

# HORNet: Hindsight Optimization Reasoning for Task-Guided Frame Selection in Visual Language Models

Anonymous CVPR submission

Paper ID 2

## Abstract

001 *Video question answering (VQA) with vision-language*  
002 *models (VLMs) depends critically on which frames are se-*  
003 *lected from the input video, yet most systems rely on uni-*  
004 *form or heuristic sampling that cannot be optimized for*  
005 *downstream answering quality. We introduce **HORNet**, a*  
006 *lightweight frame selection policy trained with Group Rel-*  
007 *ative Policy Optimization (GRPO) to learn which frames*  
008 *a frozen VLM needs to answer questions correctly. With*  
009 *fewer than 1M trainable parameters, HORNet reduces in-*  
010 *put frames by up to 99% and VLM processing time by up to*  
011 *93%, while improving answer quality on short-form bench-*  
012 *marks (+1.7% F1 on MSVD-QA) and achieving strong per-*  
013 *formance on temporal reasoning tasks (+7.3 points over*  
014 *uniform sampling on NEXt-QA). We formalize this as Select*  
015 *Any Frames (SAF), a task that decouples visual input cu-*  
016 *ration from VLM reasoning, and show that GRPO-trained*  
017 *selection generalizes better out-of-distribution than super-*  
018 *vised and PPO alternatives. HORNet’s policy further trans-*  
019 *fers across VLM answerers without retraining, yielding an*  
020 *additional 8.5% relative gain when paired with a stronger*  
021 *model. Evaluated across six benchmarks spanning 341,877*  
022 *QA pairs and 114.2 hours of video, our results demonstrate*  
023 *that optimizing what a VLM sees is a practical and comple-*  
024 *mentary alternative to optimizing what it generates while*  
025 *improving efficiency.*

## 026 1. Introduction

027 Existing state-of-the-art VLMs rely on scaling large visual-  
028 text data pairs to improve performance [8, 13, 23, 31], and  
029 while these efforts have yielded measurable gains on VQA  
030 benchmarks, the underlying mechanism—a vision encoder  
031 tokenizes image patches, a projection layer maps them into  
032 the language model’s embedding space, and an autoregres-  
033 sive LLM decodes the response—has remained largely un-  
034 changed. Videos are first sampled and transformed into  
035 visual tokens, which are then aligned with textual inputs

through cross-attention mechanisms [12]. The data-hungry 036  
nature of such architecture brings significant downfalls in 037  
“Small Data” domains, where data collection is costly, inef- 038  
ficient and sometimes facing regulations, limiting the adopt- 039  
ing in these situations [4, 27, 29]. Some approaches at- 040  
tempt to enhance spatial and temporal reasoning by apply- 041  
ing LoRA-based [15] fine-tuning on datasets specifically 042  
curated for reasoning tasks, thereby improving a model’s 043  
ability to understand and interpret video content. Most other 044  
methods rely on minor modifications to the attention archi- 045  
tecture to adapt the model to specific domains and applica- 046  
tions [17]. VLMs for video largely inherit the architecture 047  
and biases of image-based models, with temporal reasoning 048  
added only superficially. Video-LLaVA [22] samples eight 049  
frames through a shared frozen encoder with no tempo- 050  
ral module, LLaVA-OneVision [20] treats video frames as 051  
multiple images and shows that an image-only checkpoint 052  
already performs competitively on video benchmarks, and 053  
Video-ChatGPT [25] reduces temporal reasoning to mean- 054  
pooling of per-frame CLIP features. The resulting perform- 055  
ance gap is stark: InternVL2.5-78B achieves 95.1% on 056  
DocVQA [26] yet only 72.1% on Video-MME [8]; a 23- 057  
point drop that reflects the absence of genuine temporal rea- 058  
soning rather than mere task difficulty. In practice, most 059  
systems rely on pragmatic frame-sampling strategies, effec- 060  
tively reducing videos to sets of isolated images, leading to 061  
unavoidable information loss while increasing signal-noise 062  
ratio. The frames with necessary information for VLM 063  
to reason might be discarded through this sampling proce- 064  
dures, degrading answers’ quality. The broader question 065  
of how visual inputs should be structured, represented, and 066  
processed within VLMs to preserve key information while 067  
removing noise remains insufficiently examined. Although 068  
a few studies highlight the role of sampling as a form of 069  
filtering [6, 30, 41], its importance is largely overlooked. 070

Given that scaling data has improved image understand- 071  
ing far more than video understanding under this paradigm, 072  
we turn to a complementary axis of improvement: opti- 073  
mizing how models reason over their visual inputs through 074  
reinforcement-learning-based fine-tuning. Most existing 075

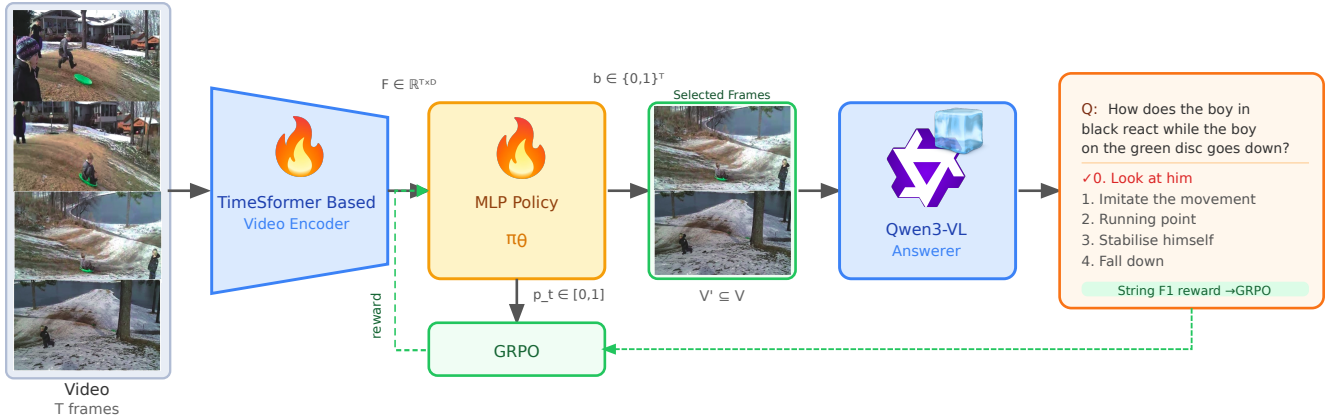


Figure 1. **HORNet pipeline.** Given a video  $\mathbf{V} = \{v_1, v_2, \dots, v_T\}$  with  $T$  uniformly sampled frames, our TimeSFormer-based video encoder  $E$  extracts per-frame features  $\mathbf{F} \in \mathbb{R}^{T \times D}$ . A lightweight trainable MLP policy  $\pi_\theta$  scores each frame independently, producing keep probabilities  $p_t \in [0, 1]$  and a binary selection mask  $\mathbf{b} \in \{0, 1\}^T$ . Only the frames selected by the mask ( $\mathbf{V}' \subseteq \mathbf{V}$ ) are passed to the frozen Qwen3-VL answerer. For example, “How does the boy in black react while the boy on the green disc goes down?”, HORNet selects only the frames capturing the key interaction moment, discarding irrelevant context; correctly predicting “Look at him”. At training time, GRPO samples  $K$  candidate subsets, evaluates each via String F1 reward against the ground-truth answer, and updates  $\pi_\theta$  through group-normalized policy gradients. The VLM answerer remain frozen throughout while the encoder and MLP policy are trainable.

076 VLMs requires supervised fine-tuning (SFT) to adapt to  
 077 new domains, which are costly and inefficient with small  
 078 data. The advent of Group Relative Policy Optimization  
 079 (GRPO) [9, 28] has opened a new avenue for end-to-end op-  
 080 timization of language model behavior via verifiable reward  
 081 signals. Inspired by its success in guiding the gradients to  
 082 improve the *outputs* of VLMs [1, 10, 39], we ask a funda-  
 083 mentally different question: can GRPO be used to optimize  
 084 *what a VLM sees (inputs)*, rather than what it says (gra-  
 085 dients)? We formalize the problem as Select Any Frames  
 086 (SAF): given a video and a question, select the subset of  
 087 frames from the full temporal sequence that maximizes the  
 088 downstream VLM’s ability to produce the correct answer.  
 089 SAF treats frame selection as a sequential decision problem  
 090 amenable to reinforcement learning, with the VLM’s QA  
 091 accuracy providing a direct, task-grounded reward signal.  
 092 This framing is intentionally simple and general; the SAF  
 093 policy is modular and can be paired with any downstream  
 094 VLM without any modifications to the architecture.

095 To address the issue, we present HORNet: Hindsight  
 096 Optimization Reasoning, a three-stage SAF pipeline that  
 097 optimize VLMs’ performance by selecting optimal frames  
 098 from video input (see Fig. 1). First, a trainable lightweight  
 099 video encoder extracts rich spatiotemporal features for each  
 100 frame independently. Second, a lightweight trainable mul-  
 101 tilayer perceptron (MLP) policy consumes these features  
 102 and outputs a per-frame keep probability. Third, GRPO  
 103 trains the video encoder and MLP by sampling multiple  
 104 candidate frame subsets per video, passing each to a frozen  
 105 VLM model for answering and computing rewards. Only  
 106 the MLP and video encoder are trained; the VLM re-

main frozen throughout. This design makes HORNet ex-  
 ceptionally parameter-efficient-suitable for small-data set-  
 tings where full fine-tuning of large models is infeasible-  
 while still benefiting from the representational capacity of  
 pretrained video and language foundations. Operating at  
 the frame-selection level rather than the token level allows  
 HORNet to substantially reduce both the memory footprint  
 and the inference latency of downstream VLM processing,  
 with these gains becoming even more pronounced as model  
 size increases.

We train HORNet on a diverse multi-dataset corpus  
 spanning MSRVTT-QA [33], MSVD-QA [34], and NExT-  
 QA [32], totaling 17,350 videos, 341,877 QA pairs, and  
 114.2 hours of content. This breadth covers descriptive,  
 causal, and temporal question types, pushing the policy  
 to discover generalizable selection strategies rather than  
 dataset-specific shortcuts. In short, our contributions are:

- We introduce **SAF (Select Any Frames)**, a task formula-  
 tion that decouples frame selection from VLM reasoning  
 and enables direct reward-based optimization of visual in-  
 puts.
- We propose **HORNet**, a GRPO-trained frame selection  
 policy built on frozen video and language foundations,  
 trainable with minimal parameters.
- We demonstrate that GRPO can be redirected from opti-  
 mizing VLM outputs to optimizing VLM inputs; a concep-  
 tual shift that is both more parameter-efficient and  
 more generally applicable.
- We provide a large-scale training benchmark combining  
 three VideoQA datasets (341,877 QA pairs, 114.2 hours)  
 for evaluating frame selection methods.

Table 1. **Research gap.** Existing frame selection methods satisfy at most two of four desirable properties simultaneously. HORNet is the first to achieve all four: learned selection, reward-based optimization, a fully frozen VLM, and parameter efficiency (<1M parameters). ✓ fully supported, ✗ not supported, ~ partial.

Method	Learned Selection	Reward Optimized	Frozen VLM	Param. Efficient
Uniform Sampling	✗	✗	✓	✓
SeViLA [36]	✓	✗	✗	✗
Frame-Voyager [37]	✓	✗	✗	✗
F2C [30]	✗	✗	✓	✓
ReFoCUS [19]	✓	✓	~	✗
ViaRL [35]	✓	✓	✗	✗
<b>HORNet (Ours)</b>	✓	✓	✓	✓

## 1.1. Small Data Statement

This work qualifies as small data research on two fronts. First, HORNet is designed for settings where annotated video-question-answer data is scarce or expensive to collect. Rather than fine-tuning a billion-parameter VLM; which typically requires hundreds of thousands of domain-specific examples, HORNet trains fewer than 1M parameters, meaning that the method can be deployed in domains where only a small number of labeled video QA pairs are available, such as medical procedures, surveillance, or industrial inspection, without risking catastrophic forgetting or overfitting of the foundation models. Second, HORNet’s training strategy is explicitly chosen to maximize sample efficiency. GRPO generates multiple candidate frame selections per training example and computes rewards from the frozen VLM’s own outputs, effectively amplifying each labeled sample by a factor of  $K = 8$  without requiring any additional annotation. Our ablation study (Table 4) confirms this advantage. Furthermore, the trained policy transfers to a different VLM answerer without retraining (Table 6), eliminating the need to recollect data when the downstream model changes. Together, these design choices—frozen foundations, reward amplification, and transferable policies—make HORNet particularly suited to the data-scarce regimes that motivate this work.

## 2. Related Work

We summarize the positioning of existing methods in Table 1 across four desirable properties: whether the method uses learned frame selection, whether it is optimized via downstream reward signals, whether the VLM remains frozen during training, and whether it is parameter-efficient. Existing approaches satisfy at most two of these properties simultaneously. HORNet is the first to satisfy all four.

**Frame selection for video understanding.** The importance of selecting the right frames, rather than sampling uni-

formly, has been recognized since Buch *et al.* [6] demonstrated that a single well-chosen frame, identified by a permutation-invariant attention module over frozen CLIP embeddings, suffices for many VideoQA benchmarks. This finding motivated a line of work on learned selection. SeViLA [36] chains a Localizer and Answerer fine-tuned from BLIP-2, using pseudo-labels from the Answerer to self-refine the Localizer. Frame-Voyager [37] enumerates frame combinations and trains a supervised selector by ranking subsets according to a Video-LLM’s prediction loss. VidF4 [21] proposes differentiable frame scoring that jointly considers question relevance and inter-frame diversity. On the training-free side, F2C [30] segments videos into temporally coherent clips using watershed-based scoring and CLIP query relevance, demonstrating that clip-level temporal coherence can outperform isolated frame selection. A.I.R. [41] employs a VLM to iteratively decompose queries and evaluate small frame batches, trading inference cost for selection accuracy. BOLT [24] and Q-Frame [40] also use CLIP similarity with different sampling strategies to balance query relevance and coverage. These methods demonstrate that the *when* of frame selection matters as much as the *how*. HORNet differs from all of these in that it *learns* the selection policy end-to-end from downstream QA rewards, without heuristics, pseudo-labels, or combinatorial enumeration.

**Reinforcement learning for frame selection.** A concurrent wave of work applies RL specifically to visual input selection. ReFoCUS [19] trains an autoregressive frame selector using reward signals from a reference VLM’s answer confidence margins. ViaRL [35] co-evolves a frame selector and answerer via iterated amplification RL, achieving strong results on temporal needle QA tasks. FrameMind [14] introduces Frame-Interleaved Chain-of-Thought with a GRPO variant for multi-turn dynamic resolution frame sampling. VideoBrain [2] trains an agent that decides when to invoke additional frame sampling using GRPO at the agent-invocation level. While these methods share our motivation, they differ in key respects: ReFoCUS uses autoregressive selection; ViaRL modifies both selector and answerer; FrameMind requires multi-turn agentic inference; and VideoBrain operates at the coarse sampling decision level rather than per-frame scoring. HORNet is simpler by design; a single forward pass through a frozen encoder followed by an MLP produces selection probabilities, and GRPO training requires no modifications to either the encoder or the VLM. This simplicity makes it particularly suited to small-data and resource-constrained settings.

**GRPO for vision-language models.** Group Relative Policy Optimization [28] was introduced to train language models on verifiable rewards without a critic network, later

224 scaled in DeepSeek-R1 [9] to incentivize emergent reason-  
 225 ing. Its application to vision-language models has since ex-  
 226 panded rapidly: Video-R1 [10] applies temporal contrastive  
 227 GRPO to video MLLMs; R1-VL [39] extends it to step-  
 228 wise multimodal reasoning; DeepVideo-R1 [1] addresses  
 229 the vanishing advantage problem specific to video GRPO;  
 230 Vision-R1 [16] demonstrates data-efficient GRPO training  
 231 for visual math reasoning; and GRPO-CARE [7] addresses  
 232 reasoning consistency degradation. Critically, all of these  
 233 works apply GRPO to improve what the VLM *generates*;  
 234 optimizing output distributions. HORNet redirects GRPO  
 235 toward optimizing what the VLM *receives*; a complemen-  
 236 tary direction that has not been explored prior to this work.

237 HORNet sits at the intersection of these three threads.  
 238 It inherits the GRPO optimization framework from the reason-  
 239 ing literature, the select-then-answer pipeline from the  
 240 frame selection literature, and the frozen foundation model  
 241 paradigm from efficient video VLMs. The key novelty is  
 242 using GRPO’s group-relative advantage estimation-critic-  
 243 free, scalable, and reward-agnostic-to directly maximize  
 244 downstream QA performance through frame selection, with  
 245 a parameter footprint small enough for low-data regimes.

### 246 3. Method

247 In this section, we first formally define the Select Any  
 248 Frames (SAF) problem. We then introduce the HORNet  
 249 architecture and detail the GRPO-based training procedure.

#### 250 3.1. Problem Formulation

251 Let  $\mathbf{V} = \{v_1, v_2, \dots, v_T\}$  denote a video represented by  
 252  $T$  uniformly sampled frames, where  $v_t \in \mathbb{R}^{H \times W \times C}$  is the  
 253  $t$ -th RGB frame. Let  $q$  be a natural language question and  
 254  $a$  the corresponding ground-truth answer. We denote by  $\mathcal{D}$   
 255 a dataset of triplets  $(\mathbf{V}, q, a)$ . A video encoder  $E$  extracts  
 256 spatiotemporal per-frame representations  $\mathbf{F} \in \mathbb{R}^{T \times D}$ , from  
 257 which a lightweight policy selects a subset  $\mathbf{V}' \subseteq \mathbf{V}$ . A  
 258 pretrained and frozen VLM  $\mathcal{M}$  then produces a predicted  
 259 answer  $\hat{a} = \mathcal{M}(\mathbf{V}', q)$ .

260 The goal of SAF is to learn a parameterized policy  $\pi_\theta$   
 261 that selects a subset  $\mathbf{V}' = \pi_\theta(\mathbf{V}, q)$  maximizing down-  
 262 stream answering performance. Formally, we seek:

$$263 \theta^* = \arg \max_{\theta} \mathbb{E}_{(\mathbf{V}, q, a) \sim \mathcal{D}} [R(\mathcal{M}(\pi_\theta(\mathbf{V}, q), q), a)], \quad (1)$$

264 where  $R(\hat{a}, a)$  is a task-specific reward function measuring  
 265 the quality of the predicted answer  $\hat{a}$  relative to the ground-  
 266 truth  $a$  (e.g., exact match accuracy). The VLM  $\mathcal{M}$  remains  
 267 frozen during training; only the policy parameters  $\theta$  are op-  
 268 timized.

269 **Policy parameterization.** We represent the policy output  
 270 as a binary selection mask  $\mathbf{b} = (b_1, \dots, b_T) \in \{0, 1\}^T$ ,

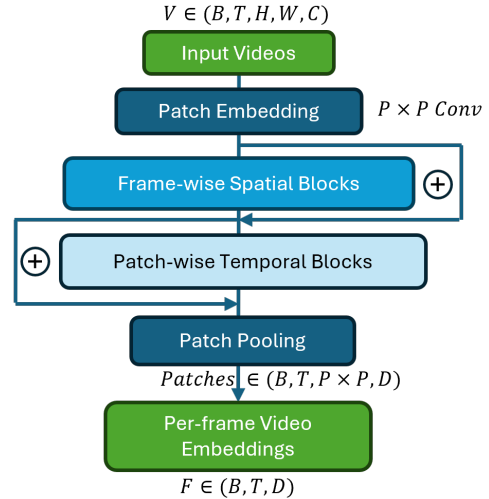


Figure 2. **HORNet encoder**  $E$ . Input frames are patchified with a  $P \times P$  convolution, processed by spatial self-attention within each frame, and then by temporal self-attention across frames at each patch location. The resulting temporally contextualized patch tokens are pooled to yield per-frame video representations used by HORNet for frame selection.  $B$  is batch size,  $T$  is frame count and  $D$  is hidden dimension. We set  $P=16$ ,  $T=32$  and  $D=768$  in our training.

where  $b_t = 1$  indicates that frame  $v_t$  is selected. The se-  
 271 lected subset is therefore  $\mathbf{V}' = \{v_t \mid b_t = 1\}$ . The policy  
 272 defines a distribution over binary masks:  
 273

$$274 \mathbf{b} \sim \pi_\theta(\mathbf{b} \mid \mathbf{V}, q), \quad (2)$$

275 which factorizes over frames via independent Bernoulli de-  
 276 cisions:

$$277 \pi_\theta(\mathbf{b} \mid \mathbf{V}, q) = \prod_{t=1}^T \text{Bernoulli}(b_t \mid p_t), \quad (3)$$

278 where  $p_t \in [0, 1]$  is the selection probability for frame  $v_t$ ,  
 279 predicted by the policy network.

280 This formulation imposes no temporal ordering or con-  
 281 tiguity constraints on frame selection, hence the name Se-  
 282 lect Any Frames (SAF). The policy may therefore learn to  
 283 select temporally sparse key events, short critical intervals,  
 284 or dense motion segments, depending solely on what maxi-  
 285 mizes the task-driven reward.

#### 286 3.2. Video Representation

287 We design HORNet to identify the most informative frames  
 288 in a video while suppressing redundant or noisy content.  
 289 This process is guided by learned video representations  
 290 rather than raw pixels, enabling the model to focus on se-  
 291 mantically meaningful cues that correlate with downstream  
 292 performance. To balance representational strength with

293 computational efficiency, HORNet employs a lightweight  
294 encoder  $E$  derived from the TimeSFormer [5] architec-  
295 ture. The encoder decouples spatial and temporal reason-  
296 ing into two separate transformer blocks to efficiently model  
297 video structure. In spatial blocks, we perform spatial self-  
298 attention independently on each of the frames to capture  
299 intra-frame relationships such as object appearance, local  
300 motion cues, and spatial layout. After spatial encoding, a  
301 second transformer stack performs temporal self-attention  
302 to capture motion patterns and temporal dependencies at  
303 each patch position. This factorized design (see Figure 2)  
304 preserves temporal modeling capacity while avoiding the  
305 prohibitive cost of joint attention over all tokens.

### 306 3.3. HORNet Architecture

307 HORNet instantiates the SAF policy using three compo-  
308 nents: a video encoder, a lightweight trainable policy net-  
309 work, and a frozen VLM answerer (Fig. 1).

310 **Video encoder.** Given a video  $\mathbf{V}$  with  $T$  frames, we ex-  
311 tract frame-level features using aforementioned encoder and  
312 obtain spatial token maps of shape  $T \times P \times P \times D$ , where  
313  $P = 16$  denotes the spatial grid resolution and  $D = 768$   
314 the feature dimension. To obtain compact per-frame repre-  
315 sentations, we apply spatial average pooling over the  $P \times P$   
316 grid:

$$317 \mathbf{F} = \text{AvgPool}_{2D}(E(\mathbf{V})) \in \mathbb{R}^{T \times D}, \quad (4)$$

318 where  $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_T]^\top$  and each  $\mathbf{f}_t \in \mathbb{R}^D$  corresponds to  
319 frame  $v_t$ .

320 **Policy network.** The SAF policy is parameterized as a  
321 frame-wise multilayer perceptron (MLP) that maps each  
322 feature vector  $\mathbf{f}_t$  to a selection probability  $p_t \in (0, 1)$ .  
323 Specifically, the network applies three linear projections  
324 with Gaussian Error Linear Unit (GELU) nonlinearities fol-  
325 lowed by a sigmoid activation:

$$326 p_t = \sigma(\mathbf{W}_2 \phi(\mathbf{W}_1 \phi(\mathbf{W}_0 \mathbf{f}_t))), \quad (5)$$

327 where  $\phi(\cdot)$  denotes GELU,  $\sigma(\cdot)$  the sigmoid function, and  
328 the weight matrices project  $D \rightarrow 512 \rightarrow 256 \rightarrow 1$ . Col-  
329 lectively, these weights define the learnable parameter set  
330  $\theta$ . The resulting probabilities  $\mathbf{p} = (p_1, \dots, p_T)$  define in-  
331 dependent Bernoulli decisions over frames, as described in  
332 the SAF formulation. This MLP constitutes the only train-  
333 able component of HORNet.

334 **Frozen VLM answerer.** For a sampled mask  $\mathbf{b}$ , the  
335 selected frames  $\mathbf{V}'$  are passed to a frozen Qwen3-VL  
336 model [3] together with the question  $q$ . The model produces  
337 a predicted answer  $\hat{a}$ , which is used to compute rewards dur-  
338 ing training and for evaluation at test time.

### 339 3.4. Training with GRPO

340 **Candidate generation.** For each training instance, we  
341 generate  $K = 8$  candidate masks  $\{\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(K)}\}$ . Candi-  
342 dates are produced using a deterministic top- $k$  sweep over  
343 sorted probabilities  $\mathbf{p}$ , progressively reducing the number  
344 of selected frames, together with one stochastic Bernoulli  
345 sample to maintain exploration.

346 **Reward computation.** Each candidate mask yields a pre-  
347 dicted answer  $\hat{a}^{(i)} = \mathcal{M}(\mathbf{V}'^{(i)}, q)$ . We define a smooth  
348 scalar reward

$$r^{(i)} = 0.1 \cdot F_1^{\text{token}}(\hat{a}^{(i)}, a) + 0.9 \cdot \text{EditSim}(\hat{a}^{(i)}, a), \quad (6) \quad 349$$

350 where  $F_1^{\text{token}}$  denotes token-level F1 after lemmatization  
351 and EditSim is normalized edit similarity in  $[0, 1]$ . This for-  
352 mulation reduces brittleness to minor lexical variations.

353 **GRPO objective.** The log-probability of candidate mask  
354  $\mathbf{b}^{(i)}$  under the current policy is

$$\log \pi_\theta(\mathbf{b}^{(i)} | \mathbf{F}) = \sum_{t=1}^T \left[ b_t^{(i)} \log p_t + (1 - b_t^{(i)}) \log(1 - p_t) \right]. \quad (7) \quad 355$$

356 Let  $\bar{r}$  and  $\sigma_r$  denote the mean and standard deviation of  
357 rewards within the group of  $K$  candidates. The normalized  
358 advantage is defined as

$$A^{(i)} = \frac{r^{(i)} - \bar{r}}{\sigma_r + \epsilon}, \quad (8) \quad 359$$

360 where  $\epsilon$  is a small constant for numerical stability. The  
361 GRPO loss is then

$$\mathcal{L}_{\text{GRPO}} = -\frac{1}{K} \sum_{i=1}^K A^{(i)} \log \pi_\theta(\mathbf{b}^{(i)} | \mathbf{F}). \quad (9) \quad 362$$

363 We optimize  $\theta$  using Adam with learning rate  $10^{-4}$ .

## 364 4. Results

365 In this section, we describe our training data and strategies,  
366 and present HORNet’s performance and efficiency gains  
367 over the baseline model through both qualitative and quanti-  
368 tative analyses. We conduct ablation studies to examine al-  
369 ternative design choices in VLM architectures, training pro-  
370 cedures, and sampling strategies. Overall, we demonstrate  
371 substantial efficiency improvements and highlight HOR-  
372 Net’s potential when scaled to larger backbone models.

#### 373 4.1. Training Data

374 HORNet is trained on a combined corpus spanning three  
375 VideoQA benchmarks: **MSRVTT-QA** [33] (10,000 videos,  
376 158,581 training QA pairs, mean 15.5s duration), **MSVD-**  
377 **QA** [34] (1,161 training videos, 30,933 QA pairs, mean  
378 9.6s), and **NExT-QA** [32] (3,870 training videos, 34,132  
379 QA pairs, mean 43.7s). In aggregate, the training set  
380 contains **223,646 QA pairs** across **15,031 videos** cover-  
381 ing 114.2 hours of content, with question types spanning  
382 descriptive (what/who), temporal (when/how), and causal  
383 (why) reasoning. This breadth ensures the selection policy  
384 generalizes across diverse temporal structures rather than  
385 overfitting to a single question distribution.

#### 386 4.2. Implementation Details

387 All experiments are conducted on a single **NVIDIA A100**  
388 **40GB** GPU. Videos are decoded and uniformly sampled to  
389  $T = 32$  frames, each resized to  $288 \times 288$  pixels. The  
390 **TimeFormer-Tiny** encoder produces spatial feature maps  
391 of shape  $16 \times 16 \times 768$ , which are spatially average-pooled  
392 to yield per-frame descriptors  $\mathbf{F} \in \mathbb{R}^{16 \times 768}$ .

393 The **MLP policy**  $\pi_\theta$  consists of a linear projection  
394 ( $768 \rightarrow 512$ ) followed by two hidden layers ( $512 \rightarrow$   
395  $1024 \rightarrow 256$ ) with GELU activations, and a final linear  
396 head ( $256 \rightarrow 1$ ) with sigmoid output. This amounts to  
397 fewer than 1M trainable parameters.

398 At each training step,  $K = 8$  candidate frame subsets  
399 are sampled per video via a top- $k$  sweep with step size  
400  $\lfloor k/K \rfloor$ . Training proceeds in two stages. In the first stage,  
401 we train on MSVD [34] and MSRVTT [33], which con-  
402 tain short videos (fewer than 100 frames) and one-word  
403 answers. Rewards are computed using an F1-Lev objec-  
404 tive: a weighted combination of token-level F1 ( $w_1=0.1$ )  
405 and normalized edit similarity ( $w_2=0.9$ ) applied to lemmat-  
406 ized predictions and ground-truth answers. In the second  
407 stage, we train on NExT-QA [32], which features MCQ-  
408 style questions and long videos (around 1,000 frames), us-  
409 ing a selection-accuracy reward tailored to the multiple-  
410 choice setting. The policy is optimized with **Adam** [18] at a  
411 learning rate of  $10^{-4}$  with batch size 8 on a total of 223,646  
412 training QA pairs. Qwen3-VL-2B stays fully frozen during  
413 training, while the video encoder and the frame-selection  
414 policy are trained jointly.

#### 415 4.3. Open-Ended QA Results

416 Table 2 reports results on three open-ended VideoQA  
417 benchmarks. On MSVD-QA, HORNet improves F1-Lev  
418 from 0.3483 to 0.3543 (+1.7%) while reducing Qwen pro-  
419 cessing time by 64% and input frames by 66%. This shows  
420 that for short videos ( $\sim 10$ s), many frames are redundant  
421 or noisy, and selecting a compact subset actually helps the  
422 VLM focus on relevant content.

On MSRVTT-QA and NExT-QA open-ended, HORNet  
trades a modest drop in F1 for substantial efficiency gains.  
MSRVTT-QA loses 5.6% F1 but reduces processing time  
by 84% and frames by 92%. NExT-QA open-ended loses  
10.1% but reduces processing time by 81% and frames by  
over 99%, compressing an average of 1,158 input frames  
down to 8. These results highlight a practical trade-off:  
HORNet enables deployment on resource-constrained set-  
tings where processing thousands of frames per video is in-  
feasible, with a bounded cost in answer quality.

#### 423 4.4. Multiple-Choice QA Results

Table 3 presents results on three MCQ benchmarks. The  
pattern mirrors the open-ended setting: HORNet consis-  
tently reduces processing time (74–93%) and frame count  
( $\geq 99\%$ ) across all datasets. On ActivityNet-QA, accuracy  
drops only 6.2% while inference becomes 93% faster. On  
NExT-QA MCQ, the gap narrows to 5.3% with 74% faster  
processing. VideoMME shows the largest accuracy gap  
(16.2%), which we attribute to its hour-scale videos where  
8 frames may be insufficient to cover the question scope.

Across both open-ended and MCQ settings, the results  
support a consistent finding: HORNet provides a con-  
trollable efficiency–accuracy trade-off, achieving order-of-  
magnitude reductions in computational cost with bounded  
quality loss. In certain cases, HORNet even improves  
the model’s predictions by discarding distracting or noisy  
frames and retaining only the most informative moments,  
producing a better answer than using VLMs alone, as illus-  
trated in Figure 3.

#### 424 4.5. Ablation Studies

**Training objective.** Table 4 compares three training  
strategies for the frame selection policy, all trained exclu-  
sively on MSVD-QA. On the in-distribution MSVD-QA  
evaluation, all three methods improve over the untrained  
baseline, with PPO achieving the highest F1 (0.3585) fol-  
lowed by GRPO (0.3543) and SFT (0.3495). However, the  
MSRVTT-QA column reveals a critical difference: since  
none of the methods were trained on MSRVTT-QA, this  
column measures out-of-distribution generalization. Here,  
all trained policies degrade relative to the untrained baseline  
(0.3209), but GRPO degrades the least (0.3029), retaining  
94% of baseline performance compared to 92% for PPO and  
90% for SFT. This suggests that GRPO’s group-relative ad-  
vantage estimation learns more transferable selection strate-  
gies, whereas PPO and SFT overfit more aggressively to the  
training distribution.

**Frame selection strategy.** Table 5 compares random,  
uniform, and HORNet selection, all restricted to exactly 4  
frames. On MSVD-QA and MSRVTT-QA, all three strate-  
gies perform within 0.01 F1 of each other. This is expected:

Table 2. Performance comparison across open-ended QA datasets. We adopt the aforementioned F1-Lev metric to measure model performance, which is a weighted combination of token-level F1 and normalized edit similarity on lemmatized texts. We also report efficiency measured in runtime and average frames passed to Qwen. Qwen’s baseline processing time includes uniform sampling, video encoding, and answer generation. In our setup, frame selection replaces the sampling step and is reported separately. Even under this accounting, the combined runtime of HORNet still yields a notable overall speedup. Additionally, when comparing generation speed, we follow Qwen’s default sampling rate (fps=2). Under the assumption of a 24-fps source video, the baseline effectively processes roughly 1/12 of all frames—still substantially more input than HORNet requires. We highlight best performance for each benchmark in **bold** and mark performance **gain** or **loss** as percentages.

Dataset	Model	F1-Lev $\uparrow$	Frame Sel. (s) $\downarrow$	Qwen Proc. (s) $\downarrow$	Avg. Frames $\downarrow$
MSVD [34]	Qwen3-VL-2B (Baseline)	0.3483	–	0.28	11.65
	HORNet+Qwen3-VL-2B (Ours)	<b>0.3543</b> +1.7%	0.12	<b>0.10</b> $\downarrow$ 64%	<b>4.00</b> $\downarrow$ 66%
MSRVTT [33]	Qwen3-VL-2B (Baseline)	<b>0.3209</b>	–	0.58	47.52
	HORNet+Qwen3-VL-2B (Ours)	0.3029 -5.6%	0.09	<b>0.09</b> $\downarrow$ 84%	<b>4.00</b> $\downarrow$ 92%
NextOE [32]	Qwen3-VL-2B (Baseline)	<b>0.3045</b>	–	1.01	1157.88
	HORNet+Qwen3-VL-2B (Ours)	0.2738 -10.1%	0.52	<b>0.19</b> $\downarrow$ 81%	<b>8.00</b> $\downarrow$ 99%

Table 3. Performance comparison across selection-based MCQ datasets. We adopt selection-accuracy as our metric and report efficiency gain of HORNet similar to Table 2. For each of these dataset we randomly sampled 1,000 QA pairs.

Dataset	Model	Accuracy(%) $\uparrow$	Frame Sel. (s) $\downarrow$	Qwen Proc. (s) $\downarrow$	Avg. Frames $\downarrow$
VideoMME [11]	Qwen3-VL-2B (Baseline)	<b>68.30</b>	–	2.53	3066.73
	HORNet+Qwen3-VL-2B (Ours)	52.10 -16.2%	1.51	<b>0.18</b> $\downarrow$ 93%	<b>8.00</b> $\downarrow$ 99%
ActivityNetQA [38]	Qwen3-VL-2B (Baseline)	<b>75.00</b>	–	2.37	3152.49
	HORNet+Qwen3-VL-2B (Ours)	68.80 -6.2%	1.64	<b>0.17</b> $\downarrow$ 93%	<b>8.00</b> $\downarrow$ 99%
NextQA [32]	Qwen3-VL-2B (Baseline)	<b>76.80</b>	–	0.98	1157.88
	HORNet+Qwen3-VL-2B (Ours)	71.50 -5.3%	0.53	<b>0.25</b> $\downarrow$ 74%	<b>8.00</b> $\downarrow$ 99%

Table 4. **Training objective ablation.** All variants use the same TimeSformer-Tiny encoder, MLP policy (<1M params), and frozen Qwen3-VL-2B answerer. Trained on MSVD-QA only. MSRVTT-QA results show out-of-distribution generalization.

Training	MSVD (F1-Lev $\uparrow$ )	MSRVTT (F1-Lev $\uparrow$ )
No training (baseline)	0.3483	<b>0.3209</b>
SFT (weighted BCE)	0.3495	0.2882
PPO (clipped surrogate)	<b>0.3585</b>	0.2948
<b>GRPO (Ours)</b>	0.3543	0.3029

Table 5. **Frame selection strategy ablation.** All methods select 4 frames and pass them to a frozen Qwen3-VL-2B answerer. MSVD-QA and MSRVTT-QA report F1-Lev; NEX-T-QA reports MCQ accuracy (%).

Strategy	MSVD (F1-Lev $\uparrow$ )	MSRVTT (F1-Lev $\uparrow$ )	NEX-T-QA (Acc. $\uparrow$ )
Random	0.3527	0.3027	65.88
Uniform	0.3493	<b>0.3058</b>	64.24
<b>HORNet</b>	<b>0.3543</b>	0.3029	<b>71.50</b>

Table 6. **VLM answerer ablation on MSVD-QA.** The same HORNet policy (trained with GRPO) selects frames. Only the frozen answerer is swapped; frame selection is identical.

VLM Answerer	Size	F1-Lev $\uparrow$
Qwen3-VL-Instruct (baseline)	2B	0.3483
Qwen3-VL-Instruct + HORNet	2B	0.3543
Qwen2.5-VL-Instruct + HORNet	3B	<b>0.3846</b>

frames in these videos carry similar visual content, and any 4-frame sample is likely to capture the relevant information. Notably, the fact that aggressive subsampling (4 out of 32 frames) does not substantially hurt performance reinforces our core premise: many frames are redundant or noisy, and discarding them does no harm.

The picture changes on NEX-T-QA, where videos average 44 seconds and questions require causal and temporal reasoning. Here, HORNet achieves 71.50% accuracy, outperforming random (65.88%) by 5.6 points and uniform (64.24%) by 7.3 points. When the temporal structure of the video matters, learned selection provides a clear advantage over blind sampling.

473 with average durations of 10s and 15s respectively, most

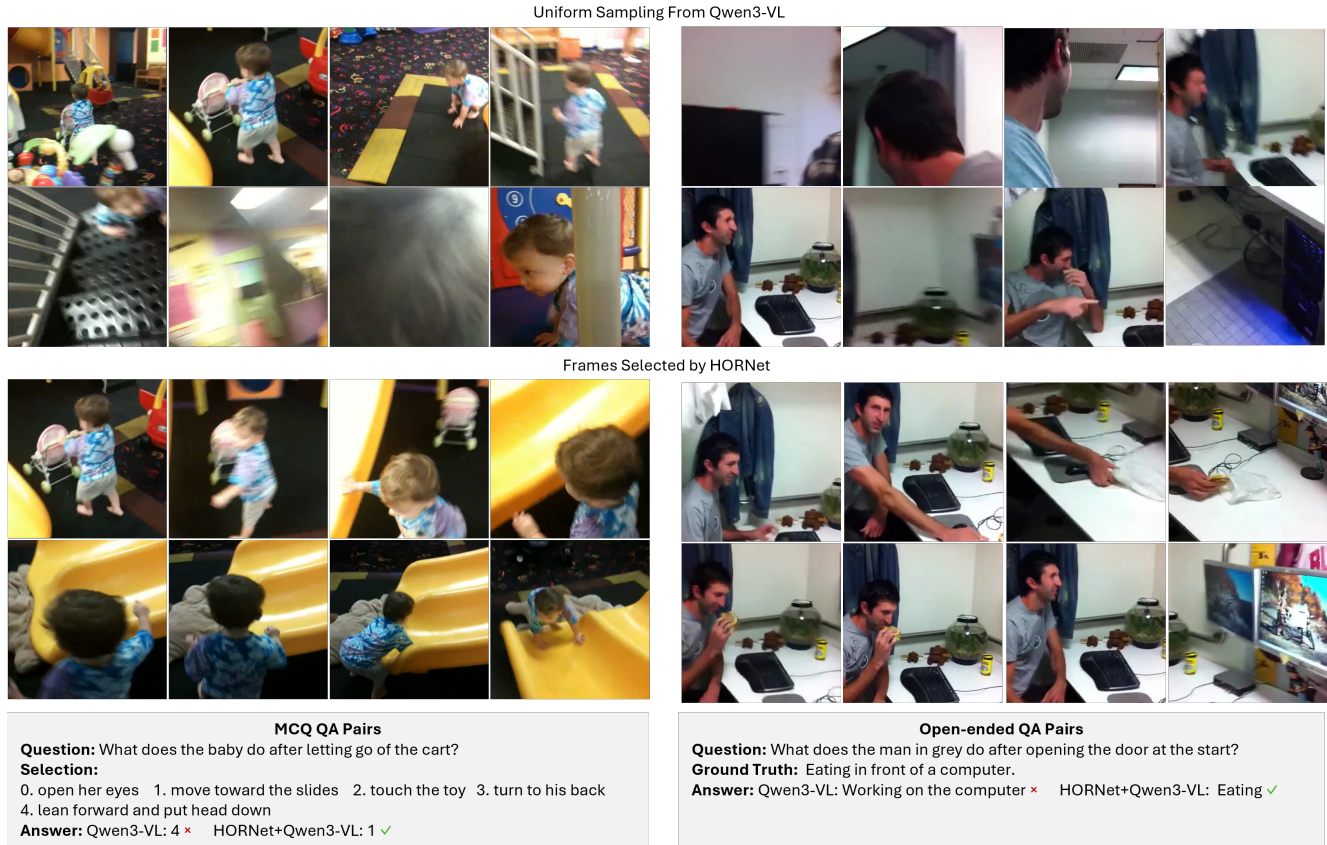


Figure 3. Qualitative example of HORNet’s frame-selection behavior on an MCQ and open-ended sample from the NExt-QA dataset [32]. Given fixed 8-frame input, uniform sampling in Qwen-VL captures frames of the child crawling instead of the slide following the discard of the cart (left), and a frame of a person working at a computer while missing the eating frames (right), leading the model to produce an incorrect answer. With a dense initial sampling ( $T=256$  frames), HORNet selects the full 8-frame sequence of action-relevant frames while discarding distractors, enabling the model to recover the correct prediction.

487 **VLM answerer.** Table 6 swaps only the frozen VLM answerer while keeping the same HORNet policy. Replacing  
488 Qwen3-VL-2B with the larger Qwen2.5-VL-3B improves  
489 F1-Lev from 0.3543 to 0.3846, a 8.5% relative gain. This  
490 confirms that HORNet’s frame selection transfers across  
491 VLM answerers without retraining, and that pairing the pol-  
492 icy with a stronger answerer amplifies the benefit of intelli-  
493 gent frame selection.  
494

## 495 5. Conclusion

496 We introduced HORNet, a lightweight frame selection pol-  
497 icy trained with GRPO that optimizes what a frozen VLM  
498 sees rather than what it generates, requiring fewer than 1M  
499 trainable parameters. Our experiments show that aggressive  
500 frame reduction (to as few as 4 frames) causes no mean-  
501 ingful quality loss on short-form videos, while on longer videos  
502 with temporal and causal questions, learned selection out-  
503 performs uniform and random baselines by up to 7.3 per-  
504 centage points. Across all benchmarks, HORNet reduces

VLM processing time by 64–93% and input frames by up to 99%. Ablation studies further confirm that GRPO generalizes better out-of-distribution than PPO and SFT, and that the learned policy transfers across VLM answerers without retraining, yielding an 8.5% relative gain when paired with a stronger model.

HORNet’s limitation is that the accuracy gap widens on hour-scale videos (e.g., VideoMME), where a fixed budget of 8 frames may be insufficient. Future work could address this through adaptive frame budgets that scale with video duration and hierarchical strategies that first localize relevant temporal segments before selecting frames within them. We also plan to incorporate visual reward signals that directly assess the perceptual quality and informativeness of selected frames, complementing the current text-based QA reward. Additionally, we aim to explore partially unfreezing the VLM answerer so that it can provide gradient-based feedback to the selection policy, enabling a tighter co-optimization loop between frame selection and answer generation.

525

**References**

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

- [1] authors. DeepVideo-R1: Video reinforcement fine-tuning via difficulty-aware regressive GRPO. *arXiv preprint arXiv:2506.07464*, 2025. 2, 4
- [2] authors. VideoBrain: Learning adaptive frame sampling for long video understanding. *arXiv preprint arXiv:2602.04094*, 2025. 3
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Tang, Zhenyu Wang, Peng Wang, et al. Qwen3-VL technical report. *arXiv preprint arXiv:2511.21631*, 2025. 5
- [4] Xiangyu Bai, Le Jiang, Yedi Luo, Aniket Gupta, Pushyami Kaveti, Hanumant Singh, and Sarah Ostadabbas. An evaluation platform to scope performance of synthetic environments in autonomous ground vehicles simulation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 1
- [5] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Icml*, page 4, 2021. 5
- [6] Shyamal Buch, Cristobal Eyzaguirre, Adrien Gaidon, Jiajun Wu, Li Fei-Fei, and Juan Carlos Niebles. Revisiting the “video” in video-language understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 3
- [7] Yi Chen et al. GRPO-CARE: Consistency-aware reinforcement learning for multimodal reasoning. *arXiv preprint arXiv:2506.16141*, 2025. 4
- [8] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 1
- [9] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 2, 4
- [10] Kaituo Feng, Kaixiong Li, Bohao Liu, Jiaming Li, Yueting Ge, Xiangyu Li, Lewei Lu, Kai Chen, and Xiangyu Wang. Video-R1: Reinforcing video reasoning in MLLMs. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2025. 2, 4
- [11] Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024. 7
- [12] Bishoy Galoaa, Xiangyu Bai, and Sarah Ostadabbas. Lang2motion: Bridging language and motion through joint embedding spaces. *arXiv preprint arXiv:2512.10617*, 2025. 1
- [13] Bishoy Galoaa, Xiangyu Bai, and Sarah Ostadabbas. Structured over scale: Learning spatial reasoning from educational video. *arXiv preprint arXiv:2601.23251*, 2026. 1
- [14] Haonan Ge et al. FrameMind: Frame-interleaved video reasoning via reinforcement learning. *arXiv preprint arXiv:2509.24008*, 2025. 3

- [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Liang Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *Iclr*, 1(2):3, 2022. 1
- [16] Wenxuan Huang et al. Vision-R1: Incentivizing reasoning capability in multimodal large language models. *arXiv preprint arXiv:2503.06749*, 2025. 4
- [17] Xiaofei Huang, Lingfei Luan, Elaheh Hatamimajoumerd, Michael Wan, Pooria Daneshvar Kakhaki, Rita Obeid, and Sarah Ostadabbas. Posture-based infant action recognition in the wild with very limited data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4912–4921, 2023. 1
- [18] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [19] Hosu Lee, Junho Kim, Hyunjun Kim, and Yong Man Ro. ReFoCUS: Reinforcement-guided frame optimization for contextual understanding. *arXiv preprint arXiv:2506.01274*, 2025. 3
- [20] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 1
- [21] Jianxin Liang et al. End-to-end video question answering with frame scoring mechanisms and adaptive sampling. *arXiv preprint arXiv:2407.15047*, 2024. 3
- [22] Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-LLaVA: Learning united visual representation by alignment before projection. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2024. 1
- [23] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 1
- [24] Shuming Liu et al. BOLT: Boost large vision-language model without training for long-form video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 3
- [25] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-ChatGPT: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, 2024. 1
- [26] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021. 1
- [27] Sarah Ostadabbas, Somaieh Amraee, Elaheh Hatamimajoumerd, and Michael Wan. Special issue 1251 editorial: computer vision with small data: a focus on human and animals transforming computer vision into equitable and impactful ai. *Multimedia Tools and Applications*, 84(21): 24515–24519, 2025. 1
- [28] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

- 639 of mathematical reasoning in open language models. *arXiv*  
640 *preprint arXiv:2402.03300*, 2024. 2, 3
- 641 [29] Liyang Song, Hardik Bishnoi, Sai Kumar Reddy Manne,  
642 Sarah Ostadabbas, Briana J Taylor, and Michael Wan. Over-  
643 coming small data limitations in video-based infant respi-  
644 ration estimation. In *Proceedings of the IEEE/CVF Win-  
645 ter Conference on Applications of Computer Vision*, pages  
646 6340–6349, 2026. 1
- 647 [30] Guangyu Sun, Archit Singhal, Burak Uz Kent, Mubarak  
648 Shah, Chen Chen, and Garin Kessler. From frames to clips:  
649 Training-free adaptive key clip selection for long-form video  
650 understanding. *arXiv preprint arXiv:2510.02262*, 2025. 1, 3
- 651 [31] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,  
652 Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin  
653 Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui  
654 Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Jun-  
655 yang Lin. Qwen2-vl: Enhancing vision-language model’s  
656 perception of the world at any resolution. *arXiv preprint*  
657 *arXiv:2409.12191*, 2024. 1
- 658 [32] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua.  
659 NEX-T-QA: Next phase of question-answering to explaining  
660 temporal actions. In *Proceedings of the IEEE/CVF Confer-  
661 ence on Computer Vision and Pattern Recognition (CVPR)*,  
662 2021. 2, 6, 7, 8
- 663 [33] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang,  
664 Xiangnan He, and Yueting Zhuang. Video question answer-  
665 ing via gradually refined attention over appearance and mo-  
666 tion. In *Proceedings of the ACM International Conference*  
667 *on Multimedia (ACM-MM)*, 2017. 2, 6, 7
- 668 [34] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang,  
669 Xiangnan He, and Yueting Zhuang. Video question answer-  
670 ing via gradually refined attention over appearance and mo-  
671 tion. In *Proceedings of the ACM International Conference*  
672 *on Multimedia (ACM-MM)*, 2017. 2, 6, 7
- 673 [35] Ziqiang Xu et al. ViaRL: Adaptive temporal grounding via  
674 visual iterated amplification reinforcement learning. *arXiv*  
675 *preprint arXiv:2505.15447*, 2025. 3
- 676 [36] Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal.  
677 Self-chained image-language model for video localization  
678 and question answering. In *Advances in Neural Information*  
679 *Processing Systems (NeurIPS)*, 2023. 3
- 680 [37] Sicheng Yu et al. Frame-voyager: Learning to query frames  
681 for video large language models. In *International Confer-  
682 ence on Learning Representations (ICLR)*, 2025. 3
- 683 [38] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting  
684 Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for  
685 understanding complex web videos via question answering.  
686 In *AAAI*, pages 9127–9134, 2019. 7
- 687 [39] Jingyi Zhang, Jiaxing Huang, Wenhao Jin, and Shijian Lu.  
688 R1-VL: Learning to reason with multimodal large language  
689 models via step-wise group relative policy optimization. In  
690 *Proceedings of the IEEE/CVF International Conference on*  
691 *Computer Vision (ICCV)*, 2025. 2, 4
- 692 [40] Shaojie Zhang et al. Q-Frame: Query-aware frame selec-  
693 tion and multi-resolution adaptation for video-LLMs. In  
694 *Proceedings of the IEEE/CVF International Conference on*  
695 *Computer Vision (ICCV)*, 2025. 3
- [41] Yuanhao Zou, Shengji Jin, Andong Deng, Youpeng Zhao, 696  
Jun Wang, and Chen Chen. A.I.R.: Enabling adaptive, itera- 697  
tive, and reasoning-based frame selection for video question 698  
answering. In *International Conference on Learning Repre- 699*  
*sentations (ICLR)*, 2026. 1, 3 700