



Wolf: Dense Video Captioning with a World Summarization Framework

Boyi Li^{1,2} Ligeng Zhu^{1,3} Ran Tian^{1,2} Shuhan Tan^{1,4}
 Yuxiao Chen¹ Yao Lu¹ Yin Cui¹ Sushant Veer¹ Max Ehrlich¹
 Jonah Philion^{1,5} Xinshuo Weng¹ Fuzhao Xue¹ Jim Fan¹ Yuke Zhu^{1,4}
 Jan Kautz¹ Andrew Tao¹ Ming-Yu Liu¹ Sanja Fidler^{1,5} Boris Ivanovic¹
 Trevor Darrell² Jitendra Malik² Song Han^{1,3} Marco Pavone^{1,6}
¹NVIDIA ²UC Berkeley ³MIT ⁴UT Austin ⁵University of Toronto ⁶Stanford University

Abstract: We propose **Wolf**, a WOrld summarization Framework for accurate video captioning. Wolf is an automated captioning framework that adopts a mixture-of-experts approach, leveraging complementary strengths of Vision Language Models (VLMs). By combining 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 (VILA-1.5, 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.

Keywords: Dataset Curation, Video Captioning

1 Introduction

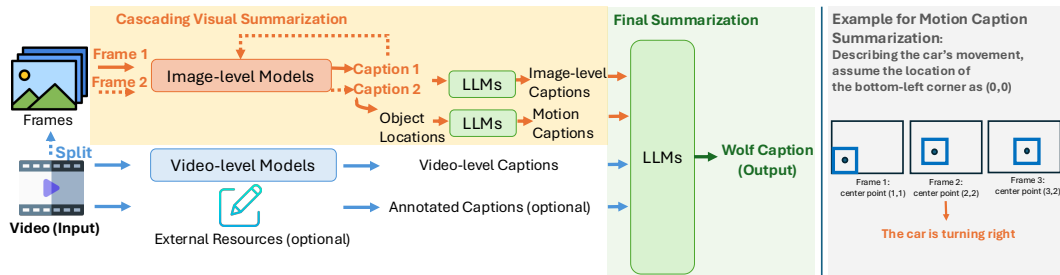


Figure 1: Overview of proposed Wolf framework. Wolf utilizes both image-level and video-level models to generate diverse and detailed captions, which are then summarized for cross-checking. On the right side, we also provide an example of how we obtain motion captions based on object locations extracted from image captions.

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 for tasks like video generation, such as Sora [1], Runaway [2] and Wan2.1 [3]. 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 internet can be faulty and misaligned and human annotation is prohibitively expensive for large datasets. *Secondly*, video captioning is inherently more challenging than image captioning due to the additional complexity of

temporal correlation and camera motion. Existing captioning models [4, 5] struggle with temporal reasoning and fail to achieve accurate scene understanding. *Thirdly*, there is no established benchmark to measure captioning progress. Existing video QA benchmarks [6] are often limited to short answers, making it difficult to measure hallucinations in detailed long captions. *Fourthly*, the correctness and completeness of the captions are crucial for safety-critical tasks. In the era of large language models (LLMs), text descriptions of scenarios used by embodied agents for planning and control become increasingly common [7, 8, 9, 10]. 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 vegetable-chopping manipulator can lead to an injury.

To handle these challenges, we introduce **WOrLd summarization Framework (Wolf)**, a novel summarization 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 [11], 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 [12], CogAgent [4]) to commercial solutions (such as Gemini-Pro-1.5 [13], GPT-4V [14]). 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 and is more effective than several widely used captioning metrics. The results show that our method improves CapScore by a large margin.
2. We introduce four 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 benchmark 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 leading performance in image captioning tasks, largely due to the success of LLMs. CLIP [15] pioneered this field by training a shared feature space for vision and language modalities on image-caption pairs. Building on CLIP, BLIP [16] and BLIP-2 [17] improved performance by aligning the pre-trained encoder with LLMs. Following the direction, LLaVA [18] and InstructBLIP [19] demonstrated that jointly training on diverse datasets as an instruction-following task leads to strong generalization across various tasks. VILA [12] highlighted the importance of pre-training with diverse data, and therefore significantly scaled up the pre-training dataset. Kosmos-2 [20] and PaLI-X [21] further introduced pseudo-labeling bounding boxes from open-vocabulary object detectors to scale up the size of pre-training dataset.

Video Captioning. As image-based VLMs are not trained with video data, they are limited in describing details present in the video data [22, 23, 24]. To improve video captioning, PLLaVa [25] builds on top of LLaVa and introduced a parameter-free pooling strategy to enhance the caption quality. Video-LLaVA [26] achieves state-of-the-art performance on several benchmarks by conducting joint training on images and videos, thereby learning a unified visual representation. Video-LLaMA [5] incorporates both video and audio into LLMs by introducing two Q-formers to extract features. Vid2seq [27] conducts large-scale pre-training with narrated videos for dense video captioning.

¹We also provide ethical statement and reproducibility in Appendix.

Meanwhile, MV-GPT [28] employs an automated speech recognition (ASR) model to provide additional labeling for the videos.

LLM-based Summarization. Recently many works have found that it is efficient to summarize useful information using LLMs. For example, LLaDA [9] can provide users with helpful instructions based on the user request and corresponding traffic rules in the desired location. OpenAI team finds re-captioning [29] via LLMs can be very helpful.

3 Wolf Framework

We propose Wolf, which is an automated dense 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 [4], GPT-4V [7] to generating image-level captions, and use VILA-1.5-7B [12], Gemini-Pro-1.5 [13] to generate video captions.

Cascading Visual Summarization. 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 **cascading visual summarizing program** to obtain video captions from image-level models. As illustrated in Figure 1, we first split the video into sequential images, sampling two key-frames every second. We **start by** feeding Image 1 into the Image-level Model to obtain Caption 1, where we require the model to generate detailed scene-level information and object locations. Given the temporal correlation between key frames in a video, we then feed both Caption 1 and Image 2 into the model to generate Caption 2. By repeating this procedure, we generate captions for all sampled frames. Finally, we use GPT-4 to summarize the information from all captions with the prompt “*Summarize all the captions to describe the video with accurate temporal information*”. We also extract the bounding box locations for each object in each frame, then feed them into LLMs to **summarize** the trajectory of the moving object. For example, in a driving video, a blue car is driving into the right lane, and the centers of the bounding boxes are (0,0), (1,1), (1,2). We provide the car’s location to the LLM, and it outputs ‘the blue car is driving to the right,’ which we refer to as a ‘*Motion Caption*’.

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

4 Benchmarking Video Captioning

To showcase the effectiveness of Wolf, we constructed four distinct datasets (please check the examples in Figure 2). These include two autonomous driving video captioning datasets based on the open-sourced NuScenes [30] dataset (Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International Public License), a general daily video captioning dataset from Pexels ², and a robot manipulation video captioning dataset from an open-source robot learning dataset [31]. These benchmark datasets are tailored to assess the caption model’s scene comprehension and its behavior understanding capabilities, both of which are vital for auto-labeling in embodied AI tasks. All captions were generated using a combination of ground truth information, rule-based heuristics, human labeling, and rewriting.

²<https://www.pexels.com/>



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. Our Wolf dataset is fully manually annotated to ensure a robust evaluation for the community. We present our dataset’s statistics in Table 1.

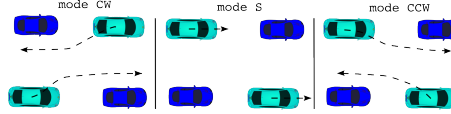


Figure 3: Illustration of homotopy types of different relative motions between a pair of vehicles.

4.1 Wolf Dataset Curation

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 understanding ability of a captioning model. Our dataset contains 500 intensely interactive video-caption pairs (≈ 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 (≈ 6 hours). Our caption generation process consists of three steps: i) agent-level motion annotation, ii) ego-centric interaction annotation, and iii) information aggregation via LLM.

Step 1: agent-level motion annotation. The NuScenes dataset provides full annotations of traffic elements in each scene, including 3D bounding boxes, element categories, and semantic map information. Similar to DriveVLM [32], we utilize this ground truth data along with lane topology information [33] to generate text descriptions of both speed and angular motion characteristics for the ego vehicle and other traffic participants within a video clip. Specifically, we classify agent actions into 11 categories, including Stopping, Accelerating, Decelerating, Lane Changes, Turns, and more, based on their observed movements and behaviors.

Step 2: egocentric interaction annotation. Beyond each agent’s dynamics information, we also aim to capture the ego vehicle’s interactions with other traffic participants (e.g., crossing pedestrians,

blocking traffic cones) depicted in the video clip. To efficiently describe interactions, we use two categorical modes: the lane relationship (*agent-ego lane mode*) and relative motion (*homotopy*) between a traffic participant and the ego vehicle [34]. At each time step t , the agent-ego lane mode encodes the topological relationship between the ego vehicle’s current lane and the traffic agent’s lane. The categories include *LEFT*, *RIGHT*, *AHEAD*, *BEHIND*, and *NOTON*, where *NOTON* indicates that the traffic agent is on a lane that cannot directly reach the ego vehicle’s lane. To compute the agent-ego lane mode, we follow [34] by identifying each agent’s lane and using a lane topology map for annotation. Homotopy describes the relative motion between agents in a video and is categorized as: [*S*, *CW*, *CCW*] (*static*, *clockwise*, *counterclockwise*), as shown in Figure 3.

Step 3: information aggregation. By combining agent-ego lane mode, homotopy, traffic agents’ ground truth dynamics, and scene context (e.g., the ego vehicle is near an intersection), we can apply heuristics to annotate interaction descriptions. For example, in a video clip, a static object’s agent-ego lane mode changes from *AHEAD*, to *LEFT*, to *BEHIND*, and the ego vehicle’s first performs *RIGHT-LANE-CHANGE*, *KEEP-LANE*, then *LEFT-LANE-CHANGE*, indicating the ego vehicle overtakes that object from the ego vehicle’s left side. We identified six interaction categories 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 lane dividers; 5) overtake traffic agent via lane-change; 6) other non-intensive interactions. With both agent-level motion annotations and ego-centric interaction annotations, we employ an LLM to aggregate this information and generate a human-like scene description. While any off-the-shelf LLM could be used for this task, we opted for the GPT-3.5 model. Additionally, we experimented with the llama 3 model and observed similar performance.

Task Type	Source	Size	Annotation Type
Normal Driving Scenes	Nuscenes	4,785	Manually
Challenging Driving Scenes	Nuscenes	500	Manually
General Daily Scenes	Pexels	473	Manually
Robot Manipulation	UCB	100	Manually

Table 1: Statistics of the Wolf dataset.

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. [31] 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.

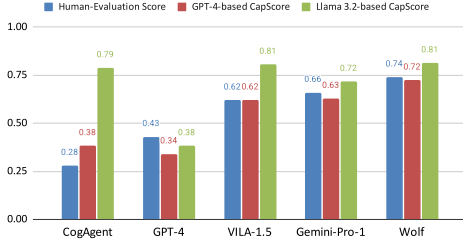
4.1.3 Pexels Dataset

To evaluate caption models in general daily environments, we further collect high quality (360p to 1080p) videos from Pexels. 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.

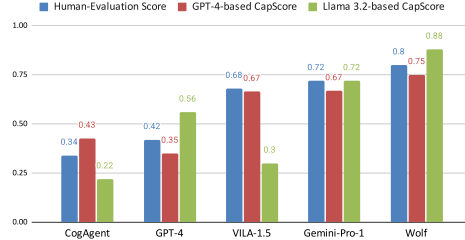
4.2 Wolf Evaluation Metric

4.2.1 CapScore: Evaluating Captions with LLMs

Video captioning has been an ill-posed problem since there is no metric to evaluate the quality of captions and the alignment between the video and the caption. Inspired by BERTScore [35], CLIPScore [36] and the stability of LLMs on evaluation [37, 38, 39], we introduce **CapScore** (Captioning Score), a quantitative metric to use LLMs to evaluate the similarity between predicted and human-annotated (ground truth) captions. We tried both GPT-4 (model=“gpt-4”) and Llama



(a) Comparison on Caption Similarity.



(b) Comparison on Caption Quality.

Figure 4: Comparisons on Human-Evaluation Score and Llama 3.2-based CapScore and GPT4-based CapScore (proposed).

3.2 [40] as our LLM to summarize the captions. We noticed that GPT-4 can always obtain stable results over 3 runs. However, for Llama 3.2, the results varied over different runs. We tried to lower the temperature (from 0.9 to 0.5) to make the inference stable, however, we noticed that the scores are not consistent with human evaluation. Therefore we select GPT-4 as our LLM to conduct the experiments. Assume we have 6 captions, we feed all the captions into GPT-4 and add the prompt “Can you give a score (two decimal places) from 0 to 1 for captions 1, 2, 3, 4 and 5, indicating which one is closer to the ground truth caption (metric 1) and which contains fewer hallucinations and less misalignment (metric 2)? Please output only the scores of each metric separated only by a semicolon. For each metric, please output only the scores of captions 1, 2, 3, 4 and 5 separated by commas, in order—no text in the output. ”. We ask GPT-4 caption similarity and caption quality scores.

We set the range [0,1] to align with several widely used NLP metrics, such as BLEU [41], ROUGE [42], and BERTScore [35]. To address the potential concern, we followed the same settings as Table 1 and used the range [0,5] to calculate CapScore. The trend remains precisely the same, with Wolf achieving scores of 3.61 for similarity and 3.70 for quality - almost five times the values shown in Table 1, demonstrating CapScore’s stability and robustness regardless of the range.

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

Through our experiments, we find that GPT-4 is very robust for calculating the scores. We have run the experiments for 1-3 times, the results appear to be stable and less than 0.05 changes. To alleviate 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 Figure 4. Beyond that, we also conduct experiments comparing CapScore with other widely used image captioning evaluation metrics, as is shown in Appendix (Sec A.5). We observe that CapScore aligns with trends observed in other metrics but highlights a larger performance gap between models as a more effective evaluation metric.

Method	Caption Similarity \uparrow			Caption Quality (eg. reduced hallucination) \uparrow		
	Nuscenes	Pexels	Robotics	Nuscenes	Pexels	Robotics
CogAgent [4]	0.18	0.68	0.38	0.24	0.72	0.43
GPT-4V [44]	0.31	0.72	0.34	0.36	0.75	0.35
VILA-1.5-7B [12]	0.21	0.85	0.62	0.25	0.86	0.67
Gemini-Pro-1.5 [13]	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. Our Wolf exhibits better performance than both open- and closed-source models.

4.2.3 Benchmarking Video Captioning

To our best knowledge, no standard evaluation benchmarks have been established for video understanding and captioning. To accelerate the advancement of this field, we have developed the first leaderboard for video captioning. As LLM evaluation has become increasingly popular [43], we realized the lack of a standard platform to evaluate VLM’s performance on video understanding. We assume this is due to the difficulty of collecting ground truth captions that accurately align with videos. We will release the initial version of our captioning leaderboard upon publication.

5 Experiments

5.1 Experimental Setup

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

Comparison Setup. We use our proposed CapScore to evaluate the similarity between predicted and ground truth captions. CogAgent and GPT-4V are image-level methods, so we upload sequential frames into the model to obtain the output. VILA-1.5-7B and Gemini-Pro 1.5 are video-based, so we directly feed a video into the model. As for the prompt for each captioning model, we use “*elaborate on the visual and narrative elements of the video in detail, particularly the motion behavior*”. We compare with four state-of-the-art image-level and video-level captioning Vision-Language Models (VLMs) CogAgent [4], GPT-4V [44], VILA-1.5 [12] and Gemini-Pro-1.5 [13]. As for CogAgent, we feed the middle frame of the video into the model to obtain captions. As for GPT-4V, we uniformly sample 16 frames from a video and feed the sequential images into the model to obtain captions. As for VILA-1.5-7B and Gemini-Pro-1.5, we feed the video into the model to obtain the captions.

5.2 Qualitative and Quantitative Results

To illustrate enhanced captioning ability by Wolf, we show the qualitative results in Section C of the Appendix. We compare Wolf with various state-of-the-art captioning models and display the results on 4 datasets in 2 and Table 2 of the Appendix. In the default setting, Wolf uses CogAgent, GPT-4V, VILA-1.5-7B, and Gemini-Pro-1.5 as Video-level models. Due to the running cost, we use Wolf (based on VILA-1.5) on the Nuscenes Normal dataset, which only uses CogAgent and VILA-1.5-7B. We notice that existing image-level models fail to capture the temporal information in detail. Video-level models perform better, while Wolf can achieve the best results compared to all state-of-the-art captioning models.

5.3 Finetuning VLMs with Wolf Captions

5.3.1 Comparison on Wolf Dataset

To further verify the effectiveness of Wolf, we finetune VILA-1.5-7B based on Wolf’s captions on 4,785 normal Nuscenes videos and evaluate it on 500 highly interactive Nuscenes videos, which

VILA-1.5-7B	Caption Similarity \uparrow	Caption Quality \uparrow
Default	0.21	0.25
Fine-tuned with Wolf annotation	0.36	0.37

Table 3: Comparison on 500 highly interactive Nuscenes videos.

VILA-1.5-13B	ActivityNet	MSRVTT
Default	54.7	60.2
Fine-tuned with Wolf annotation	55.2	60.9

Table 4: QA Accuracy comparison of the fine-Tuned Model on Activity and MSRVTT datasets.

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 Table 3, 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.

5.3.2 Comparison on Other Benchmark Datasets

To scalable measure the quality of captions, we compare the 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 [45] and MSRVTT [46] and display the results in Table 4, the improved performance effectively demonstrates the efficiency of Wolf.

Method	Caption Similarity \uparrow	Caption Quality \uparrow
CogAgent	0.18	0.24
Wolf CogAgent part (Cascading Visual Summarization)	0.26	0.32
Wolf video part (VILA-1.5-7B+Gemini-Pro-1.5+GPT-4V)	0.40	0.42
Wolf (based on VILA-1.5-7B)	0.35	0.37
Wolf (based on VILA-1.5-7B+Gemini-Pro-1.5)	0.48	0.49
Wolf (based on VILA-1.5-7B+Gemini-Pro-1.5+GPT-4V)	0.55	0.56

Table 5: Ablation study on 500 highly interactive Nuscenes videos. Note: The first row shows the results using *only image-level models*, the second row shows the results using *only video-level models*, and the last row used *both image-level models (CogAgent part) and various video-level models*.

5.4 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 5. We first compare the caption from the middle frame of CogAgent with Wolf Caption based on the visual cascading summarization approach (only using CogAgent). The visual cascading summarization procedure could largely improve the video understanding quality from an image-level model such as CogAgent. Then, we conduct an ablation using only the video-level models. Finally, we compare Wolf with various combinations of video captions. We notice that Wolf consistently shows better CapScore as its dense framework reduces hallucination and incorporate video details from different models.

Additionally, we include an ablation study on token efficiency, please see Section D of the Appendix.

6 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. We will also establish a comprehensive library that includes various types of videos with high-quality captions, regional information such as 2D and 3D bounding boxes and depth, as well as multiple object motions and interactions. For discussion and future works, please check Section E of the Appendix for details.

References

- [1] T. Brooks, B. Peebles, C. Holmes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman, E. Luhman, C. Ng, R. Wang, and A. Ramesh. Video generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>.
- [2] Runway. Gen-3 alpha. <https://runwayml.com/ai-tools/gen-3-alpha/>, 2024. Accessed on [Insert Date].
- [3] W. Team. Wan: Open and advanced large-scale video generative models. 2025.
- [4] W. Hong, W. Wang, Q. Lv, J. Xu, W. Yu, J. Ji, Y. Wang, Z. Wang, Y. Dong, M. Ding, and J. Tang. Cogagent: A visual language model for gui agents, 2024.
- [5] H. Zhang, X. Li, and L. Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.
- [6] M. Maaz, H. A. Rasheed, S. H. Khan, and F. S. 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>.
- [7] J. Mao, Y. Qian, H. Zhao, and Y. Wang. Gpt-driver: Learning to drive with gpt. *arXiv preprint arXiv:2310.01415*, 2023.
- [8] J. Mao, J. Ye, Y. Qian, M. Pavone, and Y. Wang. A language agent for autonomous driving. *arXiv preprint arXiv:2311.10813*, 2023.
- [9] B. Li, Y. Wang, J. Mao, B. Ivanovic, S. Veer, K. Leung, and M. Pavone. Driving everywhere with large language model policy adaptation. 2024.
- [10] Y. Ding, X. Zhang, C. Paxton, and S. Zhang. Task and motion planning with large language models for object rearrangement. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2086–2092. IEEE, 2023.
- [11] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l. Casas, E. B. Hanna, F. Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- [12] J. Lin, H. Yin, W. Ping, Y. Lu, P. Molchanov, A. Tao, H. Mao, J. Kautz, M. Shoenybi, and S. Han. Vila: On pre-training for visual language models, 2023.
- [13] G. Team, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [14] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. URL <https://arxiv.org/abs/2303.08774>.
- [15] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. *arXiv: Computer Vision and Pattern Recognition*, 2021. URL <https://arxiv.org/abs/2103.00020>.
- [16] J. Li, D. Li, C. Xiong, and S. Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022.
- [17] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.

- [18] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. 2023.
- [19] W. Dai, J. Li, D. Li, A. M. H. Tiong, J. Zhao, W. Wang, B. A. Li, P. Fung, and S. 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>.
- [20] Z. Peng, W. Wang, L. Dong, Y. Hao, S. Huang, S. Ma, and F. Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.
- [21] X. Chen, J. Djolonga, P. Padlewski, B. Mustafa, S. Changpinyo, J. Wu, C. R. Ruiz, S. Goodman, X. Wang, Y. Tay, et al. Pali-x: On scaling up a multilingual vision and language model. *arXiv preprint arXiv:2305.18565*, 2023.
- [22] X. Zhou, A. Arnab, S. Buch, S. Yan, A. Myers, X. Xiong, A. Nagrani, and C. Schmid. Streaming dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18243–18252, 2024.
- [23] M. Kim, H. B. Kim, J. Moon, J. Choi, and S. T. Kim. Do you remember? dense video captioning with cross-modal memory retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13894–13904, 2024.
- [24] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, and J. Carlos Niebles. Dense-captioning events in videos. In *Proceedings of the IEEE international conference on computer vision*, pages 706–715, 2017.
- [25] L. Xu, Y. Zhao, D. Zhou, Z. Lin, S. K. Ng, and J. Feng. Pllava: Parameter-free llava extension from images to videos for video dense captioning. *arXiv preprint arXiv:2404.16994*, 2024.
- [26] B. Lin, Y. Ye, B. Zhu, J. Cui, M. Ning, P. Jin, and L. Yuan. Video-llava: Learning united visual representation by alignment before projection, 2023.
- [27] A. Yang, A. Nagrani, P. H. Seo, A. Miech, J. Pont-Tuset, I. Laptev, J. Sivic, and C. Schmid. Vid2seq: Large-scale pretraining of a visual language model for dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10714–10726, 2023.
- [28] P. H. Seo, A. Nagrani, A. Arnab, and C. Schmid. End-to-end generative pretraining for multimodal video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17959–17968, 2022.
- [29] J. Betker, G. Goh, L. Jing, T. Brooks, J. Wang, L. Li, L. Ouyang, J. Zhuang, J. Lee, Y. Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- [30] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom. nuscenes: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019.
- [31] A. Padalkar, A. Pooley, A. Jain, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Singh, A. Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.
- [32] X. Tian, J. Gu, B. Li, Y. Liu, C. Hu, Y. Wang, K. Zhan, P. Jia, X. Lang, and H. Zhao. Drivevlm: The convergence of autonomous driving and large vision-language models. *arXiv preprint arXiv:2402.12289*, 2024.
- [33] A. Naumann, F. Hertlein, D. Grimm, M. Zipfl, S. Thoma, A. Rettinger, L. Halilaj, J. Luetttin, S. Schmid, and H. 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*, pages 3247–3256, June 2023.

- [34] Y. Chen, S. Tonkens, and M. Pavone. Categorical traffic transformer: Interpretable and diverse behavior prediction with tokenized latent. *arXiv preprint arXiv:2311.18307*, 2023.
- [35] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2019.
- [36] J. Hessel, A. Holtzman, M. Forbes, R. L. Bras, and Y. Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- [37] D. Chan, S. Petryk, J. E. Gonzalez, T. Darrell, and J. Canny. Clair: Evaluating image captions with large language models. *arXiv preprint arXiv:2310.12971*, 2023.
- [38] Z. Lin, D. Pathak, B. Li, J. Li, X. Xia, G. Neubig, P. Zhang, and D. Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *European Conference on Computer Vision*, pages 366–384. Springer, 2025.
- [39] Z. Lin, X. Chen, D. Pathak, P. Zhang, and D. Ramanan. Revisiting the role of language priors in vision-language models. *arXiv preprint arXiv:2306.01879*, 2023.
- [40] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [41] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- [42] C.-Y. Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [43] W.-L. Chiang, L. Zheng, Y. Sheng, A. N. Angelopoulos, T. Li, D. Li, H. Zhang, B. Zhu, M. Jordan, J. E. Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*, 2024.
- [44] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [45] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015.
- [46] J. Xu, T. Mei, T. Yao, and Y. Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016.