

ResAdapt: Adaptive Resolution for Efficient Video Reasoning

Anonymous authors

Paper under double-blind review

Abstract

1 Multimodal Large Language Models (MLLMs) achieve stronger visual
 2 understanding by scaling input fidelity, but the resulting token explosion
 3 makes it prohibitive to jointly sustain high spatial resolution and long
 4 temporal context. Existing efficiency strategies only partially resolve this
 5 tension: *model-side* compression discards fine-grained evidence after encod-
 6 ing, while *output-side* agentic reasoning adds multi-pass latency. We argue
 7 that the true bottleneck lies in the initial volume of pixels processed, and in-
 8 troduce **ResAdapt**, an **input-side adaptation** framework that dynamically
 9 allocates visual budgets *before* encoding. ResAdapt couples a lightweight
 10 Allocator with a frozen MLLM backbone, formulated as a contextual bandit.
 11 To train the Allocator, we propose **Cost-Aware Policy Optimization**
 12 (**CAPO**) to translate sparse rollout feedback into a stable accuracy–cost
 13 signal, complemented by a temporal-similarity regularizer that prevents
 14 redundant allocations across visually similar frames. Extensive evaluations
 15 demonstrate that ResAdapt establishes a new Pareto frontier in low-budget
 16 video QA and reasoning-augmented temporal grounding. By reinvesting
 17 saved spatial compute into extended temporal coverage, ResAdapt sup-
 18 ports up to $16\times$ more frames at the same visual budget while delivering
 19 over 15% performance gains. Our results highlight that learning a sparse,
 20 content-dependent, single-pass allocation policy is a highly effective route
 21 to long-context multimodal reasoning under strict efficiency constraints.
 22 Code is available at <https://anonymous.4open.science/r/ResAdapt>.

23 1 Introduction

24 Multimodal Large Language Models (MLLMs) achieve stronger visual understanding by
 25 scaling input fidelity, yet the resulting visual-token growth makes jointly sustaining high
 26 spatial resolution and long temporal context prohibitive (Guo et al., 2025a; Bai et al., 2025a;
 27 Liu et al., 2025a; Shu et al., 2025; Shao et al., 2025b). In practice, this trade-off is central to
 28 video reasoning: reducing resolution risks losing the small visual cues that determine the
 29 answer, whereas shortening the clip removes the temporal context needed for long-horizon
 30 inference. Even architecturally efficient encoders (Zhang et al., 2026; Liu et al., 2025b) do
 31 not remove this tension; they merely shift where it becomes painful.

32 Existing efficiency methods typically fall into two paradigms (Figure 1a), both of which inter-
 33 vene late in the pipeline. *Model-side* approaches prune or merge tokens post-encoding (Khaki
 34 et al., 2025; Xu et al., 2025; Bolya et al., 2022; Tao et al., 2025); however, discarded fine-grained
 35 evidence cannot be recovered, and irregular token layouts disrupt optimized attention ker-
 36 nels (Dao, 2024; Kwon et al., 2023; Zheng et al., 2024). Conversely, *output-side* agentic
 37 methods rely on iterative retrieval or zooming (Zhang et al., 2025b; Yang et al., 2025d; Shen
 38 et al., 2025b; Zheng et al., 2025b). While improving coverage, they introduce multi-turn
 39 latency and risk missing decisive cues if the initial coarse view is over-compressed.

40 To overcome these limitations, we shift focus from post-encoding compression to pre-
 41 encoding pixel allocation. We propose **Input-side adaptation**, which reallocates the visual
 42 budget *before* encoding. Our method, **ResAdapt**, employs a lightweight allocator to predict
 43 per-frame visual budgets conditioned on coarse visual features and the text query. This
 44 budget is then materialized via operations like dynamic resolution resizing or frame se-
 45 lection. Consequently, the backbone processes a standard token sequence in a single pass,

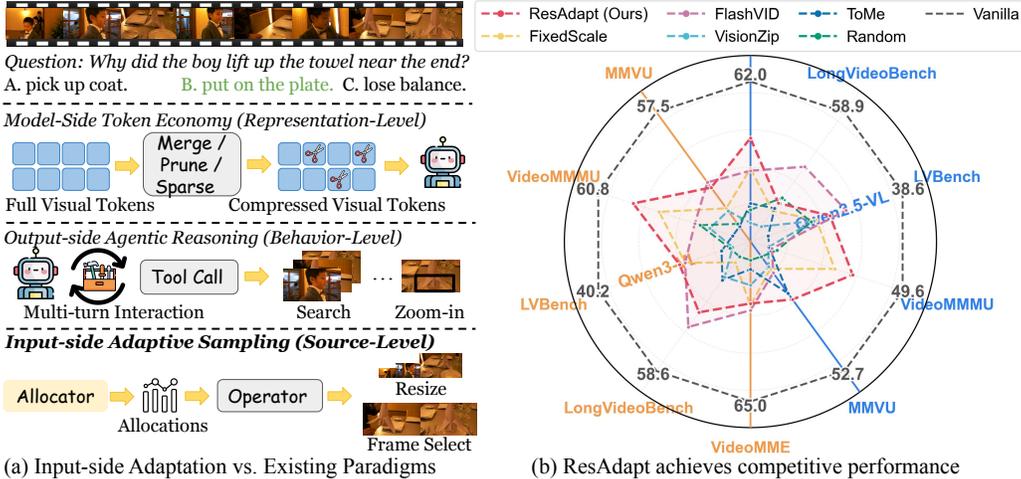


Figure 1: Input-side Adaptation improves the visual-token efficiency frontier. (a) Comparison of efficiency paradigms. While model-side methods compress tokens post-encoding and output-side agents rely on iterative retrieval, ResAdapt dynamically allocates per-frame visual budgets *before* encoding. (b) Performance of Qwen2.5-VL-7B using 32 frames at $\sim 10\%$ visual token retention. ResAdapt establishes Pareto frontier, yielding the most significant gains on reasoning benchmarks.

seamlessly integrating with modern inference engines (Dao, 2024; Kwon et al., 2023). Unlike prior slow-fast pipelines (Yang et al., 2025a; Zhang et al., 2026) that rely on query-agnostic heuristics, ResAdapt learns a query-aware allocation policy directly from task rewards.

Optimizing this allocator presents significant challenges: the allocation space is continuous, the budget operator is non-differentiable, and naive accuracy-cost trade-offs often collapse into degenerate, low-budget policies. To address these issues, we introduce **Cost-Aware Policy Optimization (CAPO)**, which transforms sparse rollout feedback into a stable, asymmetric learning signal. Coupled with a temporal-similarity regularizer that penalizes redundant high-resolution allocations across adjacent frames, ResAdapt emerges as a trainable, content-aware policy rather than a handcrafted heuristic.

Empirically, ResAdapt establishes a new efficiency-accuracy Pareto frontier for video reasoning. As shown in Figure 1b, at an aggressive 90% token reduction, ResAdapt consistently outperforms existing token economy methods across comprehensive benchmarks. By eliminating spatial redundancy, the saved compute can be seamlessly reinvested to expand temporal coverage—processing $16\times$ more frames under the same budget to unlock substantial performance gains. Furthermore, the learned policy exhibits *active perception*, autonomously allocating high resolutions to information-dense frames in a single forward pass without requiring explicit saliency supervision. Our main contributions are as follows:

1. We introduce **ResAdapt**, an *input-side adaptation* framework that formulates dynamic per-frame visual budgeting as a contextual bandit problem, fully preserving the native architecture and hardware optimizations of MLLMs.
2. We propose **CAPO** with a temporal similarity regularizer, providing a stable, asymmetric learning signal to jointly optimize accuracy and cost without hand-crafted heuristics.
3. Through extensive experiments and ablations, we show that ResAdapt achieves better efficiency-accuracy Pareto frontier across video QA and temporal grounding tasks.

2 Background and Problem Formulation

2.1 Preliminaries

Given a text query q and a video $\mathcal{V} = \{f_t\}_{t=1}^T$, a backbone policy π_ϕ encodes frames at fixed fidelity and autoregressively generates a rollout $\mathbf{y} = (y_1, \dots, y_L)$:

$$\pi_\phi(\mathbf{y} | q, \mathcal{V}) = \prod_{j=1}^L \pi_\phi(y_j | y_{<j}, q, \mathcal{V}). \quad (1)$$

75 Visual cost scales with total pixel volume, whereas answer-critical evidence is temporally
76 sparse, causing inefficiency.

77 To control pre-encoding cost, we introduce an Allocator policy π_θ that emits a per-frame
78 allocation vector $\mathbf{s} = (s_1, \dots, s_T)$ conditioned on the input $\mathbf{x} = (\mathbf{q}, \mathcal{V})$:

$$\mathbf{s} \sim \pi_\theta(\cdot | \mathbf{x}), \quad s_t \in [s_{\min}, s_{\max}], \quad (2)$$

79 and applies a *visual budget operator* \mathcal{O} to each frame: $\tilde{f}_t = \mathcal{O}(f_t, s_t)$. The backbone then
80 generates from the transformed input $\tilde{\mathbf{x}} = (\mathbf{q}, \{\tilde{f}_t\}_{t=1}^T)$:

$$\pi_\phi(\mathbf{y} | \tilde{\mathbf{x}}) = \prod_{j=1}^L \pi_\phi(y_j | y_{<j}, \tilde{\mathbf{x}}). \quad (3)$$

81 The framework is operator-agnostic, supporting resizing, frame selection, or other pre-
82 encoding budget controls.

83 2.2 Problem Formulation

84 Because the Allocator acts once before decoding, the outer problem is a *Contextual Bandit*
85 with continuous allocation vector $\mathbf{s} \in [s_{\min}, s_{\max}]^T$. For joint training, the two-stage policy
86 factorizes as $p_{\theta, \phi}(\mathbf{s}, \mathbf{y} | \mathbf{x}) = \pi_\theta(\mathbf{s} | \mathbf{x}) \pi_\phi(\mathbf{y} | \tilde{\mathbf{x}})$.

87 With $C(\mathbf{s})$ denoting visual cost and $Q(\mathbf{x}, \mathbf{y})$ the response quality, Lagrangian relaxation of a
88 budgeted objective yields the unconstrained utility:

$$\max_{\theta, \phi} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{s} \sim \pi_\theta, \mathbf{y} \sim \pi_\phi} [Q(\mathbf{x}, \mathbf{y}) - \lambda C(\mathbf{s})], \quad (4)$$

89 for trade-off coefficient $\lambda \geq 0$. Section 3 instantiates this objective with an Input-side adap-
90 tation policy, CAPO, and temporal regularization. Detailed derivations are in Appendix D.

91 3 Method

92 Figure 2 summarizes the Input-side adaptation framework. At inference, the Allocator
93 predicts one allocation per frame and applies a pre-encoding operator before the video
94 reaches the backbone in a single pass. In the experimental instantiation studied here, \mathcal{O} is
95 bilinear resizing, so the allocation becomes a resize factor s_t and $\tilde{f}_t = \mathcal{R}(f_t, s_t)$. At training,
96 rollout feedback updates the Allocator and, optionally, the backbone.

97 3.1 Joint RL Optimization Framework

98 The objective in Eq. (4) motivates a one-step expected-reward formulation. For a fixed
99 context \mathbf{x} , the ideal rollout reward is $R_{\mathbf{s}, \mathbf{y}}^{\text{ideal}} = Q(\mathbf{x}, \mathbf{y}) - \lambda C(\mathbf{s})$, yielding the target:

$$\max_{\theta, \phi} \mathcal{J}(\theta, \phi) = \mathbb{E}_{\pi_\theta(\mathbf{s} | \mathbf{x})} \left[\mathbb{E}_{\pi_\phi(\mathbf{y} | \tilde{\mathbf{x}})} \left[R_{\mathbf{s}, \mathbf{y}}^{\text{ideal}} \right] \right]. \quad (5)$$

100 Direct optimization of Eq. (5) via score-function estimators (detailed in Appendix D) is brittle
101 due to three challenges: 1) **Policy parameterization**: π_θ must efficiently predict continuous
102 allocations; 2) **Credit assignment**: the raw reward $R_{\mathbf{s}, \mathbf{y}}^{\text{ideal}}$ is high-variance and often causes
103 budget collapse; 3) **Temporal structure**: rollout-level rewards lack explicit signals to prevent
104 redundant allocations across adjacent frames. We address these sequentially below.

105 3.2 Allocator Architecture

106 To resolve the parameterization bottleneck, the Allocator processes coarse visual features
107 with negligible overhead (see Appendix C). Each frame f_t is encoded by a frozen lightweight
108 visual encoder, while the query is encoded separately. A shallow decoder alternates temporal
109 self-attention with gated cross-attention to the query, producing hidden states $\{h_t\}_{t=1}^T$.

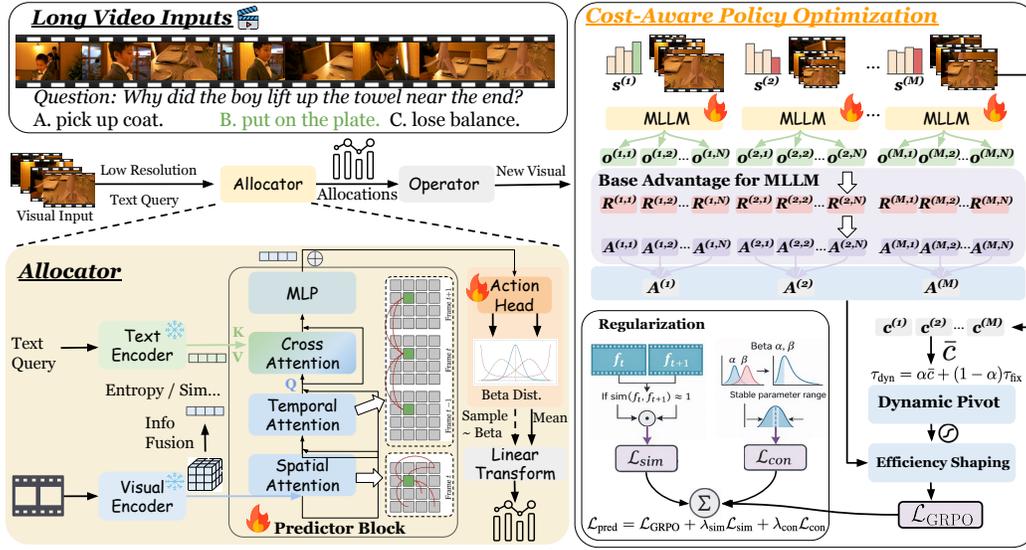


Figure 2: ResAdapt framework. (a) At inference, the Allocator π_θ maps coarse visual features and the query to latent actions $a_t \sim \text{Beta}(\alpha_t, \beta_t)$, which parameterize per-frame input allocations. In the resize instantiation used in our experiments, these allocations are realized as scales $s_t \in [s_{\min}, s_{\max}]$, and the resized frames are processed by the MLLM in a single call. (b) During training, CAPO reshapes group-relative advantages with a dynamic cost pivot τ_{dyn} , while temporal-similarity regularization suppresses redundant high-budget allocation on adjacent similar frames.

110 To preserve exploration, we parameterize latent actions with Beta distributions, mapping
 111 their bounded support to $[s_{\min}, s_{\max}]$:

$$a_t \sim \text{Beta}(\alpha_t, \beta_t), \quad s_t = s_{\min} + a_t (s_{\max} - s_{\min}). \quad (6)$$

112 Conditioned on $\{h_t\}$, the latent policy factorizes as $\log q_\theta(\mathbf{a} | \mathbf{x}) = \sum_{t=1}^T \log \text{Beta}(a_t; \alpha_t, \beta_t)$,
 113 inducing the continuous allocation policy $\pi_\theta(\mathbf{s} | \mathbf{x})$.

114 3.3 Cost-Aware Policy Optimization (CAPO)

115 A flat penalty on visual cost $C(\mathbf{s})$ often collapses the policy toward uniformly tiny budgets.
 116 To solve the credit assignment bottleneck, CAPO replaces the raw penalty with a shaped
 117 surrogate learning signal.

118 **Compute metric.** For the resize operator, frame f_t resized by s_t produces $n_t(s_t) \propto$
 119 $[s_t H_t / P][s_t W_t / P]$ tokens. Physical compute is measured by the *token retention ratio*
 120 $\rho(\mathbf{s}) = \sum n_t(s_t) / \sum n_t(1)$, well-approximated by the average quadratic scale. To reduce vari-
 121 ance during optimization, we use the smoother linear proxy $c_m = (\bar{s} - s_{\min}) / (s_{\max} - s_{\min})$
 122 for allocation m .

123 **Notation bridge.** During training, $R_{m,n}^{\text{task}}$ denotes the task reward of rollout (m, n) , $A_{m,n}^{\text{base}}$
 124 the GRPO group-normalized advantage, and $A_{m,n}$ the final CAPO-shaped advantage. Ap-
 125 pendix D details their relation to the ideal reward. We also define $u_{m,n} \in \{0, 1\}$ as a binary
 126 correctness indicator (e.g., exact-match for QA, thresholded score for generation).

127 **Dynamic cost pivot.** A fixed threshold does not track the evolving policy, whereas a purely
 128 group-dependent statistic is too noisy. CAPO therefore interpolates between a fixed target
 129 τ_{fix} and the prompt-local mean $\bar{c}_{\text{group}} = \frac{1}{M} \sum_{m=1}^M c_m$:

$$\tau_{\text{dyn}} = \kappa_{\text{mix}} \bar{c}_{\text{group}} + (1 - \kappa_{\text{mix}}) \tau_{\text{fix}}, \quad (7)$$

130 where $\kappa_{\text{mix}} \in [0, 1]$ controls adaptivity.

131 **Asymmetric shaping.** With τ_{dyn} as pivot, we apply a correctness-dependent bonus or
 132 penalty with $\lambda_- > \lambda_+ > 0$:

$$S_{m,n} = \begin{cases} \lambda_+ \sigma\left(\frac{\tau_{\text{dyn}} - c_m}{\tau_s}\right) & \text{if } u_{m,n} = 1, \\ -\lambda_- \sigma\left(\frac{c_m - \tau_{\text{dyn}}}{\tau_s}\right) & \text{if } u_{m,n} = 0, \end{cases}, \quad (8)$$

133 Efficient correct rollouts receive a moderate bonus, whereas costly incorrect rollouts receive
 134 a stronger penalty. The sigmoid temperature τ_s smooths the transition around the pivot.

135 **Final CAPO advantage.** Let $\tilde{A}_{m,n} = A_{m,n}^{\text{base}} + \lambda_{\text{capo}} S_{m,n} - \gamma c_m$. The final advantage is

$$A_{m,n} = \begin{cases} \max(\tilde{A}_{m,n}, \varepsilon_+) & \text{if } u_{m,n} = 1, \\ \tilde{A}_{m,n} & \text{if } u_{m,n} = 0, \end{cases} \quad (9)$$

136 where $\lambda_{\text{capo}} > 0$ scales CAPO shaping, $\gamma \geq 0$ applies a residual global cost penalty, and
 137 the floor $\varepsilon_+ > 0$ ensures that correct low-cost rollouts retain a positive learning signal. The
 138 dominant anti-collapse term is the pivoted asymmetric shaping in $S_{m,n}$.

139 3.4 Regularization and Training Objective

140 While CAPO balances global task reward and efficiency, it acts strictly at the rollout level.
 141 As observed empirically (Section 4.4), this can lead to the temporal structure bottleneck:
 142 the policy may allocate uniformly low-variance scales across highly redundant frames. To
 143 enforce localized selectivity, we augment the optimizer with structural regularization.

144 **Temporal similarity loss (\mathcal{L}_{sim}).** CAPO optimizes global quality–efficiency but does not
 145 penalize redundant allocations on near-duplicate frames. Reusing the coarse features f_t
 146 from Sec. 3.2, we penalize similar adjacent pairs:

$$\mathcal{L}_{\text{sim}} = \frac{1}{T-1} \sum_{t=1}^{T-1} w_t \cdot \max(0, \log s_t + \log s_{t+1} + \eta_{\text{sim}}), \quad (10)$$

147 where $w_t = \sigma((\cos(f_t, f_{t+1}) - \tau_{\text{sim}}) / \gamma_{\text{sim}})$ activates when frames exceed a similarity thresh-
 148 old τ_{sim} . No penalty is incurred when $s_t s_{t+1} \leq e^{-\eta_{\text{sim}}}$.

149 **Concentration loss (\mathcal{L}_{con}).** To prevent Beta distributions from collapsing to near-
 150 deterministic spikes, we softly cap total concentration at $\kappa_{\text{max}} > 0$: $\mathcal{L}_{\text{con}} =$
 151 $\frac{1}{T} \sum_{t=1}^T \max(0, \alpha_t + \beta_t - \kappa_{\text{max}})$. Together, \mathcal{L}_{sim} encourages differentiated allocation, while
 152 \mathcal{L}_{con} preserves exploration.

153 **Practical training objective.** We optimize both policies in a single GRPO-style loop (Zheng
 154 et al., 2025a; Yu et al., 2025). For each prompt \mathbf{x} , the Allocator samples M allocation trajec-
 155 tories $s_{1:M}$; each transformed input $\tilde{\mathbf{x}}^{(m)}$ produces N response rollouts. CAPO computes
 156 rollout advantages $A_{m,n}$, serving as the shared learning signal.

157 **Allocator objective.** We aggregate rollout advantages per allocation, $A_m^{\text{CAPO}} = \frac{1}{N} \sum_n A_{m,n}$,
 158 and optimize the per-frame clipped surrogate objective:

$$\mathcal{L}_\theta = -\frac{1}{MT} \sum_{m=1}^M \sum_{t=1}^T \min\left(r_{\theta,t}^{(m)} A_m^{\text{CAPO}}, \text{clip}\left(r_{\theta,t}^{(m)}, 1-\varepsilon, 1+\varepsilon\right) A_m^{\text{CAPO}}\right), \quad (11)$$

159 where $r_{\theta,t}^{(m)} = q_\theta(a_t^{(m)} | \mathbf{x}) / q_{\theta_{\text{old}}}(a_t^{(m)} | \mathbf{x})$. The full loss is $\mathcal{L}_{\text{alloc}} = \mathcal{L}_\theta + \lambda_{\text{sim}} \mathcal{L}_{\text{sim}} + \lambda_{\text{con}} \mathcal{L}_{\text{con}}$.

160 **Backbone update.** Conditioned on the sampled allocations, the backbone is optionally up-
 161 dated (for ResAdapt-RL) with the token-level clipped surrogate objective using advantages
 162 $A_{m,n}$ and ratios $r_{\phi,j}^{(m,n)}$ (see Appendix D). In practice, $\mathcal{L}_{\text{alloc}}$ and \mathcal{L}_ϕ are optimized alternately.

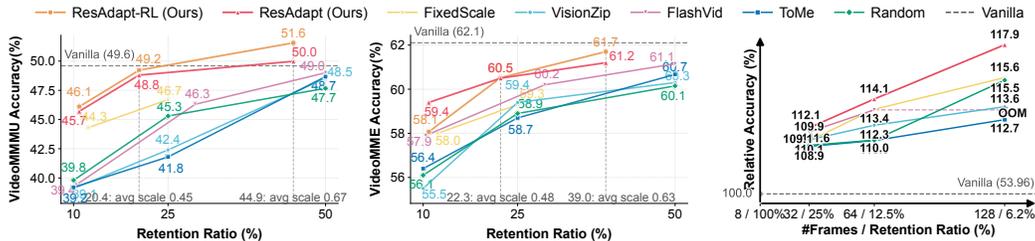


Figure 3: Efficiency-accuracy trade-offs and temporal reallocation. (a,b) VideoMMM U and VideoMME versus visual-token retention ratio R . ResAdapt is on or near the Pareto frontier, with the clearest advantage on reasoning-heavy settings at low retention. (c) Relative gain from trading spatial resolution for temporal coverage under a fixed 8-frame-equivalent budget.

163 **4 Experiments**

164 **4.1 Setup**

165 **Implementation.** The Allocator π_θ adopts the lightweight SmolVLM architecture (Marafioti
 166 et al., 2025) for high-throughput prediction. We instantiate input-side allocation via contin-
 167 uous per-frame resizing, providing the smooth action space required by our optimizer, with
 168 frame dropping naturally emerging at the zero-budget limit. The Allocator is trained on
 169 Qwen2.5-VL-7B-Instruct (Bai et al., 2025b) and additionally evaluated for zero-shot transfer
 170 on Qwen3-VL-8B-Instruct (Bai et al., 2025a). We report two paradigms: **ResAdapt-RL**,
 171 which jointly updates the Allocator and the backbone, and **ResAdapt**, which applies the
 172 trained Allocator to a frozen backbone to assess plug-and-play generalization.

173 **Baselines.** We compare against **heuristic methods** (Random Drop, FixedScale), **model-**
 174 **side compression** (ToMe (Bolya et al., 2022), FlashVid (Fan et al., 2026), VisionZip (Yang
 175 et al., 2025c)), and **inference augmentation** (VideoAuto-R1 (Liu et al., 2026)). Efficiency
 176 is measured by the visual-token retention ratio R . For inference-augmented models, R
 177 accounts strictly for visual encoder tokens. Since several baselines only support discrete
 178 operating points, we align budgets approximately and emphasize Pareto frontier trade-offs.

179 **Benchmarks.** Evaluations span *video QA* (VideoMME (Fu et al., 2025a), MMVU (Zhao et al.,
 180 2025b), LongVideoBench (Wu et al., 2024), MLVU (Zhou et al., 2025), VideoMMM U (Hu
 181 et al., 2025), LVBench (Wang et al., 2025b)), *temporal grounding* (Charades-STA (Gao et al.,
 182 2017), ActivityNet (Fabian et al., 2015), NExT-GQA (Xiao et al., 2024)), and *image under-*
 183 *standing* (MathVista (Lu et al., 2023), MMMU (Yue et al., 2024), OCRBench (Liu et al., 2024),
 184 ChartQA (Masry et al., 2022), AI2D (Kembhavi et al., 2016), TextVQA (Singh et al., 2019)).
 185 Unless otherwise stated, analyses default to Qwen2.5-VL-7B with 32 frames. All evalu-
 186 ations utilize lmms-eval (Zhang et al., 2024a). Full hyperparameters, prompts, reward
 187 formulations, and detailed configurations are deferred to Appendix B.

188 **4.2 Main Results**

189 We evaluate ResAdapt across video QA, temporal grounding, and zero-shot image transfer
 190 (Table 1). Our analysis focuses on the extreme low-budget regime ($\sim 10\%$ token retention),
 191 where standard compression typically collapses, demonstrating the robustness of input-side
 192 adaptation.

193 **Video QA: Preserving Reasoning under Extreme Compression.** Post-encoding token re-
 194 duction methods (e.g., ToMe, VisionZip) indiscriminately discard semantic features, causing
 195 precipitous drops in reasoning accuracy under high compression. In contrast, ResAdapt
 196 successfully isolates sparse but critical evidence. On the reasoning-heavy VideoMMM U
 197 benchmark, it maintains a substantial margin (up to +6.5 points) over baselines at matched
 198 budgets. This resilience highlights the core advantage of *pre-encoding* allocation: by dynami-

Table 1: Evaluation Results on Video QA Benchmarks. Retention ratio R reflects visual token count; Reasoning (\checkmark/\times) indicates chain-of-thought use; **bold** marks the best result. ResAdapt yields larger gains on the reasoning benchmark than on the perception benchmarks.

Backbone	Method	Retention Ratio R	Reasoning	Video Perception			Video Reasoning		Grounding (mIoU)		NExT-GQA			
				VideoMME	LongVideoBench	MMVU	MLVU	VideoMMU	LVBench	Charades	ActivityNet	Acc	mIoU	
32 Frames														
Owens2.5V-L7B	Vanilla	100%	\times	62.0	58.9	52.7	63.1	49.6	38.6	47.3	22.6	78.9	28.0	
	Random Drop	25.0%	\times	58.9	57.8	49.6	58.3	45.3	36.7	25.7	11.7	77.5	16.6	
	ToMe (Bolya et al., 2022)	25.0%	\times	58.7	58.0	51.0	58.7	41.8	37.7	26.0	12.1	77.8	16.3	
	VisionZip (Yang et al., 2025c)	25.0%	\times	59.4	57.1	49.8	57.9	42.4	36.5	-	-	-	-	
	FlashVid (Fan et al., 2026)	29.3%	\times	60.2	58.6	51.1	59.2	46.3	36.9	26.6	12.0	78.1	16.5	
	FixedScale	25.0%	\times	60.0	56.8	51.2	59.8	46.7	37.3	24.9	14.1	77.7	12.3	
	ResAdapt (Ours)	23.8%	\times	60.3	58.2	51.9	60.1	48.8	37.9	35.6	15.3	76.6	23.2	
	Random Drop	10.0%	\times	56.1	55.6	47.1	56.5	39.8	35.2	24.6	11.1	76.3	15.4	
	ToMe (Bolya et al., 2022)	10.0%	\times	56.4	55.2	48.9	58.0	39.2	33.6	27.4	12.2	77.3	15.7	
	VisionZip (Yang et al., 2025c)	10.0%	\times	55.5	54.5	47.6	57.3	39.1	35.3	-	-	-	-	
	FlashVid (Fan et al., 2026)	10.4%	\times	57.9	56.8	47.9	57.7	39.4	36.5	25.1	11.8	77.4	16.1	
	FixedScale	12.3%	\times	58.0	55.1	47.7	57.5	44.3	35.4	32.0	13.3	77.1	13.7	
	ResAdapt (Ours)	11.4%	\times	59.4	55.4	49.2	58.4	45.7	35.9	27.2	12.5	74.3	20.4	
	VideoAuto-RI (Liu et al., 2026)	100%	\checkmark	63.2	58.9	55.0	60.1	53.6	41.5	41.5	34.4	73.6	33.8	
	+ ResAdapt (Ours)	23.8%	\checkmark	60.4	57.1	53.2	61.1	51.2	38.7	49.1	44.7	79.3	35.3	
+ ResAdapt (Ours)	11.4%	\checkmark	59.3	56.3	51.8	59.3	49.1	36.7	30.0	24.4	74.7	24.7		
128 Frames														
Owens2.5V-L7B	Vanilla	100%	\times	65.3	60.3	53.1	66.5	47.9	42.0	52.8	34.4	79.8	29.9	
	Random Drop	25.0%	\times	64.9	61.2	50.8	64.8	48.1	41.3	20.7	18.8	80.3	10.7	
	ToMe (Bolya et al., 2022)	25.0%	\times	65.1	61.6	51.9	63.1	46.6	42.1	20.7	19.1	80.3	10.9	
	VisionZip (Yang et al., 2025c)	25.0%	\times	64.8	61.3	51.1	64.5	47.3	41.5	-	-	-	-	
	ResAdapt (Ours)	22.9%	\times	65.6	60.2	52.8	65.9	51.1	42.1	42.0	24.3	78.1	27.2	
	Random Drop	10.0%	\times	63.0	59.0	45.8	63.4	46.7	38.0	24.7	17.0	79.4	12.8	
	ToMe (Bolya et al., 2022)	10.0%	\times	60.6	56.3	44.2	63.5	41.8	39.5	17.9	16.4	79.1	11.1	
	VisionZip (Yang et al., 2025c)	10.0%	\times	61.8	56.3	44.8	63.2	42.1	38.2	-	-	-	-	
	FlashVid (Fan et al., 2026)	12.3%	\times	64.1	60.9	49.6	64.5	46.9	40.3	22.7	18.3	77.9	11.3	
	ResAdapt (Ours)	11.1%	\times	63.8	58.6	49.0	64.3	49.2	39.9	28.9	17.2	76.2	23.9	
	VideoAuto-RI (Liu et al., 2026)	100%	\checkmark	64.7	59.1	56.7	65.1	52.2	41.2	28.9	33.5	68.0	31.0	
	+ ResAdapt (Ours)	23.8%	\checkmark	66.2	60.2	53.5	66.0	52.6	41.8	49.1	44.7	79.3	35.3	
	+ ResAdapt (Ours)	11.4%	\checkmark	64.7	57.8	52.4	64.6	51.3	39.5	34.2	35.7	76.6	29.4	
	32 Frames													
	Owens-VL-8B	Vanilla	100%	\times	65.0	58.6	57.5	64.0	60.8	40.2	46.4	31.8	78.7	34.2
Random Drop		25.0%	\times	61.3	58.4	57.1	60.2	53.4	37.8	12.1	10.0	77.2	15.6	
ToMe (Bolya et al., 2022)		25.0%	\times	62.4	57.4	56.0	60.8	49.1	36.4	43.1	32.6	77.1	31.7	
VisionZip (Yang et al., 2025c)		25.0%	\times	61.8	57.2	54.4	60.6	51.5	37.3	-	-	-	-	
FlashVid (Fan et al., 2026)		30.0%	\times	63.9	59.0	54.8	61.9	55.1	38.5	47.7	36.8	77.8	33.9	
ResAdapt (Ours)		23.8%	\times	62.6	57.5	55.3	61.0	58.4	38.5	39.9	28.5	75.1	30.2	
Random Drop		10.0%	\times	58.8	54.7	53.2	56.6	47.1	35.5	4.4	5.0	74.3	11.3	
ToMe (Bolya et al., 2022)		10.0%	\times	59.2	55.5	53.1	58.5	42.7	35.8	41.8	34.1	79.2	34.0	
VisionZip (Yang et al., 2025c)		10.0%	\times	59.9	55.4	53.7	58.8	45.8	35.4	-	-	-	-	
FlashVid (Fan et al., 2026)		12.2%	\times	61.0	54.0	54.4	60.6	51.5	37.3	44.6	35.2	75.6	31.8	
FixedScale		12.3%	\times	60.8	54.9	53.8	58.4	52.6	37.1	37.9	28.4	74.2	29.9	
ResAdapt (Ours)		11.4%	\times	60.7	56.6	54.6	59.6	56.1	37.3	33.6	27.2	71.8	28.2	
128 Frames														
Owens-VL-8B		Vanilla	100%	\times	69.4	64.3	58.5	72.7	63.0	45.7	45.6	33.9	81.1	36.6
		Random Drop	25.0%	\times	67.2	61.3	56.8	67.4	53.3	27.4	26.3	29.3	79.3	22.4
	ToMe (Bolya et al., 2022)	25.0%	\times	67.2	62.0	55.9	70.4	53.5	43.1	-	-	-	-	
	VisionZip (Yang et al., 2025c)	25.0%	\times	67.1	61.3	55.7	69.2	56.8	41.2	-	-	-	-	
	ResAdapt (Ours)	22.9%	\times	67.4	61.9	56.3	70.8	59.6	43.3	39.8	30.0	76.8	33.3	
	Random Drop	10.0%	\times	64.1	58.3	55.4	62.4	55.5	36.8	21.9	24.8	76.9	19.9	
	ToMe (Bolya et al., 2022)	10.0%	\times	64.7	58.6	55.1	67.3	46.3	40.5	38.1	34.4	77.4	31.5	
	VisionZip (Yang et al., 2025c)	10.0%	\times	64.2	59.1	54.2	66.8	47.6	39.4	-	-	-	-	
	FixedScale	12.3%	\times	66.7	59.5	54.4	67.7	56.3	41.7	38.1	29.5	75.4	32.6	
	ResAdapt (Ours)	11.1%	\times	66.8	60.2	55.4	69.4	58.2	42.6	33.7	28.4	73.2	43.9	

199 cally assigning pixels based on content, ResAdapt protects high-frequency visual details
200 exactly where complex reasoning demands it.

201 *Trading spatial redundancy for temporal context.* The primary benefit of this efficiency is the
202 ability to reinvest saved spatial budgets into broader temporal coverage. By extending the
203 context window from 32 to 128 frames using the budget freed by ResAdapt, we surpass
204 the dense 128-frame baseline’s performance while remaining computationally cheaper
205 (Figure 3). This confirms a critical insight: reasoning models benefit significantly more from
206 a longer, adaptively compressed timeline than a shorter, fully dense one.

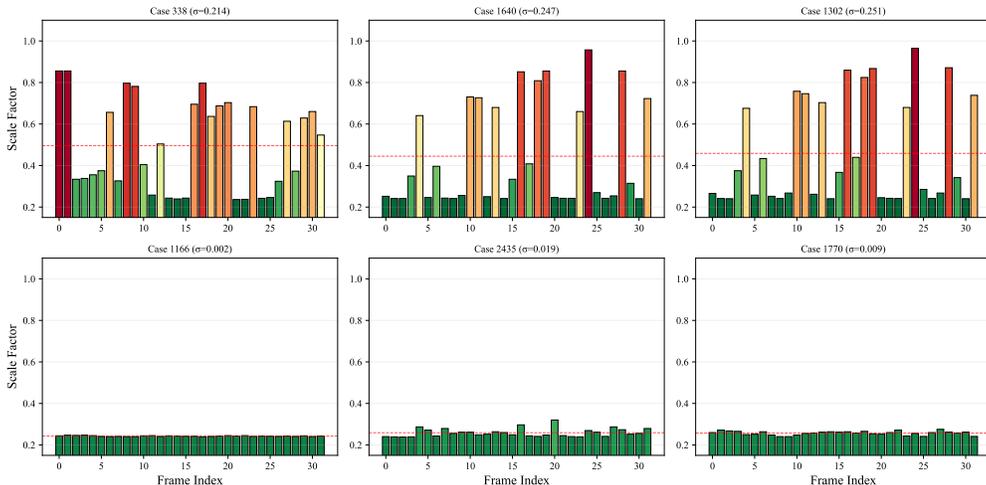
207 **Temporal Grounding: Anchoring and Denoising.** Grounding tasks are inherently
208 compression-sensitive as they demand precise spatio-temporal anchors. While naive frame
209 dropping destroys these anchors and severely degrades performance (Appendix E.6), ResAdapt
210 mitigates this by heavily compressing static segments while explicitly preserving
211 the full resolution of critical transition events.

212 *Synergy with long-context reasoning.* ResAdapt exhibits a powerful synergy with long-context
213 models. Scaling to 128 frames often causes dense models to struggle with distraction,
214 paradoxically degrading grounding accuracy. Integrating ResAdapt not only recovers
215 this loss but yields a net improvement over the short-context baseline. This reveals an
216 emergent *denoising* effect: by suppressing irrelevant visual noise, ResAdapt actively guides
217 the model’s attention toward valid temporal anchors, achieving higher accuracy with under
218 10% of the original token cost.

219 **Exploratory Image Transfer.** Zero-shot transfer to static images confirms our policy does not
220 overfit to video dynamics, though text-dense images naturally require a minimum spatial
221 resolution threshold (Appendix E.7, Table 7).

Table 2: Latency breakdown (ms, ↓) on Qwen2.5-VL-7B with single-GPU Allocator and 4-GPU vLLM engine. Averaged over 200 runs after 5 warm-up; E2E latency = Scale Time + Gen. Time.

Method	#Frames	Retention Ratio R	Scale						Inference			Total	
			TFLOPs	Text Enc.	Visual Enc.	Scale Pred.	Scale Apply	Scale Time	TFLOPs	TTFT	Gen. Time	TFLOPs	E2E Time
Vanilla	16	100%	—	—	—	—	—	—	111.4	378.9	527.9	111.4	527.9
ResAdapt	16	76.3%	1.5	19.8	94.1	85.6	6.3	205.8	77.2 (↓30.7%)	272.5 (↓28.1%)	370.7 (↓29.8%)	80.1 (↓28.1%)	576.5 (↓9.2%)
ResAdapt	16	52.8%	1.5	19.9	102.9	94.5	8.4	225.7	51.5 (↓53.8%)	261.5 (↓31.0%)	313.1 (↓40.7%)	54.4 (↓51.2%)	538.8 (↓12.1%)
ResAdapt	16	28.9%	1.5	20.4	103.4	92.2	9.0	225.0	31.0 (↓72.2%)	227.2 (↓40.0%)	237.9 (↓54.9%)	33.9 (↓69.6%)	462.9 (↓12.3%)
Vanilla	32	100%	—	—	—	—	—	—	222.5	723.3	881.9	222.5	881.9
ResAdapt	32	74.4%	2.9	19.9	204.1	97.4	14.4	335.9	153.9 (↓30.8%)	589.4 (↓18.5%)	627.6 (↓28.8%)	159.7 (↓28.2%)	963.5 (↓9.2%)
ResAdapt	32	51.5%	2.9	20.0	193.2	92.0	16.2	321.4	102.4 (↓54.0%)	505.0 (↓30.2%)	467.1 (↓47.0%)	108.2 (↓51.4%)	788.5 (↓10.6%)
ResAdapt	32	28.2%	2.9	20.3	190.4	90.3	17.3	318.3	61.4 (↓72.4%)	451.8 (↓37.5%)	332.6 (↓62.3%)	67.2 (↓69.8%)	650.9 (↓16.2%)
Vanilla	64	100%	—	—	—	—	—	—	444.6	1457.5	2059.6	444.6	2059.6
ResAdapt	64	73.2%	5.8	19.8	389.5	95.8	26.4	531.5	307.3 (↓30.9%)	1093.1 (↓25.0%)	1327.0 (↓35.6%)	318.9 (↓28.3%)	1858.5 (↓9.8%)
ResAdapt	64	50.7%	5.8	20.1	382.1	94.9	29.9	527.0	204.3 (↓54.0%)	991.8 (↓31.9%)	740.5 (↓64.0%)	215.9 (↓51.4%)	1267.5 (↓38.5%)
ResAdapt	64	27.8%	5.8	20.0	371.6	90.2	34.8	516.6	122.2 (↓72.5%)	899.2 (↓38.3%)	511.4 (↓75.2%)	133.8 (↓69.9%)	1028.0 (↓50.1%)
Vanilla	128	100%	—	—	—	—	—	—	888.9	2936.3	4877.0	888.9	4877.0
ResAdapt	128	74.2%	11.6	20.1	766.3	95.0	53.1	934.5	614.1 (↓30.9%)	2286.6 (↓22.1%)	2323.6 (↓52.4%)	637.3 (↓28.3%)	3258.1 (↓33.2%)
ResAdapt	128	51.4%	11.6	20.2	755.3	93.8	59.4	928.7	408.0 (↓54.1%)	2071.0 (↓29.5%)	1496.0 (↓69.3%)	431.2 (↓51.5%)	2424.7 (↓50.3%)
ResAdapt	128	28.2%	11.6	20.4	734.5	92.0	68.6	915.5	243.9 (↓72.6%)	1766.7 (↓39.8%)	1061.8 (↓78.2%)	267.1 (↓70.0%)	1977.3 (↓59.5%)

**Figure 4: Emergent active perception.** Per-frame scale s_t over frame index for VideoMME, grouped by intra-video scale diversity σ . High-diversity videos show localized scale spikes on scene changes or text overlays, while low-diversity videos remain uniformly compressed.

222 4.3 Runtime Overhead

223 We measure pipeline latency to determine when the Allocator’s upstream cost is amortized
 224 by downstream token savings. Table 2 reports latency using a dedicated single-GPU
 225 Allocator coupled with a 4-GPU vLLM engine. By adjusting the maximum allowed scale,
 226 ResAdapt acts as a continuous accuracy–speed dial, spanning conservative ($R \approx 74\%$) to
 227 aggressive ($R \approx 28\%$) compression.

228 The break-even point depends heavily on context length. At conservative compression,
 229 end-to-end (E2E) wall-clock savings emerge at ≥ 64 frames (-9.8%). Under aggressive
 230 compression, E2E savings begin as early as 16 frames (-12.3%) and accelerate to -59.5%
 231 at 128 frames, alongside a 78% reduction in generation time. This scaling trajectory reflects
 232 the quadratic complexity of attention: as sequence length grows, the downstream backbone
 233 savings compound much faster than the linear Allocator overhead, making ResAdapt
 234 exceptionally well-suited for long-context generation.

235 4.4 Analysis and Ablation

236 **Emergent Active Perception.** ResAdapt learns a highly sparse temporal allocation policy.
 237 As shown in Figure 4, the Allocator concentrates high resolution on brief, informative events

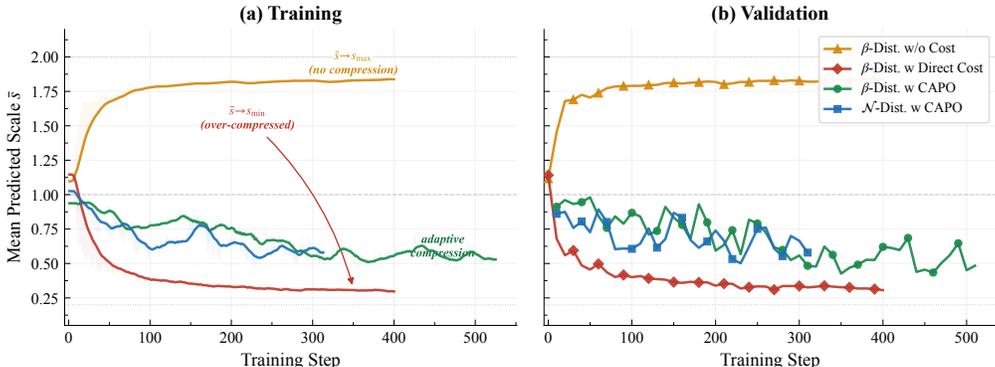


Figure 5: Reward-design ablation. Mean predicted scale \bar{s} during training and validation. Direct cost penalties collapse to the minimum scale, whereas CAPO variants converge to stable intermediate operating points.

238 (e.g., text overlays, scene transitions) while heavily compressing redundant segments. This
 239 behavior goes beyond a trivial positional prior: high-scale allocations emerge as localized,
 240 content-dependent bursts. This confirms that the model selectively expends its pixel budget
 241 precisely where reasoning demands it (see Appendix E.1.1 for detailed global allocation
 242 statistics).

243 **CAPO Ablation.** We evaluate the necessity of CAPO’s asymmetric cost shaping (Ap-
 244 pendix E.2.2). While exact parametric forms (β - vs. \mathcal{N} -CAPO) yield marginal differences,
 245 CAPO consistently outperforms direct cost penalties. As illustrated dynamically in Figure 5,
 246 direct penalties cause the policy to collapse to the minimum scale, whereas CAPO stabilizes
 247 at an intermediate operating point, effectively rewarding selective allocation.

248 **Temporal Regularization.** While CAPO shapes the cost objective, it does not prevent
 249 uniform scaling across visually redundant frames. Removing the temporal similarity loss
 250 \mathcal{L}_{sim} causes the policy to collapse into near-constant scaling (Appendix E.2.1). Reintroducing
 251 \mathcal{L}_{sim} restores sharp frame-level differentiation, demonstrating that temporal regularization
 252 is essential for breaking symmetry and enforcing sparse allocation.

253 **Further Analysis.** Although trained exclusively for continuous resizing, ResAdapt’s pre-
 254 dicted scales generalize zero-shot to discrete frame selection, outperforming uniform sam-
 255 pling at lower budgets (Appendix E.3). Furthermore, while ResAdapt effectively redistri-
 256 butes spatial budgets, we provide a qualitative analysis of its boundary cases—such as
 257 missing extremely brief, simple cues—in Appendix E.5.

258 **5 Conclusion**

259 We introduce ResAdapt, an input-side adaptation framework that shifts visual efficiency
 260 optimization from post-encoding token compression to pre-encoding budget control. By
 261 employing a lightweight Allocator trained via Cost-Aware Policy Optimization (CAPO)
 262 and temporal-similarity regularization, ResAdapt learns a sparse, content-dependent policy
 263 to dynamically resize frames before visual encoding. This approach effectively preserves
 264 critical high-frequency details while heavily compressing redundant segments. Extensive
 265 evaluations demonstrate that ResAdapt establishes a new Pareto frontier in low-budget
 266 video QA and significantly enhances reasoning-augmented long-context grounding by
 267 reinvesting spatial savings into extended temporal coverage. While the current open-loop
 268 design precludes iterative recovery of missed evidence, our findings highlight pre-encoding
 269 allocation as a highly promising and generalizable paradigm for efficient long-context
 270 multimodal reasoning.

271 **References**

- 272 Saeed Ranjbar Alvar, Gursimran Singh, Mohammad Akbari, and Yong Zhang. Divprune:
273 Diversity-based visual token pruning for large multimodal models. In *Proceedings of the*
274 *Computer Vision and Pattern Recognition Conference*, pp. 9392–9401, 2025.
- 275 Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao
276 Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang,
277 Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei
278 Li, Kaixin Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang
279 Liu, Dayiheng Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen
280 Meng, Xuancheng Ren, Xingzhang Ren, Sibao Song, Yuchong Sun, Jun Tang, Jianhong
281 Tu, Jianqiang Wan, Peng Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao
282 Xie, Yiheng Xu, Haiyang Xu, Jin Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang,
283 Bowen Yu, Fei Zhang, Hang Zhang, Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou,
284 Fan Zhou, Jing Zhou, Yuanzhi Zhu, and Ke Zhu. Qwen3-vl technical report, 2025a. URL
285 <https://arxiv.org/abs/2511.21631>.
- 286 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang,
287 Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint*
288 *arXiv:2502.13923*, 2025b.
- 289 Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and
290 Judy Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022.
- 291 Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao
292 Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration
293 for large vision-language models. In *European Conference on Computer Vision*, pp. 19–35.
294 Springer, 2024.
- 295 Yukang Chen, Wei Huang, Baifeng Shi, Qinghao Hu, Hanrong Ye, Ligeng Zhu, Zhijian Liu,
296 Pavlo Molchanov, Jan Kautz, Xiaojuan Qi, et al. Scaling rl to long videos. *arXiv preprint*
297 *arXiv:2507.07966*, 2025.
- 298 Zeyuan Chen, Kai Zhang, Zhuowen Tu, and Yuanjun Xiong. Soft tail-dropping for adaptive
299 visual tokenization. *arXiv preprint arXiv:2601.14246*, 2026.
- 300 Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning.
301 In *International Conference on Learning Representations*, 2024.
- 302 Caba Heilbron Fabian, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activi-
303 tynet: A large-scale video benchmark for human activity understanding. In *Proceedings of*
304 *the ieee conference on computer vision and pattern recognition*, pp. 961–970, 2015.
- 305 Ziyang Fan, Keyu Chen, Ruilong Xing, Yulin Li, Li Jiang, and Zhuotao Tian. Flashvid:
306 Efficient video large language models via training-free tree-based spatiotemporal token
307 merging. *arXiv preprint arXiv:2602.08024*, 2026.
- 308 Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Junfei
309 Wu, Xiaoying Zhang, Benyou Wang, and Xiangyu Yue. Video-r1: Reinforcing video
310 reasoning in mllms. *arXiv preprint arXiv:2503.21776*, 2025.
- 311 Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang,
312 Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever compre-
313 hensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the*
314 *Computer Vision and Pattern Recognition Conference*, pp. 24108–24118, 2025a.
- 315 Shenghao Fu, Qize Yang, Yuan-Ming Li, Xihan Wei, Xiaohua Xie, and Wei-Shi Zheng.
316 Love-r1: Advancing long video understanding with an adaptive zoom-in mechanism via
317 multi-step reasoning. *arXiv preprint arXiv:2509.24786*, 2025b.
- 318 Tianyu Fu, Tengxuan Liu, Qinghao Han, Guohao Dai, Shengen Yan, Huazhong Yang, Xuefei
319 Ning, and Yu Wang. Framefusion: Combining similarity and importance for video token
320 reduction on large vision language models. *arXiv preprint arXiv:2501.01986*, 2024.

- 321 Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity local-
322 ization via language query. In *Proceedings of the IEEE international conference on computer*
323 *vision*, pp. 5267–5275, 2017.
- 324 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
325 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
326 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025a.
- 327 Weiyu Guo, Ziyang Chen, Shaoguang Wang, Jianxiang He, Yijie Xu, Jinhui Ye, Ying Sun,
328 and Hui Xiong. Logic-in-frames: Dynamic keyframe search via visual semantic-logical
329 verification for long video understanding. *arXiv preprint arXiv:2503.13139*, 2025b.
- 330 Zefeng He, Xiaoye Qu, Yafu Li, Siyuan Huang, Daizong Liu, and Yu Cheng. Framethinker:
331 Learning to think with long videos via multi-turn frame spotlighting. *arXiv preprint*
332 *arXiv:2509.24304*, 2025.
- 333 Jack Hong, Chenxiao Zhao, ChengLin Zhu, Weiheng Lu, Guohai Xu, and Xing Yu. Deep-
334 eyesv2: Toward agentic multimodal model. *arXiv preprint arXiv:2511.05271*, 2025.
- 335 Kairui Hu, Penghao Wu, Fanyi Pu, Wang Xiao, Yuanhan Zhang, Xiang Yue, Bo Li, and Ziwei
336 Liu. Video-mmmu: Evaluating knowledge acquisition from multi-discipline professional
337 videos. *arXiv preprint arXiv:2501.13826*, 2025.
- 338 Xiaohu Huang, Hao Zhou, and Kai Han. Prunevid: Visual token pruning for efficient video
339 large language models. In *Findings of the Association for Computational Linguistics: ACL*
340 *2025*, pp. 19959–19973, 2025.
- 341 Jeongseok Hyun, Sukjun Hwang, Su Ho Han, Taeh Kim, Inwoong Lee, Dongyoon Wee,
342 Joon-Young Lee, Seon Joo Kim, and Minhoo Shim. Multi-granular spatio-temporal token
343 merging for training-free acceleration of video llms. In *Proceedings of the IEEE/CVF*
344 *International Conference on Computer Vision*, pp. 23990–24000, 2025.
- 345 Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and
346 Ali Farhadi. A diagram is worth a dozen images. In *European conference on computer vision*,
347 pp. 235–251. Springer, 2016.
- 348 Samir Khaki, Junxian Guo, Jiaming Tang, Shang Yang, Yukang Chen, Konstantinos N
349 Plataniotis, Yao Lu, Song Han, and Zhijian Liu. Sparsevila: Decoupling visual sparsity for
350 efficient vlm inference. In *Proceedings of the IEEE/CVF International Conference on Computer*
351 *Vision*, pp. 23784–23794, 2025.
- 352 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu,
353 Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large
354 language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th*
355 *Symposium on Operating Systems Principles*, 2023.
- 356 Jialuo Li, Bin Li, Jiahao Li, and Yan Lu. Divide, then ground: Adapting frame selection to
357 query types for long-form video understanding. *arXiv preprint arXiv:2512.04000*, 2025a.
- 358 Xinhao Li, Ziang Yan, Desen Meng, Lu Dong, Xiangyu Zeng, Yinan He, Yali Wang, Yu Qiao,
359 Yi Wang, and Limin Wang. Videochat-r1: Enhancing spatio-temporal perception via
360 reinforcement fine-tuning. *arXiv preprint arXiv:2504.06958*, 2025b.
- 361 Hao Liang, Jiapeng Li, Tianyi Bai, Xijie Huang, Linzhuang Sun, Zhengren Wang, Conghui
362 He, Bin Cui, Chong Chen, and Wentao Zhang. Keyvideollm: Towards large-scale video
363 keyframe selection. *arXiv preprint arXiv:2407.03104*, 2024.
- 364 Jiaheng Liu, Dawei Zhu, Zhiqi Bai, Yancheng He, Huanxuan Liao, Haoran Que, Zekun
365 Wang, Chenchen Zhang, Ge Zhang, Jiebin Zhang, et al. A comprehensive survey on long
366 context language modeling. *arXiv preprint arXiv:2503.17407*, 2025a.
- 367 Shuming Liu, Mingchen Zhuge, Changsheng Zhao, Jun Chen, Lemeng Wu, Zechun Liu,
368 Chenchen Zhu, Zhipeng Cai, Chong Zhou, Haozhe Liu, et al. Videoauto-r1: Video auto
369 reasoning via thinking once, answering twice. *arXiv preprint arXiv:2601.05175*, 2026.

- 370 Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng
371 Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr
372 in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024.
- 373 Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng
374 Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, et al. Nvila: Efficient frontier visual language
375 models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 4122–
376 4134, 2025b.
- 377 Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee,
378 and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint*
379 *arXiv:2503.20783*, 2025c.
- 380 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao
381 Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathe-
382 matical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*,
383 2023.
- 384 Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril
385 Zakka, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, et al. Smolvlm: Redefining
386 small and efficient multimodal models. *arXiv preprint arXiv:2504.05299*, 2025.
- 387 Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A
388 benchmark for question answering about charts with visual and logical reasoning. In
389 *Findings of the association for computational linguistics: ACL 2022*, pp. 2263–2279, 2022.
- 390 Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System
391 optimizations enable training deep learning models with over 100 billion parameters. In
392 *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data*
393 *mining*, pp. 3505–3506, 2020.
- 394 Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive
395 token reduction for efficient large multimodal models. In *Proceedings of the IEEE/CVF*
396 *International Conference on Computer Vision*, pp. 22857–22867, 2025.
- 397 Kele Shao, Keda Tao, Can Qin, Haoxuan You, Yang Sui, and Huan Wang. Holitom: Holistic
398 token merging for fast video large language models. *arXiv preprint arXiv:2505.21334*,
399 2025a.
- 400 Kele Shao, Keda Tao, Kejia Zhang, Sicheng Feng, Mu Cai, Yuzhang Shang, Haoxuan You,
401 Can Qin, Yang Sui, and Huan Wang. When tokens talk too much: A survey of multi-
402 modal long-context token compression across images, videos, and audios. *arXiv preprint*
403 *arXiv:2507.20198*, 2025b.
- 404 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
405 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathe-
406 matical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 407 Leqi Shen, Guoqiang Gong, Tao He, Yifeng Zhang, Pengzhang Liu, Sicheng Zhao, and
408 Guiguang Ding. Fastvid: Dynamic density pruning for fast video large language models.
409 *arXiv preprint arXiv:2503.11187*, 2025a.
- 410 Xiaoqian Shen, Min-Hung Chen, Yu-Chiang Frank Wang, Mohamed Elhoseiny, and Ryo
411 Hachiuma. Zoom-zero: Reinforced coarse-to-fine video understanding via temporal
412 zoom-in. *arXiv preprint arXiv:2512.14273*, 2025b.
- 413 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua
414 Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In
415 *Proceedings of the Twentieth European Conference on Computer Systems*, pp. 1279–1297, 2025.
- 416 Yan Shu, Zheng Liu, Peitian Zhang, Minghao Qin, Junjie Zhou, Zhengyang Liang, Tiejun
417 Huang, and Bo Zhao. Video-xl: Extra-long vision language model for hour-scale video
418 understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp.
419 26160–26169, 2025.

- 420 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi
421 Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the*
422 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 8317–8326, 2019.
- 423 Mingyang Song, Haoyu Sun, Jiawei Gu, Linjie Li, Luxin Xu, Ranjay Krishna, and Yu Cheng.
424 Adareasoner: Dynamic tool orchestration for iterative visual reasoning. *arXiv preprint*
425 *arXiv:2601.18631*, 2026.
- 426 Guangyu Sun, Archit Singhal, Burak Uzkent, Mubarak Shah, Chen Chen, and Garin Kessler.
427 From frames to clips: Training-free adaptive key clip selection for long-form video under-
428 standing. *arXiv preprint arXiv:2510.02262*, 2025.
- 429 Xi Tang, Jihao Qiu, Lingxi Xie, Yunjie Tian, Jianbin Jiao, and Qixiang Ye. Adaptive keyframe
430 sampling for long video understanding. *arXiv preprint arXiv:2502.21271*, 2025.
- 431 Keda Tao, Can Qin, Haoxuan You, Yang Sui, and Huan Wang. Dycoke: Dynamic compres-
432 sion of tokens for fast video large language models. In *Proceedings of the Computer Vision*
433 *and Pattern Recognition Conference*, pp. 18992–19001, 2025.
- 434 Haozhe Wang, Alex Su, Weiming Ren, Fangzhen Lin, and Wenhua Chen. Pixel reasoner:
435 Incentivizing pixel-space reasoning with curiosity-driven reinforcement learning. *arXiv*
436 *preprint arXiv:2505.15966*, 2025a.
- 437 Weihang Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Ming Ding,
438 Xiaotao Gu, Shiyu Huang, Bin Xu, et al. Lvbench: An extreme long video understanding
439 benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
440 22958–22967, 2025b.
- 441 Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji.
442 Executable code actions elicit better llm agents. In *Forty-first International Conference on*
443 *Machine Learning*, 2024.
- 444 Ye Wang, Ziheng Wang, Boshen Xu, Yang Du, Kejun Lin, Zihan Xiao, Zihao Yue, Jianzhong
445 Ju, Liang Zhang, Dingyi Yang, et al. Time-r1: Post-training large vision language model
446 for temporal video grounding. *arXiv preprint arXiv:2503.13377*, 2025c.
- 447 Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for
448 long-context interleaved video-language understanding. *Advances in Neural Information*
449 *Processing Systems*, 37:28828–28857, 2024.
- 450 Junbin Xiao, Angela Yao, Yicong Li, and Tat-Seng Chua. Can i trust your answer? visually
451 grounded video question answering. In *Proceedings of the IEEE/CVF Conference on Computer*
452 *Vision and Pattern Recognition*, pp. 13204–13214, 2024.
- 453 Long Xing, Qidong Huang, Xiaoyi Dong, Jiajie Lu, Pan Zhang, Yuhang Zang, Yuhang
454 Cao, Conghui He, Jiaqi Wang, Feng Wu, et al. Pyramiddrop: Accelerating your
455 large vision-language models via pyramid visual redundancy reduction. *arXiv preprint*
456 *arXiv:2410.17247*, 2024.
- 457 Ruyi Xu, Guangxuan Xiao, Yukang Chen, Liuning He, Kelly Peng, Yao Lu, and Song
458 Han. Streamingvlm: Real-time understanding for infinite video streams. *arXiv preprint*
459 *arXiv:2510.09608*, 2025.
- 460 Biao Yang, Bin Wen, Boyang Ding, Changyi Liu, Chenglong Chu, Chengru Song, Chongling
461 Rao, Chuan Yi, Da Li, Dunju Zang, et al. Kwai keye-v1.5 technical report. *arXiv preprint*
462 *arXiv:2509.01563*, 2025a.
- 463 Cheng Yang, Yang Sui, Jinqi Xiao, Lingyi Huang, Yu Gong, Chendi Li, Jinghua Yan, Yu Bai,
464 Ponnuswamy Sadayappan, Xia Hu, et al. Topv: Compatible token pruning with inference
465 time optimization for fast and low-memory multimodal vision language model. In
466 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 19803–19813,
467 2025b.

- 468 Senqiao Yang, Yukang Chen, Zhuotao Tian, Chengyao Wang, Jingyao Li, Bei Yu, and Jiaya
469 Jia. Visionzip: Longer is better but not necessary in vision language models. In *Proceedings*
470 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19792–19802,
471 2025c.
- 472 Zuhao Yang, Sudong Wang, Kaichen Zhang, Keming Wu, Sicong Leng, Yifan Zhang, Bo Li,
473 Chengwei Qin, Shijian Lu, Xingxuan Li, and Lidong Bing. Longvt: Incentivizing “thinking
474 with long videos” via native tool calling. *arXiv preprint arXiv:2511.20785*, 2025d.
- 475 Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai,
476 Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement
477 learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- 478 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,
479 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline
480 multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the*
481 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.
- 482 Boqiang Zhang, Lei Ke, Ruihan Yang, Qi Gao, Tianyuan Qu, Rossell Chen, Dong Yu, et al.
483 Penguin-vl: Exploring the efficiency limits of vlm with llm-based vision encoders. *arXiv*
484 *preprint arXiv:2603.06569*, 2026.
- 485 Ce Zhang, Kaixin Ma, Tianqing Fang, Wenhao Yu, Hongming Zhang, Zhisong Zhang, Yaqi
486 Xie, Katia Sycara, Haitao Mi, and Dong Yu. Vscan: Rethinking visual token reduction for
487 efficient large vision-language models. *arXiv preprint arXiv:2505.22654*, 2025a.
- 488 Congzhi Zhang, Zhibin Wang, Yinchao Ma, Jiawei Peng, Yihan Wang, Qiang Zhou, Jun Song,
489 and Bo Zheng. Rewatch-r1: Boosting complex video reasoning in large vision-language
490 models through agentic data synthesis. *arXiv preprint arXiv:2509.23652*, 2025b.
- 491 Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai
492 Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Reality
493 check on the evaluation of large multimodal models, 2024a. URL [https://arxiv.org/](https://arxiv.org/abs/2407.12772)
494 [abs/2407.12772](https://arxiv.org/abs/2407.12772).
- 495 Qizhe Zhang, Aosong Cheng, Ming Lu, Renrui Zhang, Zhiyong Zhuo, Jiajun Cao, Shaobo
496 Guo, Qi She, and Shanghang Zhang. Beyond text-visual attention: Exploiting visual cues
497 for effective token pruning in vlms. In *Proceedings of the IEEE/CVF International Conference*
498 *on Computer Vision*, pp. 20857–20867, 2025c.
- 499 Shaojie Zhang, Jiahui Yang, Jianqin Yin, Zhenbo Luo, and Jian Luan. Q-frame: Query-
500 aware frame selection and multi-resolution adaptation for video-llms. *arXiv preprint*
501 *arXiv:2506.22139*, 2025d.
- 502 Yi-Fan Zhang, Xingyu Lu, Shukang Yin, Chaoyou Fu, Wei Chen, Xiao Hu, Bin Wen, Kaiyu
503 Jiang, Changyi Liu, Tianke Zhang, et al. Thyme: Think beyond images. *arXiv preprint*
504 *arXiv:2508.11630*, 2025e.
- 505 Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng,
506 Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm:
507 Visual token sparsification for efficient vision-language model inference. *arXiv preprint*
508 *arXiv:2410.04417*, 2024b.
- 509 Shitian Zhao, Haoquan Zhang, Shaoheng Lin, Ming Li, Qilong Wu, Kaipeng Zhang, and
510 Chen Wei. Pyvision: Agentic vision with dynamic tooling. *arXiv preprint arXiv:2507.07998*,
511 2025a.
- 512 Yilun Zhao, Haowei Zhang, Lujing Xie, Tongyan Hu, Guo Gan, Yitao Long, Zhiyuan
513 Hu, Weiyuan Chen, Chuhan Li, Zhijian Xu, et al. Mmvu: Measuring expert-level multi-
514 discipline video understanding. In *Proceedings of the Computer Vision and Pattern Recognition*
515 *Conference*, pp. 8475–8489, 2025b.

- 516 Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang,
517 Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. *arXiv preprint*
518 *arXiv:2507.18071*, 2025a.
- 519 Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu,
520 Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying
521 Sheng. SGLang: Efficient execution of structured language model programs. In *Advances*
522 *in Neural Information Processing Systems*, 2024.
- 523 Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao Zhao, Guohai Xu, Le Yang, Chao Shen,
524 and Xing Yu. Deepeyes: Incentivizing “thinking with images” via reinforcement learning.
525 *arXiv preprint arXiv:2505.14362*, 2025b.
- 526 Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Zhengyang Liang, Shitao Xiao, Minghao Qin,
527 Xi Yang, Yongping Xiong, Bo Zhang, et al. Mlvu: Benchmarking multi-task long video
528 understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
529 *Recognition*, pp. 13691–13701, 2025.
- 530 Zirui Zhu, Hailun Xu, Yang Luo, Yong Liu, Kanchan Sarkar, Zhenheng Yang, and Yang
531 You. Focus: Efficient keyframe selection for long video understanding. *arXiv preprint*
532 *arXiv:2510.27280*, 2025.
- 533 Yuanhao Zou, Shengji Jin, Andong Deng, Youpeng Zhao, Jun Wang, and Chen Chen. Air:
534 Enabling adaptive, iterative, and reasoning-based frame selection for video question
535 answering. *arXiv preprint arXiv:2510.04428*, 2025.

536 Limitations and future work.

- 537 ResAdapt improves the efficiency–accuracy trade-off for long-video MLLMs, but the current
538 evidence is still bounded by four concrete design choices.
- 539 (i) *Front-end overhead is amortized only in the long-context regime.* The Allocator adds a fixed
540 pre-encoding cost—coarse visual encoding, cross-frame fusion, and distribution prediction—
541 before any backbone savings are realized. When the sequence is short ($T \leq 32$), this constant
542 overhead can offset a meaningful fraction of the downstream attention reduction, so the
543 clearest wall-clock gains appear only when temporal context is long (Sec. 4.3). Reducing this
544 fixed cost through cached video features, cheaper front-ends, or distilled allocation rules is
545 therefore an important next step.
- 546 (ii) *Allocation is limited by coarse visual evidence.* The Allocator observes frozen coarse fea-
547 tures $f_t \in \mathbb{R}^D$ rather than the full high-resolution frame. This is sufficient to detect broad
548 redundancy and scene structure, but it is weaker on small text, subtle objects, and brief
549 answer-critical cues embedded in otherwise simple frames (Figure 20). Multi-scale condi-
550 tioning, motion-aware features, or lightweight local refinement would help close this gap
551 without giving up the speed advantage of the current front-end.
- 552 (iii) *The present study validates the framework through one video-centric instantiation.* Our formu-
553 lation is general input-side adaptation, but the experiments instantiate the operator with
554 resizing and train the policy primarily on video tasks. As a result, transfer beyond this
555 regime is uneven: the learned policy sometimes recognizes image inputs that need more
556 fidelity, yet it does not deliver uniformly efficiency-preserving gains on static-image bench-
557 marks (Table 7). Extending training to mixed image–video data and alternative operators
558 such as hard frame selection remains open.
- 559 (iv) *Allocation is open-loop rather than reasoning-aware.* All budget decisions are committed
560 before the backbone processes any visual token. The policy therefore cannot revise a
561 mistaken low-resolution choice after partial reasoning or uncertainty signals emerge. A
562 natural extension is closed-loop allocation, where early backbone states trigger re-encoding,
563 budget revision, or a second visual pass only when needed.

564 Software and Data

565 The code for this paper is available at: <https://github.com/Xnhyacinth/ResAdapt>

566 A Related Work

567 **Input-side adaptation.** Reducing visual cost *before* encoding typically involves keyframe
568 selection (Liang et al., 2024; Zhu et al., 2025; Sun et al., 2025; Tang et al., 2025), recently
569 augmented by query-aware iterative search (Zou et al., 2025; Li et al., 2025a; Guo et al.,
570 2025b; He et al., 2025). Alternatively, multi-resolution encoding routes frames to varying
571 resolutions based on inter-frame similarity (Yang et al., 2025a; Zhang et al., 2026) or query
572 conditioning (Zhang et al., 2025d; Chen et al., 2026). However, these rely on handcrafted
573 heuristics, binary routing, or fixed resolution bins. In contrast, ResAdapt learns an optimal,
574 continuous input-side allocation policy directly from task rewards via RL, generalizing
575 across pre-encoding operators like resizing and selection.

576 **Model-side token economy.** Post-encoding methods compress visual tokens in the em-
577 bedding space via merging (Bolya et al., 2022), saliency-guided pruning (Chen et al., 2024;
578 Yang et al., 2025c; Shang et al., 2025; Zhang et al., 2025c), progressive dropping (Xing et al.,
579 2024; Zhang et al., 2024b), and diversity-based allocation (Alvar et al., 2025; Yang et al.,
580 2025b; Zhang et al., 2025a). Video extensions further exploit spatiotemporal redundancy
581 through token separation (Huang et al., 2025; Shen et al., 2025a), hierarchical merging (Hyun
582 et al., 2025), and segment-level fusion (Tao et al., 2025; Fu et al., 2024; Shao et al., 2025a).
583 These techniques are complementary to ResAdapt; they operate *after* encoding, whereas we
584 determine the initial pixel budget to prevent irreversible loss of high-frequency details.

585 **Output-side agentic reasoning.** To recover efficiency, iterative reasoning methods keep
586 the input fixed but re-query the model after retrieving frames or zooming into regions.
587 These employ either static cropping/clipping operators (Zheng et al., 2025b; Wang et al.,
588 2025a; Song et al., 2026) or dynamic code-generation tools (Zhang et al., 2025e; Zhao et al.,
589 2025a; Hong et al., 2025; Wang et al., 2024). While precise, they are inherently multi-pass
590 and suffer from latency and control overhead. ResAdapt demonstrates that a *single-pass*
591 pre-encoding allocation policy can achieve similar precision without iterative interaction.

592 **RL for multimodal reasoning.** RL post-training, successful in LLMs (Shao et al., 2024; Guo
593 et al., 2025a), has recently been adapted for multimodal reasoning (Liu et al., 2025c; Yu et al.,
594 2025; Zheng et al., 2025a) and video understanding via iterative evidence refinement (Feng
595 et al., 2025; Li et al., 2025b; Liu et al., 2026; Yang et al., 2025d; Chen et al., 2025; Wang
596 et al., 2025c; Fu et al., 2025b). Our approach is orthogonal: rather than training output-side
597 reasoning policies, we utilize RL for *input-side perception control* to optimize frame-level
598 visual allocations under a strict accuracy–cost trade-off. To this end, CAPO specifically
599 prevents the degenerate low-budget collapse caused by naive cost penalties.

600 B Implementation Details

601 B.1 Training Data

602 **Data Composition.** We build the training set from the difficulty-filtered data of VideoAuto-
603 R1 (Liu et al., 2026), keeping only image and video samples and discarding pure-text
604 examples. To improve coverage of visually demanding subdomains, we additionally sample
605 16,500 video instances from Video-R1 (Feng et al., 2025), focusing on OCR, free-form QA,
606 and regression-style tasks. The merged pool contains approximately 93.4K training samples.
607 We manually remove all evaluation examples from our benchmark suite to avoid leakage.

608 B.2 Training Configuration

609 Unless otherwise noted, training runs for one epoch with global batch size 128 and AdamW.
 610 The learning rate is 2×10^{-5} for the Allocator and 1×10^{-6} for the backbone, with weight
 611 decay 0.01 and gradient clipping at 1.0. We set the maximum video token budget to 8,192,
 612 use $T=128$ frames during training, and allow scales in the range $[s_{\min}, s_{\max}] = [0.2, 1.8]$,
 613 which permits both downscaling and selective upscaling. CAPO samples $M=16$ allocation
 614 trajectories per prompt and $N=1$ rollout per trajectory. Training is conducted on 32 H100
 615 GPUs with VeRL (Sheng et al., 2025), DeepSpeed (Rasley et al., 2020), and vLLM (Kwon
 616 et al., 2023). Evaluation uses lmms-eval (Zhang et al., 2024a); unless stated otherwise, we
 617 cap response length at 256 tokens and increase it to 4,096 for reasoning models.

618 B.3 Reward Design

619 We provide full details complementing Sec. 3.3. The base scalar reward $R_{m,n}^{\text{task}}$ is task-specific;
 620 efficiency enters later through CAPO advantage shaping rather than through a raw additive
 621 reward term.

622 **Base Task Reward ($R_{m,n}^{\text{task}}$).** We consider four task types:

- 623 • *Question Answering.* For math problems, we extract the numeric answer and com-
 624 pare it to the ground truth within a tolerance of 10^{-2} . For multiple-choice questions,
 625 we extract the option letter. For other QA tasks, we compare normalized strings
 626 (e.g., case-folded, whitespace-stripped). This yields the binary reward

$$R_{\text{QA}}(\hat{o}, o) \in \{0, 1\}.$$

- 627 • *Free-form Generation.* For open-ended tasks, we compute the ROUGE-L score be-
 628 tween the generated answer \hat{o} and the reference o :

$$R_{\text{Gen}}(\hat{o}, o) = \text{ROUGE-L}(\hat{o}, o) \in [0, 1].$$

- 629 • *Temporal Grounding.* Let the ground-truth segments be $\mathcal{G} = \{[s_j, e_j]\}_j$ and the pre-
 630 dicted segments be $\hat{\mathcal{G}} = \{[\hat{s}_k, \hat{e}_k]\}_k$ (each set may contain one or multiple intervals).
 631 We compute the temporal IoU and select the best-matching pair:

$$R_{\text{TG}}(\hat{\mathcal{G}}, \mathcal{G}) = \max_{[\hat{s}, \hat{e}] \in \hat{\mathcal{G}}, [s, e] \in \mathcal{G}} \text{tIoU}([\hat{s}, \hat{e}], [s, e]) \in [0, 1].$$

632 If no valid segment can be parsed from the output, we assign $R_{\text{TG}}(\hat{\mathcal{G}}, \mathcal{G}) = 0$.

- 633 • *Grounding QA.* We parse both the textual answer and the predicted temporal seg-
 634 ments from the model output, compute $R_{\text{QA}}(\hat{o}, o)$ and $R_{\text{TG}}(\hat{\mathcal{G}}, \mathcal{G})$, and sum them:

$$R_{\text{GQA}}(\hat{o}, \hat{\mathcal{G}}; o, \mathcal{G}) = R_{\text{QA}}(\hat{o}, o) + R_{\text{TG}}(\hat{\mathcal{G}}, \mathcal{G}) \in [0, 2].$$

635 These task-specific metrics define the scalar base reward $R_{m,n}^{\text{task}}$. CAPO additionally uses
 636 a binary success indicator $u_{m,n} \in \{0, 1\}$: for exact-match QA tasks we use the binary
 637 correctness outcome directly, whereas for continuous metrics (ROUGE-L, temporal IoU,
 638 and their grounding-QA combination) we threshold the scalar score at 0.35, matching the
 639 implementation. When format validation is enabled, a weighted format term is added
 640 before GRPO normalization, but $u_{m,n}$ is computed from the task metric alone.

641 **Format Reward.** We employ a binary format reward $R_{\text{fmt}}(\hat{o}) \in \{0, 1\}$ enforced via strict
 642 regex validation. The output must contain exactly one `<think>...</think>` block and one
 643 `<answer>...</answer>` block, with the final answer enclosed in `\boxed{...}` within the
 644 `<answer>` tags:

$$R_{\text{fmt}}(\hat{o}) = \begin{cases} 1 & \text{if format matches regex,} \\ 0 & \text{otherwise.} \end{cases}$$

645 In the implementation, malformed outputs receive a penalty before weighting, and the
 646 format term enters the scalar reward with weight 0.2.

Table 3: Prompt template used for CAPO training. The template presents video frames and the task question, requires intermediate reasoning inside `<think>` tags, and places the final answer in `\boxed{}` within `<answer>` tags. This structure enables automatic reward extraction from MLLM outputs.

Prompt Template for Training with Thinking
<p>System Prompt: You are a helpful assistant. You FIRST think about the reasoning process as an internal monologue and then provide the final answer. The reasoning process MUST BE enclosed within <code><think></code> <code></think></code> tags and the answer MUST BE enclosed within <code><answer></code> <code></answer></code> tags. The final answer MUST BE put in <code>\boxed{}</code> and the <code>\boxed{}</code> expression MUST BE contained entirely within the <code><answer></code> <code></answer></code> tags. Do not include any reasoning or explanations outside these tags.</p>

647 B.4 Prompt Template

648 We employ the standard prompt for GRPO training, shown in Table 3. The model generates
 649 a reasoning trace within `<think>` `</think>` tags (optional for ResAdapt since reasoning is
 650 handled by the MLLM π_ϕ , but maintained for compatibility with reasoning-based baselines),
 651 followed by the final answer enclosed in `\boxed{}`.

652 C Complexity Analysis

653 We derive formal computational bounds for ResAdapt to clarify when Allocator overhead
 654 is negligible relative to the savings induced in the backbone. For readability, we assume
 655 a standard Transformer backbone with quadratic self-attention and a uniform native reso-
 656 lution $H \times W$ over T frames; the extension to heterogeneous resolutions is immediate by
 657 replacing HW with per-frame products $H_t W_t$.

658 **Baseline cost.** Let P denote the ViT patch size. A vanilla MLLM encoding T frames at full
 659 resolution incurs a total visual token count of:

$$N_0 = T \cdot \left\lceil \frac{H}{P} \right\rceil \left\lceil \frac{W}{P} \right\rceil \approx \frac{THW}{P^2}. \quad (12)$$

660 **Adaptive cost and token retention ratio.** For the resize instantiation analyzed in this paper,
 661 frame f_t is rescaled by factor $s_t \in [s_{\min}, s_{\max}]$, producing $n_t(s_t) = \lceil s_t H/P \rceil \lceil s_t W/P \rceil \approx$
 662 $s_t^2 \cdot HW/P^2$ tokens. Summing over the sequence and normalizing by N_0 yields the *token*
 663 *retention ratio*:

$$N^{\text{adapt}} = \sum_{t=1}^T n_t(s_t) \approx \frac{HW}{P^2} \sum_{t=1}^T s_t^2, \quad \rho \triangleq \frac{N^{\text{adapt}}}{N_0} = \frac{1}{T} \sum_{t=1}^T s_t^2. \quad (13)$$

664 Because the learned Beta policy places most redundant frames near s_{\min} (Figure 8), ρ is
 665 much smaller than 1 in practice; across our evaluation suite, $\rho \in [0.06, 0.16]$.

666 **Quadratic FLOPs reduction.** For an L_{mlLM} -layer MLLM with hidden dimension D_{mlLM} , self-
 667 attention cost scales quadratically in the visual sequence length: $\Phi(N) = O(L_{\text{mlLM}} N^2 D_{\text{mlLM}})$.
 668 Substituting $N^{\text{adapt}} = \rho \cdot N_0$ gives:

$$\Phi_{\text{mlLM}}^{\text{adapt}} = O\left(L_{\text{mlLM}} \cdot \rho^2 N_0^2 \cdot D_{\text{mlLM}}\right), \quad (14)$$

669 a reduction by a factor of ρ^2 relative to full-resolution processing. At the representative
 670 operating point $\rho = 0.11$, we obtain $\rho^2 \approx 0.012$, corresponding to roughly $83\times$ fewer
 671 backbone attention FLOPs.

672 **Allocator overhead.** The Allocator processes $N_c = T \cdot \lceil H/P_c \rceil \lceil W/P_c \rceil$ coarsely pooled
 673 tokens across L_{pred} layers with dimension D_{pred} , where $P_c \gg P$ is the coarse spatial stride.
 674 Its cost and relative overhead are:

$$\Phi_{\text{pred}} = O\left(L_{\text{pred}} \cdot N_c^2 \cdot D_{\text{pred}}\right), \quad \frac{\Phi_{\text{pred}}}{\Phi_{\text{mllm}}^{\text{base}}} = O\left(\frac{L_{\text{pred}} D_{\text{pred}}}{L_{\text{mllm}} D_{\text{mllm}}} \cdot \left(\frac{P}{P_c}\right)^4\right) \ll 1. \quad (15)$$

675 Substituting our implementation parameters ($P_c=14$, $L_{\text{pred}}=4$, $D_{\text{pred}}=1,024$ versus
 676 $L_{\text{mllm}}=28$, $D_{\text{mllm}}=3,584$), the Allocator accounts for less than 3% of inference FLOPs. The
 677 decision stage is therefore small compared with the backbone computation it helps eliminate.

678 **Net speedup.** Combining the above under the first-order approximation $\Phi_{\text{mllm}}^{\text{base}} \gg \Phi_{\text{pred}}$:

$$\text{Speedup} \approx \frac{\Phi_{\text{mllm}}^{\text{base}}}{\Phi_{\text{mllm}}^{\text{adapt}} + \Phi_{\text{pred}}} \approx \frac{N_0^2}{(N^{\text{adapt}})^2} = \frac{1}{\rho^2}. \quad (16)$$

679 At $\rho = 0.11$, this again yields a theoretical reduction of roughly $83\times$ in backbone attention
 680 computation.

681 **Temporal context scaling.** The same savings admit a second interpretation in terms of
 682 *temporal coverage*. Under a fixed token budget B , a vanilla MLLM can process only $T_0 =$
 683 $BP^2 / (HW)$ full-resolution frames, whereas the resize instantiation of ResAdapt used in our
 684 experiments can process T_0/ρ adaptively resized frames. This yields an effective $1/\rho \approx 6-$
 685 $16\times$ increase in temporal horizon at comparable compute, which is exactly the trade-off
 686 exploited by the long-context experiments in Sec. 4.2.

687 **Remark (acceleration transparency).** A practical consequence of Input-side adaptation is
 688 that the backbone still receives an ordinary visual-token sequence, only shorter. As a re-
 689 sult, ResAdapt remains compatible with optimized attention stacks such as FlashAttention,
 690 vLLM (Kwon et al., 2023), and SGLang (Zheng et al., 2024) without kernel-level modifica-
 691 tions. By contrast, model-side pruning and merging often create irregular token layouts
 692 that are harder to route through the same optimized kernels and may require fallback
 693 implementations or architecture-specific engineering.

694 D Derivation of Joint RL Formulation

695 This appendix collects derivations omitted from Sec. 3 for space and clarifies how the one-
 696 step contextual MDP (Contextual Bandit) introduced in Sec. 2.2 motivates the practical
 697 surrogate objectives optimized by ResAdapt. All derivations are stated for a single context
 698 (video and query); the full objective is the expectation over the dataset \mathcal{D} .

699 *Notation.* The prompt context is $x = (q, \mathcal{V})$. The Allocator policy $\pi_\theta(s | x)$ samples a
 700 continuous allocation vector $s = (s_1, \dots, s_T)$. A deterministic transformation constructs the
 701 operator-transformed input $\tilde{x} = (q, \{\mathcal{O}(f_t, s_t)\}_{t=1}^T)$; in the experimental instantiation, \mathcal{O} is
 702 bilinear resizing. The MLLM backbone policy $\pi_\phi(y | \tilde{x})$ then samples a full response rollout
 703 $y = (r, o)$, where r is the reasoning trace and o is the final answer.

704 D.1 One-Step Contextual MDP and the Joint Objective

705 As defined in Sec. 2.2, the system is a one-step contextual MDP. In this setting, there are no
 706 sequential state transitions across time steps t ; the episode terminates after the allocation s
 707 is sampled and the corresponding rollout y is produced. Consequently, the value functions
 708 collapse to the immediate reward, and the standard Policy Gradient Theorem simplifies
 709 drastically without requiring temporal discount factors or credit assignment across Markov
 710 states.

711 The joint distribution of the allocation and the rollout factorizes conditionally:

$$p_{\theta,\phi}(\mathbf{s}, \mathbf{y} \mid \mathbf{x}) = \pi_{\theta}(\mathbf{s} \mid \mathbf{x}) \pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}}). \quad (17)$$

712 For a single context with ground-truth answer \mathbf{o}^* , the marginal answer probability under
713 the transformed input is

$$p_{\theta,\phi}(\mathbf{o}^* \mid \mathbf{x}) = \mathbb{E}_{\pi_{\theta}(\mathbf{s} \mid \mathbf{x})} \left[\mathbb{E}_{\pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}})} \left[\pi_{\phi}(\mathbf{o}^* \mid \tilde{\mathbf{x}}, \mathbf{r}) \right] \right]. \quad (18)$$

714 Because $\log(\cdot)$ is monotone, maximizing $\log p_{\theta,\phi}(\mathbf{o}^* \mid \mathbf{x})$ would be equivalent, but the RL
715 derivation below does not require introducing the logarithm. It only requires a scalar utility
716 evaluated after sampling (\mathbf{s}, \mathbf{y}) . We therefore abstract the answer-quality term as a rollout
717 utility $Q(\mathbf{x}, \mathbf{y})$, where $\mathbf{y} = (\mathbf{r}, \mathbf{o})$, and treat it as parameter-independent once the rollout
718 is sampled. This is a modeling abstraction rather than an exact reformulation: when Q is
719 chosen as an answer-aligned task score, the resulting RL problem is a surrogate to likelihood
720 maximization. This lets us define the ideal rollout reward

$$R_{\mathbf{s},\mathbf{y}}^{\text{ideal}} = Q(\mathbf{x}, \mathbf{y}) - \lambda C(\mathbf{s}), \quad (19)$$

721 and optimize the one-step expected return

$$\max_{\theta,\phi} \mathcal{J}(\theta, \phi) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathbb{E}_{\pi_{\theta}(\mathbf{s} \mid \mathbf{x})} \left[\mathbb{E}_{\pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}})} \left[R_{\mathbf{s},\mathbf{y}}^{\text{ideal}} \right] \right]. \quad (20)$$

722 D.2 Policy Gradient and Alternating Optimization

723 Because the objective involves two distinct parameterized policies, its gradients follow the
724 score-function estimator (the likelihood-ratio / REINFORCE identity). This is the underlying
725 policy-gradient structure; GRPO does not change that structure, but replaces the raw reward
726 with normalized advantages and clipped surrogates for practical optimization. Taking the
727 gradient of $\mathcal{J}(\theta, \phi)$ with respect to the backbone parameters ϕ :

$$\begin{aligned} \nabla_{\phi} \mathcal{J}(\theta, \phi) &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\pi_{\theta}(\mathbf{s} \mid \mathbf{x})} \left[\nabla_{\phi} \int \pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}}) R_{\mathbf{s},\mathbf{y}}^{\text{ideal}} d\mathbf{y} \right] \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\pi_{\theta}(\mathbf{s} \mid \mathbf{x})} \mathbb{E}_{\pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}})} \left[R_{\mathbf{s},\mathbf{y}}^{\text{ideal}} \nabla_{\phi} \log \pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}}) \right]. \end{aligned} \quad (21)$$

728 Taking the gradient of $\mathcal{J}(\theta, \phi)$ with respect to the Allocator parameters θ :

$$\begin{aligned} \nabla_{\theta} \mathcal{J}(\theta, \phi) &= \mathbb{E}_{\mathbf{x}} \left[\nabla_{\theta} \int \pi_{\theta}(\mathbf{s} \mid \mathbf{x}) \left(\int \pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}}) R_{\mathbf{s},\mathbf{y}}^{\text{ideal}} d\mathbf{y} \right) d\mathbf{s} \right] \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\pi_{\theta}(\mathbf{s} \mid \mathbf{x})} \left[\left(\int \pi_{\phi}(\mathbf{y} \mid \tilde{\mathbf{x}}) R_{\mathbf{s},\mathbf{y}}^{\text{ideal}} d\mathbf{y} \right) \nabla_{\theta} \log \pi_{\theta}(\mathbf{s} \mid \mathbf{x}) \right]. \end{aligned} \quad (22)$$

729 In practice, we substitute the ideal reward $R_{\mathbf{s},\mathbf{y}}^{\text{ideal}}$ with the shaped CAPO advantage $A_{m,n}$,
730 leading to the clipped surrogate objectives \mathcal{L}_{ϕ} and \mathcal{L}_{θ} described in the main text.

731 For the backbone update, the standard token-level clipped surrogate objective is:

$$\mathcal{L}_{\phi} = -\frac{1}{MN} \sum_{m=1}^M \sum_{n=1}^N \frac{1}{L_{m,n}} \sum_{j=1}^{L_{m,n}} \min \left(r_{\phi,j}^{(m,n)} A_{m,n}, \text{clip} \left(r_{\phi,j}^{(m,n)}, 1-\varepsilon, 1+\varepsilon \right) A_{m,n} \right), \quad (23)$$

732 where $L_{m,n}$ is the rollout length and

$$r_{\phi,j}^{(m,n)} = \frac{\pi_{\phi}(\mathbf{y}_j^{(m,n)} \mid \mathbf{y}_{<j}^{(m,n)}, \tilde{\mathbf{x}}^{(m)})}{\pi_{\phi_{\text{old}}}(\mathbf{y}_j^{(m,n)} \mid \mathbf{y}_{<j}^{(m,n)}, \tilde{\mathbf{x}}^{(m)})}. \quad (24)$$

733 The exact change-of-variables mapping $\mathbf{a} \mapsto \mathbf{s}$ yields the log-probability:

$$\log q_{\theta}(\mathbf{a} \mid \mathbf{x}) = \sum_{t=1}^T \log \text{Beta}(a_t; \alpha_t, \beta_t). \quad (25)$$

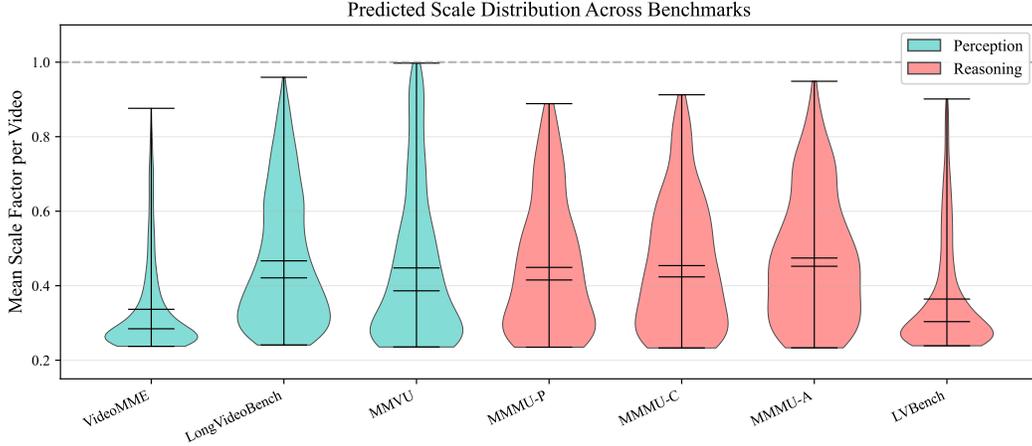


Figure 6: Per-video mean scale across benchmarks. Kernel density estimates of the per-video mean scale \bar{s} . Reasoning-heavy benchmarks shift toward larger \bar{s} than perception-heavy ones, indicating that the learned policy spends more fidelity where fine-grained evidence is more likely to matter.

734 Similarly, the gradient with respect to the Allocator parameters θ relies on the marginalized
 735 reward $R_s^{\text{ideal}} = \mathbb{E}_{\pi_{\phi}(\mathbf{y}|\tilde{\mathbf{x}})}[R_{s,\mathbf{y}}^{\text{ideal}}]$:

$$\nabla_{\theta} \mathcal{J}(\theta, \phi) = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\pi_{\theta}(s|\mathbf{x})} \left[R_s^{\text{ideal}} \nabla_{\theta} \log \pi_{\theta}(s|\mathbf{x}) \right]. \quad (26)$$

736 To optimize this objective with GRPO, we introduce importance sampling from behavior
 737 policies $\pi_{\theta_{\text{old}}}$ and $\pi_{\phi_{\text{old}}}$. A naive joint importance weight $\frac{\pi_{\theta} \pi_{\phi}}{\pi_{\theta_{\text{old}}} \pi_{\phi_{\text{old}}}}$ suffers from compounded
 738 variance. We therefore use an **alternating block-coordinate ascent** approximation. When
 739 updating the MLLM (ϕ), we fix the Allocator to its behavior policy ($\pi_{\theta} = \pi_{\theta_{\text{old}}}$), making its
 740 importance ratio exactly 1. The off-policy surrogate gradient for ϕ becomes:

$$\nabla_{\phi} \mathcal{J}_{\text{surr}}(\phi) = \mathbb{E}_{\pi_{\theta_{\text{old}}}} \mathbb{E}_{\pi_{\phi_{\text{old}}}} \left[\frac{\pi_{\phi}(\mathbf{y}|\tilde{\mathbf{x}})}{\pi_{\phi_{\text{old}}}(\mathbf{y}|\tilde{\mathbf{x}})} R_{s,\mathbf{y}}^{\text{ideal}} \nabla_{\phi} \log \pi_{\phi}(\mathbf{y}|\tilde{\mathbf{x}}) \right]. \quad (27)$$

741 Using the log-derivative identity $\nabla_{\phi} r_{\phi} = r_{\phi} \nabla_{\phi} \log \pi_{\phi}$ where $r_{\phi} = \pi_{\phi} / \pi_{\phi_{\text{old}}}$, this motivates
 742 the surrogate objective:

$$\mathcal{L}_{\phi}^{\text{ideal}} = \mathbb{E}_{\pi_{\theta_{\text{old}}}} \mathbb{E}_{\pi_{\phi_{\text{old}}}} \left[r_{\phi}(\mathbf{y}|\tilde{\mathbf{x}}) R_{s,\mathbf{y}}^{\text{ideal}} \right]. \quad (28)$$

743 Conversely, when updating the Allocator (θ), we fix the backbone to its behavior policy
 744 ($\pi_{\phi} = \pi_{\phi_{\text{old}}}$). The corresponding ideal allocator surrogate is

$$\mathcal{L}_{\theta}^{\text{ideal}} = \mathbb{E}_{\pi_{\theta_{\text{old}}}} \left[r_{\theta}(s|\mathbf{x}) R_s^{\text{ideal}} \right], \quad r_{\theta}(s|\mathbf{x}) = \frac{\pi_{\theta}(s|\mathbf{x})}{\pi_{\theta_{\text{old}}}(s|\mathbf{x})}, \quad (29)$$

745 where $R_s^{\text{ideal}} = \mathbb{E}_{\pi_{\phi_{\text{old}}}(y|\tilde{x})}[R_{s,y}^{\text{ideal}}]$. In practice, this expectation is approximated by Monte
 746 Carlo rollouts under the frozen backbone.

747 D.3 Advantage Shaping and Monte Carlo Surrogates

748 The ideal linear penalty $-\lambda C(s)$ inside R^{ideal} often causes catastrophic collapse to minimum
 749 budgets. CAPO therefore replaces the raw reward with a cost-shaped, group-normalized
 750 advantage $A_{s,y}$ (denoted $A_{m,n}$ in the main text). This replacement is *not* an unbiased baseline
 751 transformation of $R_{s,y}^{\text{ideal}}$: the CAPO signal depends on the sampled allocation, the rollout
 752 outcome, and the within-group cost statistics. Instead, it defines a deliberately biased

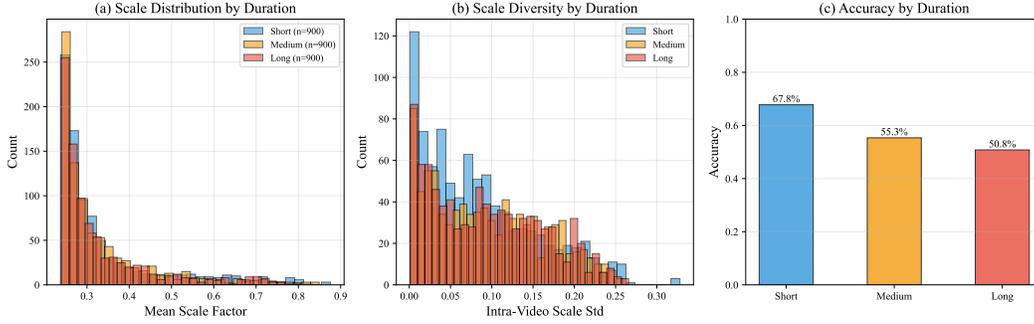


Figure 7: VideoMME broken down by video duration. As clip duration grows, the policy lowers the average scale, increases within-video scale diversity, and faces lower task accuracy. Longer clips are therefore processed more aggressively and more selectively.

753 surrogate objective that trades exact fidelity to the Lagrangian reward for lower variance
754 and stronger budget control in practice.

755 Applying surrogate clipping to the exact joint ratios would couple all frame- and token-level
756 factors, which is prohibitively noisy in practice. We therefore arrive at practical decoupled
757 objectives. For a batch of M allocations and N rollouts per allocation, the MLLM sequence-
758 level surrogate is:

$$\mathcal{L}_\phi^{\text{seq}} = -\frac{1}{MN} \sum_{m=1}^M \sum_{n=1}^N \min\left(r_\phi^{(m,n)} A_{m,n}, \text{clip}(r_\phi^{(m,n)}, 1-\varepsilon, 1+\varepsilon) A_{m,n}\right). \quad (30)$$

759 This sequence-level loss is already approximate because it uses the CAPO-shaped advantage
760 in place of the ideal reward. To achieve finer credit assignment for the autoregressive MLLM,
761 we further factorize $\pi_\phi(\mathbf{y} | \tilde{\mathbf{x}})$ into token-level probabilities, distribute the same rollout-level
762 advantage $A_{m,n}$ to all tokens, and average over the sequence length $L_{m,n}$. Equation (23)
763 should therefore be read as the standard token-level approximation to this sequence-level
764 surrogate, not as an exact decomposition of the clipped joint ratio.

765 Conversely, when updating the Allocator (θ), we fix the MLLM ($\pi_\phi = \pi_{\phi_{\text{old}}}$) and use the
766 aggregated advantage $A_m^{\text{CAPO}} = \frac{1}{N} \sum_n A_{m,n}$. Because the Allocator’s output distribution
767 factorizes conditionally across frames (Eq. 25), its score function decomposes additively:

$$\nabla_\theta \log \pi_\theta(\mathbf{s}^{(m)} | \mathbf{x}) = \sum_{t=1}^T \nabla_\theta \log \text{Beta}(a_t^{(m)}; \alpha_t, \beta_t). \quad (31)$$

768 This additive log-probability structure supports low-variance frame-level credit assignment.
769 Nevertheless, Eq. (11) remains a practical approximation to a trajectory-level clipped objec-
770 tive: conditional independence justifies decomposition of $\log \pi_\theta$, but not exact factorization
771 of the nonlinear clipping term. We use the per-frame surrogate because it is substantially
772 more stable in large-scale training.

773 E Supplementary Experiments and Analysis

774 This section first analyzes the learned allocation policy, then studies the two key ablation
775 axes, and finally reports representative qualitative cases and a boundary-case transfer test
776 beyond video. Unless otherwise noted, all plots use Qwen2.5-VL-7B with 32 uniformly
777 sampled frames.

778 E.1 Behavioral Analysis of the Learned Policy

779 E.1.1 Global Allocation Statistics

780 The global distribution of predicted scales reveals the fundamental mechanism of Re-
781 sAdapt’s efficiency. As shown in Figure 8, rather than distributing the resolution budget

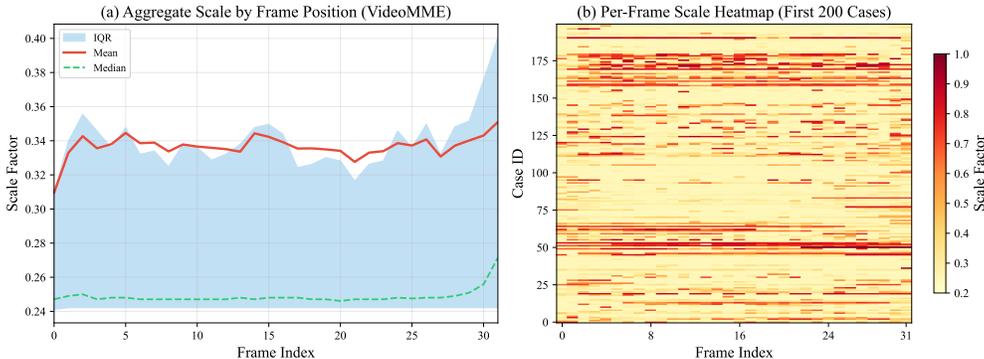


Figure 8: Global distribution of predicted scales. Across the entire dataset, the allocation policy exhibits a highly skewed distribution. The vast majority of frames are heavily compressed (near the minimum scale s_{\min}), reserving the spatial token budget for a small subset of critical, information-dense frames. This global sparsity confirms that ResAdapt achieves efficiency by aggressively reducing redundancy rather than uniformly lowering resolution.

782 evenly, the policy learns a highly skewed allocation strategy. Most frames are downscaled
 783 close to the minimum allowed bound ($s_{\min} = 0.2$), effectively discarding spatial redundancy
 784 in static or uninformative video segments. This massive compression on the majority of
 785 the video creates the budget headroom needed to process the few critical, evidence-bearing
 786 frames at much higher resolutions. This confirms that the learned behavior is genuinely
 787 selective, matching the localized scale spikes observed in the individual temporal profiles
 788 (Figure 4).

789 E.1.2 Benchmark-Level Budget Allocation

790 Figure 6 shows a clear benchmark-level ordering even though the policy never observes
 791 benchmark labels during training. Averaged across datasets, reasoning-oriented tasks use
 792 slightly higher mean scales than perception-oriented ones (0.435 vs. 0.417), with MMMU-
 793 Adaptation at the high end and VideoMME at the low end. The pattern is consistent with the
 794 main claim of the paper: the policy is not enforcing a fixed compression rule, but adapting
 795 its operating point to the expected visual difficulty of the task family.

796 E.1.3 Long-Context and Semantic Structure

797 Figure 7 is consistent with the long-context gains in the main paper. From short to long
 798 clips, the mean scale drops (0.342 \rightarrow 0.336 \rightarrow 0.332), but the within-video diversity rises (0.085 \rightarrow
 799 \sim 0.095). In other words, the policy does not merely compress longer videos more; it also
 800 becomes more selective inside them, which is exactly the regime where uniform resizing is
 801 least satisfactory.

802 Figure 9 refines the same story within a single benchmark. The policy spends the most
 803 budget on *Sports Competition* and the least on *Artistic Performance*, suggesting that even
 804 within VideoMME it distinguishes categories that are dense and spatially demanding from
 805 those that are visually simpler. This complements the main benchmark tables: the appendix
 806 focuses on *why* retained budgets differ, while the main text already reports the exact realized
 807 retention ratios.

808 E.1.4 Selectivity and Success

809 We next ask whether successful samples allocate budget *more selectively* within a clip. We
 810 quantify frame-level selectivity with the Gini coefficient of the predicted scales. High Gini
 811 means the policy concentrates budget on a small subset of frames; low Gini means the
 812 allocation is nearly uniform.

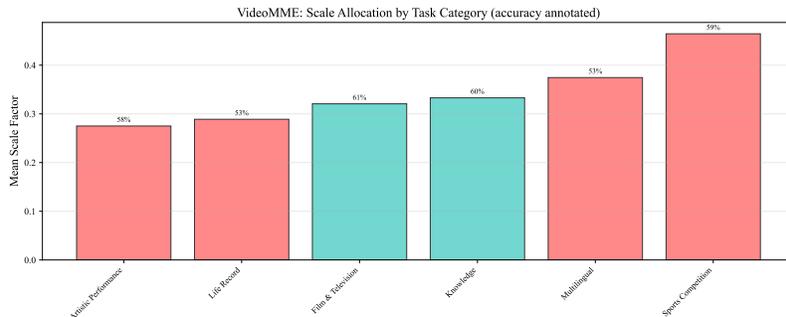


Figure 9: Scale allocation by VideoMME task category. Mean \bar{s} varies substantially across categories, with larger budgets assigned to categories that contain crowded motion or finer local evidence. Accuracy annotations show that allocation is not a trivial proxy for which category is easiest.

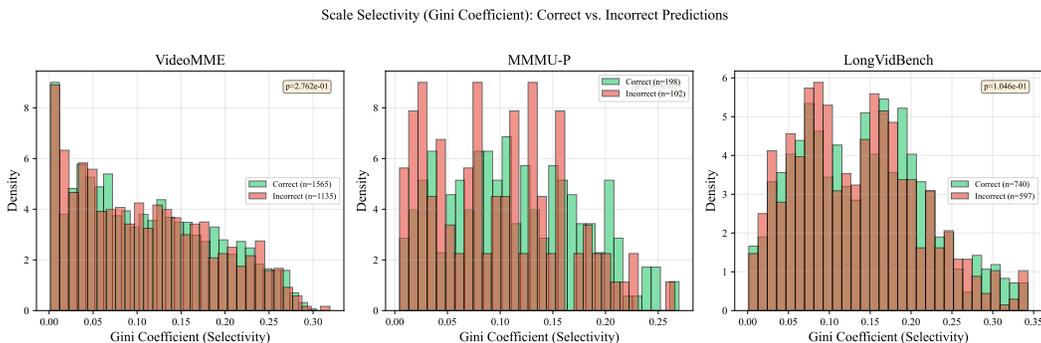


Figure 10: Selectivity versus prediction correctness on three representative benchmarks. Per-video Gini coefficients of the frame-level scales. Correct predictions tend to have higher Gini than incorrect ones, linking success to sharper concentration of resolution rather than merely larger average budgets.

813 Figure 10 shows that correct predictions consistently lie in the more selective regime, with
 814 the clearest separation on MMMU-P. This sharpens the mechanism claim of the appendix:
 815 success is associated not merely with keeping more pixels overall, but with concentrating
 816 them onto the frames that matter.

817 **Robustness and failure modes.** A final question is whether adaptive compression preserves
 818 existing correct answers or merely swaps one error pattern for another.

819 Figure 11 provides the right robustness interpretation for aggressive compression. Prediction
 820 stability remains high overall (about 89% of originally correct samples stay correct in the
 821 aggregate summary), so the policy is not helping only by randomly perturbing the answer
 822 distribution. However, error correction and error introduction are close enough that the
 823 effect should be read as *selective redistribution*: the policy repairs some failures, but it can
 824 also lose fine-grained evidence, especially when the decisive cue is brief or visually simple.

825 E.2 Ablation Studies

826 E.2.1 Temporal Similarity Ablation

827 We provide two complementary views of the temporal-similarity ablation: a cross-
 828 benchmark summary showing that the effect generalizes, and a single-benchmark diagnostic
 829 panel showing exactly how the allocation pattern changes.

830 Figure 12 makes the role of \mathcal{L}_{sim} unusually clear. Without it, the policy collapses to near-
 831 uniform scales on every benchmark ($\sigma < 0.003$); with it, the same model family recovers
 832 substantial within-video variation, with $4 \times -693 \times$ larger diversity depending on the bench-

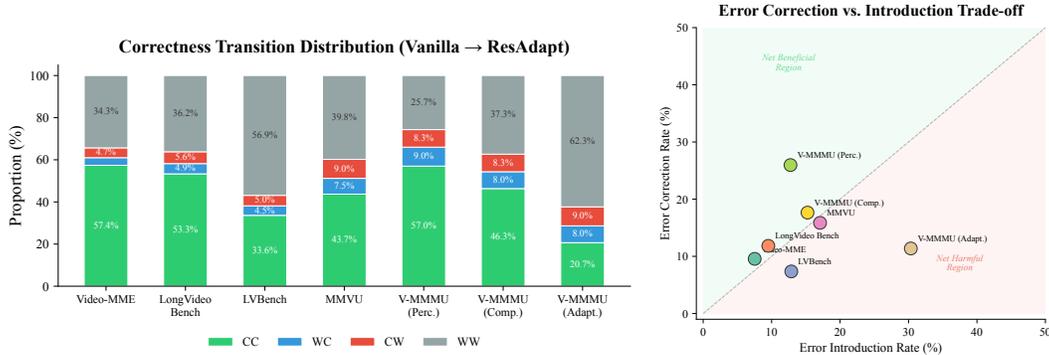


Figure 11: Sample-level robustness at 25% retention. Most originally correct predictions remain correct, but corrected and newly introduced errors are of comparable magnitude. Adaptive allocation is therefore selective rather than lossless.

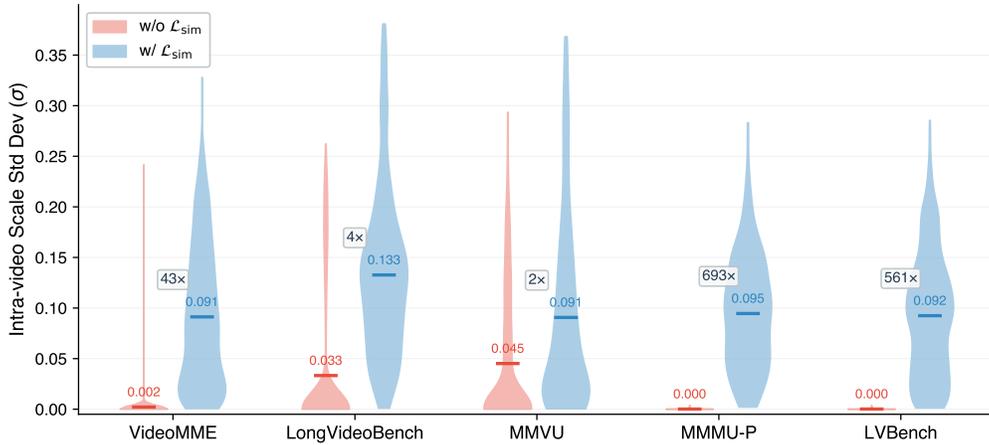


Figure 12: Cross-benchmark scale diversity with and without \mathcal{L}_{sim} . Per-video scale standard deviation σ across five benchmarks. Without the regularizer, diversity collapses toward zero; adding \mathcal{L}_{sim} restores broad within-video variation on every benchmark.

833 mark. CAPO therefore controls *where* the global budget should sit, whereas \mathcal{L}_{sim} prevents
 834 the trivial fixed-scale solution.

835 **Quantitative confirmation.** Figure 13 shows that this is not an artifact of any single statistic.
 836 The regularizer changes the global histogram, the per-video range, the frame-to-frame
 837 variation, and the Gini coefficient in the same direction, confirming that the benefit is
 838 structural rather than metric-specific.

839 **E.2.2 Reward Design Ablation**

840 We next examine whether different reward designs preserve a non-degenerate adaptive
 841 regime during training. All plots use EMA smoothing to suppress per-step noise; raw values
 842 remain visible as translucent traces.

843 **Per-sample scale adaptivity.** Figure 14 complements Figure 5 by measuring the per-sample
 844 scale range $s_{max} - s_{min}$ rather than the mean. CAPO preserves non-trivial adaptivity on
 845 validation, whereas direct cost collapses to the lower boundary and cost-free optimization
 846 drifts toward a nearly uniform high-scale policy.

847 **Convergence and stability.**

848 Figure 16 explains why CAPO works and the simpler baselines do not. The CAPO variants
 849 converge to stable interior solutions, whereas accuracy-only training saturates near s_{max} and

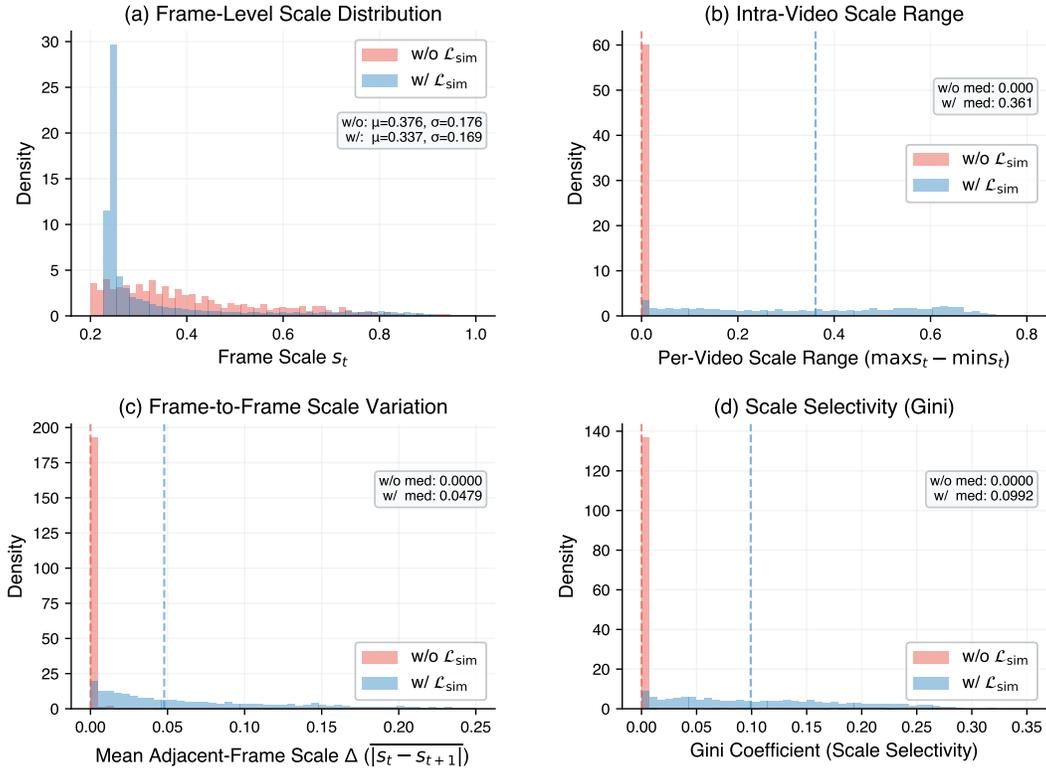


Figure 13: Four diagnostics of the \mathcal{L}_{sim} ablation on VideoMME. With the regularizer, the frame-scale histogram becomes bimodal, the per-video range expands, adjacent-frame variation increases, and the Gini coefficient rises. The policy moves from near-uniform allocation to a genuinely selective regime.

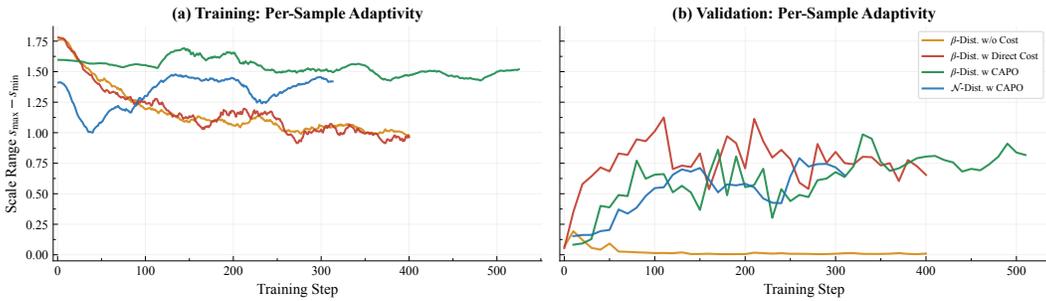


Figure 14: Per-sample scale adaptivity under different reward designs. Scale range $s_{max} - s_{min}$ over training on (a) training and (b) validation splits. CAPO keeps a non-trivial adaptive range, whereas direct cost collapses and cost-free training saturates.

850 direct cost collapses to s_{min} . This is consistent with CAPO’s intended role: balancing task
 851 reward and budget pressure without falling into either trivial boundary solution. The key
 852 result is therefore not merely convergence, but convergence to a non-degenerate operating
 853 point where content-adaptive allocation is still available.

854 **E.3 Additional Ablation Studies**

855 **E.4 Qualitative Case Studies**

856 We present four representative case studies that complement the aggregate analysis above:
 857 two task-contrast examples from Video-MMMU, one evidence-localization success from
 858 VideoMME, and one failure case. Each visualization (Figures 17–20) renders 32 uniformly

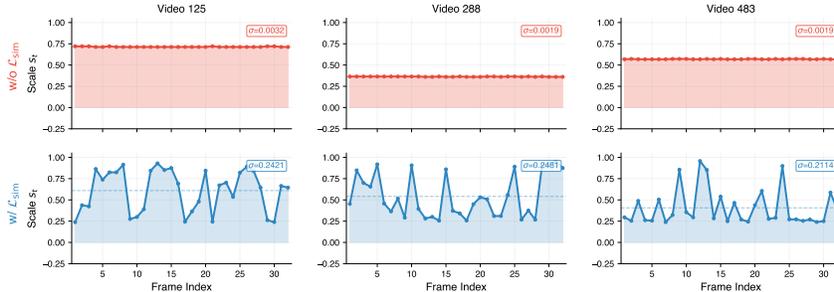


Figure 15: \mathcal{L}_{sim} ablation: per-frame scale profiles. Without temporal-similarity regularization, the Allocator approaches near-uniform scaling; with it, the policy concentrates resolution on selected frames and suppresses redundant neighbors.

Table 4: Distribution family ablation for CAPO. The two variants follow the same training protocol.

Variant	\bar{s}	VideoMME	LongVideoBench	MMVU	VideoMMMU			LVBench
					Per.	Comp.	Adap.	
β -CAPO	0.54	60.3	58.2	51.2	65.0	54.3	28.7	37.6
\mathcal{N} -CAPO	0.60	61.0	57.4	51.8	66.0	50.0	30.3	37.2

859 sampled frames at their assigned scale inside a fixed grid; warmer borders indicate larger
 860 predicted scales.

861 **Task-dependent operating regimes.** Figures 17 and 18 contrast two Video-MMMU tasks
 862 from a visually similar educational domain that nevertheless demand very different allocations.
 863 In the comprehension example, the relevant evidence is concentrated in a small set of
 864 diagram-bearing slides, so the policy adopts a sparse operating regime and suppresses the
 865 explicitly irrelevant quiz frame. In the adaptation example, the downstream reasoning depends
 866 on reading a dense numeric table, so the same policy shifts to a much higher-budget
 867 regime and preserves high fidelity much more broadly. The contrast shows that the policy
 868 responds to what the task will require, not just to generic visual clutter.

869 **Evidence localization and failure.** The VideoMME success case in Figure 19 shows a
 870 more local version of the same phenomenon: the answer depends on short text overlays
 871 embedded in otherwise repetitive footage, and the policy magnifies only those evidence-
 872 bearing moments. Figure 20 shows the failure mode that remains. The decisive cue is
 873 temporally brief and visually simple, so the policy enlarges a nearby frame but compresses
 874 the frame that actually contains the fork. This diagnosis matches the quantitative robustness
 875 analysis: ResAdapt is strong at concentrating budget, but still vulnerable when the decisive
 876 evidence is both subtle and short-lived.

877 **Summary.** Together, these case studies support the same three conclusions as the quantita-
 878 tive appendix: the policy changes its operating regime with the task, concentrates fidelity
 879 on evidence-bearing frames, and fails in interpretable ways when subtle cues are missed.
 880 The qualitative examples therefore reinforce the claim that ResAdapt learns a meaningful
 881 input-allocation strategy rather than a fixed compression heuristic.

882 **E.5 Failure Modes Analysis**

883 Adaptive allocation does not act as a lossless compression layer. In practice, ResAdapt
 884 usually preserves many originally correct predictions, but it can still miss decisive evidence,
 885 especially when the relevant cue is visually simple and appears only briefly. Because the
 886 policy is open-loop, it cannot revise allocations after reasoning begins or recover evidence
 887 that was undersampled in the initial pass. We therefore interpret its gains as selective redis-
 888 tribution of visual budget rather than as guaranteed preservation of all useful information.

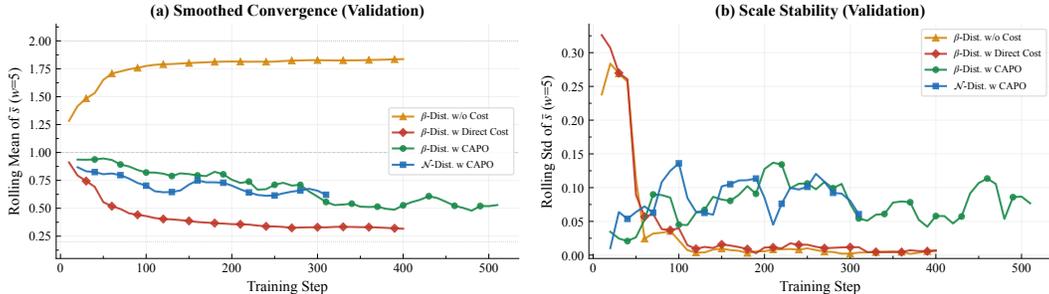


Figure 16: Validation-time convergence under different reward designs. CAPO variants converge to stable intermediate operating points, while cost-free training saturates at the upper boundary and direct cost collapses to the lower boundary. Stability alone is not sufficient; the key is where the policy stabilizes.

Table 5: Operator generalization. Zero-shot transfer of ResAdapt scores to frame selection. Combining top-K selection with adaptive resizing from 128 candidate frames outperforms uniform sampling baselines at lower token budgets.

Method	VideoMME	LongVideoBench	LVBench	MMVU
<i>Budget: 8 frames</i>				
Vanilla	54.0	53.9	33.3	48.9
Top-8 Select	52.2	51.1	32.0	49.2
<i>Budget: 16 frames</i>				
Vanilla	58.9	56.0	36.1	50.9
Threshold Select	58.0	57.4	36.4	51.0
<i>Avg. Budget (Retention Ratio)</i>	12.2f (9.5%)	23.2f (18.1%)	16.7f (13.0%)	17.2f (13.4%)
Top-32 Select + Resize	60.6	57.2	38.9	50.2
<i>Avg. Budget (Retention Ratio)</i>	11.7f (9.1%)	16.9f (13.2%)	13.7f (10.7%)	14.1f (11.0%)
<i>Budget: 32 frames</i>				
Vanilla	62.3	58.7	39.5	52.0
Top-32 Select	59.7	55.7	37.0	51.2
Top-64 Select + Resize	62.5	58.4	40.0	52.3
<i>Avg. Budget (Retention Ratio)</i>	23.8f (18.6%)	36.2f (28.3%)	24.1f (18.8%)	32.5f (25.4%)

889 **E.6 Temporal Grounding Full Results**

890 Table 6 provides the full temporal grounding evaluation results that were omitted from the
 891 main text for space.

892 **E.7 Boundary-Case Transfer Beyond Video**

893 The paper’s main claims target video QA and temporal grounding, so we place image
 894 transfer at the end of the appendix as a boundary-case analysis rather than as supporting
 895 evidence for the main contribution. Table 7 is still informative: the learned video policy
 896 sometimes identifies image inputs that warrant additional fidelity, as in ChartQA, but it does
 897 not yet yield reliable efficiency-preserving transfer on text-dense image tasks. The result is
 898 therefore best read as scope clarification. It suggests that input-side adaptation is broader
 899 than the resize-on-video setting studied here, while also showing that a video-trained policy
 900 should not be assumed to transfer cleanly to static images.

Table 6: Evaluation Results on Temporal Grounding Benchmarks. Grounding is much more compression-sensitive.

Backbone	Method	Retention Ratio R	Reasoning	Temporal Grounding Benchmark										
				Charades-STA				ActivityNet				NExT-GQA		
				0.3	0.5	0.7	mIoU	0.3	0.5	0.7	mIoU	Acc	mIoU	
<i>32 Frames</i>														
Qwen2.5-VL-7B	Vanilla	100%	✗	71.0	51.4	26.0	47.3	30.4	18.0	8.9	22.6	78.9	28.0	
	Random Drop	25.0%	✗	39.4	23.2	11.0	28.7	15.2	8.1	3.7	11.7	77.5	16.6	
	ToMe (Bolya et al., 2022)	25.0%	✗	39.5	23.9	11.4	26.0	16.0	8.4	4.0	12.1	77.8	16.3	
	FlashVid (Fan et al., 2026)	31.3%	✗	40.7	24.2	11.3	26.6	15.8	8.4	3.8	12.0	78.1	16.5	
	FixedScale	25.0%	✗	36.7	24.7	12.3	24.9	18.6	9.4	4.3	14.1	77.7	12.3	
	ResAdapt (Ours)	16.2%	✗	53.8	34.8	17.0	35.6	19.8	10.8	5.2	15.3	76.6	23.2	
	Random Drop	10.0%	✗	36.9	23.2	11.6	24.6	14.3	7.5	3.6	11.1	76.3	15.4	
	ToMe (Bolya et al., 2022)	10.0%	✗	41.3	26.9	14.1	27.4	16.0	8.4	4.0	12.2	77.3	15.7	
	FlashVid (Fan et al., 2026)	12.6%	✗	38.2	22.9	11.1	25.1	15.4	8.1	3.7	11.8	77.4	16.1	
	FixedScale	12.3%	✗	48.0	31.5	15.4	32.0	17.5	8.9	4.0	13.3	76.1	13.7	
	FixedScale	6.3%	✗	39.9	26.8	13.3	26.7	15.2	8.1	3.9	11.9	74.1	15.4	
	ResAdapt (Ours)	6.8%	✗	41.0	27.8	14.0	27.2	16.3	8.5	3.9	12.5	74.3	20.4	
	VideoAuto-RL (Liu et al., 2026)	100%	✓	60.0	48.3	27.2	41.5	50.8	34.1	17.4	34.4	73.6	33.8	
	+ ResAdapt (Ours)	6.8%	✓	43.5	30.1	15.8	30.0	35.4	21.5	10.0	24.4	74.7	24.7	
	<i>128 Frames</i>													
	Qwen2.5-VL-72B	Vanilla	100%	✗	77.5	60.3	34.1	52.8	47.9	30.9	17.5	34.4	79.8	29.9
Random Drop		25.0%	✗	32.3	19.6	7.9	20.7	26.7	13.9	6.3	18.8	80.3	10.7	
ToMe (Bolya et al., 2022)		25.0%	✗	32.4	19.8	7.9	20.7	27.2	14.4	6.4	19.1	80.3	10.9	
ResAdapt (Ours)		16.1%	✗	63.5	43.6	21.3	42.0	33.1	19.3	10.2	24.3	78.1	27.2	
Random Drop		10.0%	✗	37.8	23.8	11.2	24.7	25.8	12.0	5.3	17.0	79.4	12.8	
ToMe (Bolya et al., 2022)		10.0%	✗	27.9	16.2	7.3	17.9	22.9	11.8	5.5	16.4	79.1	11.1	
FlashVid (Fan et al., 2026)		12.3%	✗	34.7	22.3	10.5	22.7	25.0	13.8	5.9	18.3	77.9	11.3	
FixedScale		6.3%	✗	42.6	28.4	14.3	28.3	22.8	12.8	5.7	17.1	75.7	12.9	
ResAdapt (Ours)		6.8%	✗	43.5	29.8	15.0	28.9	23.5	12.9	6.1	17.2	76.2	23.9	
VideoAuto-RL (Liu et al., 2026)		100%	✓	40.3	33.7	22.1	28.9	49.4	34.3	18.5	33.5	68.0	31.0	
+ ResAdapt (Ours)		16.1%	✓	72.8	53.0	27.5	49.1	65.8	44.9	23.8	44.7	79.3	35.3	
+ ResAdapt (Ours)		6.8%	✓	50.1	33.2	16.6	34.2	53.4	34.0	16.4	35.7	76.6	29.4	
<i>32 Frames</i>														
Qwen3-VL-8B		Vanilla	100%	✗	73.0	49.0	21.4	46.4	44.6	28.3	15.5	31.8	78.7	34.2
		Random Drop	25.0%	✗	16.2	8.6	3.8	12.1	12.4	6.7	3.2	10.0	77.2	15.6
		ToMe (Bolya et al., 2022)	25.0%	✗	68.7	42.1	17.6	43.1	45.9	28.8	15.6	32.6	77.1	31.7
	FlashVid (Fan et al., 2026)	31.3%	✗	72.9	52.3	25.1	47.7	51.9	33.4	19.0	36.8	77.8	33.9	
	ResAdapt (Ours)	16.2%	✗	64.4	37.3	16.3	39.9	40.0	24.4	13.0	28.5	75.1	30.2	
	Random Drop	10.0%	✗	4.1	1.8	0.7	4.4	4.7	2.4	1.0	5.0	74.3	11.3	
	ToMe (Bolya et al., 2022)	10.0%	✗	67.6	39.3	16.6	41.8	46.3	31.0	19.2	34.1	79.2	34.0	
	FlashVid (Fan et al., 2026)	12.6%	✗	68.8	46.9	22.9	44.6	49.9	31.5	17.4	35.2	75.6	31.8	
	FixedScale	12.3%	✗	61.3	34.3	14.6	37.9	39.6	24.2	13.1	28.4	74.2	29.9	
	FixedScale	6.3%	✗	52.7	28.2	11.3	33.2	37.0	22.3	12.0	27.0	71.5	28.0	
	ResAdapt (Ours)	6.8%	✗	53.6	29.0	11.8	33.6	37.5	22.5	12.3	27.2	71.8	28.2	
	<i>128 Frames</i>													
	Qwen3-VL-88B	Vanilla	100%	✗	72.8	46.0	20.1	45.6	45.8	31.1	19.2	33.9	81.1	36.6
		Random Drop	25.0%	✗	41.6	25.2	10.6	27.4	36.1	21.1	12.7	26.3	79.3	22.4
		ResAdapt (Ours)	16.1%	✗	64.4	37.0	15.9	39.8	40.6	26.7	15.7	30.0	76.8	33.3
		Random Drop	10.0%	✗	32.6	19.0	7.8	21.9	33.5	18.6	11.5	24.8	76.9	19.9
ToMe (Bolya et al., 2022)		10.0%	✗	61.6	33.8	13.3	38.1	42.4	27.6	16.6	31.4	77.4	31.5	
FixedScale		12.3%	✗	61.7	34.9	14.7	38.1	39.9	26.2	15.3	29.5	75.4	32.6	
FixedScale		6.3%	✗	53.7	28.2	11.8	33.6	37.9	24.3	14.3	28.1	73.0	39.1	
ResAdapt (Ours)		6.8%	✗	54.3	28.0	11.7	33.7	38.3	24.5	14.4	28.4	73.2	43.9	

Q: Evaluate five statements about Urban Geography City Models (concentric zone, Hoyt sector, multiple nuclei, galactic, Latin American); identify which are correct. *Please ignore the Quiz question in last frame of the video.*

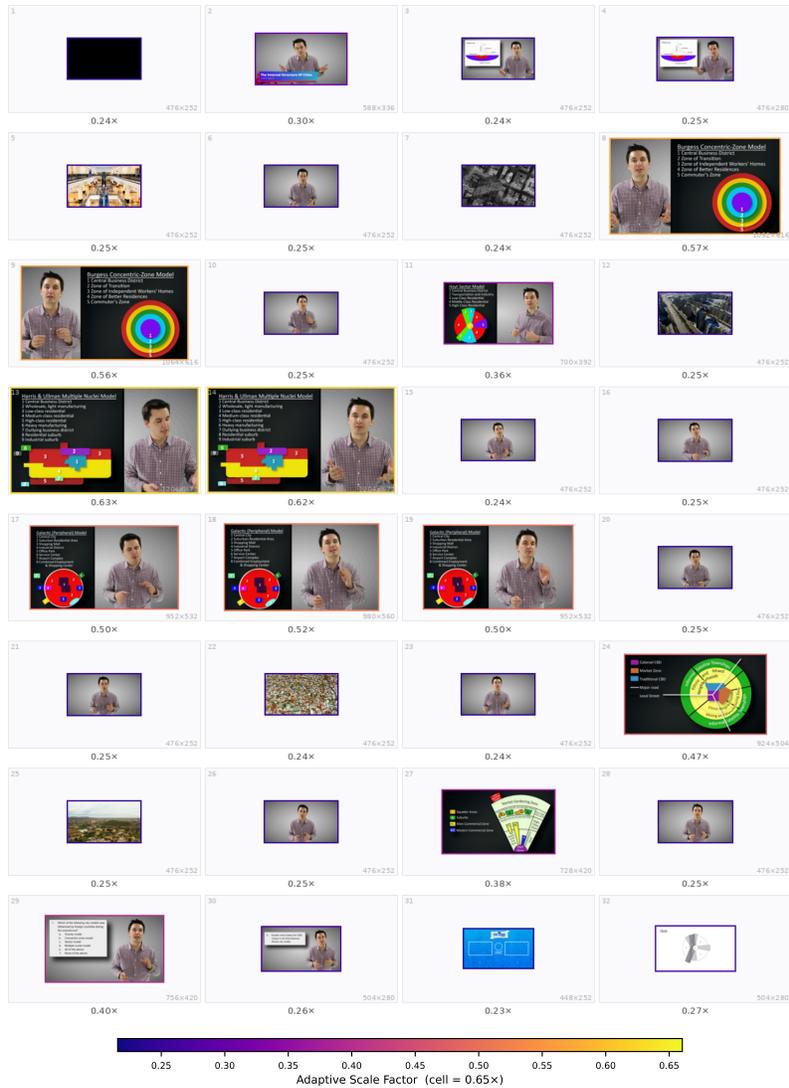


Figure 17: Case 1: Video-MMMU Comprehension (Hu et al., 2025) (Vanilla \times \rightarrow ResAdapt \checkmark). The policy concentrates resolution on diagram-bearing slide frames, compresses lecturer-only frames, and suppresses the final quiz frame that the prompt explicitly marks as irrelevant.

Q: Watch and learn the video content. Then apply what you learned to answer: Table 11.47 provides a survey of the youngest online entrepreneurs (ages 17–30) whose net worth \geq \$1M. We want to know whether ages and net worth are independent. χ^2 test statistic = -----

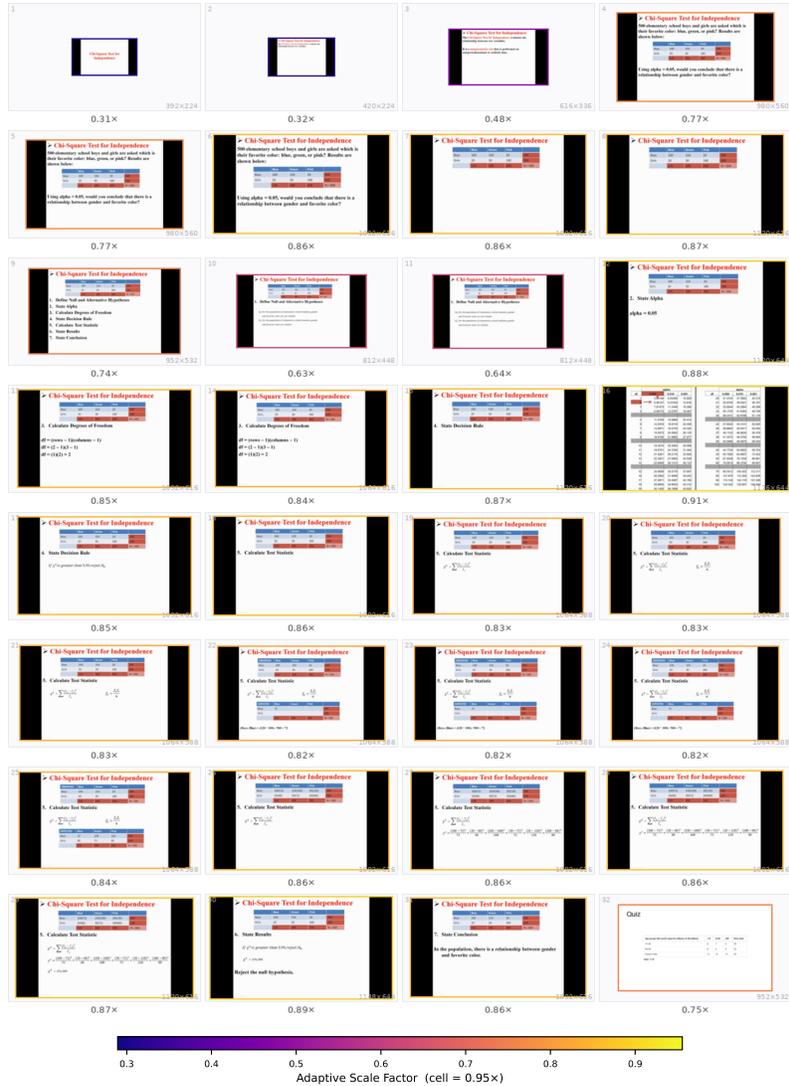


Figure 18: Case 2: Video-MMMU Adaptation (Hu et al., 2025) (Vanilla \times \rightarrow ResAdapt \checkmark). When the answer depends on reading a numeric table and performing a χ^2 computation, the policy keeps a much higher global budget and strongly upscales the table-bearing frames.

Q: When is the zodiacal light visible from the video? (A) Mar. 19, (B) Mar. 24, (C) Mar. 25, (D) Mar. 29.

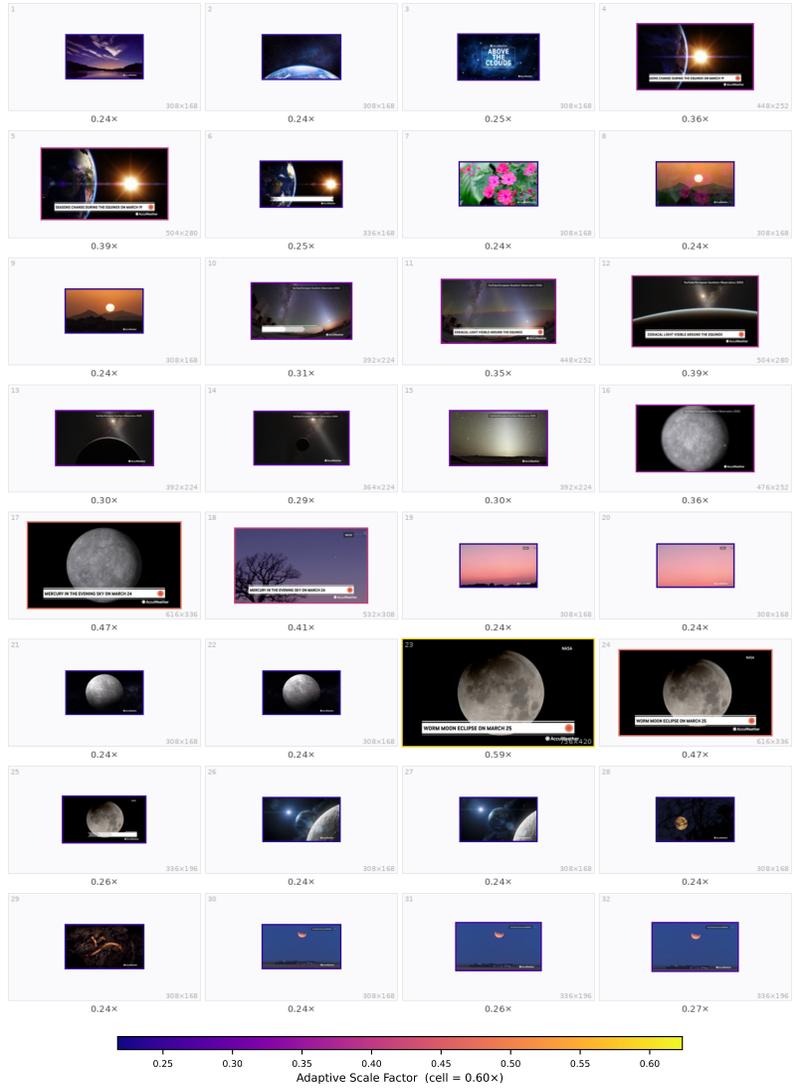


Figure 19: Case 3: VideoMME (Fu et al., 2025a) (Vanilla \times \rightarrow ResAdapt \checkmark). Frames containing the decisive date overlays are enlarged, while the largely homogeneous sky footage is compressed. The policy spends budget on answer-bearing evidence rather than on the surrounding context.

Q: Which item does the man throw into the trash at the beginning of the video? (A) A fork, (B) A pair of chopsticks, (C) A box of noodles, (D) A spoon.

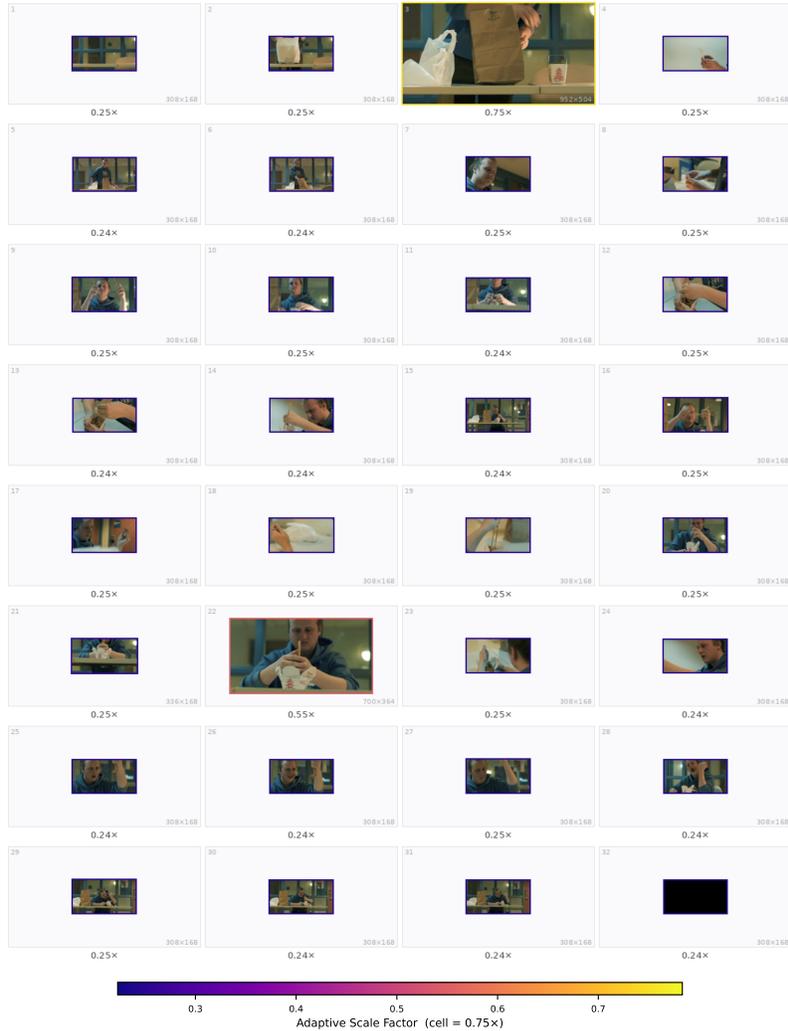


Figure 20: Case 4: VideoMME (Fu et al., 2025a) (Vanilla ✓ → ResAdapt ✗; failure case). A nearby frame is enlarged, but the actual fork-bearing frame is compressed. The decisive fine detail is therefore lost at exactly the wrong moment.

Table 7: Exploratory zero-shot transfer to image benchmarks. Parenthetical values denote per-task retention ratio R , and ResAdapt-RL additionally fine-tunes the MLLM via RL.

Model	MathVista testmini	MMMU val	OCRBench	ChartQA	AI2D	TextVQA val
Qwen2.5-VL-7B	49.1(100%)	50.9(100%)	84.2(100%)	83.9(100%)	82.5(100%)	82.9(100%)
Random Drop	44.8(50%)	49.0(50%)	74.8(50%)	71.6(50%)	80.3(50%)	78.1(50%)
ToMe (Bolya et al., 2022)	46.2(50%)	49.6(50%)	79.3(50%)	78.1(50%)	81.9(50%)	81.2(50%)
VisionZip (Yang et al., 2025c)	47.2(50%)	48.6(50%)	79.6(50%)	77.9(50%)	81.9(50%)	81.3(50%)
ResAdapt (Qwen2.5-VL-7B)	45.5(42%)	51.0(29%)	80.0(64%)	85.9(105%)	81.4(41%)	69.6(30%)
ResAdapt-RL (Qwen2.5-VL-7B)	46.7(42%)	50.9(29%)	80.8(64%)	86.6(105%)	81.1(41%)	70.1(30%)
Qwen3-VL-8B	56.1(100%)	53.4(100%)	85.0(100%)	84.0(100%)	83.5(100%)	82.1(100%)
Random Drop	47.3(50%)	48.7(50%)	62.9(50%)	70.2(50%)	79.7(50%)	76.6(50%)
VisionZip (Yang et al., 2025c)	47.8(50%)	50.3(50%)	70.5(50%)	75.0(50%)	80.5(50%)	79.3(50%)
ToMe (Bolya et al., 2022)	49.6(50%)	50.6(50%)	70.3(50%)	75.2(50%)	80.5(50%)	79.4(50%)
ResAdapt (Qwen3-VL-8B)	52.5(42%)	50.9(29%)	82.7(64%)	83.2(105%)	81.2(41%)	67.8(30%)