WOLF: ACCURATE VIDEO CAPTIONING WITH A WORLD SUMMARIZATION FRAMEWORK

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose *Wolf*, a WOrLd summarization Framework for accurate video captioning. Wolf is an automated captioning framework that adopts a mixture-ofexperts approach, leveraging complementary strengths of Vision Language Models (VLMs). By utilizing both image and video models, our framework captures different levels of information and summarizes them efficiently. Our approach can be applied to enhance video understanding, auto-labeling, and captioning. To evaluate caption quality, we introduce CapScore, an LLM-based metric to assess the similarity and quality of generated captions compared to the ground truth captions. We further build four human-annotated datasets in three domains: autonomous driving, general scenes, and robotics, to facilitate comprehensive comparisons. We show that Wolf achieves superior captioning performance compared to state-of-the-art approaches from the research community (VILA1.5, CogAgent) and commercial solutions (Gemini-Pro-1.5, GPT-4V). For instance, in comparison with GPT-4V, Wolf improves CapScore both quality-wise by 55.6% and similarity-wise by 77.4%on challenging driving videos. Finally, we establish a benchmark for video captioning and introduce a leaderboard, aiming to accelerate advancements in video understanding, captioning, and data alignment.

027 028 029

025

026

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

031 Video captioning is crucial as it facilitates content understanding and retrieval by providing accurate, searchable descriptions. It also provides pairwise data for effective training of foundation models 033 for tasks like video generation, such as Sora (Brooks et al., 2024) and Runaway (Runway, 2024). 034 However, generating descriptive, accurate, and detailed video captions remains a challenging research problem for several reasons: firstly, high-quality labeled data are scarce. Video captions from the 035 internet can be faulty and misaligned and human annotation is prohibitively expensive for large 036 datasets. Secondly, video captioning is inherently more challenging than image captioning due to the 037 additional complexity of temporal correlation and camera motion. Existing captioning models (Hong et al., 2024; Zhang et al., 2023) struggle with temporal reasoning and fail to achieve accurate scene understanding. Thirdly, there is no established benchmark to measure captioning progress. Existing 040 video QA benchmarks (Maaz et al., 2023) are often limited to short answers, making it difficult to 041 measure hallucinations in detailed long captions. Fourthly, the correctness and completeness of the 042 captions are crucial for safety-critical tasks. In the era of LLMs, text descriptions of scenarios used 043 by embodied agents for planning and control become increasingly common (Mao et al., 2023a;b; 044 Li et al., 2024; Ding et al., 2023). Consequently, a false or incomplete description of the scenario may lead to the decision-making module overlooking a critical object after training on such caption data, resulting in safety risks. For instance, missing the presence of a human in the vicinity of a 046 vegetable-chopping manipulator can lead to an injury. 047

To handle these challenges, we introduce <u>WOrL</u>d summarization Framework (*Wolf*), a novel summa rization captioning framework, along with a captioning metric CapScore, and the Wolf captioning
 benchmark with corresponding datasets. Unlike previous works that utilize a single model to generate
 captions, we propose to use multiple models to collaborate (Jiang et al., 2024), producing much more
 accurate captions. By leveraging multiple models, we can provide more fine-grained details while
 reducing hallucinations. We show that Wolf achieves superior captioning performance compared
 to state-of-the-art approaches from the research community (such as VILA (Lin et al., 2023c), Co-

gAgent (Hong et al., 2024)) and commercial solutions (such as Gemini-Pro-1.5 (Team et al., 2023), GPT-4V (OpenAI, 2023)). In summary, we have three main contributions:

- 1. We design the first world summarization framework **Wolf** for video captioning and introduce an LLM-based metric **CapScore** for evaluating the quality of captions. We have further verified that CapScore aligns with human evaluations. The results show that our method improves CapScore by a large margin.
- 2. We introduce Wolf benchmark and four human-annotated benchmark datasets. These datasets include autonomous driving, general scenes from Pexels, and robotics videos, along with human-annotated captions, referred to as the **Wolf Dataset**.
- 3. The code, data and leaderboard will be open-sourced and maintained ¹. Continuous efforts and improvements will be made to refine the Wolf Dataset, codebase, and CapScore. We hope that Wolf will raise awareness about the quality of video captioning, set a standard for the field, and boost community development.
- 2 RELATED WORKS

Image Captioning. Visual language models (VLMs) have shown rapid advancements, achieving 071 leading performance in image captioning tasks, largely due to the success of large language models. CLIP (Radford et al., 2021) pioneered this field by training a shared feature space for vision and 073 language modalities on image-caption pairs. Building on CLIP, BLIP (Li et al., 2022) and BLIP-2 (Li 074 et al., 2023) improved performance by aligning the pre-trained encoder with large language models. 075 Following the direction, LLaVA (Liu et al., 2023) and InstructBLIP (Dai et al., 2023) demonstrated 076 that jointly training on diverse datasets as an instruction-following task leads to strong generalization 077 across various tasks. VILA (Lin et al., 2023c) highlighted the importance of pre-training with diverse data, and therefore significantly scaled up the pre-training dataset. Kosmos-2 (Peng et al., 2023) and PaLI-X (Chen et al., 2023a) further introduced pseudo-labeling bounding boxes from open-vocabulary 079 object detectors to scale up the size of pre-training dataset.

081 Video Captioning. As image-based VLMs are not specifically trained with video data, they are limited in describing details present in the video data. To improve video captioning, PLLaVa (Xu 083 et al., 2024) builds on top of LLaVa and introduced a parameter-free pooling strategy to enhance the caption quality. Video-llava (Lin et al., 2023a) achieves state-of-the-art performance on several 084 benchmarks by conducting joint training on images and videos, thereby learning a unified visual 085 representation. Additionally, Video-LLama (Zhang et al., 2023) incorporates both video and audio into LLMs by introducing two Q-formers to extract features. Vid2seq (Yang et al., 2023) conducts 087 large-scale pre-training with narrated videos for dense video captioning. Meanwhile, MV-GPT (Seo 088 et al., 2022) employs an automated speech recognition (ASR) model to provide additional labeling for the videos. 090

LLM-based Summarization. Recently many works have found that it is efficient to summarize useful information using LLMs. For example, LLaDA (Li et al., 2024) can provide users with helpful instructions based on the user request and corresponding traffic rules in the desired location. OpenAI team finds re-captioning (Betker et al., 2023) via LLMs can be very helpful.

095 096

097

054

056

059

060 061

062

063 064

065

066

067

069

3 WOLF: CAPTIONING EVERYTHING WITH A WORLD SUMMARIZATION FRAMEWORK

We propose Wolf, which is an automated captioning summarization framework that adopts a mixture of experts approach to generate long, accurate, and detailed captions for videos. Figure 1 provides an overview of our framework. In this paper, we use CogAgent (Hong et al., 2024), GPT-4V (Mao et al., 2023a) to generating image-level captions, and use VILA-1.5 (Lin et al., 2023c), Gemini-Pro-1.5 (Team et al., 2023) to generate video captions.

Chain-of-thought Summarization in Image-level Models. As image-level models (image-based VLMs) have been pre-trained with a larger amount of data than video-level models (video-based VLMs), we first use image-based VLMs to generate captions. We design a Chain-of-thought program to obtain video captions from image-level models. As illustrated in Figure 1, we first split the video

¹⁰⁷

¹We also provide ethical statement and reproducibility in Appendix.



Figure 1: Overview of proposed Wolf framework. Wolf utilizes both image-level and video-level models to generate diverse and detailed captions, which are then summarized for cross-checking.

into sequential images, sampling two key-frames every second. We start by feeding Image 1 into the 122 Image-level Model to obtain Caption 1, where we require the model to generate detailed scene-level 123 information and object locations. Given the temporal correlation between key frames in a video, 124 we then feed both Caption 1 and Image 2 into the model to generate Caption 2. By repeating this 125 procedure, we generate captions for all sampled frames. Finally, we use GPT-4 to summarize the 126 information from all captions with the prompt "Summarize all the captions to describe the video with 127 accurate temporal information". Additionally, we extract the bounding box locations for each object 128 in each frame, then feed them into LLMs to summarize the trajectory of the moving object. For 129 example, in a driving video, a blue car is driving into the right lane, and the centers of the bounding 130 boxes are (0,0), (1,1), (1,2). We provide the car's location to the LLM, and it outputs 'the blue car is 131 driving to the right,' which we refer to as a **Motion Caption**.

132 **LLM-based Video Summarization.** Besides obtaining the captions from image-level models, we 133 then summarize all captions into one. We use the prompt "Please summarize on the visual and 134 narrative elements of the video in detail from descriptions from Image Models (Image-level Caption 135 and Motion Caption) and descriptions from Video Models (Video-level Caption)". Optionally, we 136 can also add the annotated caption to the summarization. Based on this simple scheme, Wolf can 137 capture a rich variety of details of the video and reduce hallucinations (in Figure 2). We assume this is because the model can compare the captions and reduce redundant and hallucinated information. 138 After obtaining the descriptions from the image-level and video-level models, we next apply the 139 prompt "Please describe the visual and narrative elements of the video in detail, particularly the 140 motion behavior". 141

141 142 143

4 WOLF BENCHMARK: BENCHMARKING VIDEO CAPTIONING

To showcase the effectiveness of Wolf, we constructed four distinct datasets. These include two 144 autonomous driving video captioning datasets based on the open-sourced NuScenes (Caesar et al., 145 2019) dataset (Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International Public 146 License), a general daily video captioning dataset from Pexels², and a robot manipulation video cap-147 tioning dataset from an open-source robot learning dataset (Padalkar et al., 2023). These benchmark 148 datasets are tailored to assess the caption model's scene comprehension and its behavior understand-149 ing capabilities, both of which are vital for auto-labeling in embodied AI tasks. All captions were 150 generated using a combination of ground truth information, rule-based heuristics, human labeling, 151 and GPT-based rewriting.

152

4.1 BENCHMARK DATASET CURATION

4.1.1 DENCIMARK DATASET CORATION 4.1.1 AUTONOMOUS DRIVING DATASET

High-quality captions of driving videos are crucial not only for training video generation models
but also for training VLMs to interpret the dynamic traffic environment. The NuScenes dataset is
a large-scale collection of driving videos designed to accelerate autonomous driving research. It
features 1,000 annotated scenes from Boston and Singapore. Each scene consists of a 20-second
driving video clip that provides an ego-centric view from the ego vehicle. We split each scene into
5-second segments and provide the corresponding captions. Our captions emphasize the high-level
driving behavior of the ego vehicle to stress-test the scene understanding ability and the behavior

²https://www.pexels.com/

162 Puntha 163 164 165 t = 1 [s]. Key features: t = 4 [s]. Key features: t = 8 [s]. Key features 166 t = 0 [s]. Key features: lane. person, beach, ocean, person, ocean, smile, person, ocean, white long hair, sunglass. sunglass, phone shirt, dark shorts 167 Caption: The video opens with a person standing on a rocky beach, holding a 168 smartphone. They are dressed in a white shirt and dark shorts, with long hair and sunglasses. In the background, there are other beachgoers and the 169 ocean. The person is seen taking a selfie, with the smartphone's screer visible in some frames. The lighting suggests it is late afternoon or early 170 evening, with the sun low on the horizon, casting a warm glow on the scene t = 2 [s]. Key features: lane traffic cones construction zone 171 The person's actions are casual and relaxed, as they pose and interact with the camera. The ocean is calm. and the sky is clear, indicating good 172 weathe 173 174 175 t = 5 [s]. Key features: lane, traffic cones, construction zone, land change 176 EA Caption: The footage captures a car navigating through an urban 177 t = 4 [s]. Key features: t = 6 [s]. Key features: t = 8 [s]. Key features: area where traffic cones are placed in its path. The cones indicate a construction zone ahead and block the vehicle's robot arm, bottle, table. push the bottle. the bottle falls onto table lane. The car maintains a steady speed and smoothly switches Caption: The video shows a yellow industrial robotic arm positioned over a 179 lanes to pass by each traffic cone on the right side. table with a green cover. The robotic arm is attempting to push a green bottle from the left side to the right side. However, during the process, the bottle falls onto the table 181

Figure 2: Wolf Dataset examples. We display the videos and corresponding human-annotated captions of autonomous driving (*Left*), Pexels (*Top-Right*), and Robot learning video dataset (*Bottom-Right*), totaling 25.7 hours for now, and the dataset size will be regularly updated and expanded.

understanding ability of the captioning model. Our dataset contains **500 intensely interactive video-caption pairs** (\approx 0.7 hours) in which the ego vehicle is involved in intense interactions with its surrounding traffic agents (such as navigating around construction zones and overtaking static obstacles) and **4785 normal driving scene video-caption pairs** (\approx 6 hours). Our caption generation process consists of three steps: i) agent-level motion annotation, ii) ego-centric interaction annotation, and iii) GPT-rewriting.

Agent-level motion annotation. The NuScenes dataset provides full annotation of the traffic elements in each scene, including the 3-D bounding box and categories of traffic elements, and semantic map information. Similar to Tian et al. (2024), we leverage this ground-truth information and the lane-topology information (Naumann et al., 2023) to annotate both the speed and angular motion characteristics of the ego vehicle and other traffic participants within a video clip. Specifically, we categorize agent actions into 11 types such as Stopping, Accelerating, Decelerating, Lane Changes, Turns, and more, based on their observed movements and behaviors.



191

192

193

194

195

196

197

199

200

201

202

203

204

Figure 3: Illustration of homotopy types of different relative motions between a pair of vehicles. **Ego-centric interaction annotation**. We are also interested in the ego vehicle's interaction with the other traffic participants (e.g., crossing pedestrians, blocking traffic cones, etc.) shown in the video clip. To efficiently caption the interaction, we leverage two types of categorical modes to describe the lane-relationship between a traffic participant and the ego vehicle (*agent*-

205 ego lane mode) and the relative motion between a traffic participant and the ego vehicle (homo-206 topy)(Chen et al., 2023b). Agent-ego lane mode at a time step t encodes the topology relationship between the ego's current lane and the traffic agent's lane, including: LEFT, RIGHT, AHEAD, 207 BEHIND, and NOTON, where NOTON describes that the traffic agent is not on any derivable lanes 208 in the scene (e.g., a parked vehicle in a parking lot). To compute the agent-ego lane mode for each 209 traffic agent, we follow (Chen et al., 2023b) to first identify the lane on which each agent is located 210 and then leverage the lane topology map to annotate the agent-ego lane mode. We project the agent's 211 center to the lane polyline and use its relative position in the local Frenet frame to determine its lane 212 association. Homotopies describe the relative motion between a pair of agents shown in the video, 213 including: [S, CW, CCW] (static, clockwise, counterclockwise), as shown in Figure 3. 214

GPT-rewriting. Combining agent-ego lane mode, homotopy, agent ground truth state information, and scene context information (e.g., ego is located near intersection) together, we can leverage

216 heuristics to annotate the interaction shown in the video clip. For example, in a video clip, a static 217 object's agent-ego lane mode changes from AHEAD, to LEFT, to BEHIND, and the ego vehicle's first 218 performs RIGHT-LANE-CHANGE, KEEP-LANE, then LEFT-LANE-CHANGE, indicating the ego 219 vehicle overtakes that object from the ego vehicle's left side. We identified 6 interaction categories 220 from the NuScenes dataset: 1) bypass blocking traffic cones to navigate around construction zone; 2) yield to crossing pedestrians; 3) yield to incoming vehicles; 4) overtake traffic agents via straddling the 221 lane dividers; 5) overtake traffic agent via lane-change; 6) other non-intensive interactions. With both 222 agent-level motion annotation and ego-centric interaction annotation, we use GPT 3.5 to summarize 223 each clip to build the final caption. 224

225 4.1.2 ROBOT MANIPULATION DATASET

In addition to the driving environment, we collect **100 robot manipulation videos** (each has a length ranging from 5 seconds to 1 minute) from Padalkar et al. (2023) that demonstrate complex robot manipulations (e.g., pick and place, push, ect.) in various environments, including kitchen, office, lab, and open world. We manually caption each video. The captions focus on the description of the scene and the interaction between the robot and the objects (see the example in Figure 2).

231 4.1.3 PEXELS DATASET

232

266

To evaluate caption models in general daily environments, we further collect high quality (360p to 1080p) videos from Pexels³. It consists of **473 high-quality videos** sourced globally, where each video has a length varying between 10 seconds and 2 minutes and the content includes 15 popular categories (details in Appendix). This diversity not only adds depth to our dataset but also provides a wide range of scenarios and contexts for our analysis.

- 238 4.2 EVALUATION METRIC AND LEADERBOARD
- 4.2.1 CAPSCORE: EVALUATING CAPTIONS WITH LLMS

241 Video captioning has been an ill-posed problem since there is no metric to evaluate the quality of 242 captions and the alignment between the video and the caption. Inspired by BERTScore (Zhang et al., 243 2019) and CLIPScore (Hessel et al., 2021), we introduce CapScore (Captioning Score), a quantitative 244 metric to use LLMs to evaluate the similarity between predicted and human-annotated (ground-truth) captions. We tried both GPT-4 and LLama 3.1 (Dubey et al., 2024) as our LLM to summarize the 245 captions. We noticed that GPT-4 can always obtain stable results over 3 runs. However, for LLama 246 3.1, the results varied over different runs. We tried to lower the temperature (from 0.9 to 0.5) to make 247 the inference stable, however, we noticed that the scores are not consistent with human evaluation. 248 Therefore we select GPT-4 as our LLM to conduct the experiments. Assume we have 6 captions, we 249 feed all the captions into GPT-4 and add the prompt "Can you give a score (two decimal places) from 250 0 to 1 for captions 1, 2, 3, 4 and 5, indicating which one is closer to the ground truth caption (metric 251 1) and which contains fewer hallucinations and less misalignment (metric 2)? Please output only the 252 scores of each metric separated only by a semicolon. For each metric, please output only the scores 253 of captions 1, 2, 3, 4 and 5 separated by commas, in order-no text in the output.". We ask GPT-4 to 254 output two scores: caption similarity and caption quality.

Caption Similarity. Caption similarity is based on how well each caption aligns with the ground truth description on a scale from 0 to 1, considering the key criteria mentioned. GPT-4 lists the requirements that affect the score: this metric measures how similar each caption is to the ground truth caption. The evaluation focuses on the content and context described in the captions, assessing whether they capture the main themes and details of the ground truth.

Caption Quality. Caption quality evaluates whether the caption contains reduced hallucination and mistakes compared to the ground truth captions on a scale from 0 to 1. GPT-4 lists the criteria that affect the score: this metric evaluates the accuracy and relevance of each caption, identifying any extraneous or incorrect details (hallucinations). Captions with fewer hallucinations and better alignment receive higher scores.

4.2.2 HUMAN-EVALUATION SCORE AND CAPSCORE

³https://www.pexels.com/



Figure 4: Comparisons on Human-Evaluation Score and CapScore.

concerns related to human alignment and correlation, we randomly selected 10 users to evaluate our set of 100 robotics videos, as detailed in Table 1 of the paper. The evaluators were presented with the videos, the generated captions, and the corresponding ground truth captions. We asked them to assign human-evaluation scores based on the CapScore standard, with the following prompt: "After reviewing the video and all the captions, please assign the caption similarity and caption quality score (floating point values) from 0 to 1 for different captions, indicating which caption is closest to the ground truth (caption similarity) and which one has fewer hallucinations and less misalignment (caption quality)." We show the results in Table 1 the corresponding visual comparison in Figure 4.

Method	Caption Similarity ↑		\mid Caption Quality (eg. reduced hallucination) \uparrow		
	Human-Evaluation Score	CapScore	Human-Evaluation Score	CapScore	
CogAgent (Hong et al., 2024)	0.28	0.38	0.34	0.43	
GPT-4V (Achiam et al., 2023)	0.43	0.34	0.42	0.35	
VILA-1.5 (Lin et al., 2023c)	0.62	0.62	0.68	0.67	
Gemini-Pro-1.5 (Team et al., 2023)	0.66	0.63	0.72	0.67	
Wolf	0.74	0.72	0.80	0.75	

Table 1: Comparison of Human-Evaluation Score and CapScore on 100 Wolf Robotics Videos.

4.2.3 BENCHMARKING VIDEO CAPTIONING

As far as we know, no standard evaluation benchmarks have been established for video understanding 301 and captioning. To accelerate the advancement of this field, we have developed the first leaderboard 302 for video captioning. As LLM evaluation has become increasingly popular (Chiang et al., 2024), 303 we realized the lack of a standard platform to evaluate VLM's performance on video understanding. 304 We assume this is due to the difficulty of collecting ground-truth captions that accurately align with 305 videos. We will release the initial version of our captioning leaderboard upon publication. 306

5 EXPERIMENTS 308

283

284

285

287

288

289

290 291

293

299

300

307

309

5.1EXPERIMENTAL SETUP

310 Data Setup. We use four sets of data to evaluate the validity of Wolf: 1) 500 Nuscences Interactive 311 Videos; 2) 4,785 Nuscences Normal Videos; 3) 473 general videos and 4) 100 robotics videos. We 312 extract 2 frames per second for autonomous driving videos. For robotics videos, we extract 1 frame 313 per second. For short videos that sample less frames, we will increase fps to capture more details.

314 **Comparison Setup.** We use our proposed CapScore to evaluate the similarity between predicted 315 and ground truth captions. CogAgent and GPT-4V are image-level methods, so we upload sequential 316 frames into the model to obtain the output. VILA-1.5 and Gemini-Pro 1.5 are video-based methods, 317 so we directly feed a video into the model. As for the prompt for each captioning model, we 318 use "elaborate on the visual and narrative elements of the video in detail, particularly the motion 319 behavior". We compare with four state-of-the-art image-level and video-level captioning method 320 CogAgent (Hong et al., 2024), GPT-4V (Achiam et al., 2023), VILA-1.5 (Lin et al., 2023c) and 321 Gemini-Pro-1.5 (Team et al., 2023). As for CogAgent, we feed the middle frame of the video into the model to obtain the captions. As for GPT-4V, we uniformly sample 16 frames from a video and feed 322 the sequential images into the model to obtain captions. As for VILA-1.5 and Gemini-Pro-1.5, we 323 feed the video into the model to obtain the captions.

324

325

326

327

328

330

331

332

333

334

335

336

337

338

339

340

341

342 343

353

354

355 356

357

Wolf 📣: The video presents a series of scenes depicting city streets t=0[s] bustling with various vehicles such as cars, trucks, motorcycles, and a person riding a bicycle. The scenes are often divided into panels, each showcasing a different vehicle in motion, including a yellow taxi, a blue car, a red car, a white truck, and a silver car. The streets are marked with white lines and are surrounded by green trees and buildings contributing to an urban landscape setting. The video captures the typical flow of traffic through a four-way intersection, with vehicles moving in different directions. The vehicles are shown in various positions in the lanes, implying movement through the area. The sequence likely shows vehicles accelerating and decelerating as they approach and leave the intersection in response to traffic signals or the flow of other vehicles. The video is a wide-angle shot of the street intersection, split into three sections, providing a panoramic view. The left section shows the road the viewer is on, the center focuses on the intersection itself, and the right section provides a view of the crossing traffic. The lighting suggests a sunny day, with a bright spot from the sun on the left section and more even lighting on the right. The video captures the typical flow of traffic at an intersection, with vehicles approaching, some waiting their turn while others proceed t=4[s] through the intersection when it's clear. The cars driving on the right side of the road suggest the location might be a country where right-hand drive is the standard. The absence of sound enhances the focus on the visual information and allows the viewer to interpret the scene without auditory cues. Overall, the video provides a snapshot of 5[s] everyday life at a busy intersection, highlighting the complexities of navigating a traffic intersection. The wide-angle perspective and focus on movement create a sense of dvnamism.

Figure 5: Wolf example for driving videos that focus on interactive operations. Wolf captions discusses the motion behavior in details and serves as a good reference for autonomous driving.

Method	Caption Similarity ↑		Caption Quality (eg. reduced hallucination) \uparrow			
	Nuscenes	Pexels	Robotics	Nuscenes	Pexels	Robotics
CogAgent (Hong et al., 2024)	0.18	0.68	0.38	0.24	0.72	0.43
GPT-4V (Achiam et al., 2023)	0.31	0.72	0.34	0.36	0.75	0.35
VILA-1.5 (Lin et al., 2023c)	0.21	0.85	0.62	0.25	0.86	0.67
Gemini-Pro-1.5 (Team et al., 2023)	0.42	0.87	0.63	0.45	0.87	0.67
Wolf	0.55	0.88	0.72	0.56	0.89	0.75

Table 2: Comparison on 500 highly interactive (difficulty and challenging) Nuscenes videos, 473 Pexels videos and 100 robotics videos. The best and second results are highlighted with **bold** and underline. Our Wolf exhibits better performance than both open- and closed-source models.

5.2 QUALITATIVE RESULTS

To illustrate enhanced captioning ability by Wolf, we show the qualitative results in Figure 5 (please check details in Appendix). We noticed that although GPT-4V is good at recognizing the scenes, capturing temporal information in a video is not ideal. Gemini-Pro-1.5 can capture video information such as "waiting their turn while others proceed through the intersection when it's clear", but it fails to describe the detailed motions. In comparison to these two state-of-the-art approaches, we observed that Wolf not only captures the motion described in Gemini-Pro-1.5 but also successfully captures "vehicles moving in different directions" and "vehicles accelerating and decelerating as they approach and leave the intersection in response to traffic signals or the flow of other vehicles".

366 5.3 QUANTITATIVE RESULTS

We compare Wolf with various state-of-the-art 367 captioning models and display the results on 4 368 datasets in Table 2 and 3. In the default setting, 369 Wolf uses CogAgent, GPT-4V, VILA-1.5, and 370 Gemini-Pro-1.5 as Video-level models. Due to 371 the running cost, we use Wolf (based on VILA-372 1.5) on the Nuscenes Normal dataset, which only 373 uses CogAgent and VILA-1.5. We notice that 374 existing image-level models fail to capture the 375 temporal information in detail. Video-level models perform better, while Wolf can achieve the 376 best results compared to all state-of-the-art cap-377 tioning models.



Figure 6: Comparison between VILA-1.5 and finetuned VILA-1.5 with Wolf provided captions. on 500 highly interactive Nuscenes videos.

Method	Caption Similarity \uparrow	Caption Quality (eg. reduced hallucination) \uparrow
CogAgent (Hong et al., 2024)	0.27	0.30
VILA-1.5 (Lin et al., 2023c)	0.35	0.39
Wolf (based on VILA-1.5)	0.56	0.60

Table 3: Comparison on 4,785 normal Nuscenes videos. The quality of Wolf is consistently better.

Method	$ $ Caption Similarity \uparrow	Caption Quality (eg. reduced hallucination) \uparrow
CogAgent	0.18	0.24
Wolf CogAgent part (chain-of-thought)	0.26	0.32
Wolf (based on VILA-1.5)	0.35	0.37
Wolf (based on VILA-1.5+Gemini-Pro-1.5)	0.48	0.49
Wolf (based on VILA-1.5+Gemini-Pro-1.5+GPT-4V)	0.55	0.56

Table 4: Ablation study on 500 highly interactive Nuscenes videos.

391 5.4 FINETUNING VIDEO CAPTIONING MODELS

To further verify the effectiveness of Wolf, we finetune VILA-1.5 based on Wolf's captions on 4,785 normal Nuscenes videos and evaluate it on 500 highly interactive Nuscenes videos, which have much more difficult captions and complex scenarios. We follow the original VILA's training setup and launch supervised-finetuning with Wolf generated video-caption pairs for one epoch. The training is performed on 8xA100 GPUs with batch size 8. We set the learning rate to 10^{-4} with warmup strategy. No weight decay is applied.

We demonstrate the results in Figure 6, corresponding to Table 2. We observe that finetuning with Wolf boosts the model performance to 71.4% on caption similarity and 48.0% on caption quality, which outperforms GPT-4V and approaches Gemini-Pro-1.5. This suggests that Wolf captions can be easily applied to push VLMs' performance to a higher level.

402 5.5 Ablation Study on Video-level Model Selection

To further evaluate how various video-level models affect the performance, we conduct an ablation study on the components of the models in Table 4. We first compare the caption from the middle frame of CogAgent with Wolf CogAgent Caption based on the chain-of-thought approach. The chainof-thought procedure could largely improve the video understanding quality from an image-level model such as CogAgent. Then we compare Wolf with various combinations of video captions. We notice that Wolf consistensly shows better CapScore as it incorporates additional video captions.

410 5.6 COMPARISON OF FINETUNED MODELS

411

378379380381382

390

While it is difficult to directly and scalable measure the quality of captions, we compare the same model (VILA-1.5-13B) trained w/ Wolf captions and w/o Wolf captions to study the effectiveness. We benchmark the WOLF-finetuned models on two widely used video datasets ActivityNet (Caba Heilbron et al., 2015) and MSRVTT (Xu et al., 2016) and display the results in Table 5.

416 5.7 Ablation Study on Token Efficiency

It is well-known that the LLMs finetuned with RLHF favor longer response (Singhal et al., 2023), a phenomenon referred to as verbosity issue. To better assess the efficiency of the captions, we performed additional evaluation using the CapScore judge. Specifically, we separate each caption result into sentences, then incrementally use more sentences to form shortened captions, starting from only using the first sentence, to the whole original caption. These shortened captions are scored via CapScore, and we plot the score against the number of tokens used. We show the results in Figure 7.

From the result, we observe that for the better performing models (Wolf, Gemini-Pro-1.5 and GPT-4V)
the similarity scores grow with token length when caption lengths are short, but quickly plateau
or even drop as the caption lengths get too long. The caption quality score demonstrates quite
diverse patterns from different models. GPT-4V maintains a relatively consistent quality score while
Gemini-Pro-1.5 and Wolf display better quality when the caption length is short.

428			
429		ActivityNet	MSRVTT
430	VILA-1.5-13B	54.7	60.2
431	VILA-1.5-13B (fine tuned with Wolf)	55.2	60.9

Table 5: QA Accuracy Comparison of fine-tuned model on Activity and MSRVTT datasets.



Figure 7: Caption Similarity / Quality evaluated by GPT-4 under varying caption length.

6 DISCUSSION AND FUTURE WORKS

443 444

445

446 Limitations and Optimization. Wolf is still significantly more cost-effective for autolabeling and 447 captioning than procuring human labels. However, there is an efficiency concern when using an 448 ensemble method like ours. This must be handled with great care to ensure that GPU resources are 449 used effectively to mitigate any throughput degradation compared to using single models, even though 450 Wolf offers a significant improvement in caption quality. Modern GPUs are based on a massively parallel pipeline, and our goal is to saturate this pipeline with meaningful work. We consider three 451 primary areas for optimization to make Wolf a unified and efficient framework: Low-Hanging Fruit, 452 Batched Inference, and Model Quantization. For example, we reduce the size of the model weights 453 for model quantization. Several recent works (Lin et al., 2023b; Dettmers et al., 2024; Ma et al., 454 2024) have noted that LLMs and VLMs can produce highly accurate results even when their weights 455 are quantized to low bit depths. Therefore, we quantize all constituent models used in Wolf to 4 456 bits to further improve efficiency. This has two benefits. First, it reduces the bandwidth required 457 for computation. These algorithms work by packing two 4-bit numbers into a single 8-bit type, so 458 when moving data on the GPU, only half the number of bits need to be moved. Since all currently 459 released GPUs support native instructions on 8-bit floating point numbers, the two 4-bit numbers 460 are extracted and expanded by each kernel. In other words, two computations can be performed for 461 every move operation. Next-generation GPUs will natively support 4-bit data types, and we expect further efficiency improvements from having dedicated 4-bit multiply and add instructions. Second, 462 it synergizes with batched inference since the model weights, which are traditionally 16-bit, now only 463 require one quarter of the GPU memory they would ordinarily use. This allows us to fit larger batch 464 sizes on each GPU and process more videos in parallel. Please check our appendix for details. 465

466 Safety Considerations. As an ensemble of captioners, Wolf mitigates the possibility of missing out on crucial information in the captions and rectifying any hallucinations that do not agree with the 467 output of most models, which is a fundamental pillar for developing safe autonomous systems, as 468 specified in the functional safety standard ISO 26262 (ROHM). Beyond the benefits of Wolf, there 469 are still various open questions pertaining to safety of VLM captioners in deployment which we aim 470 to explore more in future: (i) We need to align the captions with the task at hand; e.g., in a driving 471 scenario, a detailed description of the foliage around the road, even if correct, is irrelevant and can 472 potentially act as distractor for the decision maker. (ii) Complementary to the first point, we need to 473 *measure* how well a caption aligns with the task at hand and develop an advanced version of CapScore. 474 (iii) Finally, we need an approach to quantify the confidence we have in the captions by leveraging 475 techniques from learning theory, such as conformal prediction (Shafer & Vovk, 2008). Most prior 476 work in this direction assumes an MCQ-styled outputs or those where a unique correct answer exists 477 (Ren et al., 2023; 2024), but these approaches do not translate to free-form text descriptions.

478 479 7 CONCLUSION

In this work, we propose Wolf, a captioning framework designed to automatically and accurately annotate any video, with significant improvements in data alignment. We find out that adopting a mixture of captioning models and summarization can largely boost the quality of the captions. This enables obtaining long, detailed, and accurate video captioning. Beyond that, we set up a leaderboard to boost the development of video captioning, which preserves a guarantee for data alignment. We will also set up a thorough library that contains different types of videos with high-quality captions, regional information such as 2D or 3D bounding boxes and depth, and multiple object motions.

486 REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo
 Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang,
 Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video
 generation models as world simulators. 2024. URL https://openai.com/research/
 video-generation-models-as-world-simulators.
- Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale
 video benchmark for human activity understanding. In *Proceedings of the ieee conference on computer vision and pattern recognition*, pp. 961–970, 2015.
- Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan,
 Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving.
 arXiv preprint arXiv:1903.11027, 2019.
- Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, et al. Pali-x: On scaling up a multilingual vision and language model. *arXiv preprint arXiv:2305.18565*, 2023a.
- Yuxiao Chen, Sander Tonkens, and Marco Pavone. Categorical traffic transformer: Interpretable and diverse behavior prediction with tokenized latent. *arXiv preprint arXiv:2311.18307*, 2023b.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao
 Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*, 2024.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Albert
 Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models
 with instruction tuning. *ArXiv*, abs/2305.06500, 2023. URL https://api.semanticscholar.org/
 CorpusID:258615266.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized
 Ilms. Advances in Neural Information Processing Systems, 36, 2024.
- Yan Ding, Xiaohan Zhang, Chris Paxton, and Shiqi Zhang. Task and motion planning with large language models for object rearrangement. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2086–2092. IEEE, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang,
 Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents, 2024.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts.
 arXiv preprint arXiv:2401.04088, 2024.
- Boyi Li, Yue Wang, Jiageng Mao, Boris Ivanovic, Sushant Veer, Karen Leung, and Marco Pavone. Driving
 everywhere with large language model policy adaptation. 2024.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for
 unified vision-language understanding and generation. In *International Conference on Machine Learning*, pp. 12888–12900. PMLR, 2022.
- 539 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.

547

552

553

554

555

558

559

560

561

562

563

564

570

571

572

576

580

581

- Bin Lin, Yang Ye, Bin Zhu, Jiaxi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-Ilava: Learning united visual representation by alignment before projection, 2023a.
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation-aware weight
 quantization for llm compression and acceleration. *arXiv preprint arXiv:2306.00978*, 2023b.
- Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad
 Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023c.
- 548 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. 2023.
- Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong, Ruiping Wang, Jilong Xue, and Furu Wei. The era of 1-bit llms: All large language models are in 1.58 bits. *arXiv preprint arXiv:2402.17764*, 2024.
 - Muhammad Maaz, Hanoona Abdul Rasheed, Salman H. Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv*, abs/2306.05424, 2023. URL https://arxiv.org/abs/2306.05424.
- Jiageng Mao, Yuxi Qian, Hang Zhao, and Yue Wang. Gpt-driver: Learning to drive with gpt. *arXiv preprint arXiv:2310.01415*, 2023a.
 - Jiageng Mao, Junjie Ye, Yuxi Qian, Marco Pavone, and Yue Wang. A language agent for autonomous driving. arXiv preprint arXiv:2311.10813, 2023b.
 - Alexander Naumann, Felix Hertlein, Daniel Grimm, Maximilian Zipfl, Steffen Thoma, Achim Rettinger, Lavdim Halilaj, Juergen Luettin, Stefan Schmid, and Holger Caesar. Lanelet2 for nuscenes: Enabling spatial semantic relationships and diverse map-based anchor paths. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 3247–3256, June 2023.
- OpenAI. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023. URL https://arxiv.org/abs/2303.08774.
- Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky,
 Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x
 models. *arXiv preprint arXiv:2310.08864*, 2023.
 - Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *arXiv: Computer Vision and Pattern Recognition*, 2021. URL https://arxiv.org/abs/2103.00020.
- Allen Z Ren, Anushri Dixit, Alexandra Bodrova, Sumeet Singh, Stephen Tu, Noah Brown, Peng Xu, Leila
 Takayama, Fei Xia, Jake Varley, et al. Robots that ask for help: Uncertainty alignment for large language
 model planners. *arXiv preprint arXiv:2307.01928*, 2023.
 - Allen Z Ren, Jaden Clark, Anushri Dixit, Masha Itkina, Anirudha Majumdar, and Dorsa Sadigh. Explore until confident: Efficient exploration for embodied question answering. *arXiv preprint arXiv:2403.15941*, 2024.
- 582
 583
 584
 ROHM. ISO 26262: Functional safety standard for modern road vehicles. URL https://fscdn.rohm. com/en/products/databook/white_paper/iso26262_wp-e.pdf.
- 585 Runway. Gen-3 alpha. https://runwayml.com/ai-tools/gen-3-alpha/, 2024. Accessed on [Insert Date].
- Paul Hongsuck Seo, Arsha Nagrani, Anurag Arnab, and Cordelia Schmid. End-to-end generative pretraining for
 multimodal video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17959–17968, 2022.
- Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9 (3), 2008.
- 593 Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating length correlations in rlhf. *arXiv preprint arXiv:2310.03716*, 2023.

594	
595	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
596	Jonan Schaikwyk, Andrew M Dai, Anja Hauin, et al. Gemini: a family of nighty capable multimodal models.
507	urxiv preprint urxiv.2512.11605, 2025.
597	Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Chenxu Hu, Yang Wang, Kun Zhan, Peng Jia, Xianpeng Lang,
598	and Hang Zhao. Drivevlm: The convergence of autonomous driving and large vision-language models. arXiv
599	preprint arXiv:2402.12289, 2024.
600	Jun Yu, Tao Mai, Ting Yao, and Yang Dui, Mar utt. A large uideo decaription detect for bridging uideo and
601	language In Proceedings of the IFFF conference on computer vision and pattern recognition pp 5288–5296
602	2016.
603	
604	Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. Pllava: Parameter-free llava
605	extension from images to videos for video dense captioning. arXiv preprint arXiv:2404.16994, 2024.
606	Antoine Yang, Arsha Nagrani, Paul Hongsuck Seo, Antoine Miech, Jordi Pont-Tuset, Ivan Laptey, Josef Sivic,
607	and Cordelia Schmid. Vid2seq: Large-scale pretraining of a visual language model for dense video captioning.
608	In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10714–10726,
609	2023.
610	Hang Zhang, Yin Li, and Lidong Ring. Video llama: An instruction tuned audio visual language model for
611	video understanding. arXiv preprint arXiv:2306.02858, 2023.
612	The substanting, where propring with 200,02000, 2020.
613	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text
61/	generation with bert. In International Conference on Learning Representations, 2019.
615	
616	
010	
017	
618	
619	
620	
621	
622	
623	
624	
625	
626	
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	
638	
630	
640	
641	
041	
042	
043	
644	
645	
646	
647	