100103  
100104  
100105  
100106  
100107  
100108  
100109  
100110  
100111  
100112  
100113  
100114  
100115  
100116  
100117  
100118  
100119  
100120  
100121  
100122  
100123  
100124  
100125  
100126  
100127  
100128  
100129  
100130  
100131  
100132  
100133  
100134  
100135  
100136  
100137  
100138  
100139  
100140  
100141  
100142  
100143  
100144  
100145  
100146  
100147  
100148  
100149  
100150  
100151  
100152  
100153  
100154  
100155  
100156  
100157  
100158  
100159  
100160  
100161  
100162  
100163  
100164  
100165  
100166  
100167  
100168  
100169  
100170  
100171  
100172  
100173  
100174  
100175  
100176  
100177  
100178  
100179  
100180  
100181  
100182  
100183  
100184  
100185  
100186  
100187  
100188  
100189  
100190  
100191  
100192  
100193  
100194  
100195  
100196  
100197  
100198  
100199  
100200  
100201  
100202  
100203  
100204  
100205  
100206  
100207  
100208  
100209  
100210  
100211  
100212  
100213  
100214  
100215  
100216  
100217  
100218  
100219  
100220  
100221  
100222  
100223  
100224  
100225  
100226  
100227  
100228  
100229  
100230  
100231  
100232  
100233  
100234  
100235  
100236  
100237  
100238  
100239  
100240  
100241  
100242  
100243  
100244  
100245  
100246  
100247  
100248  
100249  
100250  
100251  
100252  
100253  
100

810  
811  
812  
C TRAINING DATA DETAILS

Modality	Task	# Samples	Dataset
Image-Text	Interleaved Pairs	20K	OBELICS
	Single-Image Captioning	558K	LCS
Video-Text	Captioning	95K	FineVideoCaptions, ANetCaptions

818  
819  
820  
821  
Table 4: Training data statistics for the alignment stage.

Modality	Task	# Samples	Dataset
Text	Instruction	93K	Evo-Instruct
Image-Text	Interleaved Pairs	20K	OBELICS
	Single-Image Captioning	50K	LCS
	Single-Image Tasks	2.8M	LLaVAOneVision
Video-Text	Captioning	5K	FineVideoCaptions, AnetCaptions
	VQA	75K	FineVideoQAs
	Dense Captioning	38K	ShareGPT4Video

832  
833  
834  
835  
Table 5: Training data statistics for the pre-training stage.

Modality	Task	# Samples	Dataset
Text	Instruction	9K	Evo-Instruct
Image-Text	Single-Image Tasks	40K	LLaVA-OneVision-SingleImage
	Multi-Images Tasks	100K	LLaVA-OneVision-MultiImages
Video-Text	Captioning	52K	TextVR, MovieChat, YouCook2
	Dense Captioning	4K	ShareGPT4Video
	Classification	1K	Kinetics-710
	VQA	354K	NExT-QA, CLEVRER, EgoQA TGIF, ShareGPTVideo, FineVideoQAs
	Instruction	188K	VideoChatGPT, VideoChat, LongVILA

848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
Table 6: Training data statistics for the supervised fine-tuning stage.

864 **D ADDITIONAL EVALUATIONS**

867 <b>Model</b>	868 <b>NExT-QA</b>		869 <b>ActivityNet-QA</b>		870 <b>Video-ChatGPT</b>			
	871 acc	872 acc	873 Correctness	874 Detail	875 Context	876 Temporal	877 Consistency	
878 LLaMA-VID (7B)	879 -	880 47.4/3.3	881 2.96	882 3.00	883 3.53	884 2.46	885 2.51	
886 Chat-UniVi (7B)	887 -	888 46.1/3.3	889 2.89	890 2.91	891 3.46	892 2.89	893 2.81	
894 Video-LLaVA (7B)	895 -	896 45.3/3.3	897 2.87	898 2.94	899 3.44	900 2.45	901 2.51	
902 VideoChat2 (7B)	903 68.6	904 49.1/3.3	905 3.16	906 3.08	907 3.69	908 2.56	909 3.14	
910 VideoLLaMA2 (7B)	911 75.6	912 50.2/3.3	913 3.30	914 33.18	915 3.78	916 2.66	917 3.12	
918 LLaVA-NeXT-Video (7B)	919 78.2	920 53.5/3.2	921 3.39	922 3.29	923 3.92	924 2.60	925 3.12	
926 Quicksviewer (8B)	927 77.5	928 47.6/2.7	929 3.10	930 3.11	931 3.09	932 2.48	933 3.04	

874 Table 7: Evaluation results on more benchmarks.

875 **E QUALITATIVE ANALYSIS**

876 Our unified perception paradigm demonstrates efficient visual understanding capabilities, successfully  
 877 processing: 1) lengthy videos, 2) high-resolution single images, and 3) multi-image contextual  
 878 reasoning tasks. Through representative qualitative analysis, we validate the model’s performance  
 879 across these diverse inputs. Specifically, we evaluate the video understanding capabilities of our  
 880 model through documentary recording and sports competitions, while the image understanding  
 881 proficiency across diverse domains including outdoor road scenes, physics/biological/historical  
 882 knowledge systems, and multi-image geographical analysis.

883 In the documentary depicting a penguin chick’s perilous encounter and subsequent escape, our model  
 884 demonstrates comprehensive video understanding by: (1) identifying the nature of the unexpected  
 885 attack, (2) precisely locating its temporal occurrence, and (3) summarizing the eventual outcome -  
 886 showcasing its advanced capabilities in long-form video understanding, including temporal action  
 887 recognition, event narrative abstraction, and exact timestamp localization.

888 In a lacrosse match video documenting a scoring play, our model precisely identifies the initiating  
 889 player when queried about the offensive sequence, subsequently describing the play development  
 890 and correctly specifying both the scoring player’s identity and jersey number. This demonstrates  
 891 the model’s dual capability of (1) recognizing individual athletes in sports footage and (2) logically  
 892 summarizing dynamic game situations.

893 We further validate our model’s image understanding capabilities across extensive scenarios. As  
 894 illustrated in Figure 6, these include: 1) traffic sign recognition in driving environments, 2) physics  
 895 problem solving, 3) biological image interpretation and association, 4) historical scene identification,  
 896 and 5) multi-image geographical reasoning. These examples demonstrate that while exhibiting strong  
 897 video understanding, our model maintains robust image understanding capabilities. Benefiting from  
 898 interleaved image-text training data, the model additionally acquires in-context learning capabilities  
 899 for complex reasoning tasks.

900  
 901  
 902  
 903  
 904  
 905  
 906  
 907  
 908  
 909  
 910  
 911  
 912  
 913  
 914  
 915  
 916  
 917

918  
919  
920  
921  
922

User: What is the main incident that happened in this video?

Quicksviewer: The main incident that happened in the video is a group of penguins standing on the snow, while a bird approaches them. The bird tries to attack one of the penguins, but it fails.

User: During what time period (HH:MM:SS - HH:MM:SS) did the attack occur? Were any penguins injured?

Quicksviewer: The attack occurred between 0:47 and 0:53, and no penguins were injured during the incident.

923  
924  
925  
926  
927  
928  
929

User: Which White team player started the last offensive sequence?

Quicksviewer: 13, Barger.

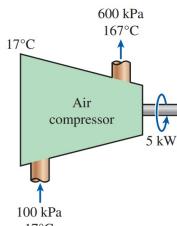
User: What happened next?

Quicksviewer: He passed the ball to his teammate, who then passed it to another teammate. White 7 Miller shot and scored a goal. The teammates celebrated by hugging and cheering.

930  
931  
932  
933

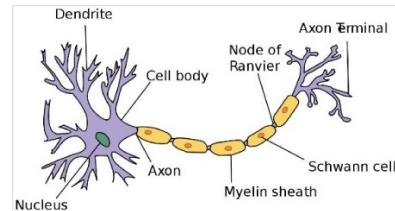
User: Where can I go by turning left?

Quicksviewer: Las Vegas.



User: Describe the process shown in the diagram including the specific values.

Quicksviewer: The thermodynamic process involves an air compressor of 5kW. The input values are 100 kPa and 17°C, while the output values are 167 kPa and 176°C.



User: Identify the biological structure and specify its computational counterpart in CS disciplines.

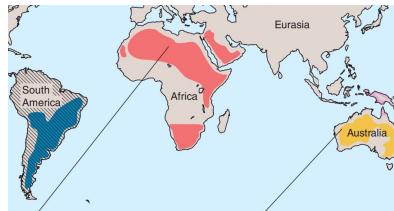
Quicksviewer: The image depicts a neuron, a fundamental structure in nervous tissue. The corresponding concept in computer science would be the neuron in artificial neural networks.

934



User: Where is it located?

Quicksviewer: This is the Longmen Grottoes, located in Luoyang, Henan, China.



User: What are the species of the birds, and which one inhabits the red area marked in the first map?

Quicksviewer: The first bird is a South American ostrich, the second is an African ostrich, and the third is a cassowary. The African ostrich inhabits the red area marked in the first map.

935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

Figure 6: Qualitative analysis shows that Quicksviewer effectively understands lengthy documentary and sports videos, as well as informative single and multiple images.