# LLMs Meet Long Video: Advancing Long Video Comprehension with An Interactive Visual Adapter in LLMs

Anonymous ACL submission

#### Abstract

Long video understanding is a significant and 001 ongoing challenge in the intersection of mul-003 timedia and artificial intelligence. Employing large language models (LLMs) for comprehending video becomes an emerging and promising method. However, this approach incurs high computational costs due to the ex-800 tensive array of video tokens, experiences reduced visual clarity as a consequence of token aggregation, and confronts challenges aris-011 ing from irrelevant visual tokens while answering video-related questions. To alleviate 012 these issues, we present an Interactive Visual 014 Adapter (IVA) within LLMs, designed to enhance interaction with fine-grained visual elements. Specifically, we first transform long videos into temporal video tokens via lever-017 018 aging a visual encoder alongside a pretrained causal transformer, then feed them into LLMs with the video instructions. Subsequently, we integrated IVA, which contains a lightweight temporal frame selector and a spatial feature interactor, within the internal blocks of LLMs to capture instruction-aware and fine-grained visual signals. Consequently, the proposed video-LLM facilitates a comprehensive understanding of long video content through appropriate long video modeling and precise visual interactions. We conducted extensive experiments on nine video understanding benchmarks and experimental results show that our interactive visual adapter significantly improves the performance of video LLMs on long video QA tasks. Ablation studies further verify the effectiveness 034 of IVA in long and short video understandings.

### 1 Introduction

041

The exponential advancement of the Internet and multimedia technologies has resulted in a significant surge in video content production by individuals and enterprises across various domains. The ability to interpret and extract meaningful content from videos is increasingly vital for meeting human demands and promoting the speed of information dissemination (Tang et al., 2023). Therefore, Video Question Answering (Yu et al., 2019; Li et al., 2023b; Castro et al., 2022) (Video OA), which allows users to ask about the content of videos through natural language and receive answers derived from their visual and auditory content, attracts tremendous research interest. Recently, large language models (LLMs) (OpenAI, 2023; Chiang et al., 2023) have demonstrated exceptional efficacy in the domains of human-machine interaction and the handling of extensive contextual information. Capitalizing on these advancements, there is a burgeoning inclination towards integrating LLMs into the realm of video information processing. This approach primarily aims to enhance the interface between users and video content through intelligent question-and-answer sessions.

043

044

045

046

047

050

051

052

053

057

058

059

060

061

062

063

064

065

067

068

069

071

073

074

075

076

077

078

079

081

The core of this innovation is a strategy that bridges the gap between the visual information in videos and the textual comprehension capabilities of LLMs. This is accomplished through a meticulously designed process that translates video data into a format comprehensible by LLMs, thereby facilitating an advanced question-answering system tailored for video content. The process involves mapping video encoding into the language space of LLMs via a learnable visual mapping network (Wu et al., 2023; Li et al., 2023d; Dai et al., 2023). Essentially, the video is converted into "video tokens", which are then fed into the LLM along with textual tokens of natural language questions. Leveraging the vast knowledge storage and natural language processing prowess of LLMs, this approach effectively handles video QA tasks. For instance, Maaz et al. (2023) performs spatial and temporal pooling for video tokens and feeds them into Vicuna (Chiang et al., 2023) to achieve the interaction between users and video content. Zhang et al. (2023b) utilizes Q-former (Li et al., 2023d) to extract question-relevant video tokens, which are

then fed into LLama (Gao et al., 2023) to generate the answer.

086

880

100

101

102

103

104

105

106

108

133

134

These LLMs-powered video understanding models (Tang et al., 2023; Song et al., 2023) mainly focus on short video modeling and have achieved a successful performance on short video captioning (Zhang et al., 2023c), question-answering (Jin et al., 2023), and summarization (Tang et al., 2023). However, the core challenges of video processing (Xu et al., 2023) stem from the need to efficiently model long video sequences and precisely respond to questions relevant to the video. Generally, using LLMs to handle long-form video often encounters the following hurdles: 1) high computational costs from a multitude of video tokens; 2) reduced visual clarity as a consequence of token aggregation such as employing average or maximum representation pooling for visual frames; 3) irrelevant visual tokens leading to incorrect answers, notably when question-relevant information is embedded within long temporal cues. Hence, previous models struggle to handle long-form videos owing to the constrained input capacity for video tokens and the challenge of distilling question-relevant, finegrained visual features during generation.

To alleviate these issues, we present a long video 109 comprehension method for LLMs, named Interac-110 tive Visual Adapter (IVA) to achieve in-depth inter-111 actions between LLMs and video content. Specif-112 ically, we first use the pretrained visual encoder 113 to obtain global and fine-grained frame represen-114 tations. We construct the temporal video tokens 115 by integrating the global features of frames with 116 temporal video embeddings, which are obtained 117 through a pretrained causal transformer. The whole 118 set of temporal video tokens is fed into the LLM to 119 attain a whole understanding of the video content. 120 Additionally, we designed a parameters-sharing 121 Interactive Visual Adapter (IVA) that contains an instruction-aware temporal frames selector and a 123 spatial feature interactor. The selector is used to ob-124 tain question-relevant frames based on contextual 125 query embeddings and global encodings of videos. 126 The selected frames are then fed into the spatial interactor to engage with the contextual query em-128 beddings, in which fine-grained representations of 129 frames are used. By doing so, LLMs could achieve 130 in-depth interaction with video content by applying 131 IVA between different layers. 132

To verify the effectiveness of our method, we conduct extensive experiments on four long video

QA and five short video understanding benchmarks. Experimental results indicate that IVA is capable of achieving effective interactions between LLMs and long or short videos. Our contributions are summarized as follows: 135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

- We analyze the challenges of modeling long videos for LLMs and propose an interactive visual adapter for LLMs to handle long videos. It realizes the in-depth interaction between LLMs and long videos based on efficient video tokens and the IVA mechanism.
- The proposed IVA is capable of selecting relevant frames and interacting with their finegrained spatial features through the internal selector and interactor, respectively. The IVA architecture is lightweight and designed to be shareable between layers of LLMs.
- Experimental results show that LLMs with IVA could achieve powerful performances in understanding long videos. Ablation studies underscore the critical role and effectiveness of IVA, confirming its significant contribution to enhanced performance.

# 2 Related Work

Traditional Video Understanding Models The rapid development of deep learning methods possesses superior task-solving capabilities for video understanding. DeepVideo (Karpathy et al., 2014) was the earliest method introducing a Convolutional Neural Network (CNN), for video understanding. Two-stream networks (Feichtenhofer et al., 2016), then integrating Convolutional Neural Networks (CNNs) (Feichtenhofer et al., 2016) and Improved Dense Trajectories (IDT) (Li et al., 2021), enhanced motion analysis in video understanding. For long-form content, Long Short-Term Memory (LSTM) (Yue-Hei Ng et al., 2015) networks were adopted, offering a robust solution for sequential data analysis over extended durations. Additionally, Temporal Segment Network (TSN) (Wang et al., 2016) advanced long-form video understanding by segmenting videos for individual analysis before aggregating insights, enabling more nuanced interpretation. Meanwhile, 3D networks started another branch by introducing 3D CNN to video understanding (C3D) (Tran et al., 2015). The introduction of Vision Transformers (ViT) (Dosovitskiy et al., 2021; Arnab et al., 2021; Fan et al., 2021) pro-

motes a series of prominent models Among the pio-183 neering efforts in this self-supervised video training 184 domain, VideoBERT (Sun et al., 2019) leverages 185 the bidirectional language model BERT (Kenton and Toutanova, 2019) for self-supervised learning from video-text data. This model, and oth-188 ers following the "pre-training and fine-tuning" 189 paradigm, such as ActBERT (Zhu and Yang, 2020), 190 SpatiotemporalMAE (Feichtenhofer et al., 2022), 191 OmniMAE (Girdhar et al., 2023), showcase the 192 diverse strategies developed to enhance video un-193 derstanding. Notably, these models have set a foun-194 dation for advanced video-language models like 195 maskViT (Gupta et al., 2022), CLIP-ViP (Xue et al., 196 2022), LF-VILA (Sun et al., 2022), further push-197 ing the boundaries of what's achievable in action classification, video captioning, and beyond. The 199 evolution from VideoBERT to more recent innovations like HiTeA (Ye et al., 2023), and CHAM-PAGNE (Han et al., 2023) underscores the rapid advancement in this field.

207

209

210

211

212

213

214

215

216

217

218

219

220

221

224

226

227

231

LLMs for Video Understanding The recent advancement in large language models (LLMs), pre-trained on expansive datasets, has ushered in groundbreaking capabilities in in-context learning (Zhang et al., 2023a) and long-form context modeling (Lyu et al., 2023). This innovation has paved the way for integrating LLMs with computer vision technologies, exemplified by initiatives like Visual-ChatGPT (Wu et al., 2023). These models transcend traditional boundaries by calling vision model APIs (Qin et al., 2023), thereby addressing complex problems within the computer vision domain. Integrating language models with video understanding technologies (Maaz et al., 2023; Zhang et al., 2023d; Li et al., 2023e; Xu et al., 2023; Song et al., 2023) enhances multimodal understanding, facilitating sophisticated processing and interpretation of the intricate interplay between visual and textual data. They leverage their extensive multimodal knowledge base and nuanced contextual understanding, mirroring a more human-like comprehension of visual content. Moreover, the exploration of LLMs in video understanding tasks (Tang et al., 2023) represents a significant stride towards harnessing their potential in analyzing and reasoning about visual data.

Multimodal Instruction Tuning for LLMs Recent advancements have significantly enhanced the performance of instruction-tuned, text-only large language models (LLMs) (Ouyang et al., 2022; Muennighoff et al., 2022; Chung et al., 2022) on various NLP tasks and human-machine interaction scenarios. SFT LLMs have demonstrated remarkable capabilities in these areas. Building on this foundation, researchers are now exploring the integration of multimodal instruction data to further refine pre-trained LLMs, aiming to elevate their multimodal human-machine interaction competencies. For instance, a study (Liu et al., 2023) employed GPT-4 to generate multimodal instruction data, which was then used to fine-tune the language model LLaMA on a synthetic dataset designed for instruction following. Similarly, another research (Zhu et al., 2023) effort constructed a well-aligned multimodal instruction-following dataset to fine-tune Vicuna, an instruction-tuned language model, which showed superior performance in open-domain multimodal dialogues. Furthermore, a lightweight adaptation method (Zhang et al., 2023d) was introduced to efficiently convert LLaMA into an instruction-following model, showcasing the potential for streamlined model enhancement. In this paper, we also explore the lightweight adapter to help LLMs understand long videos.

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

260

261

262

263

264

265

266

267

268

270

271

272

273

274

275

276

277

278

279

280

281

### 3 Methodology

# 3.1 Overview

Our work primarily introduces an interactive visual adapter for LLMs to handle long videos and answer relevant questions. The overview of workflow is shown in Figure 1. Specifically, given a video V, we first extract frames to obtain the whole sequence frame representations  $\mathbf{h}_V = (\mathbf{h}_{I_1}, ..., \mathbf{h}_{I_k}, ..., \mathbf{h}_{I_N})$ via the pretrained image encoder, where  $\mathbf{h}_{I_k}$  =  $(h_a^{I_k}, h_1^{I_k}, \dots, h_{576}^{I_k})$  refers to the representations of k th frame and N is the total number of extracted frames. Then, we use a casual transformer to acquire temporal video embeddings from the aggregated spatial representation. The overall video tokens are formed by merging temporal video embeddings and global spatial features  $[h_g^{I_1}, h_g^{I_k}, ..., h_g^{I_N}]$ , where each frame is represented by two tokens. To enhance the capability of LLMs in leveraging fine-grained visual details from videos, we have developed an Interactive Visual Adapter (IVA) that is integrated into the blocks of LLM. This integration allows LLMs to comprehend the entirety of long videos through efficient video tokens while simultaneously capturing fine-grained visual information facilitated by the IVA.



Figure 1: The overview of our framework employing LLMs to handle long video. While producing video tokens, we combine the global features and aggregated fine-grained features to represent a frame, allocating two tokens for each frame. The casual transformer is used to capture temporal relationships across frames and its output will be spliced with spatial feature sequence. The IVA will be inserted between blocks of LLMs to incorporate fine-grained visuals based on an understanding of the long video tokens, text instructions, and query tokens.

### 3.2 Producing Video Tokens

We elaborate on the detailed process employed to produce efficient tokens for long videos, characterized by the extraction of one frame per second. First, we use the self-weighted calculation on the fine-grained feature  $\mathbf{h}_{f}^{k} = (h_{1}, ..., h_{576})$  of a frame (k th) to obtain its overall representation, which will be fed into the following casual transformer. This calculation process for the k th frame is given in the following Eq. 1:

$$s_f^k = Softmax(\mathbf{W}(\mathbf{h}_f^k) + \mathbf{b}),$$
  
$$\mathbf{h}_t^k = s_f \mathbf{h}_f^k,$$
 (1)

where  $s_f \in \mathbf{R}^{1 \times 576}$  is the weight distribution and  $\mathbf{W}$  and  $\mathbf{b}$  are learnable parameters. Hence, we denote the obtained sequence-level frame representation as  $\mathbf{h}_f = (\mathbf{h}_f^1, ..., \mathbf{h}_f^N)$ .

**Casual Transformer** is employed to acquire the temporal video embeddings. Specifically, we use a four-layer transformer to facilitate interaction across frames, where a frame only attends to its previous ones. Take the first layer as an example, the specific operation of the casual transformer is presented in Eq. 2:

$$\begin{split} \mathbf{h}_s &= SelfAtten(LayerN(\mathbf{h}_f), W^{Mask}) + \mathbf{h}_f \\ \mathbf{h}_s &= LayerN(\mathbf{h}_s), \\ \mathbf{h}_o^1 &= MLP(\mathbf{h}_s) \end{split}$$

where SelfAtten and LayerN are the selfattention calculation and the feature normalization. The top output of the casual transformer will be projected into the language model by a linear layer, which is spliced with the global features  $\mathbf{h}_{g}^{V} = (h_{g}^{I_{1}}, h_{g}^{I_{k}}, ..., h_{g}^{I_{N}})$ . These global features of frames will be transferred into language models via a learnable MLP. We denote the final spliced feature to  $\mathbf{h}_{V} = (h_{V}^{1}, h_{V}^{2}, ..., h_{V}^{2N})$ . 306

307

309

310

311

312

313

314

315

316

317

318

319

321

323

324

325

327

329

331

332

333

334

335

# 3.3 Interactive Visual Adapter

After obtaining video tokens  $h_V$ , and supposing that the textual embeddings of instruction are initiated to  $h_T$  via the frozen word embedding table of LLMs, we concatenate them into a single sequence and fed it into LLMs. Considering fine-grained visual details existing in long videos, we expect that LMMs are capable of capturing the specific fine-grained visual information based on the understanding of instructions and the whole video representations. Hence, we devise a lightweight interactive visual adapter (IVA) to enable LLMs to focus on instruction-relevant fine-grained visuals during content generation.

Concretely, as the bottom part shown in Figure 1, we first introduce learnable dynamic tokens  $\mathbf{h}_D = (h_1^D, ..., h_M^D)$  as the query signals and integrate it at the end of the input token sequence. It aims to capture previous instruction and video information via the self-attention mechanism of LLMs, functioning as query tokens to engage with

(2)

305

283

287

289

296

302

303

304

the fine-grained spatial features of videos. Suppose that the output of *i* th layer of LLMs is  $h^i$ . The specific calculation process of IVA between the *i* and *i* + 1 th layers of LLMs is shown in Eq. 3 and 4 in order. Each layer of IVA consists of a selector and an interactor, which are capable of selecting relevant frames and capturing valuable fine-grained visual information. The operational process of the selector is described as follows:

$$\begin{split} \mathbf{h}_{q}^{S} &= W^{q} \mathbf{h}_{d}^{i} + b^{q}, \\ \mathbf{h}_{k}^{S} &= W^{k} \mathbf{h}_{g}^{V} + b^{k}, \\ \mathbf{M}^{S} &= \mathbf{h}_{q}^{S} (\mathbf{h}_{k}^{S})^{T}, \\ \mathbf{h}^{S} &= Softmax (\mathbf{M}^{S} / \tau) Trans([\mathbf{h}_{I_{1}}^{f}, ..., \mathbf{h}_{I_{N}}^{f}]), \end{split}$$

(3) where  $\mathbf{h}_d^i$  refers to the hidden states  $\mathbf{h}^i$  associated with the indices of dynamic tokens.  $W^q$ ,  $W^k$ ,  $b^q$ , and  $b^k$  are learnable parameters.  $\mathbf{M}^S$  signifies the distribution score on the frames, which represents the relevant attention distribution.  $\tau$  is the hyperparameter, which is set to 0.5. "Trans" refers to the transportation of feature dimension.  $[\mathbf{h}_{I_1}^f, ..., \mathbf{h}_{I_N}^f]$ represents the fine-grained features of the entire video. The output  $\mathbf{h}^S \in \mathbf{R}^{b \times M \times 576 \times d_S}$  will be fed into the following interactor as the key value, where  $d_S$  represents the dimension of the selector.

For the interactor, the specific calculation progress could be given as Eq. 4.

$$\mathbf{h}_{q}^{I} = W^{1}\mathbf{h}_{d}^{i} + b^{1},$$

$$\mathbf{h}_{k}^{I} = W^{2}\mathbf{h}^{S} + b^{2},$$

$$\mathbf{M}^{I} = \mathbf{h}_{q}^{I}(\mathbf{h}_{k}^{I})^{T},$$

$$\mathbf{h}_{c}^{S} = Softmax(\mathbf{M}^{I})(W^{3}\mathbf{h}^{S} + b^{3}),$$

$$\mathbf{h}^{S} = MLP(\mathbf{h}_{c}^{S}) + \mathbf{h}_{c}^{S}$$
(4)

where  $W^1$ ,  $W^2$ ,  $W^3$ ,  $b^1$ ,  $b^2$ , and  $b^3$  are learnable parameters. Overall, we use the same four-layer calculations of the above selector and interactor to facilitate that LLMs interact with fine-grained visual features.

### 3.4 Training

366Stage 1: Pretraining. To endow video tokens367with meaningful representation, we first train the368casual transformer, linear layers, and other learn-369able parameters during video tokens production, on370massive video-caption pairs from WebVid, a total371of 703k video-caption pairs. We freeze the other372parameters of the overall model during this process373and do not introduce the IVA module.

**Stage 2: Video Instruction Tuning**. At this stage, the model is required to generate responses that align with various instructions. These instructions often involve complex visual comprehension and reasoning, rather than merely describing visual signals. Note that the conversation data  $[Q_1, A_1, ..., Q_r, A_r]$  consists of multiple rounds.

374

375

376

377

378

379

381

383

384

385

386

387

389

391

392

393

394

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

$$X_T^r = \begin{cases} Q_1, & r = 1\\ \text{Concat}(Q_1, A_1, ..., Q_r, A_r), & r > 1 \end{cases}$$
(5)

where r represents the round count. As shown in Eq. 5, when r > 1, we concatenate the conversations from all previous rounds with the current instruction as the input for this round. The training objective remains the same as the previous stage. After this stage, the model will be capable of generating corresponding responses based on various instructions.

# 4 Experiments

#### 4.1 Data sets

While training the casual transformer, we utilize 702 thousand video-text pairs derived from Valley (Luo et al., 2023), sourced from WebVid (Bain et al., 2021). For the Video Instruction Tuning stage, we first collect instructional datasets from three sources. This includes a 100K video-text instruction dataset from Video-ChatGPT (Muhammad Maaz and Khan, 2023), a 36K short video-text instruction dataset from Valley-Instruct-73k (Luo et al., 2023), and a 34K multiple-choice QA dataset from NExT-QA (Xiao et al., 2021).

Additionally, we assessed the generalization of IVA using long and short video benchmarks. Long video benchmarks typically are characterized by videos exceeding one minute in duration. We evaluated our model using four prominent long video evaluation benchmarks: ActivityNet-QA (Yu et al., 2019), Social-IQ 2.0 (Wilf et al., 2023), LifeQA (Castro et al., 2020), WildQA (Castro et al., 2022). For short video benchmarks, the duration of the videos is often measured in seconds. We evaluated our model against three notable short video evaluation benchmarks: MSVD-QA (Xu et al., 2017), MSRVTT-QA (Xu et al., 2017) and SEED-Bench (Li et al., 2023c).

### 4.2 Baselines

We mainly compare our models with the following video LLMs that could be extended to handle long

361

336

337

341

343

346

347

354

Method	ActivityNet-QA		Social-IQ 2.0		LifeQA		WildQA	
Method	Accuracy	score	Accuracy	score	Accuracy	score	Accuracy	score
Video-LLaMA (Zhang et al., 2023b)	12.4	1.1	55.8	2.9	35.8	2.3	63.2	3.2
Video-Chat (Li et al., 2023f)	26.5	2.2	-	-	-	-	-	-
LLaMA-Adapter (Zhang et al., 2023d)	34.2	2.7	-	-	-	-	-	-
Video-ChatGPT (Maaz et al., 2023)	35.2	2.7	57.5	3.3	33.9	2.6	<u>58.0</u>	3.3
Baseline (w/o IVA)	40.8	3.0	46.8	3.0	31.7	2.3	52.2	3.1
IVA (LQ=8, NI=8)	41.6	3.0	54.0	3.6	46.5	2.8	51.2	3.1
IVA (LQ=16, NI=8)	42.1	3.0	<u>64.9</u>	<u>3.9</u>	50.5	3.0	53.5	<u>3.2</u>
IVA (LQ=32, NI=8)	41.9	3.0	57.1	3.7	51.9	3.1	53.7	3.2
IVA (LQ=16, NI=4)	42.2	3.0	63.3	3.9	50.1	3.0	52.5	3.2
IVA (LQ=16, NI=16)	42.3	3.0	55.4	3.7	50.0	3.0	55.1	3.3
IVA (LQ=16, NI=8)-272K	46.8	3.1	68.0	4.0	<u>48.1</u>	<u>2.9</u>	50.9	3.1

Table 1: **Comparison between different methods on 4 zero-shot long video QA datasets.** LLM with IVA achieves best performance on long videos compared to baselines and strong video LLMs. "LQ" refers to the length of query tokens and "NI" represents the number of interactions between LLMs and IVA. "-272K" indicates that we introduce additional training data of long video datasets like LifeQA and Social-IQ based on the original short video data.

videos. Video-ChatGPT (Muhammad Maaz and Khan, 2023) encodes frames independently and generates frame-level embeddings. Subsequently, it employs average pooling to transform these embeddings into both temporal and spatial features. These temporal and spatial features are then concatenated to derive video-level features and are fed into the LLM. Video-LLaMA (Zhang et al., 2023c) utilizes Vision-Language and Audio-Language to process video frames and audio signals separately. After fine-tuning on image instruction dataset and video instruction dataset, Video-LLaMA exhibited remarkable abilities in comprehending images and videos. Video-Chat (Li et al., 2023f) leverages perception tools to convert videos into textual descriptions in real-time, and employs a foundation model named InternVideo to encode videos into embeddings. These textual descriptions and video embeddings are then processed by an LLM for multimodal understanding. LLaMA-Adapter (Zhang et al., 2023d) is a lightweight adapter injected into the attention calculation of LLM, which could be used to handle videos, text, and image tasks.

### 4.3 Evaluation Metrics

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

For open-ended video QA tasks, we employ 444 ChatGPT-Assistant to evaluate the performance 445 following Video-ChatGPT (Muhammad Maaz and 446 Khan, 2023). First, we input the question, the pre-447 dicted answer, and the correct answer into Chat-448 GPT. Second, we request ChatGPT to verify the 449 accuracy of the predicted answer, expecting a bi-450 nary response of 'yes' for correct predictions or 451 'no' for incorrect ones. Additionally, we require 452 ChatGPT to rate the quality of the predicted an-453

swer on a scale from 0 to 5, where 5 indicates a perfect match. Finally, we determine the overall accuracy by counting the number of 'yes' responses and calculate the overall score by averaging all quality scores. This evaluation employs the "gpt-3.5-turbo" version of ChatGPT.

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

#### 4.4 Implementation Details

We employ the AdamW optimizer (Kingma and Ba, 2014) in conjunction with a cosine learning rate scheduler to train our model. We first utilize 2 A100 GPUs to train visual-language MLP with 2 million image-text pairs with a global batch size of 256 and a base learning rate of 2e-4. Subsequently, we train the causal transformers using 703K videotext pairs data on the same two GPUs, employing a global batch size of 24 and a base learning rate of 3e-4. Transitioning to the video instruction tuning stage, we scale up to 8 A100 GPUs with a global batch size of 64. Here, we leverage LoRA to efficiently fine-tune the language model LLaMA. In our implementation, we set the rank to 128 and alpha to 256, maintaining a learning rate of 1e-4 for both LoRA and IVA parameters. Given the pretraining visual-language MLP and causal transformers, we adopt a smaller learning rate of 2e-5.

#### 4.5 Main Results

We present the performance of the models on four long video QA benchmarks and five short video QA benchmarks. In zero-shot long video QA benchmarks, our model achieved state-of-theart (SOTA) results compared to the previous pure video LLMs, except WildQA. Especially on the LifeQA and Social-IQ 2.0 evaluation datasets,

Mathad	MSVD-QA		MSRVTT-QA		SEED <sup>AR</sup>	SEED <sup>AP</sup>	SEED <sup>PU</sup>	
Method	Accuracy	score	Accuracy	score	Accuracy	Accuracy	Accuracy	
Valley	-	-	-	-	31.3	23.2	20.7	
Video-LLaMA	51.6	2.5	29.6	1.8	-	-	-	
LLaMA-Adapter	54.9	3.1	43.8	2.7	-	-	-	
Video-Chat	56.3	2.8	45.0	2.5	34.9	36.4	27.3	
Video-ChatGPT	64.9	3.3	49.3	2.8	27.6	21.3	21.1	
Baseline (w/o IVA)	54.5	3.2	49.6	2.9	22.5	23.5	24.8	
IVA (LQ=8, NI=8)	53.2	3.2	47.6	2.9	32.0	31.8	27.5	
IVA (LQ=16, NI=8)	55.7	3.2	49.1	2.9	35.2	32.0	34.2	
IVA (LQ=32, NI=8)	53.0	3.2	47.2	2.9	32.2	32.1	28.8	
IVA (LQ=16, NI=4)	55.0	3.2	47.8	2.9	32.5	31.7	26.0	
IVA (LQ=16, NI=16)	53.3	3.1	47.1	2.8	31.8	29.4	31.0	
IVA (LQ=16, NI-8)-272K	58.6	3.2	50.2	2.9	32.2	30.0	31.6	

Table 2: **Comparison between different methods on 5 zero-shot short video QA benchmarks.** Benchmark names are abbreviated due to space limits. MSVD-QA(Xu et al., 2017); MSRVTT-QA(Xu et al., 2017); SEED<sup>AR</sup>: SEED-Bench Action Recognition(Li et al., 2023c); SEED<sup>AP</sup>: SEED-Bench Action Prediction(Li et al., 2023c); SEED<sup>PU</sup>: SEED-Bench Procedure Understanding(Li et al., 2023c).

our model achieved significantly higher results, surpassing the previous SOTA accuracy by **18.0** and **7.4** percentage points, respectively. In zeroshot short video QA benchmarks, our model also demonstrated strong capabilities across some evaluation datasets, especially in procedure understanding. Overall, IVA significantly enhances the capability of LLMs to analyze and interpret long videos, maintaining high-performance levels without compromising the understanding and reasoning abilities of short videos.

#### 4.6 Ablation Study

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

503

506

**Effect of IVA**. From the results in Tables 1 and 2, we found that introducing the IVA module improved the overall visual understanding of the long video datasets of Social IQ2, LifeQA, and ActivityNet-QA, as well as the short video datasets. Among them, our model achieved an improvement of over 20% on LifeQA compared to the baseline, notably suggesting the effectiveness of IVA.

Length of Query Tokens. Comparing the experimental results of IVA (LQ=8, NI=8) and IVA 508 (LQ=16, NI=8) in Tables 1 and 2, we observed 509 a significant decrease in evaluation results across 510 various benchmarks when reducing the length of 511 query tokens (16  $\rightarrow$  8). Regarding the compari-512 son between IVA (LQ=16, NI=8) and IVA (LQ=32, 513 NI=8) in long video benchmarks, we noted a slight 514 decrease in performance on the first two bench-515 marks when increasing the length. However, while 516 there was a slight improvement in LifeQA, it did 517 not conclude an overall performance enhancement. 518

In contrast, in the short video benchmarks, there was a downward trend in results across all benchmarks. Overall, increasing the length of query tokens may not lead to performance improvement. Moreover, reducing the length of query tokens may result in the loss of crucial visual information, consequently leading to performance degradation. 519

520

521

522

523

524

525

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

Impact of long video instruction data. Furthermore, we attempted to further enhance the model's performance during the Video Instruction Tuning phase by introducing more long video data. Therefore, we included the training sets of Social IQ2 and LifeQA, the video instruction part of the MIMIC-IT dataset(Li et al., 2023a), along with the open-ended question-answering training data from NExT-QA, to the existing 170K training data, forming a new 272K training dataset. From the results of IVA(lQ=16, NI=8)-272K in the Tables 1 and 2, we observed a significant improvement on the long video dataset Social IQ2 with the inclusion of more training data. However, there was little difference in the results on the remaining datasets, and in some cases, a certain degree of decline was observed. This may be attributed to the somewhat indiscriminate addition of new datasets, leading to a certain imbalance in the proportions of different data. Additionally, the training set of LifeQA only consisted of 1,383 instances, which is relatively small in proportion to the total data, thus not providing sufficient improvement.

Number of Interaction between IVA and LLMs.

Similar to the analysis of Query Tokens Length, we also conducted experiments with both doubled and



Figure 2: Five cases illustrate the comparative performances of our IVA Model and Baseline. The bottom part shows the detailed description of a video. Red words represent the inaccurate statement and the green words indicate the accurate statement.

halved number of interaction layers. The detailed injection layers are shown in Appendix A. Upon analyzing the results of IVA(LQ=16, NI=8) and IVA(LQ=16, NI=4) in Tables 1 and 2, we observe that this reduction resulted in a significant decrease in its performance on most long video datasets, especially on the Action Prediction and Procedure Understanding of the SEED-Bench. Moreover, the experimental results also indicate that increasing the number of layers ( $8 \rightarrow 16$ ) in the IVA interaction likewise caused a slight degradation in the model's performance. Given that there was no significant improvement observed when increasing the interaction times between IVA and LLMs, we set it to 8 as the standard for experimentation.

### 4.7 Case Study

553

555

557

559

562

563

564

565

570

571

574

575

577

We present four open-ended question-answering cases and one detailed description example in Figure 2. Upon examining the initial two examples, we observe that the model augmented with IVA exhibits enhanced proficiency in recognizing particular actions associated with specific frames. In response to specific queries, it could discern objects such as the 'basketball-shaped cake', which solely appears towards the video's conclusion, and the 'glass bowl,' present solely in the video's opening segment. Furthermore, the fourth questionanswering example illustrates that IVA augments the model's reasoning ability, enabling it to deduce the prevailing weather conditions based on the lighting conditions within the video. These indicate the effectiveness of IVA in incorporating fine-grained visuals of long videos. Meanwhile, the bottom detailed description example reveals that when confronted with lengthy video descriptions, IVA could refine the perceptual acuity of LLMs, resulting in more precise recognition of elements such as the environment and color. 579

580

582

583

584

586

587

588

589

592

593

594

595

596

598

600

601

602

604

# 5 Conclusion

In this study, our primary goal is to enhance the capacity of LLMs to process and interpret long video content effectively. We identified the principal obstacles in this area and introduced an Interactive Visual Adapter (IVA) designed to facilitate dynamic interaction between LLMs and extended video sequences. The IVA incorporates a selector module for identifying relevant temporal frames within long videos based on specific instructions and tokens, along with an interactor module that isolates detailed spatial visual features within long videos. The empirical results demonstrate that our IVA significantly improves LLMs' ability to comprehend and reason about long video content.

8

700

701

703

704

705

652

653

654

655

656

# Limitations

Our work, while contributing valuable insights into video understanding through LLMs, is subject to several limitations that warrant further investigation:

• Optimization for Longer Videos: Our current methodology demonstrates proficient per-611 formance in processing videos ranging from 612 a few seconds to two minutes. However, the 613 challenge of comprehensively understanding 614 longer videos remains. Specifically, the optimization of video token length and the integra-616 tion method of the Interactive Visual Adapter 617 (IVA) within LLMs require further refinement 618 to enhance their effectiveness and efficiency in handling extended content.

Impact of Interaction Frequency and Query Token Length: The stability of the IVA can be influenced by the frequency of interactions and the length of query tokens. These factors often occur in the development of multimodal large models, where a delicate balance must be struck between achieving high performance and maintaining operational efficiency, particularly in the context of long video interaction and encoding.

• Accuracy and Appropriateness of Generated Responses: Another limitation is the potential for LLMs to generate responses that may be inaccurate, contain harmful content, or be factually incorrect. This issue stems from the inherent unpredictability in the response generation process of LLMs, underscoring the need for mechanisms that can ensure the reliability and appropriateness of the output.

# References

632

633

634

638

639

641

643

646

647

650

 Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. 2021.
 Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846.

- 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 *IEEE International Conference on Computer Vision*.
- Santiago Castro, Mahmoud Azab, Jonathan Stroud, Cristina Noujaim, Ruoyao Wang, Jia Deng, and Rada

Mihalcea. 2020. LifeQA: A real-life dataset for video question answering. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4352–4358, Marseille, France. European Language Resources Association.

- Santiago Castro, Naihao Deng, Pingxuan Huang, Mihai G. Burzo, and Rada Mihalcea. 2022. In-the-wild video question answering. In *COLING*, pages 5613– 5635, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- 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 opensource chatbot impressing gpt-4 with 90%\* chatgpt quality.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- 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*.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*.
- Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. 2021. Multiscale vision transformers. In Proceedings of the IEEE/CVF international conference on computer vision, pages 6824–6835.
- Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. 2022. Masked autoencoders as spatiotemporal learners. *Advances in neural information processing systems*, 35:35946–35958.
- Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. 2016. Convolutional two-stream network fusion for video action recognition. In *Proceedings* of the IEEE conference on computer vision and pattern recognition, pages 1933–1941.
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. 2023. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*.

706

- 752

- Rohit Girdhar, Alaaeldin El-Nouby, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Omnimae: Single model masked pretraining on images and videos. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 10406–10417.
  - Agrim Gupta, Stephen Tian, Yunzhi Zhang, Jiajun Wu, Roberto Martín-Martín, and Li Fei-Fei. 2022. Maskvit: Masked visual pre-training for video prediction. In The Eleventh International Conference on Learning Representations.
  - Seungju Han, Jack Hessel, Nouha Dziri, Yejin Choi, and Youngjae Yu. 2023. Champagne: Learning real-world conversation from large-scale web videos. arXiv preprint arXiv:2303.09713.
  - Peng Jin, Jinfa Huang, Pengfei Xiong, Shangxuan Tian, Chang Liu, Xiangyang Ji, Li Yuan, and Jie Chen. 2023. Video-text as game players: Hierarchical banzhaf interaction for cross-modal representation learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2472-2482.
  - Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. 2014. Large-scale video classification with convolutional neural networks. In CVPR.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of NAACL-HLT, pages 4171–4186.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. 2023a. Mimic-it: Multi-modal in-context instruction tuning.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023b. Seed-bench: Benchmarking multimodal llms with generative comprehension. arXiv preprint arXiv:2307.16125.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023c. Seed-bench: Benchmarking multimodal llms with generative comprehension. arXiv preprint arXiv:2307.16125.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023d. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. ICML.
- KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023e. Videochat: Chat-centric video understanding. arXiv preprint arXiv:2305.06355.

Kunchang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023f. Videochat: Chat-centric video understanding. arXiv preprint arXiv:2305.06355.

759

760

761

765

766

767

768

769

770

775

781

782

783

784

785

786

787

791

792

793

796

798

799

800

801

802

803

804

805

806

807

809

810

811

- Xinyu Li, Yanyi Zhang, Chunhui Liu, Bing Shuai, Yi Zhu, Biagio Brattoli, Hao Chen, Ivan Marsic, and Joseph Tighe. 2021. Vidtr: Video transformer without convolutions. arXiv e-prints, pages arXiv-2104.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. arXiv preprint arXiv:2304.08485.
- Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. 2023. Valley: Video assistant with large language model enhanced ability. arXiv preprint arXiv:2306.07207.
- Hanjia Lyu, Jinfa Huang, Daoan Zhang, Yongsheng Yu, Xinyi Mou, Jinsheng Pan, Zhengyuan Yang, Zhongyu Wei, and Jiebo Luo. 2023. Gpt-4v (ision) as a social media analysis engine. arXiv preprint arXiv:2311.07547.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. arXiv preprint arXiv:2306.05424.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. arXiv preprint arXiv:2211.01786.
- Salman Khan Muhammad Maaz, Hanoona Rasheed and Fahad Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. ArXiv 2306.05424.

OpenAI. 2023. Chatgpt. OpenAI Blog.

- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. arXiv:2203.02155.
- Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, et al. 2023. Tool learning with foundation models. arXiv preprint arXiv:2304.08354.
- Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, et al. 2023. Moviechat: From dense token to sparse memory for long video understanding. arXiv preprint arXiv:2307.16449.

Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. 2019. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7464–7473.

812

813

814

816

817

818

819

820

821

824

825

826

827

831

832

837

838

842

843

844

845

847

851

853

860

- Yuchong Sun, Hongwei Xue, Ruihua Song, Bei Liu, Huan Yang, and Jianlong Fu. 2022. Long-form videolanguage pre-training with multimodal temporal contrastive learning. *Advances in neural information processing systems*, 35:38032–38045.
- Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, et al. 2023. Video understanding with large language models: A survey. *arXiv preprint arXiv:2312.17432*.
- Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. 2015. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497.
- Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. 2016. Temporal segment networks: Towards good practices for deep action recognition. In *European conference on computer vision*, pages 20–36. Springer.
- Alex Wilf, Leena Mathur, Sheryl Mathew, Claire Ko, Youssouf Kebe, Paul Pu Liang, and Louis-Philippe Morency. 2023. Social-iq 2.0 challenge: Benchmarking multimodal social understanding. https://github.com/abwilf/ Social-IQ-2.0-Challenge.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*.
- 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)*, 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 *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653.
- Jiaqi Xu, Cuiling Lan, Wenxuan Xie, Xuejin Chen, and Yan Lu. 2023. Retrieval-based video language model for efficient long video question answering. *arXiv preprint arXiv:2312.04931*.
- Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. 2022. Clipvip: Adapting pre-trained image-text model to videolanguage alignment. In *The Eleventh International Conference on Learning Representations*.

Qinghao Ye, Guohai Xu, Ming Yan, Haiyang Xu, Qi Qian, Ji Zhang, and Fei Huang. 2023. Hitea: Hierarchical temporal-aware video-language pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15405–15416. 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

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

- Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *AAAI*, pages 9127–9134.
- Joe Yue-Hei Ng, Matthew Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat Monga, and George Toderici. 2015. Beyond short snippets: Deep networks for video classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4694–4702.
- Daoan Zhang, Weitong Zhang, Bing He, Jianguo Zhang, Chenchen Qin, and Jianhua Yao. 2023a. Dnagpt: A generalized pretrained tool for multiple dna sequence analysis tasks. *bioRxiv*, pages 2023–07.
- Hang Zhang, Xin Li, and Lidong Bing. 2023b. Videollama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.
- Hang Zhang, Xin Li, and Lidong Bing. 2023c. Videollama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.
- Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. 2023d. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.
- Linchao Zhu and Yi Yang. 2020. Actbert: Learning global-local video-text representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8746–8755.
- A Inserting IVA in Different Layers

NI	<b>Corresponding Decoder Layers</b>
4 8	0, 8, 16, 24 0, 4, 8, 12, 16, 20, 24, 28
16	0, 2, 4, 6, 8, 10,, 22, 24, 26, 28, 30

Table 3: Ablation Study on Injection Layers for IVA. NI: Number of Inserting Layers. The incorporated inserting layers were positioned before the respective decoder layers.

In this section, we detail the methodology behind our ablation studies focusing on the variation in the

912	Number of Injection Layers. Our experiments were
913	structured around three different setups, where the
914	injection layers were configured to be 4, 8, and 16
915	in number. To ensure a uniform distribution, these
916	Injection Layers were interspersed throughout the
917	decoder layers of the language model evenly. We
918	utilized the Vicuna-7B model as our experimental
919	framework, which is equipped with 32 decoder lay-
920	ers. The specific layers of the decoder that received
921	the Injection Layers are outlined in Table 3, pro-
922	viding a clear reference to how the integration was
923	achieved in each experimental setup