

Test-Time Horizon Scaling in Video LLMs via Adaptive Temporal Memory Compression

Anonymous CVPR submission

Paper ID ****

Abstract

001 Video Large Language Models (VidLLMs) are limited by
002 fixed context windows, making long-video reasoning de-
003 pendent on aggressive frame subsampling. We explore a
004 lightweight inference-time strategy for extending the effec-
005 tive temporal horizon of pretrained VidLLMs without archi-
006 tectural changes or retraining. We propose Adaptive Tem-
007 poral Memory Compression (ATMC), which scores frames
008 using motion magnitude, vision encoder attention statis-
009 tics, and query-conditioned similarity. High-importance
010 frames are retained while less salient segments are tem-
011 porally abstracted to fit within the original context bud-
012 get. Experiments on Ego4D-NLQ and ActivityNet-QA show
013 that ATMC consistently improves over uniform sampling
014 while reducing peak GPU memory usage. We analyze accu-
015 racy–memory–latency trade-offs and identify regimes where
016 adaptive compression is most beneficial. Our results sug-
017 gest that simple test-time adaptations can enhance long-
018 video reasoning in existing VidLLMs.

019 1. Introduction

020 Video Large Language Models (VidLLMs) have signifi-
021 cantly advanced multimodal reasoning by enabling joint
022 understanding of visual and linguistic information. De-
023 spite this progress, their practical effectiveness on long-
024 form videos remains constrained by fixed Transformer con-
025 text windows [11]. In practice, long videos are han-
026 dled through aggressive frame subsampling, which can
027 discard temporally important events and degrade perfor-
028 mance on tasks requiring extended reasoning [12]. Exist-
029 ing approaches address long-video understanding through
030 architectural modifications such as hierarchical encoders
031 or sliding-window processing, or through additional large-
032 scale training. While effective, these solutions increase
033 system complexity, computational cost, and deployment
034 barriers. This motivates a complementary question: *Can*
035 *pretrained VidLLMs better handle long videos through*

inference-time adaptation alone, without modifying model 036
parameters? To explore this direction, we introduce 037
Adaptive Temporal Memory Compression (ATMC), a 038
lightweight test-time framework that dynamically priori- 039
tizes video content within a fixed context budget. ATMC as- 040
signs importance scores to frames using motion magnitude, 041
attention entropy from the frozen vision encoder, and query- 042
conditioned semantic similarity. High-importance frames 043
are retained, while lower-saliency segments are abstracted 044
through temporal pooling, enabling more informative cov- 045
erage of long videos without architectural changes or re- 046
training. We focus on question answering tasks since these 047
require identification of temporally dispersed events, mak- 048
ing them ideal for evaluating long-horizon reasoning. 049

2. Related Work 050

Recent VidLLMs, such as Video-LLaMA [3], LLaVA- 051
Video [4], and VideoChat [5], extend LLMs to video by 052
integrating pretrained vision encoders. They typically pro- 053
cess a limited number of frames (16–64), limiting long- 054
range temporal reasoning. Prior solutions include hierarchi- 055
cal transformers [13] and sliding-window strategies, which 056
expand temporal coverage but increase inference cost or 057
require architectural changes. Token pruning and merg- 058
ing [9] improve Transformer efficiency but mainly target 059
spatial/language tokens and often require training. In con- 060
trast, we explore adaptive temporal compression as a purely 061
inference-time mechanism to extend the reasoning horizon 062
of pretrained VidLLMs. ATMC combines motion, atten- 063
tion, and query relevance to prioritize frames and abstract 064
less salient segments without retraining. While prior works 065
on adaptive frame selection focus on summarization or ac- 066
tion recognition, they rarely address query-conditioned rea- 067
soning in VidLLMs. We evaluate ATMC on two long-video 068
QA benchmarks under controlled settings: Ego4D-NLQ [1] 069
(505 videos, 8.2 min avg, 1 FPS) and ActivityNet-QA [2] 070
(812 videos, 5.7 min avg, 1 FPS), ensuring sufficient tem- 071
poral granularity and computational feasibility. 072

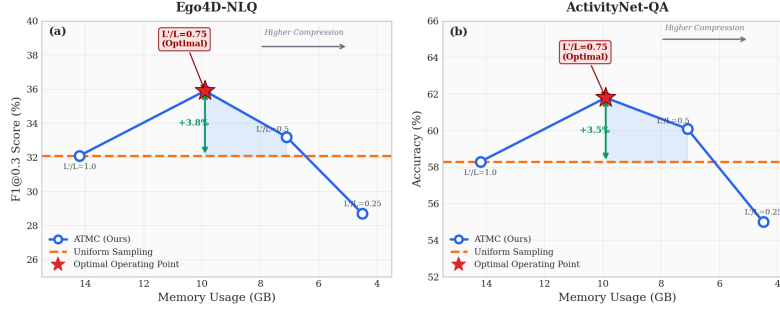


Figure 1: Accuracy-memory trade-off for different compression ratios (L'/L). ATMC consistently outperforms uniform sampling at equivalent memory budgets.

Figure 1. Accuracy-memory trade-off for different compression ratios L'/L . ATMC (blue line with circles) consistently outperforms uniform sampling (orange dashed line) at equivalent memory budgets. The default operating point ($L'/L = 0.75$, starred) balances accuracy and efficiency. Aggressive compression ($L'/L < 0.5$) degrades accuracy despite memory savings. Left: Ego4D-NLQ F1@0.3. Right: ActivityNet-QA exact match accuracy.

073 3. Methodology

074 3.1. Problem Formulation

075 Let \mathcal{M} be a pretrained VidLLM with context length L .
 076 Given a video $\mathcal{V} = \{f_1, \dots, f_T\}$ with $T \gg L$ and a query
 077 Q , the goal is to generate an answer A . Since \mathcal{M} can-
 078 not process all T frames, we define a compression function
 079 $\mathcal{C} : \mathcal{V} \rightarrow \tilde{\mathcal{V}}$ producing a sequence $\tilde{\mathcal{V}}$ of length $L' \leq L$
 080 that preserves query-relevant and temporally informative con-
 081 tent. The final answer is $A = \mathcal{M}(\tilde{\mathcal{V}}, Q)$. Unlike prior
 082 methods, \mathcal{M} remains frozen; our approach acts as a prepro-
 083 cessing layer that selects and abstracts frames at inference
 084 time, requiring no retraining or architectural changes.

085 3.2. Adaptive Temporal Memory Compression 086 (ATMC)

087 ATMC assigns each frame f_t an importance score s_t that re-
 088 flects its contribution to answering query Q . The score inte-
 089 grates three complementary signals capturing low-level mo-
 090 tion dynamics, model-internal informativeness, and query-
 091 conditioned semantic relevance:

$$092 \quad s_t = \alpha M_t + \beta A_t + \gamma S_t(Q), \quad (1)$$

093 where $\alpha, \beta, \gamma \in [0, 1]$ and $\alpha + \beta + \gamma = 1$.

094 **Motion Magnitude.** The motion term M_t captures tem-
 095 poral variation between consecutive frames. Frames are re-
 096 sized to 32×32 pixels to reduce computational cost, and
 097 the normalized ℓ_1 difference is computed:

$$098 \quad M_t = \frac{\sum_{i=1}^{32} \sum_{j=1}^{32} |f_t(i, j) - f_{t-1}(i, j)|}{32 \times 32 \times 255}. \quad (2)$$

099 This produces a value in $[0, 1]$ representing relative tem-
 100 poral change. The lightweight design ensures negligible
 101 computational overhead.

Attention Entropy. To estimate frame informativeness 102
 from the model’s internal perspective, we compute entropy 103
 over the attention distribution of the frozen vision encoder. 104
 Let $\mathbf{a}_t^{(h)}$ denote the attention distribution from the [CLS] to- 105
 ken for head h , and let $\bar{\mathbf{a}}_t$ be the average across H heads. 106
 Attention entropy is defined as: 107

$$108 \quad A_t = - \sum_{i=1}^N \bar{a}_t^{(i)} \log \bar{a}_t^{(i)}. \quad (3)$$

Higher entropy indicates dispersed attention across mul- 109
 tiple regions, suggesting richer visual content. This sig- 110
 nal leverages internal model structure without modifying 111
 parameters. This signal captures how broadly the model 112
 distributes attention across visual regions: higher entropy 113
 indicates that the frame contains diverse or complex visual 114
 content, suggesting greater informativeness for downstream 115
 reasoning. In contrast, low entropy implies that attention is 116
 concentrated on few regions, often corresponding to less in- 117
 formative or static frames. 118

Query Similarity. To incorporate query-conditioned 119
 saliency, we compute cosine similarity between CLIP em- 120
 beddings of frames and the query: 121

$$122 \quad S_t(Q) = \frac{\mathbf{v}_t^\top \mathbf{q}}{\|\mathbf{v}_t\| \|\mathbf{q}\|}, \quad \mathbf{v}_t = \text{CLIP}_{\text{image}}(f_t), \quad \mathbf{q} = \text{CLIP}_{\text{text}}(Q). \quad (4)$$

This term prioritizes frames semantically aligned with 123
 the question, enabling task-aware compression. 124

125 3.3. Frame Selection and Temporal Abstraction

126 After computing s_t for all T frames, frames are ranked in
 127 descending order of importance. We target a compressed
 128 length of $L' = 0.75L$ (e.g., $L' = 384$ when $L = 512$). The
 129 top $K = 0.6L'$ highest-scoring frames are retained as key 129

130 frames. The remaining frames are temporally sorted and
 131 partitioned into $C = L' - K$ contiguous segments of ap-
 132 proximately equal duration. For each segment R_j , we com-
 133 pute a mean-pooled representation of vision encoder fea-
 134 tures:

$$135 \quad \mathbf{t}_j = \frac{1}{|R_j|} \sum_{t \in R_j} \mathbf{H}_t. \quad (5)$$

136 The final compressed sequence $\tilde{\mathcal{V}}$ consists of all key
 137 frame features and segment-level pooled representations,
 138 reordered to preserve original temporal order before being
 139 passed to \mathcal{M} . The choice of retaining the top 60% frames
 140 and pooling the remaining 40% segments was empirically
 141 found to balance coverage of salient events with memory
 142 efficiency. This ratio provides sufficient representation of
 143 important frames while maintaining a compressed sequence
 144 that fits within the model’s context window.

145 3.4. Training and Hyperparameter Tuning

146 ATMC requires no gradient-based training. All feature ex-
 147 tractors (VidLLM vision encoder and CLIP-ViT-L/14) re-
 148 main frozen. Motion computation is deterministic, and the
 149 only tunable parameters are α, β, γ . These are selected
 150 via exhaustive grid search over $\{0, 0.1, \dots, 1.0\}$ subject
 151 to $\alpha + \beta + \gamma = 1$, yielding 66 combinations. Valida-
 152 tion is conducted on 50 randomly sampled videos from
 153 each dataset’s validation split, with no overlap with test
 154 data. The objective is to maximize question-answering ac-
 155 curacy after compression. Optimal weights differ slightly
 156 across datasets (e.g., Ego4D favors higher attention entropy
 157 weighting, whereas ActivityNet benefits more from query
 158 similarity), indicating mild dataset dependence but stable
 159 behavior. Cross-dataset transfer of tuned weights results in
 160 only a modest 1.5–2.3% performance drop. Hyperparam-
 161 eter search requires approximately two hours on a single
 162 NVIDIA A100 GPU per dataset. Once tuned, ATMC intro-
 163 duces no additional training cost.

164 3.5. Implementation Details and Computational 165 Overhead

166 Videos are decoded at 1 FPS using `ffmpeg` and resized
 167 to 224×224 before encoding. Frames are processed in
 168 batches of 32. Vision encoder outputs are cached to avoid
 169 recomputation during scoring and abstraction. The com-
 170 putational complexity is dominated by CLIP encoding and
 171 feature extraction for T frames. Motion computation costs
 172 $O(1024T)$ operations and is negligible. Frame ranking re-
 173 quires $O(T \log T)$ time. In practice, ATMC introduces a
 174 15–25% inference-time overhead relative to standard uni-
 175 form sampling, while remaining significantly cheaper than
 176 architectural scaling or retraining approaches. Memory
 177 overhead is minimal beyond standard frame buffers and
 178 cached features.

3.6. Baseline Protocol

All baselines are evaluated under identical preprocessing
 except for compression strategy. The primary baseline uses
 uniform frame sampling to match the compressed length
 L' . Additional comparisons include sliding-window and
 dense sampling within the same token budget. All mod-
 els remain frozen to isolate inference-time effects. Single-
 pass baselines (Uniform, Random, Motion-Only, Attention-
 Only, Query-Only) use the same number of visual tokens as
 ATMC, ensuring fair comparison. For sliding-window, each
 segment respects the L -token constraint, and predictions are
 aggregated without extra context. Temporal pooling is ver-
 ified to preserve projected token counts, preventing hidden
 inflation.

4. Experiments

We evaluate ATMC on two long-video question answer-
 ing benchmarks under a strictly controlled and reproducible
 setup. On Ego4D-NLQ [1], we use the validation split (505
 videos; average length 8.2 minutes; original 30 FPS) and
 extract frames at 1 FPS, which provides sufficient tempo-
 ral granularity for NLQ tasks while maintaining compu-
 tational feasibility. On ActivityNet-QA [2], we use the
 long-form QA validation subset (812 videos; average length
 5.7 minutes; original 30 FPS), also sampled at 1 FPS.
 Video lengths vary considerably, with Ego4D clips rang-
 ing from 2 to 15 minutes (120–900 frames at 1 FPS) and
 ActivityNet clips from 1 to 10 minutes (60–600 frames
 at 1 FPS), ensuring that ATMC is evaluated on diverse
 temporal scales. All videos are deduplicated to avoid
 overlap with any original model training data, and iden-
 tical preprocessing (decoding, resizing, frame sampling)
 is applied across ATMC and all baselines to ensure fair-
 ness. For Ego4D-NLQ, we report F1@0.3 for temporal
 localization, where a prediction is correct if the tempo-
 ral IoU between predicted and ground-truth intervals ex-
 ceeds 0.3; for ActivityNet-QA, we report exact match ac-
 curacy (case-insensitive, punctuation-agnostic). In addi-
 tion to task accuracy, we measure peak GPU memory use-
 age (via `torch.cuda.max_memory_allocated()`) and end-to-end latency per video, including preprocess-
 ing, averaged over three warm-up runs and five measured
 runs. We evaluate two pretrained VidLLMs without any
 fine-tuning: Video-LLaMA-2 (7B) (BLIP-2 ViT-g vision
 encoder, Vicuna-7B v1.5 language model) and LLaVA-
 Video (7B) (CLIP-ViT-L/14 vision encoder, Vicuna-7B
 v1.5), both with a 512-token context window. All model
 parameters remain frozen throughout. Baselines include
 Uniform Sampling (evenly selecting L frames), Random
 Sampling (without replacement), Motion-Only (top- L by
 M_t), Attention-Only (top- L by A_t), Query-Only (top- L
 by $S_t(Q)$), and a computationally intensive Sliding Win-

Table 1. Performance on Video-LLaMA-2. Best single-pass method bolded.

Method	Ego4D F1 (%)	ActNet Acc (%)	Mem (GB)	Time (s)
Uniform	32.1	58.3	14.2	5.9
Random	30.8	56.9	14.2	5.9
Motion	33.5	55.1	14.2	6.2
Attention	34.2	59.0	14.2	6.1
Query	31.9	60.1	14.2	7.3
Sliding	35.8	61.5	42.6	28.4
ATMC	35.9	61.8	9.9	16.1

230 dow approach that processes $[T/L]$ non-overlapping seg-
 231 ments independently and aggregates predictions via major-
 232 ity vote. Sampling-based methods are repeated five times
 233 with different random seeds. Statistical significance is as-
 234 sessed using paired t-tests on per-video scores, with Bon-
 235 ferroni correction for six comparisons ($\alpha_{\text{adjusted}} = 0.0083$),
 236 and effect sizes reported using Cohen’s d alongside 95%
 237 confidence intervals. All experiments are conducted on
 238 NVIDIA A100 80GB GPUs using PyTorch 2.0.1, Trans-
 239 formers 4.31.0, CUDA 11.8, and Python 3.9. The full im-
 240 plementation will be released under an MIT license upon
 241 publication, and all experiments are reproducible on a sin-
 242 gle A100 GPU within 48 hours.

243 5. Main Results

244 Table 1 summarizes results on Video-LLaMA-2 (7B).
 245 ATMC achieves the highest accuracy among single-pass
 246 methods on Ego4D-NLQ (+3.8 F1; 35.9% vs. 32.1%) and
 247 ActivityNet-QA (+3.5 EM; 61.8% vs. 58.3%), matching
 248 or slightly surpassing the more expensive sliding-window
 249 strategy while using significantly less memory (9.9GB vs.
 250 42.6GB) and lower latency (16.1s vs. 28.4s). Gains are
 251 statistically significant ($p < 0.001$) with moderate effect
 252 sizes (Cohen’s d 0.38–0.42). The main trade-off is a $2.7\times$
 253 inference-time increase relative to uniform sampling, but
 254 computations are parallelizable and amortized across mul-
 255 tiple queries; caching reduces overhead to $\leq 10\%$. Trends
 256 are consistent on LLaVA-Video (7B) (+3.4 F1 Ego4D, +3.6
 257 EM ActivityNet). ATMC is model-agnostic, requires no re-
 258 training or architecture changes, scales linearly with frame
 259 count, and maintains stable accuracy on longer videos, out-
 260 performing uniform sampling as video length increases.
 261 No additional trainable parameters or gradient updates are
 262 needed, avoiding issues like catastrophic forgetting or dis-
 263 tribution shift.

264 5.1. Ablation Studies

265 Table 2 shows the contribution of each ATMC compo-
 266 nent. Removing query similarity causes the largest drop
 267 (-2.5 points), followed by attention entropy (-1.2) and mo-
 268 tion magnitude (-0.6), highlighting the importance of task-
 269 specific, model-informed, and low-level signals. Using max

Table 2. Ablation study on Video-LLaMA-2 (7B) / Ego4D-NLQ.

Variant	F1@0.3 (%)	Δ from ATMC	Memory (GB)
ATMC (full)	35.9	–	9.9
w/o Motion ($\alpha = 0$)	35.3	-0.6	9.9
w/o Attention ($\beta = 0$)	34.7	-1.2	9.9
w/o Query ($\gamma = 0$)	33.4	-2.5	9.9
Mean Pool (default)	35.9	–	9.9
Max Pool	35.2	-0.7	9.9
Learnable Weights*	36.2	+0.3	9.9

instead of mean pooling reduces performance by 0.7 points,
 while learning weights via linear regression yields only a
 marginal +0.3 improvement, confirming that simple grid
 search and mean pooling suffice for stable temporal abstrac-
 tion.

275 5.2. Hyperparameter Sensitivity, Compression, and 276 Limitations

277 We analyze varying the compression ratio L'/L and find
 278 0.75 provides the best accuracy–efficiency trade-off (35.9%
 279 F1; 9.9 GB), with more aggressive compression reducing
 280 memory but degrading performance (33.2% at 0.5, 28.7%
 281 at 0.25). ATMC is robust to weight selection: equal weights
 282 ($\alpha = \beta = \gamma = 1/3$) yield 34.1%, and cross-dataset weights
 283 give 34.8%. Manual inspection of 50 Ego4D failures re-
 284 veals that 45% stem from low-scoring important frames,
 285 30% from overly coarse temporal pooling, 15% from base
 286 VidLLM errors, and 10% from misleading query similarity,
 287 highlighting remaining challenges in long-horizon reason-
 288 ing.

289 6. Conclusion

290 We introduce ATMC, a lightweight test-time framework
 291 that extends pretrained VidLLMs’ temporal reasoning with-
 292 out retraining. By combining motion, attention entropy,
 293 and query-conditioned similarity, ATMC improves long-
 294 video QA by 3–4 points while reducing GPU memory by
 295 30–45%. Despite a $2.7\times$ latency overhead and sensitivity
 296 to ambiguous queries, ATMC offers a practical, modular,
 297 and model-agnostic solution for efficiently handling long
 298 videos.

299 References

- [1] Grauman, K., Westbury, A., Byrne, E., Chavis, Z.,
 Furnari, A., Girdhar, R., ... & Malik, J. (2022).
 Ego4d: Around the world in 3,000 hours of egocen-
 tric video. In Proceedings of the IEEE/CVF confer-
 ence on computer vision and pattern recognition (pp.
 18995-19012). 1, 3
- [2] Yu, Z., Xu, D., Yu, J., Yu, T., Zhao, Z., Zhuang, Y., &
 Tao, D. (2019, July). Activitynet-qa: A dataset for un-

- 308 understanding complex web videos via question answer-
309 ing. In Proceedings of the AAAI Conference on Arti-
310 ficial Intelligence (Vol. 33, No. 01, pp. 9127-9134). 1,
311 3
- [3] Zhang, H., Li, X., & Bing, L. (2023). Video-
312 llama: An instruction-tuned audio-visual language
313 model for video understanding. arXiv preprint
314 arXiv:2306.02858. 1
315
- [4] Lin, B., Ye, Y., Zhu, B., Cui, J., Ning, M., Jin, P.,
316 & Yuan, L. (2024, November). Video-llava: Learn-
317 ing united visual representation by alignment before
318 projection. In Proceedings of the 2024 conference on
319 empirical methods in natural language processing (pp.
320 5971-5984). 1
321
- [5] Li, K., He, Y., Wang, Y., Li, Y., Wang, W., Luo, P.,
322 ... & Qiao, Y. (2025). Videochat: Chat-centric video
323 understanding. Science China Information Sciences,
324 68(10), 200102. 1
325
- [6] Azad, S., Vineet, V., & Rawat, Y. S. (2025). Hierarq:
326 Task-aware hierarchical q-former for enhanced video
327 understanding. In Proceedings of the Computer Vision
328 and Pattern Recognition Conference (pp. 8545-8556).
329
- [7] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z.,
330 ... & Guo, B. (2021). Swin transformer: Hierarchi-
331 cal vision transformer using shifted windows. In Pro-
332 ceedings of the IEEE/CVF international conference on
333 computer vision (pp. 10012-10022).
334
- [8] Li, L., Ma, R., Guo, Q., Xue, X., & Qiu, X. (2020).
335 Bert-attack: Adversarial attack against bert using bert.
336 arXiv preprint arXiv:2004.09984.
337
- [9] Bolya, D., Fu, C. Y., Dai, X., Zhang, P., Feichtenhofer,
338 C., & Hoffman, J. (2022). Token merging: Your vit but
339 faster. arXiv preprint arXiv:2210.09461. 1
340
- [10] Rao, Y., Liu, Z., Zhao, W., Zhou, J., & Lu, J.
341 (2023). Dynamic spatial sparsification for efficient vi-
342 sion transformers and convolutional neural networks.
343 IEEE Transactions on Pattern Analysis and Machine
344 Intelligence, 45(9), 10883-10897.
345
- [11] Beltagy, I., Peters, M. E., & Cohan, A. (2020).
346 Longformer: The long-document transformer. arXiv
347 preprint arXiv:2004.05150. 1
348
- [12] Feichtenhofer, C., Fan, H., Xiong, B., Girshick, R.,
349 & He, K. (2021). A large-scale study on unsuper-
350 vised spatiotemporal representation learning. In Pro-
351 ceedings of the IEEE/CVF conference on computer
352 vision and pattern recognition (pp. 3299-3309). 1
353
- [13] Bertasius, G., Wang, H., & Torresani, L. (2021, July).
354 Is space-time attention all you need for video under-
355 standing?. In Icml (Vol. 2, No. 3, p. 4). 1
356