WOLF: ACCURATE VIDEO CAPTIONING 000 001 WITH A WORLD SUMMARIZATION FRAMEWORK 002 003 (APPENDIX) 004

Anonymous authors

Paper under double-blind review

800 009

006

010

013

011 012

CONTRIBUTIONS A

014 We would like to list Wolf Contributions:

015 1) Framework and Evaluation Metric. We designed a novel world summarization framework, 016 Wolf, for video captioning and introduced an LLM-based metric, CapScore, to evaluate the quality 017 of captions. The results show that our method significantly improves CapScore. 018

019 2) Datasets and Benchmark. We introduce the Wolf benchmark (leaderboard) and four humanannotated benchmark datasets. These datasets include autonomous driving, general scenes from 020 Pexels, robotics videos, and human-annotated captions, collectively referred to as the Wolf Dataset. 021

022 3) Intended Uses. We believe Wolf can serve as one of the best practices (auto-labeling tool) for creating and curating paired datasets and benchmarks.

4) Hosting, licensing, and maintenance plan. The code, data, and leaderboard will be open-sourced 025 and maintained. Continuous efforts will be made to refine the Wolf Dataset, Wolf codebase, and 026 CapScore. We hope that Wolf will raise awareness about the quality of video captioning, set a 027 standard for the field, and boost community development. 028

В PEXEL DATASET CATEGORIES

We categorize videos from pexel into the following types: Travel & Events, Sports, Education, Pets & Animals, People & Blogs, Nonprofits & Activism, News & Politics, Music, Science & Technology, 033 Comedy, Entertainment, Film & Animation, Gaming, Robotics, How to Styles. 034

- 037

041

029

031

REPRODUCIBILITY CHECKLIST С

038 Datasets Curation. We have provided dataset curation details in "Sec 4.1 Benchmark Dataset 039 Curation" of the main paper. The main contribution of the Wolf dataset is its human-annotated 040 captions. Wolf is built upon existing data from Nuscenes, videos from the Pexels website (publicly available, downloaded, and used for free), and robot manipulation videos collected from the Open 042 X-Embodiment dataset (Padalkar et al., 2023) (under a CC License). Since all the videos have already 043 been or will be released publicly, we can directly follow the instructions to download and use them. 044 We will continually increase the size and update the data information in these folders.

Evaluation Procedures. We have provided evaluation procedures in "Sec 4.2.1 CapScore: Evaluating 046 Captions with LLMs" of the main paper. Since CapScore compares captions from different methods 047 simultaneously, it can guarantee the quality of the evaluation. In our paper, we ran all the experiments 048 three times and took the average score in the reported tables and figures of the paper.

- 049
- ETHICS STATEMENT D
- 051 052
- We commit to the ICLR Code of Ethics and affirm that our work follows the rules for experimentation. We welcome any related discussions and feedback.

E QUALITATIVE COMPARISON ON INTERACTIVE NUSCENES DRIVING VIDEOS

We display the details of Figure 5 of the paper (Wolf example for driving videos that focus on interactive operations) in Figure 1.

058 059

054

055 056

057

060 061

062

063

064

065

F WOLF EFFICIENCY OPTIMIZATION

We consider three primary areas: **Low-Hanging Fruit**, **Batched Inference**, and **Model Quantization** as optimizations which make Wolf a unified and efficient framework. Using the optimizations detailed in this section we were able to increase the speed of CogVLM by a factor of approximately 10x (450s/video to 41s/video), VILA throughput was similarly improved to only about 3s per video.

Low-Hanging Fruit. These are primarily systems concerns and work arounds for simplistically 066 written APIs. For example, the off-the-shelf CogVLM (Hong et al., 2024) and VILA (Lin et al., 2023b) 067 supporting code is heavily based on loading PIL images to present to a huggingface pipeline (Wolf 068 et al., 2019). In the naive pipeline, videos would need to be decoded and then converted to PIL images 069 before input to the respective pipelines, which in turn convert them to GPU PyTorch (Ansel et al., 070 2024) tensors. This is extremely inefficient. Instead, we can leverage the hardware video decoder 071 present in modern GPUs to decode the videos directly to GPU tensors and rewrite the preprocessing 072 pipelines to operate on these tensors directly. This has the additional benefit of shifting preprocessing 073 transform work from CPU to GPU. 074

Batched Inference. Simplifying Wolf into the simplest terms, we are essentially performing repeated 075 neural network inference. Surprisingly, most VLM supporting code is designed to run inference on 076 only a single example at a time. However, just as in other deep-learning problems, there fundamentally 077 no reason why we cannot processes multiple videos at a single time in batches. This step is crucial to 078 maximizing the use of GPU resources. Processing a single example may only use as little as 25% of a modern datacenter GPU which would either increase the time to process a dataset or the number of 080 GPUs required to complete a task in a fixed time budget. We can reimplement more of the supporting 081 code to enable processing batches of as many videos as will fit in GPU memory at a single time yielding a linear speedup in processing. For example, if we can fit batches of 4 in GPU memory we 082 observe a speedup of 4x over processing single examples. 083

084 **Model Quantization.** The final optimization we consider is to reduce the size of the model weights. 085 Several recent works (Lin et al., 2023a; Dettmers et al., 2024; Ma et al., 2024) have noted that 086 LLMs and VLMs can produce highly accurate results even when their weights are quantized to low 087 bit-depths. Therefore, we quantize all constituent models used in Wolf to 4-bits to further improve 088 efficiency. This has two benefits. First, it reduces the bandwidth required for computation. These 089 algorithms work by packing two 4-bit numbers into a single 8-bit type, so when moving data on the GPU only half the number of bits need to be moved. Since all currently released GPUs support native 090 instructions on 8-bit floating point numbers, the two 4-bit numbers are extracted and expanded by 091 each kernel. In other words, two computations can be performed for every move operation. Next 092 generation GPUs will natively support 4-bit datatypes and we expect further efficiency improvements from having dedicated 4-bit multiply and add instructions. Next, it synergizes with batched inference 094 since the model weights, which are traditionally 16-bit, now only require one quarter of the GPU 095 memory they would ordinarily use. This allows us to fit larger batch sizes on each GPU and process 096 more videos in parallel.

097 098

099 100

101

102 103

104

G UPDATED RESULTS AND DOCUMENTATION

We will release our webpage and the initial version of our captioning leaderboard upon publication. We will regularly update Wolf results and documentation on our webpage.

References

Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter
Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Anjali Chourdia, Will Constable, Alban Desmaison,
Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind, Brian Hirsh, Sherlock Huang,
Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, Mario Lezcano, Yanbo Liang, Jason Liang, Yinghai



Figure 1: Comparison of CogAgent, VILA-1.5, GPT-4, Gemini-Pro-1.5, and Wolf on Interactive Nuscenes Driving Videos, Corresponding to Figure 4 of the Paper.

162	Lu CK Luk Bert Maher Vunije Pan Christian Puhrsch Matthias Resa Mark Saroufim Marcos Vukio
163	Siraichi, Helen Suk, Michael Suo, Phil Tillet, Eikan Wang, Xiaodong Wang, William Wen, Shunting
164	Zhang, Xu Zhao, Keren Zhou, Richard Zou, Ajit Mathews, Gregory Chanan, Peng Wu, and Soumith
165	Chintala. PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph
166	Compilation. In 29th ACM International Conference on Architectural Support for Programming Languages
167	and Operating Systems, Volume 2 (ASPLOS '24). ACM, April 2024. doi: 10.1145/3620665.3640366. URL
168	https://pytorch.org/assets/pytorch2-2.pdf.
169	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized
170	Ilms. Advances in Neural Information Processing Systems, 36, 2024.
171	Wenyi Hong, Weihan Wang, Oingsong Ly, Jiazheng Yu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang
172	Yuxiao Dong Ming Ding and Jie Tang Cogagent: A visual language model for gui agents 2024
173	Tukluo Dong, hing Ding, and the Tung. Cogagent. IT tisual language model for gal agents, 2021.
174	Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation-aware weight
175	quantization for llm compression and acceleration. arXiv preprint arXiv:2306.00978, 2023a.
176	Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad
177	Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023b.
178	Shuming Ma Hangun Wang Lingvigo Ma Lai Wang Wankui Wang Charley Hunge Li Dage Dai i W
179	Shuming Ma, nongyu wang, Lingxiao Ma, Lei wang, Wennui Wang, Shaonan Huang, Li Dong, Ruiping Wang, Jilong Xue, and Furu Wei. The era of 1-hit llms: All large language models are in 1.58 hits. arXiv preprint
180	arXiv:2402.17764, 2024.
181	
182	Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky,
183	Anant Kai, Anikait Singn, Anthony Bronan, et al. Open x-embodiment: Robotic learning datasets and rt-x
184	models. <i>urxiv preprint urxiv.2510.00004</i> , 2025.
185	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric
186	Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transformers: State-of-the-art natural
187	language processing. arXiv preprint arXiv:1910.03//1, 2019.
188	
189	
100	
101	
102	
102	
10/	
105	
195	
107	
102	
100	
200	
200	
201	
202	
203	
205	
205	
200	
208	
200	
205	
210	
211	
212	
213	
214	
210	