

ATOMS TO EVENTS: CATEGORICAL EVIDENCE COMPOSITION FOR VIDEO ANOMALY DETECTION

005 **Anonymous authors**

006 Paper under double-blind review

ABSTRACT

011 Video anomaly detection (VAD) seeks to identify events that deviate from learned
 012 normality. Current Vision-Language Models (VLMs) face significant challenges:
 013 anomalies are rare, labels are weak, and visual appearance varies drastically.
 014 Mainstream VLMs directly map visual features to events, they overfit to inter-
 015 mediate incidental cues which are present during training and generalize poorly.
 016 To address this issue, we propose a categorical view of anomaly understanding.
 017 Firstly, an Unsupervised Anomalous Period Detector (UAPD) is proposed to iden-
 018 tify abnormal periods. Next, a Category-based Atom Miner (CAM) is proposed
 019 to map visual features to learned atoms in video segments, and learn the roles of
 020 atoms. In inference, CAM provides role-aware indications to VLM which maps
 021 meaningful atoms and visual features to event predictions. This framework har-
 022 nesses meaningful evidence and preserves the generalization capacity of VLMs.
 023 Extensive experiments and ablations show consistent gains over strong vision-only
 024 and fine-tuned VLM baselines.

1 INTRODUCTION

028 Video anomaly detection (VAD) seeks to localize the time spans in long videos where scenes de-
 029 viate from regular patterns—e.g., violence, accidents, explosions, or other unexpected events. The
 030 problem is challenging because anomalies are rare in contrast to the massive scale of surveillance
 031 streams, visual conditions vary substantially across cameras and over time due to viewpoint, light-
 032 ing and other factors, and labels are typically video-level or even absent. Existing approaches have
 033 explored vision-only strategies that rely purely on visual features: prediction-based models that
 034 forecast future sequences and flag deviations, reconstruction-based models that assume anomalies
 035 reconstruct poorly, representations combining multiple feature types or Multiple Instance Learning
 036 (MIL)-based approaches Li et al. (2022); Georgescu et al. (2021); Park et al. (2020); Noghre et al.
 037 (2024); Yang et al. (2023); Liu et al. (2021); Huang et al. (2025); Georgescu et al. (2021); Cho et al.
 038 (2023); Guo et al. (2023); Liu et al. (2022). Although these approaches are effective, some of them
 039 cannot perform well under domain shifts, and may mistake irrelevant variations for true anomalies.

040 To improve interpretability and zero-shot generalization, recent works leverage Vision-Language
 041 Models (VLMs). They use captions, prompts, or pseudo-labels to score abnormal semantics within
 042 sliding windows Cao et al. (2024); Chen et al. (2024); Yang et al. (2024a); Micorek et al. (2024); Zhu
 043 & Pang (2024); Li et al. (2024); Yang et al. (2024b); Tang et al. (2025). Training-free approaches
 044 Zanella et al. (2024) directly caption frames with VLM and summarize how captions change over
 045 time. Guidance-driven methods focus VLMs on anomaly evidences by leveraging verbalized ques-
 046 tions Ye et al. (2025) or sampling the most suspicious snippets Zhang et al. (2025). These methods
 047 are interpretable. However, VLM-based detectors often encode an event as one entangled concept,
 048 letting nuisance factors, such as viewpoints and backgrounds, swamp model representations. As a
 049 result, the visual features of events are not well distinguished.

050 To stably attend VLMs to anomaly-relevant semantics without being distracted by variations, chal-
 051 lenges lie in the set-theoretic fashion of VLMs. They learn a direct mapping from visual features to
 052 event labels by accumulating “relevant” cues, a video is treated as a set of features whose member-
 053 ships vote for the label. Such membership often admits incidental cues and spurious contexts. We
 054 find that meaningful features can be divided into directly relational, inherently relational and counter
 055 evidences. In this regard, we propose a categorical, role-aware view: we treat atoms as objects with

054 roles in events — Direct(support), Indirect (synergy), and Counter(inhibition). The "visual features
 055 to events" mapping in set-theoretic VLMs changes to "visual features to atoms to events". The cues
 056 without valid roles have no effect. Robust reasoning and auditable explanations are yielded.
 057

058 **REFERENCES**
 059

060 Chentao Cao, Zhun Zhong, Zhanke Zhou, Yang Liu, Tongliang Liu, and Bo Han. Envisioning
 061 outlier exposure by large language models for out-of-distribution detection. *arXiv preprint*
 062 *arXiv:2406.00806*, 2024.

063 Jiankang Chen, Tong Zhang, Wei-Shi Zheng, and Ruixuan Wang. Tagfog: Textual anchor guidance
 064 and fake outlier generation for visual out-of-distribution detection. In *Proceedings of the AAAI*
 065 *Conference on Artificial Intelligence*, volume 38, pp. 1100–1109, 2024.

066 MyeongAh Cho, Minjung Kim, Sangwon Hwang, Chaewon Park, Kyungjae Lee, and Sangyoun
 067 Lee. Look around for anomalies: Weakly-supervised anomaly detection via context-motion re-
 068 lational learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
 069 *Recognition*, pp. 12137–12146, 2023.

070 Mariana-Iuliana Georgescu, Antonio Barbalau, Radu Tudor Ionescu, Fahad Shahbaz Khan, Marius
 071 Popescu, and Mubarak Shah. Anomaly detection in video via self-supervised and multi-task
 072 learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-
 073 tion*, pp. 12742–12752. IEEE, 2021.

074 Chongye Guo, Hongbo Wang, Yingjie Xia, and Guorui Feng. Learning appearance-motion synergy
 075 via memory-guided event prediction for video anomaly detection. *IEEE Transactions on Circuits*
 076 *and Systems for Video Technology*, 2023.

077 Yuzhi Huang, Chenxin Li, Haitao Zhang, Zixu Lin, Yunlong Lin, Hengyu Liu, Wuyang Li, Xinyu
 078 Liu, Jiechao Gao, Yue Huang, et al. Track any anomalous object: A granular video anomaly
 079 detection pipeline. In *Proceedings of the Computer Vision and Pattern Recognition Conference*,
 080 pp. 8689–8699, 2025.

081 Shuo Li, Fang Liu, and Licheng Jiao. Self-training multi-sequence learning with transformer for
 082 weakly supervised video anomaly detection. In *Proceedings of the AAAI Conference on Artificial*
 083 *Intelligence*, volume 24. AAAI Press, 2022.

084 Xiaofan Li, Zhizhong Zhang, Xin Tan, Chengwei Chen, Yanyun Qu, Yuan Xie, and Lizhuang Ma.
 085 Promptad: Learning prompts with only normal samples for few-shot anomaly detection. In *Pro-
 086 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16838–
 087 16848, 2024.

088 Yuejiang Liu, Riccardo Cadei, Jonas Schweizer, Sherwin Bahmani, and Alexandre Alahi. Towards
 089 robust and adaptive motion forecasting: A causal representation perspective. In *Proceedings of*
 090 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17081–17092. IEEE,
 091 2022.

092 Zhan Liu, Yongwei Nie, Chengjiang Long, Qing Zhang, and Guiqing Li. A hybrid video anomaly
 093 detection framework via memory-augmented flow reconstruction and flow-guided frame predic-
 094 tion. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 13588–
 095 13597, 2021.

096 Jakub Micorek, Horst Possegger, Dominik Narnhofer, Horst Bischof, and Mateusz Kozinski. Mulde:
 097 Multiscale log-density estimation via denoising score matching for video anomaly detection.
 098 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
 099 18868–18877, 2024.

100 Ghazal Alinezhad Noghre, Armin Danesh Pazho, and Hamed Tabkhi. An exploratory study on
 101 human-centric video anomaly detection through variational autoencoders and trajectory predic-
 102 tion. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*,
 103 pp. 995–1004, 2024.

108 Hyunjong Park, Jongyoun Noh, and Bumsub Ham. Learning memory-guided normality for anomaly
109 detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-*
110 *nition*, pp. 14372–14381. IEEE, 2020.

111

112 Jiaqi Tang, Hao Lu, Ruizheng Wu, Xiaogang Xu, Ke Ma, Cheng Fang, Bin Guo, Jiangbo Lu, Qifeng
113 Chen, and Yingcong Chen. Hawk: Learning to understand open-world video anomalies. *Advances*
114 *in Neural Information Processing Systems*, 37:139751–139785, 2025.

115 Yuchen Yang, Kwonjoon Lee, Behzad Dariush, Yinzhi Cao, and Shao-Yuan Lo. Follow the
116 rules: Reasoning for video anomaly detection with large language models. *arXiv preprint*
117 *arXiv:2407.10299*, 2024a.

118

119 Zhiwei Yang, Jing Liu, Zhaoyang Wu, Peng Wu, and Xiaotao Liu. Video event restoration based
120 on keyframes for video anomaly detection. In *Proceedings of the IEEE/CVF Conference on*
121 *Computer Vision and Pattern Recognition*, pp. 14592–14601, 2023.

122 Zhiwei Yang, Jing Liu, and Peng Wu. Text prompt with normality guidance for weakly supervised
123 video anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
124 *Pattern Recognition*, pp. 18899–18908, 2024b.

125

126 Muchao Ye, Weiyang Liu, and Pan He. Vera: Explainable video anomaly detection via verbalized
127 learning of vision-language models. In *Proceedings of the Computer Vision and Pattern Recog-*
128 *nition Conference*, pp. 8679–8688, 2025.

129

130 Luca Zanella, Willi Menapace, Massimiliano Mancini, Yiming Wang, and Elisa Ricci. Harnessing
131 large language models for training-free video anomaly detection. In *Proceedings of the IEEE/CVF*
132 *Conference on Computer Vision and Pattern Recognition*, pp. 18527–18536, 2024.

133

134 Huaxin Zhang, Xiaohao Xu, Xiang Wang, Jialong Zuo, Xiaonan Huang, Changxin Gao, Shanjun
135 Zhang, Li Yu, and Nong Sang. Holmes-vau: Towards long-term video anomaly understanding at
136 any granularity. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp.
137 13843–13853, 2025.

138

139 Jiawen Zhu and Guansong Pang. Toward generalist anomaly detection via in-context residual learn-
140 ing with few-shot sample prompts. In *Proceedings of the IEEE/CVF Conference on Computer*
141 *Vision and Pattern Recognition*, pp. 17826–17836, 2024.

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161