VideoGPT+ **Solution:** Integrating Image and Video Encoders for Enhanced Video Understanding

Anonymous ACL submission

Abstract

Building on advances in language models, Large Multimodal Models (LMMs) have significantly improved video understanding. However, current video LMMs rely on either image or video encoders, each with limitations; image encoders capture rich spatial details but lack temporal context, while video encoders provide temporal understanding but process sparse frames at lower resolutions. To this end, we introduce VideoGPT+, which integrates image and video encoders for detailed spatial understanding and global temporal modeling. The model processes videos in segments and applies adaptive pooling on extracted features, achieving state-of-the-art results on VCGBench, MVBench, Zero-shot QA, and Video-MME. Additionally, we develop a 112K video-instruction dataset using a novel semiautomatic annotation pipeline, further enhancing performance. To comprehensively evaluate video LMMs, we present VCGBench-Diverse, a benchmark covering 18 diverse video categories, including lifestyle, sports, and surveillance. With 4,354 QA pairs, it assesses dense video captioning, spatio-temporal understanding, and complex reasoning, ensuring a robust evaluation across video types. Our code, dataset, and models will be released publicly.

1 Introduction

011

014

017

027

037

038

041

Existing methods for video understanding often rely solely on either image encoders or video encoders (Maaz et al., 2024; Jin et al., 2024; Liu et al., 2024c). Most works focus on image encoders, which encode multiple frames and either fuse the information or concatenate the embeddings before passing them to the LLM. When fusing the information, spatial or temporal pooling is typically used (Maaz et al., 2024). Spatial pooling has shown minimal effectiveness in capturing video information, whereas temporal pooling retains some



Figure 1: VideoGPT+ versus various SoTA models. VideoGPT+ performs better compared to various models (Li et al., 2023c; Jin et al., 2024; Lin et al., 2023; Maaz et al., 2024) on video conversation benchmarks: VCGBench (Maaz et al., 2024), Video-MME (Fu et al., 2024), MVBench (Li et al., 2023c), Zero-shot video QA: MSVD-QA, MSRVTT-QA, ActivityNet-QA and VCGBench-Diverse (across dense captioning, spatial understanding, and reasoning).

spatial information but lacks explicit temporal context. On the other hand, concatenating embeddings without pooling (Jin et al., 2024; Liu et al., 2024c; Zhang et al., 2024b) can rapidly increase computational complexity due to the extended context length required by the LLM, limiting the number of frames that can be processed. While this approach provides better spatial representation, the overall context is still limited to few frames. The limited context results in a poor understanding of the video, especially if a uniform sampling strategy is employed, as it only captures small segments of the video, missing important temporal dynamics.

In order to address these challenges, we propose VideoGPT+ which effectively combines the merits of both image and video encoders (see Fig. 2). By leveraging an image encoder for rich spatial details and a video encoder for global temporal context, our model achieves improved video understanding. To model finegrained temporal dynamics in VideoGPT+, we use a segment-wise sampling strategy. Unlike uniform sampling used in existing

video LMMs (Maaz et al., 2024; Li et al., 2023c), which may miss important temporal dynamics, our approach divides the video into smaller segments and applies segment-wise sampling. This ensures that the model captures representative information from different segments of the video, enabling a more comprehensive understanding.

065

066

077

078

087

880

097

100

101

102

103

104

105

106

107

108

110

111

112

113

114

To facilitate the integration of image and video features, VideoGPT+ introduces a visual adapter module that combines their complimentary benefits. It performs projection and pooling operations, mapping both image and video features to a common space while reducing computational complexity. By aligning the features in this manner, the model can utilize the combined spatial and temporal information for improved video understanding.

We demonstrate the effectiveness of VideoGPT+ across five standard video-conversation benchmarks, including VCGBench (Maaz et al., 2024), MVBench (Li et al., 2024), Zero-shot question-answering (Maaz et al., 2024) and Video-MME (Fu et al., 2024), where it performs better than previous SoTA (see Fig. 1). Further, we develop VCG+112K using a novel semi-automatic annotation pipeline (see Fig. 3), which provides dense video captions along with spatial understanding and reasoning-based question-answer (QA) pairs, further enhancing the model's performance. We also propose VCGBench-Diverse, extending VCGBench (Maaz et al., 2024) by including videos from 18 different domains to extensively evaluate the video-based conversation models in diverse domains (see Fig. 4).

Our work has three main contributions:

- We present VideoGPT+, the first videoconversation model that benefits from a dualencoding scheme based on both image and video features. These complimentary sets of features offer rich spatiotemporal details for improved video understanding (Sec. 2).
- Addressing the limitations of existing VideoInstruct100K dataset (Maaz et al., 2024), we develop VCG+112K with a novel semi-automatic annotation pipeline, offering dense video captions along with spatial understanding and reasoning-based QA pairs, improving model performance (Sec. 3).
- Recognizing the lack of diverse benchmarks for video-conversation task, we propose VCGBench-Diverse, which provides 4,354 human annotated QA pairs across 18 video cate-

gories to extensively evaluate the performance of a video-conversation model (Sec. 4).

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

156

157

158

159

160

161

162

163

164

2 Method

For effective video understanding, combining detailed spatial information with explicit temporal context is crucial. To achieve this, we propose VideoGPT+, which features a dual encoder design that leverages the complementary strengths of an image encoder and a video encoder.

Overall Architecture: The overall architecture consists of (i) segment-wise sampling, (ii) dual visual encoder, (iii) vision-language adapters that project vision features to the language domain and (iv) a large language model. Frames selected through a segment-wise sampling strategy are encoded through a dual encoder consisting of an image and a video encoder. Both sets of features are projected to language space using vision-language (V-L) adapters, and the resulting tokens are pooled through adaptive token pooling and concatenated before being fed to the LLM (see Fig. 2).

Segment-wise Sampling: To extract fine-grained temporal cues, we use a segment-wise frame sampling strategy. Given an input video $\mathbf{V} \in \mathbb{R}^{T \times H \times W \times C}$, we divide it into K segments, where each segment consists of $n = \frac{T}{K}$ frames. Thus, the video can be represented as $\mathbf{V} = [\mathbf{V}_k]_{k=1}^K$. Each segment $\mathbf{V}_k \in \mathbb{R}^{n \times H \times W \times C}$ can be described as a sequence of frames, \mathbf{X}_i , where $\mathbf{V}_k = [\mathbf{X}_{i,j}]_{j=1}^n$. The video segments are downsampled to a lower resolution of $n \times h \times w \times c$ for video encoding.

Compared to a uniform sampling, segment-wise sampling better aligns with our dual encoder design. Video encoders often face computational constraints, limiting them to processing only sparse frames. Uniform sampling increases the selfattention computation complexity as it requires attending to features of all frames. Additionally, video encoders are typically trained with sparse frames, and providing more frames can hinder their ability to accurately capture temporal information. In contrast, the segment-wise sampling strategy divides the video into smaller, manageable segments, enabling the video encoder to efficiently capture rich temporal cues within each segment.

Dual Vision Encoder: Our design leverages the complementary strengths of an image encoder that captures detailed spatial features and a video encoder that provides explicit temporal context. The image encoder g, processes T frames, $g(\mathbf{X}) \in$



Figure 2: **Overview of VideoGPT+**. VideoGPT+ is a large multimodal model for video understanding. It uses a dual-encoder design that combines the complementary strengths of an image encoder and a video encoder. The image encoder captures detailed spatial features, while the video encoder captures temporal dynamics across multiple frames. To retain fine-grained temporal details while ensuring efficiency, we use segment-wise frame sampling instead of random sparse sampling. Both sets of features are then projected into a unified space through Vision-Language (V-L) projection layers and the resulting tokens are pooled and concatenated before being processed by a Large Language Model to generate comprehensive responses to video-based questions. Symbols indicates frozen components, indicates trainable components, and the indicates LoRA-training.

 $\mathbb{R}^{T \times H_g \times W_g \times D_g}$, producing local features that provide frame-level context. Meanwhile, the video encoder h, operates on low-resolution video segments \mathbf{V}_k , yielding global features that provide segment-wise context, $h(\mathbf{V}_k) \in \mathbb{R}^{n \times h_h \times w_h \times D_h}$.

166

167

168

169

171

172

173

174

176

177

179

180

181

The primary goal of VideoGPT+ is to leverage the capabilities of a pre-trained LLM alongside visual modalities from both a pre-trained image encoder and a pre-trained video encoder. Specifically, we utilize the pre-trained CLIP model, ViT-L/14 (336×336) (Radford et al., 2021) as the image encoder, and InternVideo-v2 (224×224) (Wang et al., 2024) as the video encoder. These models are selected for their robust performance and their ability to complement each other in capturing both spatial and temporal information. Both encoders are pre-trained on large-scale datasets in a multimodal setting using contrastive loss, facilitating their integration within our architecture.

Visual Adapter: The output embeddings from the 184 second last layer of both image and video encoders 185 are passed through separate V-L projection layers, W_q and W_h , respectively. These Multi-Layer per-187 ceptrons (MLPs) project the visual features into the language space. The projection layers are trainable, while the visual encoders remain frozen, preserving 191 the rich, pre-trained representations. The projected embeddings are reshaped back into their grid forms 192 and subjected to a 2×2 adaptive token pooling, 193 which operates on the spatial dimensions of the local and global features. This pooling reduces the 195

token length by a factor of 4, thereby allowing to fit in larger visual context within the same LLM context window. The pooled embeddings from the local features form $\mathbf{E}^{img} \in \mathbb{R}^{T \times h_g \times w_g \times D_t}$, while the pooled embeddings from the global features of each segment form $\mathbf{E}^{vid} \in \mathbb{R}^{n \times h_h \times w_h \times D_t}$.

197

198

199

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

Large Language Model: We obtain the final representation by concatenating the embeddings \mathbf{E}^{img} with K segment-wise embeddings \mathbf{E}^{vid} , such that we have detailed spatial representation across all segments followed by their global temporal context. We then concatenate the text embeddings $\mathbf{E}^{text} \in \mathbb{R}^{L \times D_t}$ of the user text query with the visual embeddings,

$$\mathbf{E} = [\mathbf{E}^{img}, \mathbf{E}_1^{vid}, \dots, \mathbf{E}_K^{vid}, \mathbf{E}^{text}].$$
(1)

This integration ensures that the LLM receives a sequence of embeddings that include detailed spatial features from the image encoder and comprehensive temporal context from the video encoder, allowing for robust video understanding. The LLM is fine-tuned using LoRA (Hu et al., 2021) in an auto-regressive manner with a next-token prediction loss. Refer to Fig. 2 for detailed illustration.

3 Dataset

Video-ChatGPT (Maaz et al., 2024) introduces the VideoInstruct100K, which employs a semiautomatic annotation pipeline to generate 100K instruction-tuning QA pairs. To address the limitations of this annotation process, we present



Figure 3: **Illustration of the semi-automatic annotation process in VCG+112K**. The figure shows how we use ground-truth video captions and frame-level descriptions to generate a detailed video description. GPT-4 is used to remove irrelevant and conflicting noisy information in the frame-level descriptions to produce a high-quality video description. The semi-automatic annotation process integrates spatial, temporal and event, and reasoning details into the brief information we start with. This dense video description is then used to generate instruction-tuning QA pairs using GPT-3.5. We provide detailed prompts used in both stages in Appendix F (see Figs. 8 and 9). We also compare the video description in the VideoInstruct100K (Maaz et al., 2024) dataset to show the improvement in quality achieved by our new annotation pipeline.

 VCG+ 112K dataset developed through an improved annotation pipeline. Our approach improves the accuracy and quality of instruction tuning pairs by improving keyframe extraction, leveraging SoTA LMMs for detailed descriptions, and refining the instruction generation strategy.

227

235

240

241

242

Keyframe Extraction: VideoInstruct100K uses a fixed number of video keyframes, regardless of video length or dynamics, to generate frame-level dense captions. This often results in both insufficient and redundant information. We address this by first extracting scenes from videos (Castellano, 2022), and then selecting one keyframe/scene. Consequently, we obtain detailed information for videos with rich content and reduce redundancy for videos with less content. It provides better visual context by extracting more stable keyframes, thus offering a more accurate video representation.

Frame-Level Descriptions: After extracting 243 keyframes, we use an image LMM, LLaVA-244 v1.6 (Liu et al., 2024a), to generate dense descrip-245 tions for each keyframe. These descriptions en-246 compass visual details, including spatial attributes, 247 scene context, and object characteristics, which are often absent in concise ground truth captions. While ground truth captions are precise, they lack 251 the granularity to capture intricate visual and spatial information. To address this, we augment them with detailed but noisy information from the framelevel descriptions, thus enhancing the quality and accuracy of the subsequent video descriptions. 255

Detailed Video Descriptions: VideoInstruct100K prompts GPT-3.5 directly with frame-level descriptions and concise ground truth captions to generate QA pairs, imposing a significant cognitive load on the model to verify frame-level descriptions with the ground truth. We improve this process by first creating a coherent and detailed video description. We prompt GPT-4 to integrate the detailed framelevel descriptions with the ground truth captions by comparing information and removing any inconsistencies. The resulting detailed descriptions include a timeline of events, actions, object attributes, and scene settings, providing a thorough representation of the video content. This structured input simplifies the task for LLM, thereby enhancing the generated QA pairs quality.

256

257

258

259

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

282

Improved Instruction Tuning Data: Using the ground truth captions and detailed video descriptions, we generate two types of QA pairs using GPT-3.5: descriptive and concise. For **descriptive** instruction pairs, we focus on three categories: (i) *dense captioning*, which provides descriptions of the video covering the entire sequence of events and visual details; (ii) *detailed temporal information*, which addresses the sequence of events and their dependency to learn temporal relationships; and (iii) *generic question answering*, which involves in-depth questions about different actions, their consequences, and other detailed aspects of the video. For **concise** instruction pairs, we target (i) *spatial reasoning*, focusing on understanding



Figure 4: **Illustration of VCGBench-Diverse video conversational benchmark**. VCGBench-Diverse comprehensive benchmark is designed to evaluate video LMMs across 18 broad video categories. With 4,354 QA pairs, VCGBench-Diverse tests generalization on dense video captioning, spatial and temporal understanding, and complex reasoning. It covers five video-capturing methods, ensuring diversity and robust generalization and six reasoning complexities assessing various analytical and comprehension skills.

and describing spatial details such as scene settings, number of objects, attire, and locations; (ii) *reasoning* of events, covering the causal relationships between events; and (iii) *short temporal questions*, addressing specific moments or sequences, such as what happened at the beginning or end.

4 Proposed Benchmark

Recognizing the limited diversity in existing video conversation benchmarks, we introduce VCGBench-Diverse to comprehensively evaluate generalization ability of video LMMs. While VCG-Bench (Maaz et al., 2024) provides an extensive evaluation protocol, it is limited to videos from the ActivityNet200 (Fabian Caba Heilbron and Niebles, 2015) dataset. Our benchmark comprises a total of 877 videos, 18 broad video categories and 4,354 QA pairs, ensuring a robust evaluation framework. The detailed breakdown of VCGBench-Diverse is illustrated in Fig. 4, showcasing the distribution of videos across content domains, video capturing methods, and reasoning complexities.

We collect videos from *18 distinct domains*, including lifestyle, how-to, science and technology, news, travel, entertainment, film, sports, comedy, activism, gaming, education, surveillance, pets, cooking, music, automobile, and traffic These categories encompass a broad spectrum of real-world scenarios, ensuring that models are evaluated on a diverse set of challenges. In addition to content diversity, VCGBench-Diverse includes a variety of *video capture methods*, which ensures a comprehensive assessment of robustness to different filming techniques, camera movements, quality levels and lighting. The benchmark covers *five* video capture methods including static and controlled settings, dynamic and unpredictable settings, fixed camera perspectives, professional and high-quality videos, and uncontrolled and variable quality. Further, the benchmark evaluates models across *six reasoning complexities*, including sequential understanding, complex action and predictive reasoning, contextual and world knowledge reasoning, causal reasoning, narrative and emotional reasoning, and analytical and critical reasoning, which is crucial for understanding diverse video content. 321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

337

338

339

340

341

342

343

344

345

349

350

351

352

355

The videos in VCGBench-Diverse are sourced from HDVILA (Xue et al., 2022), MPII (Andriluka et al., 2014), YouCook2 (Zhou et al., 2018), UCF Crime (Sultani et al., 2018), and STUD Traffic (Xu et al., 2021). The video durations range from 29 sec to 471 sec, with an average of 217 sec. Human annotators are tasked with writing detailed descriptions based on their understanding of audio and visual elements of the videos. This comprehensive annotation process involves a set of annotators who are provided with an initial set of ten videos each. These annotations undergo a meta-review stage where feedback is provided, and necessary corrections are made to meet the required standards. Following this, annotators receive additional batches, with random samples being selected for quality checks by the meta-reviewer. The final human annotations are utilized to generate QA pairs using GPT-3.5, based on prompts detailed in Fig. 10.

Following VCG-Bench (Maaz et al., 2024), the evaluation is performed over five different aspects: (i) correctness of information (ii) detail orientation (iii) contextual understanding (iv) temporal understanding and (v) consistency. Additionally,

310

312

313

316

317

320

287

CI	DO	CU	TU	CO	Avg.
2.40	2.52	2.62	1.98	2.37	2.38
2.68	2.69	3.27	2.34	2.46	2.69
2.78	<u>3.10</u>	3.40	2.49	2.47	2.85
2.89	2.91	3.46	2.89	2.81	2.99
2.96	3.00	<u>3.53</u>	2.46	2.51	2.89
2.84	2.86	3.44	2.46	2.57	2.81
3.02	2.88	3.51	2.66	2.81	2.98
<u>3.11</u>	2.78	3.51	2.44	<u>3.29</u>	<u>3.03</u>
3.27	3.18	3.74	<u>2.83</u>	3.39	3.28
	CI 2.40 2.68 2.78 2.89 2.96 2.84 3.02 3.11 3.27	CI DO 2.40 2.52 2.68 2.69 2.78 3.10 2.89 2.91 2.96 3.00 2.84 2.86 3.02 2.88 3.11 2.78 3.27 3.18	CI DO CU 2.40 2.52 2.62 2.68 2.69 3.27 2.78 3.10 3.40 2.89 2.91 3.46 2.96 3.00 3.53 2.84 2.86 3.44 3.02 2.88 3.51 3.11 2.78 3.51 3.12 3.18 3.74	CI DO CU TU 2.40 2.52 2.62 1.98 2.68 2.69 3.27 2.34 2.78 3.10 3.40 2.49 2.89 2.91 3.46 2.89 2.96 3.00 3.53 2.46 2.84 2.86 3.44 2.46 3.02 2.88 3.51 2.66 3.11 2.78 3.51 2.44 3.27 3.18 3.74 2.83	CIDOCUTUCO2.402.522.621.982.372.682.693.272.342.462.783.103.402.492.472.892.913.462.892.812.963.003.532.462.512.842.863.442.462.573.022.883.512.662.813.112.783.512.443.29 3.273.183.74 2.83 3.39

Table 1: VideoGPT+ on VCGBench (Maaz et al., 2024). All models use 16 frames except Video-ChatGPT and Chat-UniVi which use 100 and 64 frames respectively.

VCGBench-Diverse provides a performance breakdown across three key aspects: (i) dense video captioning, which assesses the ability to generate detailed and accurate video descriptions, (ii) spatial understanding, which evaluates the capability to understand and describe the spatial relationships, and (iii) reasoning, which tests the adeptness in inferring and explaining causal relationships and actions within the video.

Experiments 5

357

363

371

374

384

perform quantitative evaluation We of VideoGPT+ on five standard benchmarks: i) VCG-Bench (Maaz et al., 2024), ii) VCGBench-Diverse, iii) MVBench (Li et al., 2024), iv) Video-MME (Fu et al., 2024) and v) Zero-shot QA.

Implementation Details: We use CLIP-L/14 (Radford et al., 2021) as our image encoder, InternVideo-v2 (Wang et al., 2024) stage-2 1B model as our video encoder in conjunction with Phi-3-Mini-3.8B (Abdin et al., 2024) LLM with 4K context in our experiments. The image encoder operates at 336×336 , while the video encoder operates at 224×224 resolution. Our training consists of two pretraining stages and one instruction-tuning stage. In the pretraining stage, we train with only the image encoder and only the video encoder on the CC-595K dataset (Liu et al., 2023a), with only 382 the visual adapters being learned while the rest of the model is kept frozen. During the instructiontuning stage, we use LoRA (Hu et al., 2022) with r = 64 for LLM, while visual adapters are fully trained and vision encoders are kept frozen. The LR is set to $1e^{-3}$ during pretraining and $2e^{-4}$ dur-388 ing instruction tuning. Please refer to Appendix. B for additional implementation details. 390

VCGBench: The benchmark consists of around 391 3000 QA pairs generated from 500 humanannotated videos. It evaluates responses based on five aspects: i) CI (Correctness of Information) -394

accuracy of the response with video content, ii) DO (Detail Orientation) - depth of the response, iii) CU (Contextual Understanding) - alignment with video context, iv) TU (Temporal Understanding) - accuracy in identifying temporal sequences, and v) CO (Consistency) - response consistency to similar questions. Table 1 compares our model with previous SoTA approaches. VideoGPT+ achieves an average score of 3.28 surpassing previous best method by a margin of 0.25 (5%).

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

VCGBench-Diverse: We provide a quantitative comparison of VideoGPT+ against previous SoTA approaches on VCGBench-Diverse, which contains 4,354 QA pairs from 877 videos. Following (Maaz et al., 2024), we evaluate the Correctness of Information (CI), Detail Orientation (DO), Contextual Understanding (CU), Temporal Understanding (TU), and Consistency (CO). Additionally, we provide results for dense captioning, spatial understanding, and visual reasoning abilities. The results are presented in Table 2. VideoGPT+ achieves an average score of 2.47 surpassing all previous methods. Further, we achieves a score of 1.38, 2.80, and 3.63 on dense captioning, spatial understanding, and visual reasoning, respectively. Notably, VideoGPT+ achieves improvements in spatial and temporal understanding, surpassing previous best models by 0.37 (7.4%) and 0.23 (4.6%), respectively. This is attributed to the dual encoder architecture, where the high-resolution image encoder enhances spatial understanding and the video encoder improves temporal accuracy.

To further validate the alignment of GPT scores with human preferences, we conduct a study involving human annotators. Four annotators given the same GPT scoring guidelines, each reviewed 50 questions from a pool of 200 randomly selected questions. They scored responses from three models: VideoGPT+, VideoChat2, and Chat-UniV. Their respective scores, 2.0, 1.9, and 2.3, closely matched the GPT averages of 2.3, 2.2, and 2.5 for each model. This comparison confirms that GPT scores align well with human preferences, supporting the reliability of our evaluation method.

Table 2 also shows the results of closed-source models in gray for reference. Note that the comparison between open-source and significantly larger, closed-source models is not fair due to the vast differences in scale, parameters, and training data. We compare VideoGPT+ (3.8B-scale) with similarly scaled open-source models (7B-scale), where our model demonstrates superior performance.

Method	CI	DO	CU	TU	СО	Avg.	Caption	Spatial	Reasoning
GPT4o-mini-2024-07-18	3.06	3.05	3.43	2.67	3.47	3.14	1.82	3.16	4.19
Gemini-Pro-1.5-Flash-001	3.15	3.24	3.40	2.68	3.32	3.16	2.30	3.48	3.82
Video-ChatGPT (ACL 2024) (Maaz et al., 2024)	2.07	2.42	2.46	1.39	2.06	2.08	0.89	2.25	3.60
BT-Adapter (CVPR 2024) (Liu et al., 2024b)	2.20	2.62	2.59	1.29	2.27	2.19	1.03	2.35	<u>3.62</u>
VTimeLLM (CVPR 2024) (Huang et al., 2024a)	2.16	2.41	2.48	1.46	2.35	2.17	1.13	2.29	3.45
Chat-UniVi (CVPR 2024) (Jin et al., 2024)	2.29	2.56	<u>2.66</u>	1.56	<u>2.36</u>	<u>2.29</u>	<u>1.33</u>	2.36	3.59
VideoChat2 (CVPR 2024) (Li et al., 2024)	2.13	2.42	2.51	1.66	2.27	2.20	1.26	2.43	3.13
VideoGPT+ (ours)	2.46	2.73	2.81	1.78	2.59	2.47	1.38	2.80	3.63

Table 2: Performance of VideoGPT+ on VCGBench-Diverse. All open-source models use 16 frames except Video-ChatGPT and Chat-UniVi, which use 100 and 64 frames, respectively. The good performance of our VideoGPT+ model on VCGBench-Diverse shows its generalization to diverse scenarios.

Model	Avg.	Table 3: MVBench.
Random	27.3	Comparison
GP1-4V (OpenAI, 2023)	45.5	of VideoGPT+
Otter-V (Li et al., 2023a) mPLUG Owl V (Ve et al. 2023)	26.8	methods. See
Video-ChatGPT (Maaz et al., 2024)	32.7	Tab. 6 in App. C
VideoLLaMA (Zhang et al., 2023)	34.1	for complete
VideoChat (Li et al., 2023c)	35.5	results on 20
VideoChat2 (Li et al., 2024)	51.1	sub-categories.
VideoGPI+ (ours)	58.7	

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

MVBench: We evaluate VideoGPT+ on MVBench (Li et al., 2024), which provides 4,000 QA pairs from 11 video datasets covering a broad spectrum of scenes, ranging from first-person to third-person and from indoor to outdoor environments. The tasks are categorized into 20 finegrained temporal understanding tasks. The results presented in Table 3 compare VideoGPT+ with previous methods, indicating an overall improvement of 7.6% compared to the previous best, VideoChat2. Please refer to Table 6 in Appendix C for complete results on 20 categories.

Video-MME: We evaluate the performance of our model on Video-MME, a more comprehensive benchmark that assesses video understanding across six domains and 30 subfields through 2700 multiple-choice-ga pairs from 900 videos. It covers a diverse range of video durations, from short, medium, and long videos (11 sec to 1 hour). Table 4 shows that our model achieves better performance compared to prior SoTA approaches. Specifically, our model performs well across the short, medium, and long video categories, demonstrating strong temporal understanding and effectively capturing long-range dependencies.

Zero-shot Question-Answering: We provide a 472 473 quantitative comparison of our method on the zeroshot QA task across four open-ended QA datasets, 474 including MSVD-QA (Xu et al., 2017), MSRVTT-475 QA (Xu et al., 2017), TGIF-QA (Jang et al., 2019), 476 and ActivityNet-QA (Fabian Caba Heilbron and 477

Model	Short	Med	Long	Avg
Video-LLaVA	45.3	38.0	36.2	39.9
Qwen-VL-Chat	46.9	38.7	37.8	41.1
ChatUniVi	45.7	40.3	35.8	40.6
VideoChat2	48.3	37.0	33.2	39.5
VideoGPT+	56.4	47.2	42.5	48. 7

Table 4:	Performance	comparison	of different	models on
short, med	lium, and long	g video segme	ents in Video	-MME.

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

Niebles, 2015). Results presented in Table 5 show VideoGPT+ achieves superior performance compared to previous methods, indicating its ability to adapt effectively to unseen videos and generate accurate contextually relevant responses in challenging settings.

Vision Encoder Type: We ablate our dual visual encoder design in VideoGPT+. We ablate three settings: using only the image encoder, only the video encoder, and both encoders. The results shows that our dual encoder design effectively combines both spatial and temporal information and achieves the highest score on both VCGBench and VCGBench-Diverse.

Note that the image encoder operates at a higher resolution of 336×336 , while the video encoder operates at 224×224 . The image encoder captures better spatial information and fine-grained details, while the video encoder contributes to understanding motion and action sequences. We further verify this on MVBench action categories including action sequence (+3.6%), action antonym (+1.5%), fine-grained action (+1.5%) and unexpected action (+4.0%), where video-only model performs better than the image-only model.

For completeness, we use a best response selection method with GPT4-as-a-judge to evaluate different model designs. Responses from three model variants: image encoder, video encoder and our dual encoder design are presented anonymously to GPT4 alongside the ground truth. The model selects the best response among the three and excludes cases with no clear winner. For

Model	MSVD	-QA	MSRVT	T-QA	TGIF-	QA	ActivityNet-QA		
	Accuracy	Score	Accuracy	Score	Accuracy	Score	Accuracy	Score	
FrozenBiLM (Yang et al., 2022)	32.2	-	16.8	-	41.0	-	24.7	-	
VideoChat (Li et al., 2023c)	56.3	2.8	45.0	2.5	34.4	2.3	26.5	2.2	
LLaMA Adapter (Zhang et al., 2024a)	54.9	3.1	43.8	2.7	-	-	34.2	2.7	
Video-LLaMA (Zhang et al., 2023)	51.6	2.5	29.6	1.8	-	-	12.4	1.1	
Video-ChatGPT (Maaz et al., 2024)	64.9	3.3	49.3	2.8	51.4	3.0	35.2	2.8	
ChatUniVi (Jin et al., 2024)	65.0	3.6	54.6	3.1	60.3	3.4	45.8	3.2	
LLaMA-VID (Li et al., 2023d)	70.0	<u>3.7</u>	58.9	3.3	-	-	47.5	<u>3.3</u>	
Video-LLaVA (Lin et al., 2023)	<u>70.7</u>	3.9	<u>59.2</u>	<u>3.5</u>	<u>70.0</u>	<u>4.0</u>	45.3	<u>3.3</u>	
VideChat2 (Li et al., 2024)	70.0	3.9	54.1	3.3	-	-	<u>49.1</u>	<u>3.3</u>	
VideoGPT+ (ours)	72.4	3.9	60.6	3.6	74.6	4.1	50.6	3.6	

Table 5: **Performance of VideoGPT+ on Zero-shot QA.** All the models are evaluated in zero-shot setting where none of the videos were included in the training set. VideoGPT+ achieves good results on all datasets.

511 VCGBench (VCG), 732 out of 2000 samples were scored, where the dual encoder design was pre-512 ferred in 51% of cases, compared to 22% for the 513 514 image encoder and 27% for the video encoder. For VCGBench-Diverse (VCG-Div), 792 out of 4354 515 samples were scored, with the dual encoder pre-516 ferred in 42% of cases, compared to 28% for the 517 image encoder and 30% for the video encoder, in-518 519 dicating that our dual encoding design as a clear winner among other uni-encoder alternatives. 520

521

522

524

526

528

530

531

532

533

534

535

536

537

538

541

Vision			Temporal	Spatial	GPT4 as Judge			
Encoder	VCG	VCG-Div	Score	Score	VCG	VCG-Div		
Image-only	3.17	2.36	1.61	2.70	22	28		
Video-only	3.20	2.38	1.69	2.64	27	30		
Dual (ours)	3.28	2.47	1.78	2.80	51	42		

Frame-level and Video-level Feature Fusion: Though our design uses some known components, their meticulous combination to develop an efficient pipeline for video understanding in MLLMs has not been demonstrated. We ablate our approach with two alternatives: i) Without segment-wise sampling, resulting in less effective temporal information captured by the video encoder impacting performance; ii) Without adaptive token pooling, which limits the model's ability to utilize the LLM context length effectively, restricting the model to fewer frames. The performance on VCGBench and VCGBench-Diverse benchmarks indicates the effectiveness of our proposed fusion strategy.

Setting	VCG	VCG-Div
w/o Segment-wise Sampling	3.21	2.40
w/o Adaptive Pooling	3.08	2.31
Video-GPT+ (ours)	3.28	2.47

VCG+ 112K: To demonstrate the effectiveness of VCG+ 112K, we train VideoGPT+ with and without it and report its impact on the performance across multiple benchmarks, including VCGBench, MVBench, VCGBench-Diverse and VideoMME. On VCGBench, our data improves performance, particularly in detail orientation (DO) and tem-

poral understanding (TU). The performance on MVBench shows minimal gains when incorporating the VCG+112k data. This is attributed to the distribution differences, as MVBench predominantly includes short videos averaging 5-40 seconds, whereas the VCG+112k dataset comprises videos from ActivityNet with an average duration of 3 minutes. However VCGBench-Diverse and VideoMME, do not include data from ActivityNet, ensuring a fair evaluation. The results shows improvement on both VCGBench-Diverse and VideoMME. This improvement can be attributed to our novel semi-automatic annotation pipeline and the enhanced instruction tuning data, which focuses on generating both detailed and concise instruction pairs. Refer to Fig. 3 for qualitative visualization of the data.

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

VCG+ 112K	VCG	MVBench	VCG-Div	VideoMME
\checkmark	3.17	58.7	2.4	46.2
×	3.28	58.8	2.5	48.7

6 Conclusion

In this work, we introduce VideoGPT+, a novel video conversation model that leverages the complementary benefits of image and video encoders to achieve enhanced video understanding. VideoGPT+ demonstrates better performance across multiple video benchmarks, owing to its dual-encoder design, lightweight visual adapters that map image/video features to a common space and a segment-wise sampling strategy that retains fine-grained temporal information. We also develop VCG+112K, a 112K video-instruction set using a resource-efficient semi-automated annotation pipeline that delivers further gains. Lastly, we propose VCGBench-Diverse, a diverse benchmark covering 18 video categories, to comprehensively evaluate video LMMs.

576

7 Limitations

Despite reported improvements, video LMMs still find challenges in precise action localization, un-578 derstanding very long videos, and navigating long 579 paths; areas where major improvements can un-580 lock new applications. Further, the use of closedsource LLMs (e.g., GPT-3.5 and GPT-4) for open-582 ended evaluation of video conversations limits re-583 producibility. Although we have designed our evaluation prompts to be model-agnostic, switching to a different LLM for evaluation (e.g., in case a proprietary model is discontinued) may lead to minor variations in the results. Therefore, designing a comprehensive and reliable evaluation metric for open-ended video conversation evaluation is highly desirable to ensure consistency and reproducibility 591 in future research.

References

594

595

611

612

613

615

616

617

618

619

621

623

627

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, and 1 others. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Meta AI. 2024. Llama 3. https://llama.meta.com/ llama3.
- Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2014. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.*
- Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. 2021. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Brandon Castellano. 2022. Pyscenedetect: Automated video scene detection. https://github. com/Breakthrough/PySceneDetect.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. https://lmsys.org/blog/ 2023-03-30-vicuna.
- Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, and 1 others. 2023. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. arXiv preprint arXiv:2312.16886.

Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, and 1 others. 2024. Mobilevlm v2: Faster and stronger baseline for vision language model. *arXiv preprint arXiv:2402.03766*.

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

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

- Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Niebles. 2015. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.*
- Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, and 1 others. 2024. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*.
- Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, and 1 others. 2017. The" something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Machine Learning*.
- Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu Zhu. 2024a. Vtimellm: Empower llm to grasp video moments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.*
- De-An Huang, Shijia Liao, Subhashree Radhakrishnan, Hongxu Yin, Pavlo Molchanov, Zhiding Yu, and Jan Kautz. 2024b. Lita: Language instructed temporal-localization assistant. *arXiv preprint arXiv:2403.19046*.
- Yunseok Jang, Yale Song, Chris Dongjoo Kim, Youngjae Yu, Youngjin Kim, and Gunhee Kim. 2019. Video Question Answering with Spatio-Temporal Reasoning. *International Journal of Computer Vision*.
- Peng Jin, Ryuichi Takanobu, Caiwan Zhang, Xiaochun Cao, and Li Yuan. 2024. Chat-univi: Unified visual representation empowers large language models with image and video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.*
- Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio

- 695 698 700 701 703 704 710 711 712 715 717 718 719 721 726 727 728 729 731 733 734

737

Viola, Tim Green, Trevor Back, Paul Natsev, and 1 others. 2017. The kinetics human action video dataset. arXiv preprint arXiv:1705.06950.

- Wonkyun Kim, Changin Choi, Wonseok Lee, and Wonjong Rhee. 2024. An image grid can be worth a video: Zero-shot video question answering using a vlm. arXiv preprint arXiv:2403.18406.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multi-modal model with in-context instruction tuning. arXiv preprint arXiv:2305.03726.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. BLIP-2: bootstrapping language-image pretraining with frozen image encoders and large language models. In International Conference on Machine Learning.
- Kunchang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023c. Videochat: Chat-centric video understanding. arXiv preprint arXiv:2305.06355.
- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, and 1 others. 2024. Mybench: A comprehensive multi-modal video understanding benchmark. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Yanwei Li, Chengyao Wang, and Jiaya Jia. 2023d. Llama-vid: An image is worth 2 tokens in large language models. arXiv preprint arXiv:2311.17043.
- Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. Video-llava: Learning united visual representation by alignment before projection. arXiv preprint arXiv:2311.10122.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning. arXiv preprint arXiv:2310.03744.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023b. Improved baselines with visual instruction tuning. arXiv:2310.03744.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024a. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning. In Advances in Neural Information Processing Systems.
- Ruyang Liu, Chen Li, Yixiao Ge, Ying Shan, Thomas H Li, and Ge Li. 2024b. One for all: Video conversation is feasible without video instruction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Ruyang Liu, Chen Li, Haoran Tang, Yixiao Ge, Ying Shan, and Ge Li. 2024c. St-llm: Large language models are effective temporal learners. arXiv preprint arXiv:2404.00308.

Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2024. Video-chatgpt: Towards detailed video understanding via large vision and language models. In Association for Computational Linguistics.

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

758

759

760

761

763

764

765

766

767

769

770

771

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

- Shehan Munasinghe, Rusiru Thushara, Muhammad Maaz, Hanoona Abdul Rasheed, Salman Khan, Mubarak Shah, and Fahad Khan. 2023. Pg-videollava: Pixel grounding large video-language models. ArXiv 2311.13435.
- OpenAI. 2023. Gpt-4v(ision) system card. https://api.semanticscholar.org/CorpusID: 263218031.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. ArXiv, abs/2306.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In International Conference on Machine Learning.
- Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M. Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S. Khan. 2024. Glamm: Pixel grounding large multimodal model. The IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, and 1 others. 2024. Moviechat: From dense token to sparse memory for long video understanding. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Waqas Sultani, Chen Chen, and Mubarak Shah. 2018. Real-world anomaly detection in surveillance videos. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 6479-6488.
- Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Jilan Xu, Zun Wang, and 1 others. 2024. Internvideo2: Scaling video foundation models for multimodal video understanding. arXiv preprint arXiv:2403.15377.
- 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, pages 9777-9786.
- 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 International Conference on Multimedia.

Li Xu, He Huang, and Jun Liu. 2021. Sutd-trafficqa: A question answering benchmark and an efficient network for video reasoning over traffic events. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9878–9888.

793

794

796

797

799

806

807

810

811

813

815

816

817

818

820

822

825

826

827

829

834

839

840

841

842

846

- Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and Baining Guo. 2022. Advancing high-resolution videolanguage representation with large-scale video transcriptions. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Zero-shot video question answering via frozen bidirectional language models. In Advances in Neural Information Processing Systems.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, and 1 others. 2023. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*.
- Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B Tenenbaum. 2019. Clevrer: Collision events for video representation and reasoning. *arXiv preprint arXiv:1910.01442*.
- Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. 2023. Ferret: Refer and ground anything anywhere at any granularity. arXiv preprint arXiv:2310.07704.
- Hang Zhang, Xin Li, and Lidong Bing. 2023. Videollama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.
- Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. 2024a. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. In *International Conference on Learning Representations*.
- Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. 2024b. Llava-next: A strong zero-shot video understanding model.
- Luowei Zhou, Chenliang Xu, and Jason Corso. 2018. Towards automatic learning of procedures from web instructional videos. In *AAAI*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2024. Minigpt-4: Enhancing vision-language understanding with advanced large language models. In *International Conference on Learning Representations*.

Supplemental Material

We provide supplementary material for a deeper understanding and more analysis related to the main paper, arranged as follows: 1. Related Works (Appendix A)

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

- 2. Additional Implementation Details (Ap-
- pendix B)
- 3. Complete MVBench Results (Appendix C)
- 4. Qualitative results (Appendix D)
- 5. Additional ablations (Appendix E)
- 6. GPT Prompts (Appendix F)
- 7. Ethics and societal impact (Appendix G)

A Related Works

Building on advances in language models, LLMs offer a flexible interface for various multimodal applications. Early efforts in image-based conversation models such as BLIP-2 (Li et al., 2023b), MiniGPT-4 (Zhu et al., 2024) and LLaVA (Liu et al., 2023c,b) project image features into the language space through a learnable module and perform instruction tuning for visual conversations capabilities. Other efforts extend these models to visual grounding tasks (Peng et al., 2023; Rasheed et al., 2024; You et al., 2023), exploring the potential of LLMs in complex vision tasks.

Video Conversation Models: Initial works like Video-ChatGPT (Maaz et al., 2024) and Video-LLaMA (Zhang et al., 2023) extend image-based models to the video domain by introducing components to encode temporal features, where framelevel visual features are fed to the LLM. However, this is computationally expensive and quickly fills its context window. To address this issue, Video-ChatGPT (Maaz et al., 2024) employs spatial and temporal pooling. LLaMA-Vid (Li et al., 2023d) proposes representing a single image with two tokens, context and content. IG-VLM (Kim et al., 2024) treats a video as a grid of images, while LITA (Huang et al., 2024b) employs slow-fast token pooling to reduce the number of visual features. Chat-UniVi (Jin et al., 2024) uses clustering in both spatial and temporal dimensions to merge tokens, and VideoChat (Li et al., 2023c) uses Q-Former (Li et al., 2023b) to learn a fixed number of queries by cross-attending to the visual features. MobileVLM (Chu et al., 2023, 2024) utilize a lightweight CNN to reduce the spatial dimensions. Other notable methods include (Liu et al., 2024b; Lin et al., 2023; Munasinghe et al., 2023; Song et al., 2024; Huang et al., 2024a).

Alternatively, methods such as VideoChat2 (Li et al., 2024) use pretrained video encoders. Although video encoders provide temporal context, they are limited by computational constraints, operating with limited frames at lower resolutions, restricting temporal context and spatial understanding. Our VideoGPT+ model addresses these issues by using segment-wise sampling and effectively combining image and video encoders to capture rich spatial and temporal details (see Fig. 2).

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

919

921

923

925

927

Video Instruction Tuning **Datasets:** VideoChat (Li et al., 2023c) builds a videoinstruction tuning dataset consisting of 7K instructions using videos from WebVid-10M (Bain et al., 2021). Video-ChatGPT (Maaz et al., 2024) introduces a semi-automatic annotation pipeline to generate VideoInstruct100K using videos from ActivityNet (Fabian Caba Heilbron and Niebles, 2015). VideoChat2 (Li et al., 2024) combines multiple existing image and video datasets to develop a 1.9M joint image-video instruction tuning dataset. In our experiments, we use VideoInstruct100K and a subset of the dataset from VideoChat2. Additionally, addressing the limitations of the VideoInstruct100K dataset (Maaz et al., 2024), we develop VCG+112K through a novel semi-automatic annotation pipeline, which provides dense video captions along with 112K QA pairs targeting reasoning, spatial and temporal understanding, which further improves model's understanding of video content (see Fig. 3).

Video Conversation Benchmarks: Video-ChatGPT (Maaz et al., 2024) introduces VCG-929 Bench and zero-shot QA benchmarks, where 930 931 VCGBench includes 500 videos with 3000 QA pairs, evaluated using GPT-3.5 across various met-932 rics. Despite its comprehensive evaluation, it 933 only contains videos from the ActivityNet dataset. The Zero-shot evaluation covers MSVD-QA (Xu 935 et al., 2017), MSR-VTT-QA (Xu et al., 2017), 936 TGIF-QA (Jang et al., 2019), and ActivityNet-937 QA (Fabian Caba Heilbron and Niebles, 2015). 938 MVBench (Li et al., 2024) consists of 4K QA pairs evaluating 20 temporal tasks, though it mostly in-940 cludes short videos averaging 5-40 seconds. Another recent benchmark, Video-MME (Fu et al., 2024), addresses the issue of diversity by incor-944 porating a wide range of videos. However, both MVBench and Video-MME are limited to MCQs, 945 which, while straightforward for evaluation, restrict the range of questions that can be asked and reduce the depth of understanding the model

can demonstrate. By confining to predefined 949 choices, MCQs introduce bias and fail to cap-950 ture the model's true understanding. Considering 951 the limitation of existing benchmarks, which of-952 ten lack focus on generalization and diversity, we 953 propose VCGBench-Diverse, featuring 4,354 QA 954 pairs from 877 videos across 18 domains, evaluated 955 using open-ended questions (see Fig. 4). 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

B Additional Implementation Details

In this section, we provide additional implementation details regarding our training setup and compute requirements. For experiments on VCG-Bench, VCGBench-Diverse and Zero-shot QA, we sample 16 frames from videos, while for MVBench which consists of relatively shorter videos, we sample 8 frames. We keep the same sampling strategy during inference. For VCGBench and VCGBench-Diverse, the model is trained on VideoInstruct100K (Maaz et al., VCG+112K, conversation and caption 2024), data from VideoChat (Li et al., 2023c) and VQA dataset from WebVid (Bain et al., 2021), that combines to approximately 260K single turn conversations. For MVBench, the model is trained on Kinetics-710 (Kay et al., 2017), Something-Something-v2 (Goyal et al., 2017), conversations from VideoChat (Li et al., 2023c), CLEVRER (Yi et al., 2019), VQA dataset from WebVid (Bain et al., 2021) and NExT-QA (Xiao et al., 2021) datasets, which combines to approximately 330K single turn conversations. We run all trainings for one epoch. Following previous approaches (Maaz et al., 2024; Jin et al., 2024; Liu et al., 2024c), we employ GPT-3.5-Turbo-0613 for VCGBench and Zero-shot QA evaluation. However, for our proposed VCGBench-Diverse, we employ the latest GPT-3.5-Turbo-0125 for evaluation for better reproducibility purposes.

All of our experiments are conducted using 8xA100 40GB GPUs. The training for VCG-Bench experiments takes around 12 hours to complete, while the training for MVBench experiments finishes in around 10 hours. We use the model trained for the VCGBench task to evaluate on VCGBench-Diverse and zero-shot questionanswering benchmarks. All of our training and evaluation codes, pretrained models and dataset will be publicly released.

Model	AS	AP	AA	FA	UA	OE	OI	os	MD	AL	ST	AC	MC	MA	SC	FP	со	EN	ER	CI	Avg.
Random	25.0	25.0	33.3	25.0	25.0	33.3	25.0	33.3	25.0	25.0	25.0	33.3	25.0	33.3	33.3	25.0	33.3	25.0	20.0	30.9	27.3
GPT-4V (OpenAI, 2023)	55.5	63.5	72.0	46.5	73.5	18.5	59.0	29.5	12.0	40.5	83.5	39.0	12.0	22.5	45.0	47.5	52.0	31.0	59.0	11.0	43.5
Otter-V (Li et al., 2023a)	23.0	23.0	27.5	27.0	29.5	53.0	28.0	33.0	24.5	23.5	27.5	26.0	28.5	18.0	38.5	22.0	22.0	23.5	19.0	19.5	26.8
mPLUG-Owl-V (Ye et al., 2023)	22.0	28.0	34.0	29.0	29.0	40.5	27.0	31.5	27.0	23.0	29.0	31.5	27.0	40.0	44.0	24.0	31.0	26.0	20.5	29.5	29.7
Video-ChatGPT (Maaz et al., 2024)	23.5	26.0	62.0	22.5	26.5	54.0	28.0	<u>40.0</u>	23.0	20.0	31.0	30.5	25.5	39.5	48.5	29.0	33.0	<u>29.5</u>	26.0	35.5	32.7
VideoLLaMA (Zhang et al., 2023)	27.5	25.5	51.0	29.0	39.0	48.0	40.5	38.0	22.5	22.5	43.0	34.0	22.5	32.5	45.5	32.5	40.0	30.0	21.0	37.0	34.1
VideoChat (Li et al., 2023c)	33.5	26.5	56.0	33.5	40.5	53.0	40.5	30.0	25.5	27.0	48.5	35.0	20.5	42.5	46.0	26.5	<u>41.0</u>	23.5	23.5	36.0	35.5
VideoChat2 (Li et al., 2024)	<u>66.0</u>	<u>47.5</u>	83.5	49.5	60.0	<u>58.0</u>	71.5	42.5	23.0	23.0	88.5	<u>39.0</u>	<u>42.0</u>	<u>58.5</u>	44.0	49.0	36.5	35.0	<u>40.5</u>	65.5	<u>51.1</u>
VideoGPT+ (ours)	69.0	60.0	<u>83.0</u>	<u>48.5</u>	66.5	85.5	75.5	36.0	44.0	34.0	89.5	39.5	71.0	90.5	45.0	53.0	50.0	<u>29.5</u>	44.0	<u>60.0</u>	58.7

Table 6: **Performance of VideoGPT+ on MVBench.** Following (Li et al., 2024), we evaluate on 20 tasks including <u>AS</u>: Action Sequence, <u>AP</u>: Action Prediction, <u>AA</u>: Action Antonym, <u>FA</u>: Fine-grained Action, <u>UA</u>: Unexpected Action, <u>OE</u>: Object Existence, <u>OI</u>: Object Interaction, <u>OS</u>: Object Shuffle, <u>MD</u>: Moving Direction, <u>AL</u>: Action Localization, <u>ST</u>: Scene Transition, <u>AC</u>: Action Count, <u>MC</u>: Moving Count, <u>MA</u>: Moving Attribute, <u>SC</u>: State Change, <u>FP</u>: Fine-grained Pose, <u>CO</u>: Character Order, <u>EN</u>: Egocentric Navigation, <u>ER</u>: Episodic Reasoning and <u>CI</u>: Counterfactual Inference.

C Complete MVBench Results

997

999

1000

1001

1003

1004

1005

1006

1007

1008

1009

1011

1012

1013

1015

1016

1017

1018

1019

1020

1021

1022

1023

1025

1026

1027

1028

1029

1030

1031

1033

We provide complete results on 20 sub-categories in MVBench in Table. 6. Specifically, VideoGPT+ achieves SoTA results in 14 out of 20 tasks and comes second in 4 out of 20 tasks, obtaining an average score of 58.7% across the 20 tasks. Additionally, VideoGPT+ shows significant improvements in the Action Prediction (+12.5%), Object Existence (OE) (+27.5%), Moving Direction (MD) (+17%), Moving Count (MC) (+29%) and Moving Attributes (MA) (+32%) indicating the rich spatial information and temporal context achieved by our model.

D Qualitative Results

We provide a qualitative comparison of our VideoGPT+ with the previous state-of-the-art approach, VideoChat2 (Li et al., 2024), in Fig. 5. The example shows an advertisement video for sunscreen, where multiple scene changes are present. The video starts with a close-up view of the sunscreen, followed by a woman applying sunscreen on her hand, then applying sunscreen near a beach. The woman is then seen applying sunscreen on her arms, and finally, the video shows the key ingredients of the sunscreen and ends with the cover of the sunscreen.

As shown in Fig. 5, our VideoGPT+ correctly identifies the events present in the video and provides a detailed and accurate description. On the other hand, VideoChat2 struggles to accurately capture all the events. Further, our model generates an advertisement post highlighting one of the unique features of the sunscreen shown in the video, namely that it functions as both sunscreen and moisturizer. Lastly, our VideoGPT+ correctly identifies the SPF value and brand name of the sunscreen, while VideoChat2 struggles to correctly

	In	nage P	ooling	Video	Pooling
	CNN	$4 \times$	$4 2 \times 2$	Time	Space
	3.25	3.25	3.28	3.23	3.28
	Tab	le 7: A	blation on p	pooling st	rategy.
,	Training	Data	MVBench	VCG	VCG-Div
,	Task-spec	ific	58.7	3.28	2.47
(Combined	1	58.3	3.27	2.45

Table 8: Ablation on generalization across different tasks. Our model trained on combined instruction tuning dataset performs comparable to the task-specific variants.

identify the brand name. We present further comparison in Fig. 6-7.

E Additional Ablations

Pooling Strategy: We ablate different pooling strategies for the image and video encoders. The image encoder outputs a 24×24 feature map from a 336×336 input. We compare two downsampling methods: a learnable lightweight CNN (LDPv2 from (Chu et al., 2024)) and a non-learnable adaptive average pooling with a 2×2 kernel. Results indicate that adaptive pooling performs better than CNN. A 4×4 adaptive pooling was also tested but showed inferior performance.

Similarly, we ablate the pooling choice for the video encoder, which takes an input of size $T \times 224 \times 224 \times C$ and outputs a feature map of $T \times 16 \times 16 \times d$. We compare two pooling strategies: time pooling across the temporal dimension to reduce the feature map to $1 \times 16 \times 16 \times d$, and space pooling across the spatial dimension with a 2×2 kernel. Results shows that space pooling effectively preserves temporal information and yields better results.

Generalization across video conversation datasets: To ensure a fair comparison with existing methods (Li et al., 2024; Liu et al., 2024c), 1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058



Describe the video in detail.

2

The video shows a woman applying a cream to her face and then rubbing it in. She then applies a serum to her face and rubs it in. The woman then applies a sunscreen to her face and rubs it in. The video also shows the woman applying a cream to her arm and rubbing it in.

The video starts with a close-up shot of a bottle of sunscreen, followed by a woman applying it to her skin. The camera then pans out to show the woman sitting on the beach, applying the sunscreen to her arms. The sunscreen is shown again in a close-up shot, highlighting its texture and consistency. The video then shows a series of close-up shots of the sunscreen, with text appearing on the screen to highlight its key ingredients. Overall, the video showcases the importance of using sunscreen to protect the skin from harmful UV rays, with a focus on the key ingredients that make up the product.

> Generate a creative advertisement feature about this product that can be used to post on the product's instagram page.



The creative advertisement feature for this product could be a video showcasing the product's benefits and how it can be used. The video could include a woman demonstrating how to use the product and the various ways it can be applied.

Introducing the ultimate sun protection solution for your skin! Our new sunscreen moisturizer is the perfect blend of protection and hydration for your skin. With its broadspectrum protection, you can enjoy the sun without worrying about harmful UV rays. Plus, our moisturizing formula keeps your skin looking and feeling healthy and radiant. Say goodbye to dry, flaky skin and hello to a smooth, supple complexion. Don't forget to apply it generously and reapply every two hours for maximum protection. Your skin deserves the best, and our sunscreen moisturizer is the perfect way to give it the care it needs. Try it today and see the difference for yourself!

What is the SPF value of the sunscreen and what is the name of the sunscreen?

The sunscreen in the video has an SPF value of 30. The name of the sunscreen shown in the video is Key Benefits.



1061

1062

1063

1066

1068

1070

The sunscreen has an SPF value of 30 and the sunscreen is called "Skin Authority Sunscreen Moisturizer"



Figure 5: Qualitative comparison of VideoGPT+ with VideoChat2. Our VideoGPT+ demonstrates superior temporal understanding by correctly identifying multiple events in the video, effective reasoning in generating a creative advertisement, and accurate spatial understanding by identifying the SPF value and brand name of the sunscreen.

we train our model on different combination of datasets for evaluation on MVBench and VCGBench. To further clarify the generalization capability of our model, we provide results on three benchmarks, VCGBench, MVBench and VCGBench-Diverse, using a single model trained on a combined dataset. The results demonstrate that our model maintains performance across all benchmarks, indicating its ability to generalize effectively across diverse video conversation datasets.

Feature concatenation strategy: We conduct anablation study to determine the optimal order in

which image and video features should be input to the LLM. Specifically, we perform two exper-1074 iments. In the first experiment, image and video 1075 features are extracted for each video segment and concatenated in an interleaved manner before sending as input to the LLM. For example, the video is divided into segments of equal size, and then 1079 the image and video features from each segment 1080 are concatenated and input to the LLM. In the sec-1081 ond experiment, we first place all the image features followed by all the video features. The re-1083 sults shown in Table 9, indicate that the sequential design, where the image features are placed first

followed by the video features, yields better performance. This can be justified by the fact that we use different visual adapters for image and video features, so interleaving the features from both modalities can create a larger distribution shift, hindering the learning process.

Feature		Avg.				
Concatenation	CI	DO	CU	TU	СО	0
Interleaved	3.25	3.17	3.72	2.78	3.39	3.26
Sequential	3.27	3.18	3.74	2.83	3.39	3.28

Table 9: Ablation on Feature Concatenation Strategy. Performance comparison between interleaved and sequential feature concatenation strategies. The sequential feature concatenation performs better.

Generalization of VideoGPT+ to other LLMs: We train VideoGPT+ with different LLMs including Vicuna 7B and 13B (Chiang et al., 2023) and LLaMA-3 8B (AI, 2024). We observe slight improvements in VCGBench scores when training using better LLMs, including Vicuna 13B and LLaMA-3 8B models.

F GPT Prompts

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

In this section, we provide the GPT prompts used for the following tasks: (i) Dense video description generation for VCG+112K, (ii) Question-answer generation for VCG+112K and (iii) Questionanswer generation for VCGBench-Diverse.

Dense Video Description Generation for VCG+ 112K: To generate dense video captions, we provide GPT-4 with a concise ground truth caption of the video and detailed frame-level captions of the key-frames generated from LLaVAv1.6 (Liu et al., 2024a). GPT-4 is then prompted to combine this information into a detailed caption for the entire video. As illustrated in Fig. 8, the prompt includes clear instructions to eliminate any conflicting information, ensuring an accurate and detailed caption.

1116Question-answer generation for VCG+ 112K:1117After generating detailed video descriptions us-1118ing GPT-4, we use GPT-3.5 to create question-1119answer pairs for instruction tuning. Fig. 9 shows1120the prompt to generate detailed summary question-1121answer pair using the ground truth caption and the1122dense description of the video.

1123Question-Answer Generation for VCGBench-1124Diverse: We provide prompts used to gen-1125erate comprehensive question-answer pairs for1126VCGBench-Diverse. As illustrated in Fig. 10, the1127questions are generated in three categories: tem-

LLM	VCGBench					Avg.
	CI	DO	CU	TU	СО	8
Phi3-Mini-3.8B	3.27	3.18	<u>3.74</u>	2.83	3.39	3.28
Vicuna-7B	3.22	3.14	3.69	2.65	3.46	3.23
Vicuna-13B	3.30	3.20	3.75	2.77	3.48	3.30
LLaMA3-8B	<u>3.29</u>	3.21	3.73	2.86	3.38	<u>3.29</u>

Table 10: **Ablation on LLM type.** We train and evaluate VideoGPT+ with different LLMs, including vicuna (Chiang et al., 2023) and LLaMA3 (AI, 2024), which further improves accuracy.

poral, spatial, and reasoning. Similar prompts are used to generate consistency and summary questions, offering an extensive evaluation protocol for VCGBench-Diverse. 1128

1129

1130

1131

1132

G Ethics and societal impact

We use multiple open-source video datasets in-1133 cluding ActivityNet (Fabian Caba Heilbron and 1134 Niebles, 2015), WebWid (Bain et al., 2021), 1135 MSVD-QA (Xu et al., 2017), MSRVTT-QA (Xu 1136 et al., 2017), TGIF-QA (Jang et al., 2019), 1137 HDVILA (Xue et al., 2022), MPII (Andriluka 1138 et al., 2014), YouCook2 (Zhou et al., 2018), UCF 1139 Crime (Sultani et al., 2018), and STUD Traffic (Xu 1140 et al., 2021) in our work. To the best of our knowl-1141 edge, the dataset does not portray any strong biases 1142 or discrimination. We urge for the responsible use 1143 of VideoGPT+ and VCG+ 112K, promoting research 1144 progress while safeguarding privacy. 1145

	DDS TV
	GT Answer: The car with the dashboard camera shifted out of the road because the white
	pick-up truck tried to change lanes, causing the car to avoid a potential collision.
	The car with the dashboard camera shifted out of the road because it was being driven by a person who was not paying attention.
	The car with the dashboard camera shifted out of the road to avoid a collision with another vehicle.
\$	The car with the dashboard camera shifted out of the road due to a collision with a white pickup truck.
+	The car did not move out of its lane. A white pickup truck attempts to move into the lane, but then brakes to avoid a collision with a large truck in the lane, staying in its lane.
	Category: Traffic
THIS SAN FRANCISCO CAN ASSPER & BOOTSIN	
	Q: Why does the video highlight these specific bodybuilders' cases?
	GT Answer: The video highlights these specific bodybuilders' cases to underscore the potential dangers and severe health repercussions associated with the misuse of steroids and hormones within the bodybuilding community, serving as cautionary tales.
	The video highlights these specific bodybuilders' cases to showcase their achievements and success in their respective fields.
	The video highlights these cases to showcase the impact of substance abuse on bodybuilders' health and careers, emphasizing the consequences of such actions.
\$	The video highlights these specific cases to discuss the extreme measures some individuals have taken for bodybuilding, which have led to significant health risks, reputational damage, and unfortunate events.
+	The video highlights the cases of bodybuilders who suffered adverse effects from the use of steroids.
	Category: Education
	Q: How many men are initially seen walking along the road in the CCIV footage?
	GT Answer: Four men are initially seen walking along the road in the CCTV footage.
	There are two men walking along the road in the CCTV footage.
	Initially, four men are seen walking along the road in the CCTV footage.
\$	In the CCTV footage provided, initially, there are three men seen walking along the road.
$ \blacklozenge $	There are five men initially seen walking along the road.
	Category: Surveillance
	R User 😻 VideoChat2 🔗 VideoGPT+ 🚱 GPT-4V 🔶 Gemini-Pro-V

Figure 6: Qualitative comparison from VCGBench-Diverselof VideoGPT+. We show qualitative comparison of VideoGPT+ with VideoChat2 and propriety models GPT-4V and Gemini-1.5-Pro-V from three different categories including traffic, education and surveillance from VCGBench-Diverse.

	Q: Which team is participating in the featured match?
	GT Answer: The teams participating in the featured match are Valencia CF and Real Valladolid.
3	The teams participating in the featured match are Real Madrid and Atletico Madrid.
	The video shows that the teams playing are Real Valladolid and Valencia CF.
\$	The teams participating in the featured match are Valencia CF and Real Valladolid.
•	The teams participating are Valencia CF and Real Valladolid.
	Category: Sports
C The v	
	Q: What significant declaration was made by President Rodrigo Duterte a year ago?
	GT Answer: President Rodrigo Duterte declared martial law in Mindanao in response to the Marawi siege.
	Last year, President Rodrigo Duterte made a significant declaration regarding the Philippines' territorial claims in the South China Sea.
	A year ago, President Rodrigo Duterte declared martial law in Mindanao.
\$	A year before, President Rodrigo Duterte declared martial law in Mindanao in response to the Marawi siege.
•	A year ago, Philippine President Rodrigo Duterte declared martial law in Mindanao, in response to the Marawi siege.
	Category: News
Ī	
	Q: What is the color of the self-parking car?
	GT Answer: The self-parking car is black.



Figure 7: Qualitative comparison from VCGBench-Diverse of VideoGPT+. We show qualitative comparison of VideoGPT+ with VideoChat2 and propriety models GPT-4V and Gemini-1.3-Pro-V from three different categories including sports, news and automobiles videos from VCGBench-Diverse.

Gemini-Pro-V

VideoChat2

User

Generate a detailed and accurate description of a video based on the given ground-truth video caption and multiple frame-level captions. Use the following details to create a clear and complete narrative: Ground-truth Video Caption: [Ground-truth caption here] Frame-level Captions: [Frame-level caption 1]; [Frame-level caption 2]; [Frame-level caption 3]; ... Instructions for writing the detailed description: 1. Focus on describing key visual details such as appearance, motion, sequence of actions, objects involved, and interactions between elements in the video. 2. Check for consistency between the ground-truth caption and frame-level captions, and prioritize details that match the ground-truth caption. Ignore any conflicting or irrelevant details from the frame-level captions. 3. Leave out any descriptions about the atmosphere, mood, style, aesthetics, proficiency, or emotional tone of the video. 4. Make sure the description is no more than 20 sentences. 5. Combine and organize information from all captions into one clear and detailed description, removing any repeated or conflicting details. 6. Emphasize important points like the order of events, appearance and actions of people or objects, and any significant changes or movements. 7. Do not mention that the information comes from ground-truth captions or frame-level captions. 8. Give a brief yet thorough description, highlighting the key visual and temporal details while keeping it clear and easy to understand. Use your intelligence to combine and refine the captions into a brief yet informative description of the entire video.

Figure 8: **Prompt for Dense Video Captions Generation for VCG+ 112K.** We use GPT-4 to generate detailed video captions using concise ground truth and frame-level detailed captions.

System Prompt You are an AI assistant tasked with generating questions and answers about video content to create a video instruction tuning dataset. Your goal is to extract detailed visual and temporal information from the video, ensuring the explanations are comprehensive enough for someone to understand the entire sequence of events in the video.
<pre>##TASK: 1. Users provide a video ground truth caption and a detailed description. 2. Generate three questions that effectively prompt a detailed description of the entire video content and sequence of events.</pre>
<pre>##INSTRUCTIONS: - Ensure each question targets the goal of generating a detailed description of the entire video from start to end. - Avoid questions that focus on small parts, less relevant details, or abstract concepts such as logical reasoning, attention to subtle details, overall aesthetic. - Every answer must include all the details from the ground truth caption and integrate additional specifics from the detailed description. - Focus on visual and temporal details.</pre>
<pre>##SAMPLE QUESTIONS: - Can you describe the entire video in detail from start to finish? - What happens throughout the entire video, including all key actions and events? - Could you provide a detailed walkthrough of the entire video?</pre>
<pre># User Prompt: The video ground truth caption is: [Ground-truth caption here]. The noisy detailed description is: [Dense description here].</pre>
Generate three questions and answers about the entire content and sequence of events in the video. Each question should aim to elicit a comprehensive description of the full sequence of events in the video from start to finish. Each answer must include all the details from the ground truth caption and integrate additional specifics from the detailed description. Format the output as a list of dictionaries in JSON style, with each dictionary containing a 'Q' key for the question and an 'A' key for the answer. For example:
[{'Q': 'Your first question here', 'A': 'Your first answer here'}, {'Q': 'Your second question here', 'A': 'Your second answer here'}, {'Q': 'Your third question here', 'A': 'Your third answer here'}].
Most importantly, every answer must provide a full understanding of the video by incorporating ALL the details from the ground truth caption and additional specifics from the detailed description.

Figure 9: Prompt for Question-answer generation for VCG+ 112K. We use GPT-3.5 to generate question-answer pairs for instruction tuning using the concise video ground truths and detailed video descriptions.

System Prompt:

You are an AI assistant tasked with generating questions and detailed answers based on a video description. Your goal is to extract important information from the video content, focusing on temporal events, visual details, and reasoning behind actions.

##TASK:

You will receive a video description, and based on it, you must generate a set of questions and answers in three distinct categories:

1. Temporal - These questions should focus on the sequence and timing of events. Use approximate time references where necessary.

2. Spatial - These questions should address visual aspects such as appearance, objects, colors, attire, displayed texts, number of objects or people, location, and other significant visual details.

3. Reasoning - These questions should delve into the actions, motivations, and consequences as depicted in the video description.

##INSTRUCTIONS:

- Each question must directly relate to and be answerable by the provided video description. Avoid assumptions and fabrication of details not present in the description.

- Provide clear, unambiguous questions that allow for definitive answers based on the description. - If the video description does not contain enough information to formulate a question in any category, do not include a question for that category.

##SAMPLE QUESTIONS:

Temporal: Describe the entire process the person goes through from start to finish or What happens at the beginning of the video? or What does the person do right after the dog appears?
Spatial: Can you provide a detailed description of the appearance and activities of all

individuals or What is the color of the main character's shirt? or What is the name of the drink on the bottle? How many people are at the table?

- Reasoning: What action does the coach take after the whistle blows? or Why did the player throw the ball? or Who is John Davis in the video?

User Prompt:

The video description is: [Dense description here].

Format the output as a dictionary in JSON style, with each key representing a question category and containing a sub-dictionary with 'Q' for the question and 'A' for the answer. Example output with all three categories filled:

{'temporal': {'Q': 'Temporal question here...', 'A': 'Answer here...'},
'spatial': {'Q': 'Spatial question here...', 'A': 'Answer here...'},
'reasoning': {'Q': 'Reasoning question here...', 'A': 'Answer here...'}}.

If a category cannot be filled:

{'temporal': {'Q': 'Describe the sequence of events in the video.', 'A': 'The video starts with...'},'spatial': {'Q': 'What is the main character wearing?', 'A': 'The main character is dressed in...'}} # reasoning omitted due to lack of information

Importantly, the answers MUST extract information DIRECTLY from the given description. Do not include categories that cannot be filled based on the video description alone.

Figure 10: **Prompt for Question-Answer Generation for VCGBench-Diverse.** We use GPT-3.5 to generate temporal, spatial, and reasoning question-answer pairs.