Too Many Frames, Not All Useful: Efficient Strategies for Long-Form Video QA

Anonymous ACL submission

Abstract

Long-form videos that span across wide temporal intervals are highly information redundant and contain multiple distinct events or en-004 tities that are often loosely related. Therefore, 005 when performing long-form video question answering (LVQA), all information necessary to generate a correct response can often be contained within a small subset of frames. Recent literature leverage large language models (LLMs) in LVQA benchmarks, achieving exceptional performance, while relying on vision 011 012 language models (VLMs) to convert all visual content within videos into natural language. Such VLMs often independently caption a large number of frames uniformly sampled from long videos, which is not efficient and can mostly be redundant. Motivated by this inefficiency, 017 we propose LVNet, a modular and trainingfree framework featuring a novel Hierarchical 019 Keyframe Selector (HKS) that efficiently selects a minimal set of informative frames tailored to each question. LVNet's modularity al-022 lows easy integration with existing approaches for more efficient LVQA. We achieve stateof-the-art performance among similarly configured models across four benchmark LVQA datasets: EgoSchema, NExT-OA, IntentOA, VideoMME. Code will be released publicly.

1 Introduction

034

041

043

Video understanding is a long-standing vision problem (Aggarwal and Ryoo, 2011) with numerous real-world applications. It has been traditionally studied even before the era of differentiable representation learning, with hierarchical approaches focusing on longer videos (Allen and Ferguson, 1994; Ivanov and Bobick, 2000; Shi et al., 2004; Hongeng et al., 2004; Ryoo and Aggarwal, 2006). Today, video understanding research involving the language modality is particularly popular, with tasks such as video question answering (QnA) that involve generating human-style conversations using large language models (LLMs) (Tapaswi et al., 2016; Zeng et al., 2017; Xu et al., 2017).



Figure 1: **High Accuracy with Low Compute:** LVNet achieves state-of-the-art performance on EgoSchema (subset) while processing only a fraction of frame captions (below 12 per video) with the expensive LLM. More detailed analysis presented in Section 4.3.

044

045

046

047

051

056

057

059

060

061

062

063

065

066

067

069

Motivated by successful image-language models (Liu et al., 2023; Dai et al., 2023), several works build video LLMs (Yu et al., 2023; Papalampidi et al., 2023; Maaz et al., 2023; Wang et al., 2024b) that directly process video frames in one stage to perform QnA. However, for long videos (containing 1000s of frames) these models require processing large visual token sequences with the LLM, making inference computationally expensive or even infeasible (see Table 5). An alternate line of works (Zhang et al., 2023; Wang et al., 2023; Ranasinghe et al., 2024) follow multi-stage pipelines that first extract frame-level information followed by temporal modeling with an LLM. While these multi-stage works similarly encounter compute bottlenecks with increased frame processing for longer videos, it is possible to feed the LLM with descriptors of non-uniformly sampled frames (Kahatapitiya et al., 2024; Wang et al., 2024d,f), motivating our exploration into keyframe selection, i.e. identifying a minimal set of frames most useful for correctly answering a given video-question pair.

Therein, we propose LVNet, a framework containing a novel Hierarchical Keyframe Selector (HKS) that performs efficient key-frame selection followed by VLM for caption generation and LLM



Figure 2: (top) **Overview:** LVNet uses a Hierarchical Keyframe Selector (HKS) module to select keyframes, followed by VLM & LLM for caption and answer generation. (below) **HKS Module** processes dense frames with lighter modules and progressively exploits heavier, more performance-oriented modules on smaller subsets of frames to ensure efficient computation.

for answer generation as illustrated in Figure 2. Aligned with prior work (Zhang et al., 2023; Wang et al., 2024f,d), per-frame captions are processed with a powerful LLM to generate correct answers for a given video-question pair. As shown in Fig. 1, LVNet achieves strong performance efficiently, processing only a fraction of frame descriptors (captions) with the LLM. Compared to feeding all captions or all frames directly to a powerful model (e.g. GPT-40), our LVNet performs inference at a fraction of the cost (see Table 6).

070

071

074

086

094

100

101

We summarize our key contributions as follows:

- (a) Efficient Frame Selection: Hierarchical Keyframe Selector (HKS) efficiently extracts keyframes from sequences up to 1800 frames (hour long videos) with a filter rate over 98%.
- (b) **Video Training Free:** Our framework requires no video level training and simply uses existing off-the-shelf modules.
- (c) **Versatile Framework:** Existing methods can easily be integrated with LVNet to further boost their performance (details in Appendix A.6).

Proposed LVNet achieves state-of-the-art results (at common inference compute) on multiple long-form video question answering benchmarks demonstrating strong performance and generality.

2 Related Work

Video Question Answering: Visual question answering (VQA) involves generating open-ended textual content conditioned on an image and natural language query (Agrawal et al., 2015). Its video variant, Video-VQA (Yu et al., 2019a) replaces im-

Feature	Ours	VA	Tr.	VT	VC	FV
(effective selection) Uses non-uniform sampling Scene continuity-based selection Robust to initial frames Fine-grained visual refinement	\$ \$ \$ \$ \$ \$	✓ × ×	✓ × ✓ ×	√ √ ✓ ×	√ √ ✓ ×	√ √ √ √ ×
(compute efficient) Lightweight feature extraction Single pass inference Video Training Free	\ \ \	× × ✓	× × ✓	x x x	ン ン メ	、 ×

Table 1: LVNet exhibits unique features compared to prior work VideoAgent (VA) (Fan et al., 2024), Traveler (Tr.) (Shang et al., 2024), VideoTree (VT) (Wang et al., 2024g), VideoChat-T (VC) (Zeng et al., 2024), and Frame-Voyager (FV) (Yu et al., 2024b). See Appendix A.7 for details.

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

134

135

136

137

138

139

140

141

ages with videos. Multiple early datasets focus on querying objects or events based on referential and spatial relations (Xu et al., 2017; Zeng et al., 2017; Yu et al., 2019a). Later tasks require explicit temporal understanding of sequential events (Lei et al., 2018, 2020; Yu et al., 2019b). More recent datasets focus on longer videos containing multiple actions and scenes spread over wide time intervals (termed long-form videos) (Xiao et al., 2021; Li et al., 2022). Referred to as long-form video question answering (LVQA), these benchmarks are constructed to specifically test strong causal and temporal reasoning (Xiao et al., 2021) over long temporal windows (Mangalam et al., 2023). Some works tackling such video VQA tasks leverage graph networks to model cross object / event relations (Hosseini et al., 2022; Xiao et al., 2022a,b). A more recent line of works integrate LLMs to tackle this task (Zhang et al., 2023; Wang et al., 2023; Kahatapitiya et al., 2024; Wang et al., 2024d; Ranasinghe et al., 2024; Wang et al., 2024f; Fan et al., 2024) utilizing the strong reasoning skills of LLMs. A common aspect is the use of a vision language model (VLM) to convert frame level visual information into natural language. This in turn is input to the LLM which makes a final prediction.

Unlike these methods, LVNet incorporates a unique Hierarchical Keyframe Selector that progressively reduces the number of keyframe candidates. Lighter modules are applied to dense frames, while heavier, more performance-focused modules are applied to a small subset of filtered frames. Additionally, LVNet does not require video-level training, unlike earlier supervised approaches.

Frame Selection in Videos: The task of frame selection in videos has been long explored in video (Davis and Bobick, 1997; Zhao et al., 2017) with more recent works focused directly on long-form video question answering (Buch et al., 2022; Wang et al., 2024g; Fan et al., 2024; Zeng et al., 2024; Yu

195

196

197

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

et al., 2024b). Most similar to our work is Wang et al. (2024d) which employs an LLM based strategy for video frame selection. However, our LVNet differs with several unique features as summarized in Table 1.

Method 3

142

143

144

145

146

147

148

149

151

154

155

156

157

158

159

160

161

162

163

164

165

167

169

170

171

172

173

174

175

177

178

187

191

In this section, we present our training-free (i.e. zero-shot) framework for long-form video QA, LVNet. Videos are a dense form of data with even a few seconds long clip being composed of 100s of frames (individual images). In the case of longform videos, this frame count is even greater. However, the information necessary to answer a given question is often contained in a handful of those frames. Our framework tackles this challenge of selecting an optimal and minimal set of informative frames. We refer to this as keyframe selection. Given such a set of useful frames, we also establish optimal strategies for extracting their information using modern large language models (LLMs), taking into account their sequential nature.

Our proposed LVNet comprises of three components: a Hierarchical Keyframe Selector (HKS), a Vision Language Model (VLM), and a Large Language Model (LLM) as illustrated in Figure 2. The HKS, an efficient, hierarchical keyframe selector, is the core contribution of our work. First, the model processes 900 uniformly sampled frames and clusters them into distinct scenes Next, it extracts keywords from a given natural language query via LLM and selects the frames most relevant to those keywords. Finally, the selected frames are described in natural language by a more powerful and computationally intensive VLM. Finally, an LLM processes the language descriptions of the selected frames to answer a given query.

Background 3.1

Recent approaches utilizing LLMs for long video 179 question answering (LVQA) (Zhang et al., 2023; Wang et al., 2023; Kahatapitiya et al., 2024; Ranasinghe et al., 2024; Wang et al., 2024d) can be viewed as a composition of three sequential stages: 183 a) frame selection, b) VLM based frame captioning, and c) LLM based answer generation. Note that the complexity of each stage varies across methods given their focus on different aspects of the LVQA task (e.g. frame selection in some is simply uniform 188 sampling). In our work, we also follow this structure, but we focus on improving the frame selection 190 stage. Under such a framework, our proposed HKS can serve as plug-in modules to replace the frame 192

selection stage and the later two stages are similar to these prior works.

3.2 Architecture

Consider a video, $\mathbf{x} \in \mathbb{R}^{T \times C \times H \times W}$ with T. C. H, W for frames, channels, height, width respectively and its paired natural language query q. Also consider a frame in \mathbf{x} at timestamp t as $\mathbf{x}[t] \in \mathbb{R}^{C \times H \times W}$. Our goal is to output a response, referred as r, suitable for the given query q based on information contained in the video x.

Our LVNet processes a given video-query (x, q)pair to output a response, $\hat{\mathbf{r}}$. The HKS module initially processes this video-query pair, selects T' keyframes, and outputs a deterministically subsampled video $\mathbf{x}' \in \mathbb{R}^{T' \times C \times H \times W}$. Each of these T' frames is then passed through the captioning stage of our VLM to generate a set of natural language descriptions, $D = \{d_1, d_2, \dots, d_{T'}\}$ where d_i describes the frame $\mathbf{x}'[i]$. Finally, the LLM processes all descriptions D and the query q to generate response $\hat{\mathbf{r}}$. We illustrate this overall architecture in Figure 2.

Hierarchical Keyframe Selector 3.3

We now describe our proposed Hierarchical Keyframe Selector (HKS) module. As illustrated in Figure 2, our proposed HKS comprises of three sequential submodules, each reducing the frame count to T_a , T_b , and $T_c = T'$ respectively.

Temporal Scene Clustering (TSC): The role of TSC is to perform visual content aware preliminary frame sampling. The established approach for preliminary frame selection is uniform sampling (limited to at most 200 frames). In contrast, TSC processes 900 to 1800 uniformly sampled frames to extract per-frame visual features using a lightweight deep neural network (ResNet-18) followed by a clustering procedure to identify n non-overlapping frame sets. Within each of the n sets, we uniformly sample $\leq \tau$ frames obtaining a total of $T_a \leq \tau \times n$. Our iterative clustering procedure is outlined in Algorithm 1. It calculates pairwise distances between all frames accounting for intra-frame local information using the extracted per-frame features, followed by n iterative frame similarity based clustering operations. A single cluster could contain just one frame or significantly more based on frame feature similarities, leading to a non-uniform sam*pling of frames* across the entire video. This allows more frames to be sampled from the information heavy temporal regions of videos.



Figure 3: **Qualitative example**: We illustrate a challenging long-video QA scenario from EgoSchema (Mangalam et al., 2023). We consider an input of 900 frames, which first get clustered into scenes and subsampled to retain around 390 frames. Next, the Coarse Keyframe Detector selects only 32 frames out of them, based on the alignment with keywords (Here, keywords are extracted based on answer options, via an LLM). Such coarse keyframes are then ranked based on the combination of confidence value and temporal span, and grouped into four sets, each containing eight frames. These sets are then processed through visual templating (*i.e.* simple concatenation across space) and fed into a VLM for Fine Keyframe Detection, resulting in just 12 frames.

Coarse Keyframe Detector (CKD): Unlike TSC in the prior stage, CKD reasons across both visual and language modalities (using the paired textual query, q) to further sub-sample T_a into T_b frames. CKD contains three elements: keyword generation strategy, dual-encoder image-text model, and similarity based confidence assignment algorithm. Keyword generation utilizes the given query, \mathbf{q} , alongside hand-crafted templating operations or an LLM to select or generate suitable keywords. The dualencoder image-text model uses a spatially aware contrastive language image pre-training (CLIP) network from (Ranasinghe et al., 2023). For confidence assignment, we construct an algorithm as outlined in Algorithm 2 which processes two lists, one of frames and one of keywords, and then calculates their pairwise likelihood of occurrence to assign each frame a confidence value (that reflects its usefulness to answer the query, q). See Appendix A.4 for more details.

243

246

247 248

251

257

261

266

267

269

270

271

For a single query, there can be multiple regions in a video that are highly informative but not useful or relevant in answering that query. A single query can also contain multiple different concepts and attributes that must be given attention to construct a correct answer: the keyword generation attempts to capture each of these distinct attributes. On the visual modality, a single frame will also encode multiple concepts and attributes. Our design choice for the spatially aware CLIPpy dual-encoder VLM from (Ranasinghe et al., 2023) is motivated by this nature of individual frames. Finally, confidence assignment takes into account these multiple modes of information within each frame and the query to suitably assign confidence scores to each frame that reflects its query relevance. We also highlight how the confidence scores are directly linked to the related keyword (i.e. reason that makes the frame relevant), leading to better interpretability and the ability to perform further keyword-based refinement in later stages.

Fine Keyframe Detector (FKD): In the prior CKD stage, cross-modal selection utilizes a dual-encoder VLM that is constrained by the set of keywords provided and performs limited reasoning at frame level. In contrast, FKD uses a *visual templating module* to combine multiple frames and uses VLM to generate open-ended natural language output through higher-level reasoning. The input in this stage is the set of F_b frames, with each frame having an assigned confidence score and keyword.

Our visual templating module partitions the T_b frames into sets of 8 ordered by their confidence scores, arranges frame sets as grids to form a collage-style image, and annotates that image with visually identifiable tags corresponding to each frame. We further illustrate this process in Figure 3 (see Visual Templating column). Each of

300

272

273

274

275

276

277

278

279

280

these visual templated images also contain a sub-301 set of keywords that correspond to their 8 images. 302 These resulting visual templated images along with 303 a prompt containing their associated keywords and 304 instructions to select a frame subset based on valid association between keywords and images (see Appendix A.5 for details) are input to the VLM. The output of the VLM is used to select a subset of each 8 image group. These frames are collected as output of the FKD stage, overall resulting in T_c frames. 311

312

317

331

332

336

337

338

339

342

The purpose of the initial visual templating module is to allow reasoning across a set of frames using the image-text VLM (which is trained to pro-314 cess a single image at time). This partitioning of the input T_b frames is performed based on confidence scores from the prior stage and timestamps. The eight frames with top confidence scores are grouped into the first visual template, followed by the next eight and so forth. This ensures the VLM selects both high confidence concepts and low confidence concepts, accounting for biases and weaknesses in our CKD stage. After that, we temporally reorder some image sets with low confidence scores to cover keyframes distributed across 326 long-range segments, while the sets with high confidence scores concentrate on keyframes in shortrange segments. A total of 16 low-score frames are temporally reordered in this process. The algorithm is described in Algorithm 3 and the prompting technique is explained in Appendix A.5. Our intuition is that such a mechanism allows one to best utilize the complementary strengths of two different 334 VLMs from CKD and FKD stages for better frame selection overall.

4 Experiments

In this section, we first discuss our experimental setup followed by quantitative evaluations comparing to existing baselines and ablations of our proposed components. We then present qualitative results for our method and outline some limitations of our approach.

Experimental Setup 4.1

Datasets: Given the training free nature of our framework, we do not utilize any video datasets for training. Datasets are used purely for evaluation. We select three benchmark video visual question 347 answering datasets focused on long-form videos for this purpose: EgoSchema (Mangalam et al., 2023), NExT-QA (Xiao et al., 2021), and IntentQA (Li et al., 2023). In addition, to further highlight 351

the strength of our approach on longer videos, we include results on VideoMME's long split (Fu et al., 2024). These datasets are public available and can be used freely for academic research. The first dataset, EgoSchema, consists of 5031 questions and each video lasts three-minute and have multiple choice question. The second dataset, NExT-QA, is another rigorously designed video question answering benchmark containing questions that require causal & temporal action reasoning, and common scene comprehension to correctly answer. These questions are further classified as Causal (Cau.), Temporal (Tem.), and Descriptive (Des.) and we evaluate on its validation set containing 4996 questions over 570 videos. The third dataset, IntentQA, is based on NExT-QA videos corresponding to temporal and causal reasoning quetions. It consists of 16k multiple-choice questions which are classified as Why?, How? or Before/After (B./A.). The fourth dataset, VideoMME, consists of very long videos—some up to one hour long, with an average duration of 44 minutes, and provides 900 Q&A.

352

353

354

355

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

379

380

381

382

383

384

385

386

387

390

391

392

393

394

395

396

397

398

399

400

401

402

403

Model Choices & Hyperparameters: For the HKS module, we use the ResNet-18 (He et al., 2016a) for the TSC, CLIP-B/16 (Ranasinghe et al., 2023) for the CKD and GPT-40 for the FKD. We select ResNet-18 and CLIP-B/16 due to their smaller models sizes—0.01B and 0.12B, respectivelywhich are significantly lighter compared to LLMs that are on a billion parameter scale. In line with previous state-of-the-art work (Wang et al., 2024f; Zhang et al., 2023; Wang et al., 2023), we use up to 4 NVIDIA RTX A5000 GPUs and GPT or DeepSeek APIs for running baselines and our setup in all evaluations. Also following prior work, we report results over single evaluation runs.

4.2 Evaluation

Quantitative Results: We evaluate LVNet on the EgoSchema, NExT-QA, and IntentQA dataset and present our results in Table 2. Models with videocaption pretraining are de-emphasized in grey to ensure fairness with image-level pertaining. Models utilizing significantly more captions than the 12 frames are downplayed in light green to consider caption efficiency. We reiterate how number of captions input to LLM affects inference compute cost of methods the most.

For EgoSchema(fullset), LVNet achieves 61.1%, the highest among the models utilizing approximately 12 captions. This result outperforms VideoAgent, the next best model using 8.4 captions, by +7%, is on par with VideoTree while using only

Model	Ego	Schema	NE	ExT-QA	In	VT Free	
model	Cap.↓	Acc.† (%)	Cap.↓	Acc.† (%)	Cap.↓	Acc.† (%)	• 1 1100
VideoLLaMA 2 (Cheng et al., 2024)	-	53.3	-	-	-	-	no
InternVideo2 (Wang et al., 2024e)	-	60.2	-	-	-	-	no
Tarsier (Wang et al., 2024a)	-	61.7	-	79.2	-	-	no
Vamos (Wang et al., 2023)	-	48.3	-	-	-	-	yes
IG-VLM (Kim et al., 2024)	-	59.8	-	68.6	-	65.3	yes
VIOLET (Fu et al., 2023)	5	19.9	-	-	-	-	yes
mPLUG-Owl (Ye et al., 2023)	5	31.1	-	-	-	-	yes
VideoAgent (Wang et al., 2024d)	8.4	54.1	8.2	71.3	-	-	yes
MVU (Ranasinghe et al., 2024)	16	37.6	16	55.2	-	-	yes
MoReVQA (Min et al., 2024)	30	51.7	16	69.2	-	-	yes
VFC (Momeni et al., 2023)	-	-	32	51.5	-	-	yes
ProViQ (Choudhury et al., 2023)	50	57.1	50	64.6	-	-	yes
VideoTree (Wang et al., 2024g)	62.4	61.1	(56)	73.5	(56)	66.9	yes
FrozenBiLM (Yang et al., 2022)	90	26.9	-	-	-	-	yes
LifelongMemory (Wang et al., 2024f)	90	62.1	-	-	-	-	yes
TraveLER (Shang et al., 2024)	(101)	53.3	(65)	68.2	-	-	yes
LangRepo (Kahatapitiya et al., 2024)	180	41.2	90	60.9	90	59.1	yes
LLoVi (Zhang et al., 2023)	180	50.3	90	67.7	90	64.0	yes
LVNet (ours)	12	61.1	12	72.9	12	71.7	yes

Table 2: **Long Video Evalation:** LVNet achieves state-of-the-art accuracies of 61.1%, 72.9%, and 71.7% on EgoSchema, NExT-QA, and IntentQA datasets respectively using just 12 frames compared to models using a similar number of captions. Models are ordered based on number of captions processed per video. Models with video-level training (VT) or utilizing significantly more captions than 12 frames used by LVNet are de-emphasized in grey or downplayed in light green to ensure fair comparison. Numbers in parentheses () indicate the maximum number of frames used. See Sec. A.2 in appendix for detailed results.

Method	LLM Param/Type	VT Free	TS Cap.↓	Acc. \uparrow
VideoChat-T	7B/OS	no	N/A	43.8
Frame-Voyager	34B/OS	no	N/A	51.2
LVNet (DS-V3)	37B/OS	yes	24	53.1
VideoTree+GPT-4o	<1.8T/PP	yes	98.0	53.1
VideoAgent+GPT-40	<1.8T/PP	yes	24.6	46.4
LVNet (GPT-4o)	<1.8T/PP	yes	24	53.9

Table 3: **Comparison to Keyframe Selection Methods on VideoMME:** We compare LVNet with both single-stage methods that rely on video-level training (VideoChat-T, Frame-Voyager) and two-stage methods (VideoTree, VideoAgent). Each model uses the maximum feasible frames or default settings. *Legend: DS-V3=DeepSeek-V3, OS=open-source, PP=proprietary, VT Free=no video-level training, TS Cap.=two-stage caption numbers.*. Notably, LVNet outperforms both single-stage and two-stage methods without requiring video-level training.

1/5 of the captions, and outperforms TraveLER by +7.8% while utilizing only 12% of the captions.

We next evaluate on NExT-QA dataset, which has a particular focus on both temporal and casual reasoning based question-answer pairs. LVNet achieves state-of-the-art performance on this benchmark outperforming prior work among the models utilizing approximately 12 captions. In fact, our LVNet outperforms VideoAgent by +1.6%.

In IntentQA dataset, LVNet outperforms all prior work, including de-emphasized models with videocaption pretraining and downplayed models utilizing significantly compute (captions input to LLM) than ours. In fact, LVNet shows a substantial improvement of +4.8% over the next best VideoTree, while using only 13% of captions (12 vs 90).

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

Lastly, Table 3 shows the LVNet's performance of on VideoMME's long split (Fu et al., 2024), which features videos up to an hour long-significantly exceeding the 12-minute average in MLVU (Zhou et al., 2024) or the 17minute average in the overall VideoMME.In combination with an open-source LLM (DeepSeek-V3), LVNet outperforms the single-stage keyframe selection methods VideoChat-T (Zeng et al., 2024) and Frame-Voyager (Yu et al., 2024b) by 9.3% and 1.9%, respectively, without any video-level training. Moreover, LVNet with GPT-40 surpasses VideoAgent and VideoTree by +7.5% when using a similar number of captions, and by +0.8% with only onequarter of the frames, underscoring its efficiency in extracting keyframes for very long videos. See Appendix A.3 for the detail result.

Given the generative nature of VQA tasks as well as the limited availability and noisy nature of fully-annotated video VQA corpora, building generalizable fully-supervised models are challenging for these tasks. Nevertheless, we highlight how our zero-shot and video level training-free framework is competitive with the best supervised approaches on this dataset. This indicates the promise of utilizing pretrained models, especially those equipped with extensive world knowledge and reasoning

Model	A	Avg. Frame Captions \downarrow						
1110401	8	11 12		16 70				
VideoAgent+GPT-40	63.2	60.8	_	_	_			
VideoTree+GPT-40	_	_	_	_	67.0			
LVNet (GPT-40)	64.4	-	68.2	67.8	-			

		TSC	CVD	EVD	1 22 1
Templating Order	Acc ↑	ISC	CKD	гкd	Acc.
Templating Older	1100.	X	X	X	62.6
Temporal	65.2	1	X	X	64.5
Confidence	67.6	1	1	X	65.8
Hybrid (both)	68.2	1	1	1	68.2

(a) **Frame Caption Count**: LVNet ourperforms VideoAgent and (b) VideoTree using fewer captions, with all models leveraging GPT- Contact 40 and using the default settings to ensure a fair comparison.

(b) **Visual Templating:** (c) **HKS Ablation**: LVNet Combination of confidence- accuracy consistently imbased & temporal ordering proves with each HKS subgives the best performance. module.

Table 4: **Ablation study** on EgoSchema (Mangalam et al., 2023): We evaluate different design decisions of our framework on EgoSchema 500-video subset for zero-shot VQA.

Method	LLM AP	Frames \uparrow	LLM frames \downarrow	Acc. \uparrow
Qwen-VL	7B	128	128	37.8
Qwen-VL	7B	256	256	OOM
Qwen-VL	7B	1800	1800	OOM
LVNet (DS-V3)	37B	1800	24	53.1

Table 5: **Single-Stage Method Comparison:** We report Accuracy on VideoMME long split (average video length 41 mins) along with LLM active parameters (LLM AP), total frames processed, and frames input to LLM as tokens / captions (LLM frames). LVNet DeepSeek-V3 (DS-V3) variant is used. Processing lengthy frame sequences with single-stage models at fixed compute becomes infeasible. Inference tested on 4 x 24GB A5000 GPUs. Similar findings in Zeng et al. (2024).

skills from alternate modality specific learning (i.e. in our cases image domain VLMs and language domain LLMs).

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

465

466

467

468

469

470

Qualitative Analysis of Hierarhical Keyframe Selector: We compare open-ended responses of LVNet and the uniform sampling method in Figure 4 to understand the effectiveness of the hierarchical keyframe selector in LVNet. The frames chosen by LVNet and the naive uniform sampling method are indicated by blue and red checkmarks in the images, respectively. LVNet selects frames at 5, 69, and 135 seconds by executing the hierarchical keyframe selector and generates captions based on those frames. When we feed the concatenated captions to the LLM to answer the given question: "Based on the video, what are the three main types of tools that C uses..." in an open-ended manner, the output identifies two main activities: welding torches and measuring tapes, among the three main activities described in Option 3 (welding handle, hammer, measuring tape), which is the correct answer, leading LVNet to choose the correct option.

In contrast, the uniform sampling method selects frames at 0, 16, and 32 seconds and generates

Model	$\mathrm{FV}\uparrow$	$FL\downarrow$	$VC\downarrow$	$\mathrm{TC}\downarrow$	Acc. \uparrow
GPT-40	384	384	\$2.88	\sim \$2592	65.3
LVNet (GPT-40)	1800	24	\$0.19	\$171	53.9

Table 6: **API Comparison:** We perform cost comparison with using GPT-40 directly on frames vs with LVNet. We report accuracy on VideoMME long split (Acc), frames processed per video (FV), frames / captions input to LLM (FL), per video cost (VC), and total evaluation cost (TC). Our LVNet achieves competitive performance at over 10x less inference cost. GPT-40 accuracy result from official benchmark leaderboard.

captions based on those frames. Similarly, when we feed the concatenated captions to the LLM to answer the same question, the output identifies only one activity—welding tools—resulting in the selection of the incorrect option. This example highlights the importance of keyframe selection and demonstrates the effectiveness of hierarchical keyframe selection in LVNet.

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

4.3 Ablations

In this section, we present ablations on key design decisions such as the sorting order in FKD, the number of frames for captions, and the effect of different components in HKS. In all ablations, we use a subset of EgoSchema (Mangalam et al., 2023), composed of 500 videos. Additional ablations about *Choice of LLM* and *Effect of Patch Size on Keyword Matching in CKD* are in Appendix A.1

Number of Frame Captions: We conduct an ablation on the number of frame captions, comparing our approach with VideoAgent (Wang et al., 2024d) and VideoTree (Wang et al., 2024g), all using GPT-40 for a fair comparison. As shown in Table 4a, LVNet achieves the highest accuracy of 68.2% with 12 captions, outperforming VideoAgent (63.2% at 8 captions) by +5% in a similar low caption settings and surpassing VideoTree (67.0% at 70 captions) by +1.2% while using only one-sixth as many cap-

••••						
Time(s):	0	5	16	32	69	135
Caption:	A person is using a hacksaw	holding a welding torch	A person is welding or soldering	C is using a power tool	A person is working on unwrapping	They are manipulating a measuring tape

Prompt: {Concatenated Captions from $\sqrt{}$ or $\sqrt{}$ }. I request you to answer the following question based on the preceding descriptions in less than 50 words. Question: Based on the video, what are the three main types of tools that c uses, and how do their roles in shaping the iron differ from one another?

LLM Output:

498

499

500

502

503

505

506

507

508

511

512

513

514

- LVNet (Ours): Based on the video, the three main types of tools used are welding torches, pliers, and measuring tapes...
- **Uniform Sampling**: The three main types of tools used are hacksaws, welding tools, and power tools...

GT (option 3): The three main types of tools that c uses are a welding handle, a hammer, and a measuring tape...

Figure 4: **Open-ended Responses from LVNet vs Uniform Sampling**: The frames chosen by LVNet and the naive uniform sampling method are indicated with blue and red checkmarks, respectively. LVNet identifies both welding torches and measuring tapes, choosing the correct option, whereas uniform sampling only detects welding tools and selects the incorrect answer. The blue, red, and purple highlights correspond to the three main activities in the video—welding a handle, using a hammer, and using a measuring tape, respectively.

tions. Additionally, LVNet exceeds VideoAgent by +1.2% even at VideoAgent's optimal 8 caption setting. More details can be found in Appendix A.1 and Appendix A.2.

Visual Templating Order: In visual templating, prioritizing frames by keyword confidence scores followed by reordering low-confidence frames based on timestamp proves more effective than using confidence scores or temporal order alone, as shown in Table 4b. In this hybrid approach, highconfidence frames capture short but important segments of videos, while low-confidence keyframes, which are crucial but visually challenging for keyword matching, are temporally ordered to cover broader segments. This hybrid approach outperforms solely temporal ordering or confidence-based ordering by +3% and +0.6%, respectively.

Effect of Hierarchical Keyframe Modules: Ta-515 ble 4c demonstrates the impact of incrementally 516 adding the temporal scene clustering (TSC), coarse 517 keyframe detector (CKD), and fine keyframe de-518 tector (FKD) modules. Without any of these modules, the model relies on uniform sampling and achieves 62.6%. When TSC is added and 12 frames 521 are selected uniformly, the accuracy increases to 522 64.5%. Adding both TSC and CKD raises the ac-523 curacy to 65.8%. Finally, incorporating all three modules-TSC, CKD, and FKD-into the model, 525

which is LVNet, results in an accuracy of 68.2%. This demonstrates the importance of including all modules in LVNet for optimal performance.

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

Singe-Stage Comparison: Table 5 highlights the scalability limitations of single-stage models on long videos. Models like Qwen-VL (Bai et al., 2023) run out of memory beyond 128 frames, while LVNet achieves significantly higher accuracy using just 24 keyframes as LLM inputs.

API Comparison: Table 6 compares inference cost and performance against purely GPT-40. LVNet achieves competitive accuracy while reducing per-video LLM cost by over 10×, demonstrating exceptional compute efficiency.

5 Conclusion

We proposed a novel approach for Long-form Video Question Answering (LVQA) that achieves state-of-the-art performance compared to the model using the similar-scale captions across 3 benchmarks datasets. Our Hierarchical Keyframe Selector demonstrates the effectiveness of keyframe selection in understanding a very long-form video QA. Additionally, we highlight the zero-shot capability for long-form video comprehension of our LVNet framework, which requires no video-level training. Our experiments showcase its significant advantage over previous methods.

Limitations

which we discuss below.

Despite the effectiveness of LVNet, as demon-

strated by benchmark experiments and comprehen-

sive ablations, our study has certain limitations,

• First, we acknowledge that we are unable to eval-

uate LVNet and other models with all available

VLMs or LLMs due to computational constraints and high costs. However, we carefully select

GPT-40 and DeepSeek-v3, the LLMs commonly

used in video understanding research, for our

main experiments and provide ablation studies

comparing various LLMs (e.g. GPT-3.5, GPT-4,

and GPT-40) to ensure a fair performance com-

parison, as presented in Table 4a and Table A.7a.

three components: TSC, CKD, and FKD. While

we demonstrated the effectiveness of each com-

ponent in Table 4c, we did not have the time or

resources to develop a unified module that could

replace all three. Although this is beyond the

scope of this paper, exploring a more efficient

implementation that integrates these three mod-

ules into a single model would be an interesting

• Like any LLM-based approach, LVNet is sensi-

tive to prompting. To ensure the transparency,

we provide examples of these prompts in Fig-

ure 4 and Figure A.6. We also plan to release

the code to enable further exploration by other

• Finally, we acknowledge that, as our approach

is zero-shot, any inherent limitations or biases in

the pretrained models may persist in the outputs

Jake K. Aggarwal and Michael S. Ryoo. 2011. Human

Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Devi Parikh,

and Dhruv Batra. 2015. Vqa: Visual question an-

swering. International Journal of Computer Vision,

James F Allen and George Ferguson. 1994. Actions and

and computation, 4(5):531-579.

events in interval temporal logic. Journal of logic

activity analysis. ACM Computing Surveys (CSUR),

direction for future research.

researchers.

of LVNet.

References

43:1 - 43.

123:4 - 31.

· Our hierarchical keyframe selector consists of

- 555
- 556
- 557
- 559
- F
- 562

564

- 565
- Ę

566

Ę

- į
- ļ
- 5
- 5
- 5

5

- 5
- 5
- Ę

F

Į

5

588

- 589
- 59

593 594

595

- 59
- 598 599

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile visionlanguage model for understanding, localization, text reading, and beyond. *Preprint*, arXiv:2308.12966. 600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

- S. Buch, Cristobal Eyzaguirre, Adrien Gaidon, Jiajun Wu, Li Fei-Fei, and Juan Carlos Niebles. 2022. Revisiting the "video" in video-language understanding. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 2907–2917.
- Yukang Chen, Fuzhao Xue, Dacheng Li, Qinghao Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian Tang, Shang Yang, Zhijian Liu, Ethan He, Hongxu Yin, Pavlo Molchanov, Jan Kautz, Linxi Fan, Yuke Zhu, Yao Lu, and Song Han. 2024a. Longvila: Scaling long-context visual language models for long videos. *Preprint*, arXiv:2408.10188.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024b. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024c. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67(12):220101.
- Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. 2024. Videollama
 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv*:2406.07476.
- Rohan Choudhury, Koichiro Niinuma, Kris M Kitani, and László A Jeni. 2023. Zero-shot video question answering with procedural programs. *arXiv preprint arXiv:2312.00937*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*.
- James Davis and Aaron Bobick. 1997. The representation and recognition of action using temporal templates. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2736–2744.
- Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. 2024. Videoagent: A memory-augmented multimodal agent for video understanding. *arXiv preprint arXiv:2403.11481*.
- Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. 2024.
- 9

Video-mme: The first-ever comprehensive evaluation

benchmark of multi-modal llms in video analysis.

Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin,

William Yang Wang, Lijuan Wang, and Zicheng Liu. 2023. An empirical study of end-to-end video-

language transformers with masked visual model-

ing. In Proceedings of the IEEE/CVF Conference

on Computer Vision and Pattern Recognition, pages

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016a. Deep residual learning for image recog-

nition. In Proceedings of the IEEE conference on

computer vision and pattern recognition, pages 770-

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016b. Deep residual learning for image recog-

nition. In Proceedings of the IEEE conference on

computer vision and pattern recognition, pages 770-

Somboon Hongeng, Ram Nevatia, and Francois Bre-

mond. 2004. Video-based event recognition: activity

representation and probabilistic recognition meth-

ods. Computer Vision and Image Understanding,

Pedram Hosseini, David A. Broniatowski, and Mona

Diab. 2022. Knowledge-augmented language models

for cause-effect relation classification. In Proceed-

ings of the First Workshop on Commonsense Representation and Reasoning (CSRR 2022), pages 43-48,

Dublin, Ireland. Association for Computational Lin-

Yuri A. Ivanov and Aaron F. Bobick. 2000. Recogni-

Kumara Kahatapitiya, Kanchana Ranasinghe, Jongwoo

Wonkyun Kim, Changin Choi, Wonseok Lee, and Wonjong Rhee. 2024. An image grid can be worth a

Jie Lei, Licheng Yu, Mohit Bansal, and Tamara Berg.

Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal.

2020. TVQA+: Spatio-temporal grounding for video

question answering. In Proceedings of the 58th An-

nual Meeting of the Association for Computational

Linguistics, pages 8211-8225, Online. Association

2018. TVQA: Localized, compositional video ques-

tion answering. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language

video: Zero-shot video question answering using a

Park, and Michael S Ryoo. 2024. Language reposi-

Machine Intelligence, 22(8):852–872.

tory for long video understanding.

vlm. arXiv preprint arXiv:2403.18406.

Processing (EMNLP).

for Computational Linguistics.

tion of visual activities and interactions by stochastic

parsing. IEEE Transactions on Pattern Analysis and

arXiv preprint arXiv:2405.21075.

22898-22909.

778.

778.

96(2):129-162.

guistics.

- 673
- 674
- 675
- 677
- 679

694

- 700
- 701 702 703

704

Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llava-onevision: Easy visual task transfer. Preprint, arXiv:2408.03326.

710

711

712

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

752

756

757

758

759

760

761

- Jiangtong Li, Li Niu, and Liqing Zhang. 2022. From representation to reasoning: Towards both evidence and commonsense reasoning for video questionanswering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).
- Jiapeng Li, Ping Wei, Wenjuan Han, and Lifeng Fan. 2023. Intentqa: Context-aware video intent reasoning. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 11963-11974.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. arXiv preprint arXiv:2304.08485.
- Muhammad Maaz, Hanoona Abdul Rasheed, Salman H. Khan, and Fahad Shahbaz Khan. 2023. Videochatgpt: Towards detailed video understanding via large vision and language models. In Annual Meeting of the Association for Computational Linguistics.
- Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. 2023. Egoschema: A diagnostic benchmark for very long-form video language understanding. ArXiv, abs/2308.09126.
- Juhong Min, Shyamal Buch, Arsha Nagrani, Minsu Cho, and Cordelia Schmid. 2024. Morevga: Exploring modular reasoning models for video question answering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13235-13245.
- Liliane Momeni, Mathilde Caron, Arsha Nagrani, Andrew Zisserman, and Cordelia Schmid. 2023. Verbs in action: Improving verb understanding in videolanguage models. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 15579-15591.
- Pinelopi Papalampidi, Skanda Koppula, Shreya Pathak, Justin Chiu, Joe Heyward, Viorica Patraucean, Jiajun Shen, Antoine Miech, Andrew Zisserman, and Aida Nematzdeh. 2023. A simple recipe for contrastively pre-training video-first encoders beyond 16 frames. arXiv preprint arXiv:2312.07395.
- Kanchana Ranasinghe, Xiang Li, Kumara Kahatapitiya, and Michael Ryoo. 2024. Understanding long videos in one multimodal language model pass.
- Kanchana Ranasinghe, Brandon McKinzie, Sachin Ravi, Yinfei Yang, Alexander Toshev, and Jonathon Shlens. 2023. Perceptual grouping in contrastive visionlanguage models. In ICCV.

- 763 764 765 766 767
- 768 769 770 771
- 772 773 774 775 776 776 777 778 779
- 780 781 782
- 784 785
- . 7 7
- 789 790
- 791 792
- 793
- 794 795 796
- 798 799

- 80
- 803 804
- 00
- 806 807
- 808
- 810 811

812 813

- 814 815
- 815

817

- Michael S. Ryoo and Jake K. Aggarwal. 2006. Recognition of composite human activities through contextfree grammar based representation. 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), 2:1709–1718.
- Chuyi Shang, Amos You, Sanjay Subramanian, Trevor Darrell, and Roei Herzig. 2024. Traveler: A modular multi-Imm agent framework for video questionanswering. *arXiv preprint arXiv:2404.01476*.
- Yifan Shi, Yan Huang, David Minnen, Aaron Bobick, and Irfan Essa. 2004. Propagation networks for recognition of partially ordered sequential action. In *CVPR*.
- Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2016. MovieQA: Understanding stories in movies through question-answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).*
- Jiawei Wang, Liping Yuan, and Yuchen Zhang. 2024a. Tarsier: Recipes for training and evaluating large video description models. *arXiv preprint arXiv:2407.00634*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Ke-Yang Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024b. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *ArXiv*, abs/2409.12191.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024c. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Shijie Wang, Qi Zhao, Minh Quan Do, Nakul Agarwal, Kwonjoon Lee, and Chen Sun. 2023. Vamos: Versatile action models for video understanding. *arXiv preprint arXiv:2311.13627*.
- Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. 2024d. Videoagent: Long-form video understanding with large language model as agent.
- Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Jilan Xu, Zun Wang, et al. 2024e. Internvideo2: Scaling video foundation models for multimodal video understanding. *arXiv preprint arXiv:2403.15377*.
- Ying Wang, Yanlai Yang, and Mengye Ren. 2024f. Lifelongmemory: Leveraging llms for answering queries in long-form egocentric videos. *Preprint*, arXiv:2312.05269.

Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jaehong Yoon, Feng Cheng, Gedas Bertasius, and Mohit Bansal. 2024g. Videotree: Adaptive tree-based video representation for llm reasoning on long videos. *arXiv preprint arXiv:2405.19209*. 818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

- Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. NExT-QA: Next phase of questionanswering to explaining temporal actions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).*
- Junbin Xiao, Angela Yao, Zhiyuan Liu, Yicong Li, Wei Ji, and Tat-Seng Chua. 2022a. Video as conditional graph hierarchy for multi-granular question answering. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence (AAAI)*, pages 2804–2812.
- Junbin Xiao, Pan Zhou, Tat-Seng Chua, and Shuicheng Yan. 2022b. Video graph transformer for video question answering. In *European Conference on Computer Vision*, pages 39–58. Springer.
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *ACM Multimedia*.
- Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Zero-shot video question answering via frozen bidirectional language models. *Advances in Neural Information Processing Systems*, 35:124–141.
- Zeyuan Yang, Delin Chen, Xueyang Yu, Maohao Shen, and Chuang Gan. 2024. Vca: Video curious agent for long video understanding. *arXiv preprint arXiv:2412.10471*.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*.
- Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. 2023. Self-chained image-language model for video localization and question answering. *ArXiv*, abs/2305.06988.
- Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. 2024a. Self-chained image-language model for video localization and question answering. *Advances in Neural Information Processing Systems*, 36.
- Sicheng Yu, Chengkai Jin, Huanyu Wang, Zhenghao Chen, Sheng Jin, Zhongrong Zuo, Xiaolei Xu, Zhenbang Sun, Bingni Zhang, Jiawei Wu, et al. 2024b. Frame-voyager: Learning to query frames for video large language models. *arXiv preprint arXiv:2410.03226*.

Zhou Yu, D. Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019a. Activitynet-qa: A dataset for understanding complex web videos via question answering. ArXiv, abs/1906.02467.

872

873

875

887

891

892

894

896

897

901

903

904

905

906 907

908

909

910

911

- Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019b. ActivityNet-QA: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*.
 - Kuo-Hao Zeng, Tseng-Hung Chen, Ching-Yao Chuang, Yuan-Hong Liao, Juan Carlos Niebles, and Min Sun. 2017. Leveraging video descriptions to learn video question answering. In AAAI Conference on Artificial Intelligence.
 - Xiangyu Zeng, Kunchang Li, Chenting Wang, Xinhao Li, Tianxiang Jiang, Ziang Yan, Songze Li, Yansong Shi, Zhengrong Yue, Yi Wang, et al. 2024. Timesuite: Improving mllms for long video understanding via grounded tuning. arXiv preprint arXiv:2410.19702.
 - Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas Bertasius. 2023. A simple llm framework for longrange video question-answering. *arXiv preprint arXiv:2312.17235*.
 - Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024a. Video instruction tuning with synthetic data. *Preprint*, arXiv:2410.02713.
 - Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024b. Video instruction tuning with synthetic data. *Preprint*, arXiv:2410.02713.
 - Zhichen Zhao, Huimin Ma, and Shaodi You. 2017. Single image action recognition using semantic body part actions. In *The IEEE International Conference on Computer Vision (ICCV)*.
- Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. 2024. Mlvu: A comprehensive benchmark for multi-task long video understanding. *arXiv preprint arXiv:2406.04264*.

914 915 916

917

918

919

920

924

926

930

931

932

933

934

936

937

938

941

947

A.1	Additional	Ablations

Appendix

In this section, we present additional experiments conducted to inform the LVNet's design. We have tested different LLMs and experimented with various scales of the visual feature map.

LLM	Acc. (%)	Patch Size	Acc. (%)
GPT-3.5	61.0	1x1	63.6
GPT-4	65.4	7x7	66.2
GPT-40	68.2	14x14	68.2

(a) **Choice of LLM**: We consider different options for our LLM for video QA. GPT-40 performs the best

(b) **Effect of Patch Size in CKD**: A larger patch size in Keyword Matching performs better.

Table A.7: **Additional ablations experiments** on EgoSchema (Mangalam et al., 2023): We evaluate different design decisions of our framework on EgoSchema 500-video subset for zero-shot VQA. Default setting is highlighted.

Choice of LLM: Table A.7a shows that GPT-4o outperforms both GPT-4 and GPT-3.5 by +2.8% and +7.2%, respectively. Given that GPT-4o is more cost-effective and lightweight compared to GPT-4, we have selected it as our default LLM.

Effect of Patch Size on Keyword Matching in CKD: Table A.7b shows the effect of the scales of the patch sizes in the CKD. Since keywords can represent activities spanning the entire image or confined to a small region, we adjust the resolution of the visual feature map output from the spatially aware contrastive image pre-training (CLIP) network (Ranasinghe et al., 2023) to match keywords. Our findings show that higher resolutions lead to better accuracy. In LVNet, we use a 14×14 feature map and determine the confidence level of the keyword by selecting the maximum value between the 14×14 patches and the keyword's text embedding.

Exact Frame Caption Counts: For completeness, Table A.8 lists the original (unrounded) frame caption counts and corresponding accuracies for VideoAgent(Wang et al., 2024d), VideoTree(Wang et al., 2024g), and LVNet (ours). These values supplement the rounded numbers presented in Table 4a of the main text.

A.2 Extended results on NExT-QA and IntentQA

We present extended zero-shot evaluation results on NExT-QA in Table A.9, comparing LVNet with prior zero-shot models across different task categories: causal, temporal, and descriptive reasoning. Models are ordered based on the number of captions processed per video, highlighting the tradeoffs between caption efficiency and performance.

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

LVNet achieves state-of-the-art performance with an overall accuracy of 72.9%, outperforming most models while using only 12 captions per video. Notably, it attains 75.0% on causal reasoning, which is the highest among all models evaluated. For temporal reasoning, LVNet achieves 65.5%, remaining competitive despite using significantly fewer captions than models like VideoTree (56 captions) and LangRepo (90 captions). In descriptive reasoning, LVNet reaches 81.5%, matching VideoTree while processing significantly fewer captions.

Compared to VideoAgent, the closest competing model in terms of caption efficiency (8.4 captions), LVNet demonstrates a substantial performance gain across all categories, with a +2.8% improvement in overall accuracy. While models like VideoTree and TraveLER show strong performance, they process significantly more captions (56 and 65, respectively), indicating that LVNet achieves a superior balance between efficiency and accuracy.

We present extended zero-shot evaluation results on IntentQA in Table A.10, comparing LVNet with prior zero-shot models across different reasoning categories: *Why?*, *How?*, and *B.A.* (Before/After). Models are ordered based on the number of captions processed per video, highlighting the balance between caption efficiency and performance.

LVNet achieves an overall accuracy of 71.7%, outperforming all models while using only 12 captions per video. It achieves 75.0% on the *Why*? category, 74.4% on the *How*? category, and 62.1% on the *B.A.* category. Compared to VideoTree, which processes 56 captions and achieves an overall accuracy of 66.9%, LVNet outperforms it by +4.8% while using significantly fewer captions. Similarly, LangRepo and LLoVi, which process 90 captions, achieve overall scores of 59.1% and 64.0%, respectively, further demonstrating LVNet's caption efficiency.

To ensure fairness, models that utilize videocaption pretraining or process substantially more captions than LVNet are de-emphasized in grey or downplayed in light green in Table A.9 and Table A.10. We adopt the reported GPT-40 results for VideoAgent and VideoTree in Table 3 and Table 4a from the VCA (Yang et al., 2024), but do not compare directly against VCA for two rea-



Question: Identify a recurring action in the video

Figure A.5: **Comparison of Keyframe Selection**: Comparison of LVNet and VideoAgent in keyframe selection for video question answering. LVNet refines frames through a multi-stage process (TSC, CKD, FKD) to form a non-uniform keyframe distribution, capturing relevant moments tied to the query. In contrast, VideoAgent relies on uniform sampling and LLM-based frame selection, which fails to focus on crucial keyframes, leading to incorrect predictions. The final keyframe distributions illustrate LVNet's ability to retrieve meaningful frames directly related to the answer, while VideoAgent selects irrelevant frames.

Table A.8: **Exact (Unrounded) Frame Caption Counts.** These values supplement the rounded numbers in Table 4a All models are based on either GPT-40 or <u>GPT-4</u>.

Model						Avg	. Frame	Captio	ons↓					
	6.4	7.1	8	8.1	8.4	9.7	10.7	11	11.3	12	16	32	62.4	69.5
VideoAgent	<u>58.4</u>	-	-	63.2	60.2	-	60.8	<u>57.4</u>	-	-	-	-	-	-
VideoTree	-	<u>61.0</u>	-	-	-	<u>61.6</u>	-	-	<u>62.2</u>	<u>62.5</u>	<u>66.2</u>	<u>64.4</u>	<u>66.2</u>	67.0
LVNet (Ours)	-	-	64.4	-	-	-	-	-	-	68.2	67.8	-	-	-

1001 sons: (1) the VCA paper does not provide code or implementation details (e.g., inference speed), 1002 making replication infeasible; and (2) its reported results cover only a subset of Egoschema, prevent-1004 1005 ing a fair comparison to our approach on the fullscale EgoSchema, NExT-QA, and IntentQA. Ove-1006 reall, these clarifications, alongside the results in 1007 Table A.9 and Table A.10, underscore LVNet's ef-1008 fectiveness in achieving high accuracy while maintaining computational efficiency. 1010

1011

1013

1014

1016

1019

1020

1022

1024

1025

1027

1028

1030

1033

1034

1036

1038

1039

1040

1041

1042

1043

A.3 Extended results on VideoMME

Table A.11 provides an extended comparison of both *single-stage* and *two-stage* approaches on the VideoMME long-video benchmark. In contrast to single-stage methods, which typically require costly video-level training on large datasets which has the risk of containing evaluation datasets. LVNet is video-level training free and thus retains broader generality. Critically, single-stage pipelines often become infeasible for long videos because they must process entire sequences—potentially exceeding 1,000 frames—through a heavy vision-language model (VLM) or LLM.

Keyframe-selection approaches mitigate this challenge by filtering a minimal set of relevant frames before any large-model processing. As shown in Table A.11, LVNet can handle up to 1,800 frames by reducing them to just a few dozen keyframes, thereby remaining both memory- and compute-efficient while still achieving strong accuracy. For Qwen-VL and GPT40, we demonstrate an ablation with 1,800 frames to highlight how easily single-stage methods can run out of memory (OOM) under 4 NVIDIA RTX A5000 GPUs or encounter GPT API errors.

Moreover, our approach achieves strong performance with both open-source and proprietary LLM backbones, all without requiring any videolevel training. LVNet+DeepSeek-V3 outperforms these single-stage keyframe selection models (VideoChat-T, Frame-Voyager) with equal or less frame selection complexity and outperforms singlestage VideoLLMs with a similarly sized open-
source LLM. LVNet surpasses two-stage keyframe1044selection models (VideoTree, VideoAgent) under
the same GPT40, underscoring its effectiveness1045over both single-stage and two-stage video meth-
ods. The accuracy comes from each model's rec-
ommended configuration.1046

1051

1075

A.4 Algorithms in Detail

Our algorithms are presented in full detail in Al-1052 gorithm 1, Algorithm 2, and Algorithm 3. TSC in 1053 Algorithm 1 extracts per-frame visual features us-1054 ing ResNet-18, followed by an iterative clustering 1055 procedure to identify n non-overlapping frame sets. 1056 Within each of the *n* sets, we uniformly sample $\leq \tau$ 1057 frames, obtaining a total of $T_a \leq \tau \times n$ frames. For example, LVNet sets $\psi = 5, \lambda = 12, \tau = 18$, resulting in approximately $n \sim 25$ and $T_a \sim 390$ 1060 on the EgoSchema dataset. CKD in Algorithm 2 1061 selects top L frames based on similarity/confidence 1062 scores, which are calculated using cosine similar-1063 ity between frames and keywords with CLIP-B/16. 1064 LVNet employs L = 32, $len(K) \le 25$ on the 1065 EgoSchema dataset. FKD in Algorithm 3 sorts 1066 frames and their corresponding keywords by con-1067 fidence scores, and reorder the K frames with the 1068 lowest scores temporally. It groups frames sequen-1069 tially into visual templates, each consisting of N1070 frames. From each template, the M frames and 1071 keywords most relevant among the N pairs are 1072 selected using GPT-40. We set L = 32, K =1073 16, N = 8, M = 3.1074

A.5 Prompting: Fine Keyframe Detector

We prompt the VLM to select frames that are most
compatible with the list of given keywords. Each
template image contains 8 images, and their order
is described in language (e.g. top left to right, bot-
tom left to right) and the VLM outputs the selected
images according to our prompting as described in
1081
Figure A.6.1076
1078
1079

Model	Cap.	Cau. (%)	Tem. (%)	Des. (%)	All (%)
IG-VLM (Kim et al., 2024)	-	69.8	63.6	74.7	68.6
Tarsier (Wang et al., 2024a)	-	-	-	-	79.2
VideoAgent (Wang et al., 2024d)	8.2	72.7	64.5	81.1	71.3
MVU (Ranasinghe et al., 2024)	16	55.4	48.1	64.1	55.2
MoReVQA (Min et al., 2024)	16	70.2	64.6	-	69.2
VFC (Momeni et al., 2023)	32	45.4	51.6	64.1	51.5
SeViLA ^{\dagger} (Yu et al., 2024a)	32	61.3	61.5	75.6	63.6
VideoTree (Wang et al., 2024g)	(56)	75.2	67.0	81.3	73.5
ProViQ (Choudhury et al., 2023)	60	-	-	-	64.6
TraveLER (Shang et al., 2024)	(65)	70.0	60.5	78.2	68.2
LangRepo (Kahatapitiya et al., 2024)	90	64.4	51.4	69.1	60.9
LLoVi (Zhang et al., 2023)	90	69.5	61.0	75.6	67.7
LVNet (ours)	12	75.0	65.5	81.5	72.9

Table A.9: Extended results on NExT-QA (Xiao et al., 2021). We compare LVNet against prior zero-shot models across different reasoning categories: causal, temporal, and descriptive. LVNet achieves an overall accuracy of 72.9% while using only 12 captions per video, demonstrating strong performance across all reasoning types. Notably, it outperforms all models in causal reasoning (75.0%) and matches the best performance in descriptive reasoning (81.5%), despite processing significantly fewer captions than models like VideoTree (56 captions) and TraveLER (65 captions). Models that utilize video-caption pretraining or process substantially more captions than LVNet are de-emphasized in gray or downplayed in light green to ensure fairness in comparison. Numbers in parentheses () indicate the maximum number of frames used.

Model	Cap.	Why? (%)	How? (%)	B./A. (%)	All (%)
IG-VLM (Kim et al., 2024)	-	-	-	-	65.3
SeViLA ^{\dagger} (Yu et al., 2024a)	32	-	-	-	60.9
VideoTree (Wang et al., 2024g)	(56)	-	-	-	66.9
LangRepo (Kahatapitiya et al., 2024)	90	62.8	62.4	47.8	59.1
LLoVi (Zhang et al., 2023)	90	68.4	67.4	51.1	64.0
LVNet (ours)	12	75.0	74.4	62.1	71.7

Table A.10: **Extended results on IntentQA (Li et al., 2023).** We compare LVNet against prior zero-shot models across different reasoning categories: *Why?, How?*, and *B.A.* (Belief/Action). LVNet achieves an overall accuracy of 71.7%, surpassing all models while using only 12 captions per video. It reaches 75.0% in the *Why?* category, 74.4% in the *How?* category, and 62.1% in the *B.A.* category. Compared to VideoTree, which processes 56 captions and achieves 66.9% accuracy, LVNet outperforms it by +4.8% while using significantly fewer captions. Additionally, LVNet demonstrates superior reasoning-based performance compared to LangRepo (90 captions, 59.1%) and LLoVi (90 captions, 64.0%). Models with video-caption pretraining or utilizing significantly more captions than 12 frames used by LVNet are de-emphasized in grey or downplayed in light green to ensure fairness with image-level pretraining or highlight caption efficiency. Numbers in parentheses () indicate the maximum number of frames used.

A.6 Integration to Existing Methods

1084

1085

1086

1087

1088

1090

1091

1092

1093

1094

1095

1096

Our LVNet has been successfully integrated with other works for evaluation on long video benchmarks. For example, in Ranasinghe et al. (2024) LVNet is integrated with their proposed MVU to gain further performance boosts on the EgoSchema and NextQA benchmarks.

A.7 Comparison with Other Keyframe Selection Methods

We aim to highlight the main advantage of the Hierarchical Keyframe Selector over other existing keyframe selection methods. Models like VideoAgent, VideoTree, and TraveLER provide useful comparisons, as they utilize keyframe selection mechanism with similar or different scale of frames. VideoAgent and TraveLER rely on uniform frame 1098 selection in the first iteration without analyzing 1099 the entire video even though they perform non-1100 uniform sampling in the next iterations. They iden-1101 tify important segments based solely on these ini-1102 tial frames and the LLM's response, which can be 1103 problematic if the initial uniformly selected frames 1104 are not representative of the entire video or if the 1105 LLM misinterprets the captions and prompts. In 1106 such cases, the LLM might incorrectly identify 1107 segments for further analysis. If the LLM fails 1108 to pinpoint the correct segment initially, the entire 1109 process can break down because subsequent frames 1110 will be similar to the first set, leading the LLM to 1111 continuously select frames within or near the ini-1112

Method	LLM Active Params/Type	TS	VT Free	# Frames (n) \uparrow	TS Cap. \downarrow	$FR\uparrow$	Complexity	Accuracy \uparrow
Qwen-VL (Bai et al., 2023)	7B/OS	no	no	4	N/A	-	37.8	
Qwen-VL	7B/OS	no	no	1800	N/A	-	_	OOM
LongVILA (Chen et al., 2024a)	8B/OS	no	no	128	N/A	-	38.8	
LongVILA	8B/OS	no	no	256	N/A	-	_	39.7
LLaVA-OneVision (Li et al., 2024)	7B/OS	no	no	8	N/A	-	_	43.8
VideoChat-T (Zeng et al., 2024)	7B/OS	no	no	128	N/A	87.5	$\mathcal{O}(n)$	41.9
Frame-Voyager (Yu et al., 2024b)	7B/OS	no	no	128	N/A	93.7	$O(n^r)$	48.9
LLaVA-NexT-Video (Zhang et al., 2024a)	34B/OS	no	no	32	N/A	-	_	44.3
Frame-Voyager	34B/OS	no	no	128	N/A	93.7	$O(n^r)$	51.2
InternVL2 (Chen et al., 2024c)	34B/OS	no	no	16	N/A	-	_	52.6
LVNet (DeepSeek-v3)	37B/OS	yes	yes	1800	24	-	$\mathcal{O}(n)$	53.1
VideoAgent+GPT40	<1.8T/PP	yes	yes	-	24.6	-	$\mathcal{O}(m)$	46.4
VideoTree+GPT40	<1.8T/PP	yes	yes	300	98	67.3	$O(n^r)$	53.1
GPT40 direct	<1.8T/PP	yes	yes	1800	1800	-	_	API Error
LVNet (GPT-40)	<1.8T/PP	no	yes	1800	24	98.7	$\mathcal{O}(n)$	53.9
LLaVA-Video (Zhang et al., 2024b)	72B/OS	no	no	64	N/A	-	-	61.5
Qwen2-VL (Wang et al., 2024c)	72B/OS	no	no	768	N/A	-	-	62.2
InternVL2.5 (Chen et al., 2024b)	72B/OS	no	no	64	N/A	-	-	62.6

Table A.11: **Detailed Comparison on VideoMME (Fu et al., 2024)**. This table expands on Table 3 by including two-stage (TS), video-level training free (VT Free), number of input frames used, two-stage captions fed to LLM (TS Cap.), ratio of frames/captions provided to the LLM relative to input frames (FR)., and frame selection complexity. "OS" denotes open-source, "OOM" stands for out-of-memory. O(m) denotes the complexity proportional to the number LLM calls to predict keyframe timelines. The bottom 3 results are from benchmark leaderboard; we are unable to replicate these on our compute resources. Our LVNet variants require no video-level training yet achieve competitive results on very long videos.

tial segment. Additionally, for videos that are as challenging or more difficult than EgoSchema in terms of temporal complexity and activities, existing keyframe selection models such as VideoAgent, VideoTree, and TraveLER may require numerous iterations by running heavy visual/language models to finalize keyframes selection. This results in higher computational and latency costs, as it necessitates numerous runs of resource-intensive VLM and LLM models.

1113

1114

1115

1116

1117

1118

1119 1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138 1139

1140

1141

1142

1143

1144

In contrast, our method analyzes the entire video with high frame rates using a lightweight ResNet-18 (He et al., 2016a) and segments the video nonuniformly based on scene continuity. We then select several frames in each segment by measuring feature similarity between frame features and keywords using the CLIP-B/16 (0.12B) (Ranasinghe et al., 2023) which is lighted than VideoAgent's EVA-CLIP-8Bplus (8B). By reviewing the entire video and non-uniformly selecting keyframes based on scene continuity and similarity scores, these keyframes accurately represent the question-based important frames distribution in the entire video. Furthermore, we use VLM for a fine-grained selection of keyframes, improving keyframe selection when CLIP-B/16 struggles to understand detailed atomic activities in the frames. By hierarchically segmenting the video with different modules, the resulting segments and keyframes are more reliable than those from VideoAgent. Even with more challenging videos, our process only needs to go through the video once to collect keyframes, maintaining computational efficiency.

Figure A.5 visualizes the differences of 1146 the keyframe selection mechanism bewtween 1147 LVNet and VideoAgent. On the left, LVNet begins 1148 with uniformly sampled frames and filters them 1149 through multiple stages, resulting in a non-uniform 1150 distribution of frames over time. First, the tem-1151 poral scene clustering (TSC) selects some frames 1152 that represent temporally distinct activities. Next, 1153 the coarse keyframe detector (CKD) targets frames 1154 most relevant to the question. Finally, the fine 1155 keyframe detector (FKD) further refines this selec-1156 tion to ensure the keyframes accurately capture the 1157 activity in question. As a result, LVNet produces 1158 12 frames, with 8 of them (67%) directly depict-1159 ing "usage of phones," which is the correct answer 1160 and leads the model to select the right option. On 1161 the right, VideoAgent also starts with the uniform 1162 frames but relies on a LLM to request additional 1163 frames. Since the initial frames do not capture 1164 enough relevant content, the LLM again selects 1165 frames uniformly, adding more irrelevant samples 1166 that lack the crucial information about "usage of 1167 phones." As a result, VideoAgent ultimately selects 1168 the wrong option. 1169

Algorithm 1: Temporal Scene Clustering	Algorithm 1	: Temporal	Scene	Clustering
--	-------------	------------	-------	------------

1: Require: ResNet-18 (He et al., 2016b) pretrained on imagenet dataset f, frame list List *frame*, image index list $List_{index} \in \{1, \ldots, N\}$, minimum number of list length ψ , temperature λ , number of sample τ , function to find index of x in list w index(x, w), and function to sort list sort(List) 2: for all img^i in \mathbf{List}_{frame} do $\mathbf{F}^i \leftarrow f(img^i)$ 3: $List_{feat}.insert(\mathbf{F}^i)$ 4: 5: end for 6: for all \mathbf{F}^i in \mathbf{List}_{feat} do $\textbf{List}_{dist} \leftarrow \frac{\sum_{y}\sum_{x}\sqrt{(\textbf{F}_{i}-\textbf{List}_{feat})^{2}}}{x}$ 7: \mathbf{M}_{dist} . insert(\mathbf{List}_{dist}) 8. 9: end for 10: while length of $\mathbf{List}_{index} > \psi \ \mathbf{do}$ 11: $\text{List}_{sample} \leftarrow \emptyset$ 12: $\text{List}_{\delta} \leftarrow \emptyset$ 13: $i \leftarrow \text{List}_{index} . \text{pop}(0)$ 14: $\mathbf{p}^{i} \leftarrow \mathsf{softmax}(\mathbf{M}_{dist}^{i})$ $\mu_{\mathbf{p}^{i}}, \sigma_{\mathbf{p}^{i}} \leftarrow \mathsf{mean}(\mathbf{p}^{i}), \mathsf{std}(\mathbf{p}^{i})$ 15: 16: $\beta \leftarrow \mu_{\mathbf{p}^{\mathbf{i}}} - \sigma_{\mathbf{p}^{\mathbf{i}}} \sum_{i=0} e^{1-i/\lambda}$ for all prob in pⁱ do 17. **if** prob < β **then** 18: 19: $List_{selected}$.insert(index(prob, p^i)) 20: end if 21: end for 22: for all γ in List_{selected} do $\delta \leftarrow \gamma \ th$ value in \mathbf{List}_{index} 23: **List** $_{\delta}$.insert(δ) 24: $List_{index}$.pop(γ) 25 end for 26: $List_{\delta}.insert(i)$ 27: $\mathbf{List}_{sample} \leftarrow \mathtt{sample} \ au \ \mathtt{items} \ \mathtt{from} \ \mathbf{List}_{\delta}$ 28: sort(List_{sample}) 29: for all $frame^j$ in \mathbf{List}_{frame} do 30: 31: if j in List_{sample} then 32: **Outputs**.insert(*frame^j*) 33: end if end for 34. 35: end while

Algorithm 2: Keyword-Image Matching Process in CKD

- 1: **Require:** keyword set **K**, image set **I**, total length of selected image set *L*, function to calculate similarity matrix **sim**(**K**, **I**), function to sort similarity matrix and return indices **sort**(**S**)
- 2: $\mathbf{S} \leftarrow \mathbf{sim}(\mathbf{K}, \mathbf{I})$
- 3: $\mathbf{S}_{sorted}, \mathbf{idx}_{sorted} \leftarrow \mathbf{sort}(\mathbf{S})$
- 4: Initialize \mathbf{P}_{best} as an empty list
- 5: Initialize $\mathbf{I}_{selected}$ as an empty set
- 6: while length of $\mathbf{I}_{selected} < L \, \mathbf{do}$
- 7: for $k \in \mathbf{K}$ do
- 8: for $i \in \mathbf{I}$ do
- 9: $i_{index} \leftarrow \mathbf{idx}_{sorted}[k][i]$
- 10: **if** i_{index} not in **I**_{selected} **then**
- 11: $\mathbf{P}_{best}.insert(k, i_{index})$
- 12: $\mathbf{I}_{selected}.insert(i_{index})$
- 13: break
- 14: end if
- 15: end for
- 16: **if** length of $\mathbf{I}_{selected} \geq L$ then
- 17: break
- 18: end if
- 19: **end for**
- 20: end while
- 21: return \mathbf{P}_{best}

Algorithm 3: Fine Keyframe Detection Pro-	
cess (FKD)	

- Require: keyword set K, image set I, similarity score list S, total length L, number of low similarity indices K, number of frames per visual template N, number of keyframes selected per visual template M, function to sort by similarity sort(S), function to order indices temporally temporal_order()
- 2: $idx_{sorted} \leftarrow sort(S)$
- 3: $\mathbf{idx}_{low_sim} \leftarrow \mathbf{idx}_{sorted}[-K:]$
- 4: $\mathbf{idx}_{temporal} \leftarrow \mathbf{temporal_order}(\mathbf{idx}_{low_sim})$
- 5: $\mathbf{idx}_{final} \leftarrow \mathsf{concatenate}(\mathbf{idx}_{sorted}[: -K], \mathbf{idx}_{temporal})$
- 6: $\mathbf{I}_{ordered}, \mathbf{K}_{ordered} \leftarrow \mathbf{I}[\mathbf{idx}_{final}], \mathbf{K}[\mathbf{idx}_{final}]$ 7: sets \leftarrow
- $create_sets(\mathbf{I}_{ordered}, \mathbf{K}_{ordered}, L//N)$
- 8: for each set \in sets do
- 9: $\mathbf{I}_{selected} \leftarrow \texttt{select_top_M}(\mathbf{set}, M)$
- 10: end for
- 11: return $I_{selected}$



Figure A.6: **Prompt for Fine Keyframe Detection**: The figure illustrates the input image, the prompt provided to the VLM, and the output. The input image represents a visual template composed of eight frames, and the prompt requests the three best frames along with their corresponding keywords. The output displays the top three selected frames and their associated keywords.