

Holmes-VAU: Towards Long-term Video Anomaly Understanding at Any Granularity

Supplementary Material

A. Details of the Data Engine.

To construct a dataset with hierarchical annotations with both short-term and long-term anomalies, we developed a semi-automated annotation engine that combines manual efforts with the generative capabilities of LLM. In the main paper, we present the complete annotation workflow. Below, we provide additional details about the data engine.

A.1. Hierarchical Video Decoupling

Before annotation, we collected videos from the training sets of the UCF-Crime [5] and XD-Violence [7] datasets. From UCF-Crime, we selected 758 normal videos and 735 anomaly videos, while from XD-Violence, we selected 1,904 normal videos and 2,046 anomaly videos. The anomaly videos included their original video-level labels, e.g., *Abuse*, *Explosion*. For the anomaly videos, we organized a team of five annotators to label each anomaly event within the videos. The annotation process took approximately 20 hours to complete. For the normal videos, we considered all segments to be normal and randomly cropped segments of varying lengths to serve as normal event-level video segments. These anomaly and normal event-level videos were further divided into shorter clip-level segments. For UCF-Crime, we adopted the clip-level divisions from UCA [9]. For XD-Violence, we performed uniform division.

A.2. Hierarchical Free-text Annotation

Clip Captioning. For videos in UCF-Crime, we fully utilized the manually annotated captions from UCA [9]. For videos in XD-Violence, we used LLaVA-Next-Video-7B [11] as our captioner to generate textual descriptions for clip-level videos. The specific prompt is as follows:

'Please provide a short and brief description of the video clip, focusing on the main subjects and their actions.'

Event Summary. We combined all captions and video-level category labels to generate anomaly-related summaries for each event using an LLM. We selected LLaMA3-70B [1] as our LLM due to its strong summarization capabilities. The specific prompt is as follows:

'The dense caption of the video is: {clip captions}. There are (is no) abnormal events ({video-level label}) in the video. Your response should include the following three parts: 1. Whether the anomaly exists and the specific name of the anomaly. 2. A summary of the anomaly events. 3.

Brief explanation of the basis for judging the anomaly.'

Video Summary. Similar to generating event summaries, we generated video-level summaries by analyzing the event-level summaries. The specific prompt is as follows:

'Below is a summary of all the events in the video: {event summaries}. There are (is no) abnormal events ({video-level label}) in the video. Your response should include the following three parts: 1. Whether the anomaly exists and the specific name of the anomaly. 2. Detailed description of the video anomaly event from start to end. 3. Brief analysis of the basis for judging the anomaly.'

Annotation Format. In Fig.A, we present an example of the hierarchical free-text annotations for a video.

A.3. Hierarchical Instruction Data Construction

To construct the instruction dataset, we designed question prompts tailored to different tasks, including **Caption**, **Judgment**, **Description**, and **Analysis**. For each instruction item, we randomly selected one prompt from the pool and matched it with the corresponding content from the free-text annotations as the answer.

Caption.

1. "Describe the video briefly."
2. "Describe the main events that take place in this video."
3. "Give a short description of the video."
4. "What happened in this video?"
5. "Generate a brief caption for the video."
6. "Can you provide a brief description of the video?"
7. "Briefly describe the main subjects and their actions in the video."
8. "Provide a short overview of what happens in the video?"
9. "Describe the key moments that showcase the subjects' activities in the video."
10. "Describe the sequence of events involving the main subjects in the video."
11. "What activities happen throughout the video?"
12. "Describe the main subjects and their roles in the video."
13. "What key moments stand out in the video?"
14. "What are the primary activities showcased in the video?"
15. "What happens to the main subjects as the video progresses?"
16. "What is a brief overview of what happens in the video?"
17. "Describe the main subjects and their contributions to the video."
18. "Describe the key events in the video."
19. "Describe the video's main activities."
20. "Can you describe the main action in this video briefly?"
21. "Describe the video clip concisely."
22. "Provide a brief description of the given video clip."
23. "Summarize the visual content of the video clip."
24. "Give a short and clear explanation of the subsequent video clip."

Judgement.

1. "What types of anomalies are shown in the video clip?"
2. "Are there any anomaly events detected in the video?"
3. "Detect and classify the anomaly events in the video."
4. "Identify any abnormal behaviors depicted in the video."
5. "Determine whether there are anomaly events in the video and the specific name of the anomaly."
6. "What anomalies can be identified in the video?"
7. "What categories of anomalies can be found in the video?"
8. "Could you point out any abnormal actions in the video?"
9. "Point out the abnormal actions in the video."



```

1 {
2   "video": "v=2rfyeR-YaJw__#l_label_G-0-0",
3   "n_frames": 1940,
4   "fps": 24.0,
5   "label": ["Explosion"],
6   "clips": [[[5.583, 11.903], [11.903, 18.222], [18.222, 24.542]],
7             [[36.167, 43.48], [43.48, 50.792]]],
8   "clip_captions": [
9     [
10      "A military tank moving across a barren landscape with low-rise buildings and sparse vegetation. the sky
11      ↪ is overcast, and the overall color palette is muted with earthy tones.",
12      "A series of images depicting a barren landscape with a few buildings in the background. the foreground
13      ↪ consists of a rocky terrain with sparse vegetation. the sky is overcast, and there are no visible
14      ↪ people or moving objects.",
15      "A silhouette of a person operating a large, mounted weapon on a rocky terrain under a clear sky. the
16      ↪ individual appears to be adjusting or aiming the weapon."
17    ],
18    [
19      "A series of images depicting a hazy, dusty environment with buildings in the background. the focus is on
20      ↪ a red light source that appears to be moving or flashing, possibly indicating some form of activity or
21      ↪ event occurring in the area.",
22      "A series of images depicting a military tank moving across a landscape with buildings in the background.
23      ↪ the tank is seen from different angles, and there are visible explosions occurring in the distance."
24    ]
25  ],
26   "events": [[5.583, 24.542], [36.167, 50.792]],
27   "event_summary": [
28     "The anomaly exists, specifically an Explosion anomaly. The basis for judging the anomaly is the presence of an
29     ↪ unexpected and sudden release of energy, likely accompanied by a loud noise and visible effects such as smoke,
30     ↪ fire, or debris.",
31     "The anomaly exists and its specific name is Explosion. The anomaly event is a series of visible explosions
32     ↪ occurring in the distance as a military tank moves across a landscape with buildings in the background. The
33     ↪ basis for judging the anomaly is the sudden and intense release of energy and light in the form of explosions,
34     ↪ which is an unusual and notable event in the otherwise desolate and quiet war-torn environment depicted in the
35     ↪ video",
36     "The anomaly exists, specifically named as Explosion. The anomaly event depicts a series of
37     ↪ visible explosions occurring in the distance as a military tank moves across a landscape with buildings in the
38     ↪ background, characterized by sudden and intense releases of energy and light, accompanied by loud noise and
39     ↪ visible effects such as smoke, fire, or debris. The basis for judging the anomaly lies in the unusual and
40     ↪ notable nature of these explosions, which stand out against the otherwise desolate and quiet war-torn
41     ↪ environment depicted in the video, making them an unexpected and sudden release of energy that grabs
42     ↪ attention."
43   ]
44 }

```

Figure A. **An example of hierarchical free-text annotations.** For each labeled video, the hierarchical free-text annotations include clip-level captions, event-level, and video-level anomaly analysis. Additionally, the temporal boundaries for each event and clip are annotated.

10. "Are there anomalies observed in the video clip?"

Description.

1. "Describe the anomaly events observed in the video."
2. "Could you describe the anomaly events observed in the video?"
3. "Could you specify the anomaly events present in the video?"
4. "Give a description of the detected anomaly events in this video."
5. "Could you give a description of the anomaly events in the video?"
6. "Provide a summary of the anomaly events in the video."
7. "Could you provide a summary of the anomaly events in this video?"
8. "What details can you provide about the anomaly in the video?"
9. "How would you detail the anomaly events found in the video?"
10. "How would you describe the particular anomaly events in the video?"

Analysis.

1. "Why do you judge this event to be anomalous?"
2. "Can you provide the reasons for considering it anomalous?"
3. "Can you give the basis for your judgment of this event as an anomaly?"
4. "What led you to classify this event as an anomaly?"
5. "Could you provide the reasons for considering this event as abnormal?"
6. "What evidence do you have to support your judgment of this event as an anomaly?"
7. "Can you analyze the factors contributing to this anomalous event?"
8. "Could you share your analysis of the anomalous event?"
9. "What patterns did you observe that contributed to your conclusion about this event being an anomaly?"
10. "How do the characteristics of this event support its classification as an anomaly?"

A.4. Data Samples.

To facilitate understanding, we provide the final constructed instruction data at various temporal granularities, including **clip-level**, **event-level**, and **video-level**, as shown in Fig. B, Fig. C, and Fig. D.

B. Details of the Anomaly Scorer

B.1. Model Architecture

We use UR-DMU [12] as the anomaly scorer in our Anomaly-focused Temporal Sampler. As shown in Fig. E, UR-DMU utilizes a Global and Local Multi-Head Self Attention (GL-MHSA) module to capture both long-range and short-range temporal relationships among video snippets. Furthermore, UR-DMU introduces two memory banks to store and differentiate abnormal and normal prototypes, thereby maximizing the margins between these two representations. In order to learn discriminative representations, UR-DMU uses triplet loss to increase the feature distance



```

1 {
2   "id": "Abuse007_x264_E0C2",
3   "type": "clip",
4   "video":
5     ↪ "ucf-crime/clips/train/Abuse007_x264_E0C2.mp4",
6   "conversations": [
7     {"from": "human",
8      "value": "Describe the key events in the video."},
9     {"from": "gpt",
10      "value": "The man in shorts was slapped to the
11      ↪ ground by the strong man opposite."}
12   ]
13 }

```

Figure B. An example of clip-level instruction data.



```

1 {
2   "id": "Abuse018_x264_E0",
3   "type": "event",
4   "video":
5     ↪ "ucf-crime/events/train/Abuse018_x264_E0.mp4",
6   "conversations": [
7     {"from": "human",
8      "value": "Identify any abnormal behaviors depicted
9      ↪ in the video."},
10    {"from": "gpt",
11     "value": "An anomaly exists, specifically an
12     ↪ instance of Police Brutality."},
13    {"from": "human",
14     "value": "Could you provide a summary of the anomaly
15     ↪ events in this video?"},
16    {"from": "gpt",
17     "value": "The anomaly event involves the police
18     ↪ officer kicking the man who rushed towards him,
19     ↪ causing the man to fall heavily to the ground
20     ↪ with his head on his back."},
21    {"from": "human",
22     "value": "Could you share your analysis of the
23     ↪ anomalous event?"},
24    {"from": "gpt",
25     "value": "The excessive use of force by the police
26     ↪ officer, which deviates from standard protocol
27     ↪ and may be considered an abuse of authority."}
28   ]
29 }

```

Figure C. An example of event-level instruction data.

after interacting with different memories. Simultaneously, it utilizes KL loss to constrain the normal memory to follow a Gaussian distribution, accounting for the variance introduced by noise. Furthermore, We leveraged the event-level anomaly boundaries obtained during the annotation phase to generate frame-level labels and computed a binary cross-entropy loss, *i.e.*, \mathcal{L}_{AS} , which is simple yet effective. Thus, the loss function for the anomaly scorer is defined as follows:

$$\mathcal{L} = \mathcal{L}_{AS} + \mathcal{L}_{triplet} + \mathcal{L}_{kl} \quad (1)$$



```

1 {
2   "id": "v=xLFbqPzSI2Y__#1_label_B4-0-0",
3   "type": "video",
4   "video": "xd-violence/videos/train/
5   v=xLFbqPzSI2Y__#1_label_B4-0-0.mp4",
6   "conversations": [
7     {"from": "human",
8      "value": "Are there any anomaly events detected in
9      ↪ the video?"},
10    {"from": "gpt",
11     "value": "The anomaly exists and its specific name
12     ↪ is the Riot anomaly."},
13    {"from": "human",
14     "value": "Could you give a description of the
15     ↪ anomaly events in the video?"},
16    {"from": "gpt",
17     "value": "The anomaly event depicts a chaotic and
18     ↪ intense urban unrest scene, starting with a
19     ↪ chaotic street scene where individuals are
20     ↪ running away from an area of fire and smoke,
21     ↪ with people appearing to be in a state of panic
22     ↪ or urgency, and some being sprayed with water by
23     ↪ individuals in uniform. The scene escalates to
24     ↪ show a growing fire in the background, police
25     ↪ officers in riot gear standing on the street,
26     ↪ and a formation of police officers spread out
27     ↪ across the road, suggesting an ongoing riot or
28     ↪ civil unrest. "},
29    {"from": "human",
30     "value": "Can you give the basis for your
31     ↪ judgment?"},
32    {"from": "gpt",
33     "value": "The presence of multiple indicators of a
34     ↪ riot or emergency situation, including fire,
35     ↪ smoke, and people fleeing in panic, as well as
36     ↪ the sudden appearance of flames, smoke, and
37     ↪ police in riot gear, and the presence of violent
38     ↪ and destructive events, such as fires and
39     ↪ explosions, and the response of authorities,
40     ↪ including the use of tear gas or pepper spray,
41     ↪ which deviate significantly from normal,
42     ↪ peaceful urban activity."}
43   ]
44 }

```

Figure D. An example of video-level instruction data.

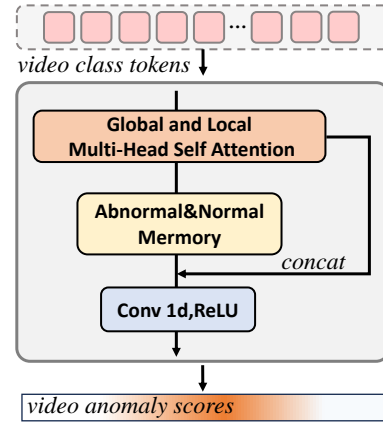


Figure E. Architecture of the Anomaly Scorer (UR-DMU [12]).

B.2. Visualization Results

In Fig. F, we present visualized results of anomaly scores and sampled frames on the UCF-Crime and XD-Violence test sets. These results demonstrate the accuracy of our

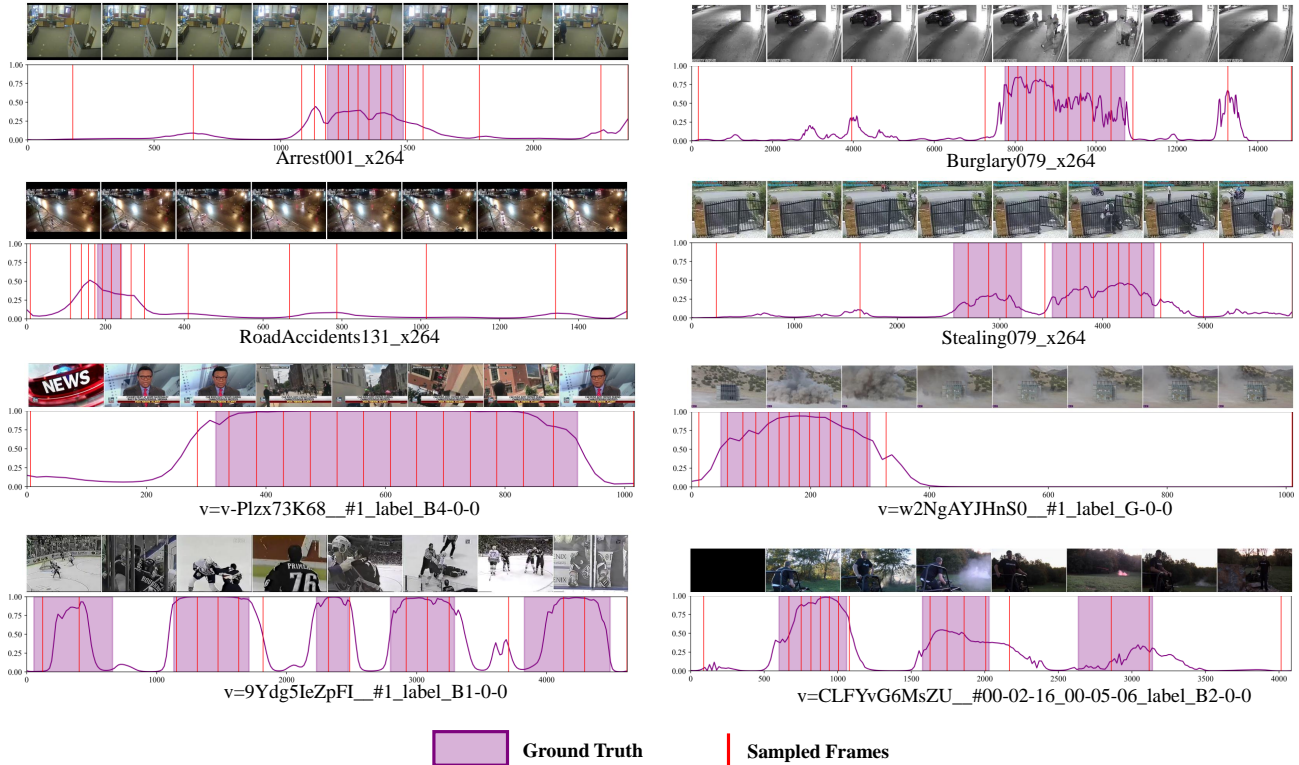


Figure F. Visualization results of anomaly scores and sampled frames output by the Anomaly-focused Temporal Sampler.

Table A. Comparison with related multimodal/explainable VAU methods and benchmarks. HIVAU-70k provides accurate temporal annotations and hierarchical anomaly-related free-text annotations.

Methods	#Categories	#Samples	Text			Temp. Anno.	MLLM tuning
			clip-level	event-level	video-level		
UCA [8]	13	23,542	✓	✗	✗	✓	✗
LAVAD [10]	N/A	N/A	✓	✗	✓	✗	✗
VAD-VideoLLama [4]	13/7	2,400	✗	✗	✓	✗	projection
CUVA [3]	11	6,000	✗	✗	✓	✗	✗
Hawk [6]	-	16,000	✗	✗	✓	✗	projection
HIVAU-70k (Ours)	19	70,000	✓	✓	✓	✓	LoRA

method in anomaly detection within complex real-world scenarios, with the sampled frames being concentrated in anomalous regions.

C. Discussion with related works.

In Table A, we provide a comprehensive comparison with related works in terms of benchmarks and methods.

Summary of related works: Recently, there has been substantial research on multi-modal Video Anomaly Understanding, making significant contributions to advancing open-world anomaly understanding. LAVAD [10] utilized several pre-trained foundational models to offer a training-free explainable VAD process. VAD-VideoLLaMA [4],

designed a three-phase training method to finetune VideoLLaMA in the VAD domain. CUVA [3] introduced a dataset and metric for evaluating causation understanding of video anomalies. Hawk [6] constructed an instruction dataset and finetuned a video-language framework that incorporates both motion and video information.

Difference and Advantages of our proposed benchmark and method:

- We develop a semi-automated annotation engine that scales hierarchical anomaly annotation efficiently, combining manual refinement with LLM-based annotation to maintain high-quality data across multiple granularities, resulting in over **70,000** annotations at clip, event, and video levels, which significantly surpasses previous

datasets in scale.

- UCA [9] only provides clip-level captions, overlooking the understanding of anomalies across longer time spans. CUVA [3] and Hawk [6], on the other hand, only offer video-level instruction data, neglecting finer-grained visual perception and anomaly analysis. In contrast, our proposed HIVA-70k takes a multi-temporal granularity perspective, offering more comprehensive and diverse anomaly annotations for open-world anomaly detection. It enables progressive and comprehensive learning, from short-term visual perception to long-term anomaly reasoning.
- We propose the **Anomaly-focused Temporal Sampler (ATS)**, integrated with a multi-modal visual-language model. Benefiting from the precise temporal annotations we provide, the ATS is able to focus on anomaly-dense video segments. This integration significantly improves efficiency and accuracy in long-video anomaly detection.

D. More Qualitative Results.

As shown from Fig. G to Fig. J, we present the output of explainable text generated by Holmes-VAU compared with the base model, *i.e.*, InternVL-2B [2]. The results demonstrate significant improvements in the model’s visual perception and anomaly analysis capabilities after fine-tuning on HIVA-70k.

E. Limitations and Future Work.

While our work demonstrates significant strides in multi-granular video anomaly understanding, several limitations present opportunities for future enhancement. First, optimizing for real-time streaming remains a challenge. Our sparse sampling approach improves efficiency, but further refinement is necessary for seamless deployment in streaming contexts. Additionally, our work has so far focused on surveillance data, extending our framework to other domains, such as industrial monitoring and medical diagnostics, will help validate its generalization capabilities. Lastly, integrating additional sensory data, like audio, and scalable hierarchical annotation could enhance anomaly detection and broaden applicability.

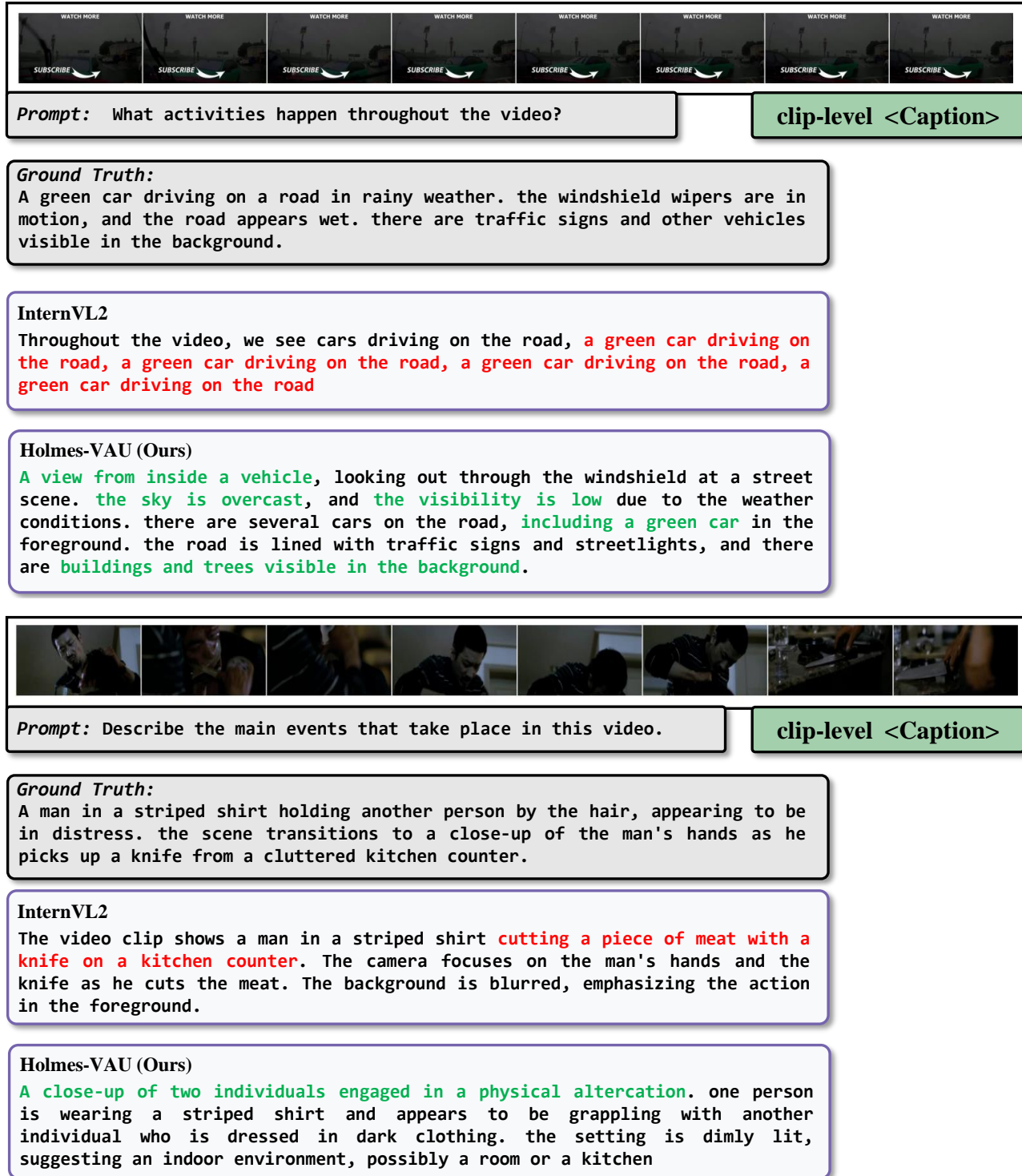


Figure G. Qualitative comparison of anomaly understanding explanation with our baseline model, i.e., InternVL-2B. Correct and wrong explanations are highlighted in green and red, respectively.



Figure H. Qualitative comparison of anomaly understanding explanation with our baseline model, i.e., InternVL-2B. Correct and wrong explanations are highlighted in green and red, respectively.

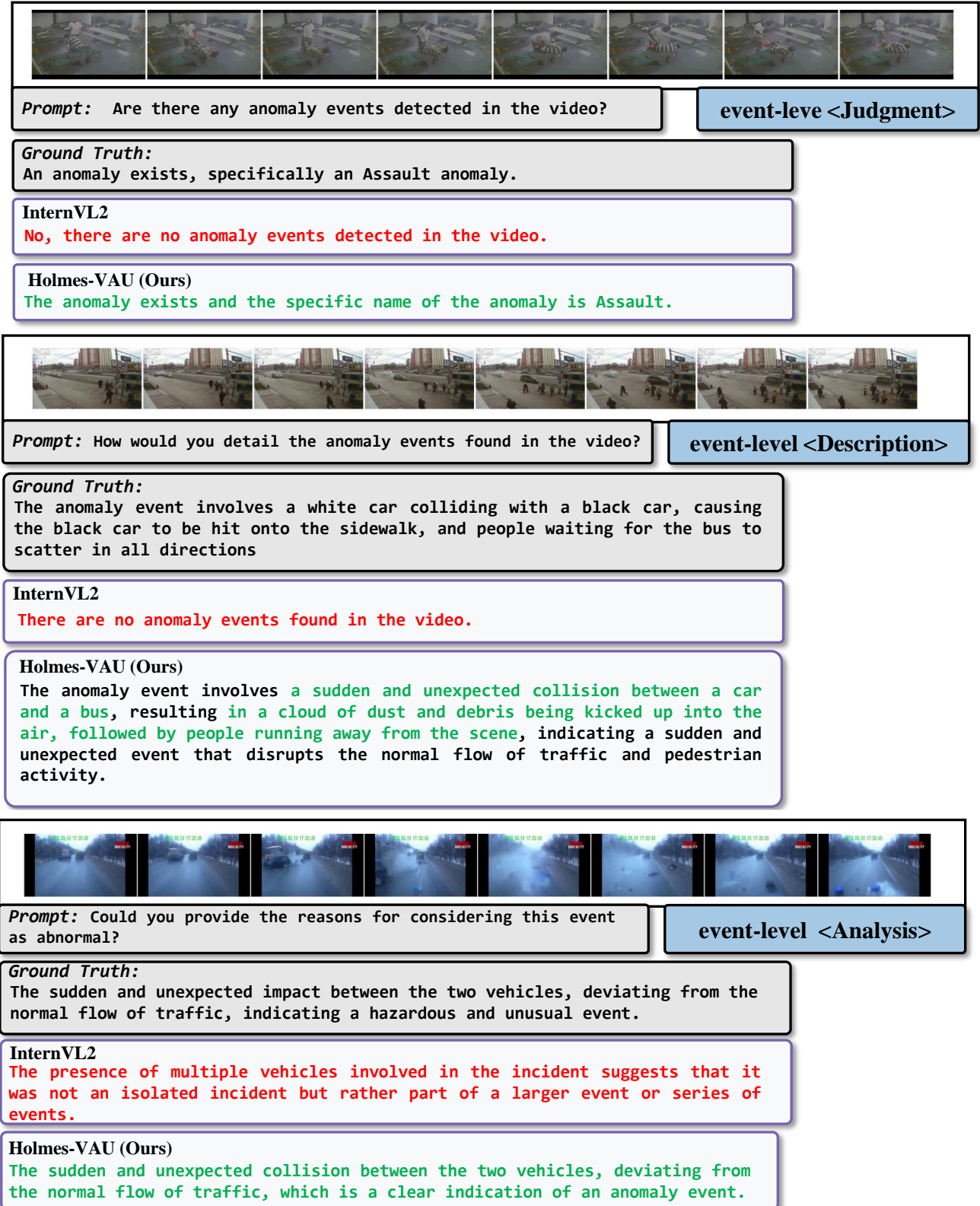


Figure I. Qualitative comparison of anomaly understanding explanation with our baseline model, i.e., InternVL-2B. Correct and wrong explanations are highlighted in green and red, respectively.



Figure J. Qualitative comparison of anomaly understanding explanation with our baseline model, i.e., InternVL-2B. Correct and wrong explanations are highlighted in green and red, respectively.

References

- [1] AI@Meta. Llama 3 model card. 2024. [1](#)
- [2] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. [5](#)
- [3] Hang Du, Sicheng Zhang, Binzhu Xie, Guoshun Nan, Jiayang Zhang, Junrui Xu, Hangyu Liu, Sicong Leng, Jiangming Liu, Hehe Fan, et al. Uncovering what why and how: A comprehensive benchmark for causation understanding of video anomaly. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18793–18803, 2024. [4](#), [5](#)
- [4] Hui Lv and Qianru Sun. Video anomaly detection and explanation via large language models. *arXiv preprint arXiv:2401.05702*, 2024. [4](#)
- [5] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6479–6488, 2018. [1](#)
- [6] Jiaqi Tang, Hao Lu, Ruizheng Wu, Xiaogang Xu, Ke Ma, Cheng Fang, Bin Guo, Jiangbo Lu, Qifeng Chen, and Yingcong Chen. Hawk: Learning to understand open-world video anomalies. *arXiv preprint arXiv:2405.16886*, 2024. [4](#), [5](#)
- [7] Peng Wu, Jing Liu, Yujia Shi, Yujia Sun, Fangtao Shao, Zhaoyang Wu, and Zhiwei Yang. Not only look, but also listen: Learning multimodal violence detection under weak supervision. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages 322–339. Springer, 2020. [1](#)
- [8] Tongtong Yuan, Xuange Zhang, Kun Liu, Bo Liu, Chen Chen, Jian Jin, and Zhenzhen Jiao. Towards surveillance video-and-language understanding: New dataset, baselines, and challenges, 2023. [4](#)
- [9] Tongtong Yuan, Xuange Zhang, Kun Liu, Bo Liu, Chen Chen, Jian Jin, and Zhenzhen Jiao. Towards surveillance video-and-language understanding: New dataset baselines and challenges. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22052–22061, 2024. [1](#), [5](#)
- [10] Luca Zanella, Willi Menapace, Massimiliano Mancini, Yiming Wang, and Elisa Ricci. Harnessing large language models for training-free video anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18527–18536, 2024. [4](#)
- [11] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. [1](#)
- [12] Hang Zhou, Junqing Yu, and Wei Yang. Dual memory units with uncertainty regulation for weakly supervised video anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3769–3777, 2023. [2](#), [3](#)