

Video Anomaly Detection Using Multi-level Text Prompts and Context-Aware Frame Sampling for Long Duration Videos

Anonymous WACV **Algorithms Track** submission

Paper ID ****

Abstract

001 Weakly supervised video anomaly detection (WSVAD) is an
002 essential yet challenging problem, particularly for long-
003 duration surveillance videos. Recent approaches have
004 leveraged contrastive language-image pretraining (CLIP)
005 to enhance frame-level feature representations, yielding sig-
006 nificant improvements. However, most existing methods
007 overlook the challenges associated with the duration of the
008 videos. These methods employ uniform or stepwise frame
009 sampling to manage video length, which risks discarding
010 crucial fine-grained temporal information and leads to in-
011 formation bottlenecks. In this paper, we propose MLP-
012 VAD, a novel framework that addresses these limitations
013 through text-guided frame selection using a frozen Image-
014 Text Matching CLIP model. By leveraging class-level se-
015 mantic descriptions, our approach adaptively filters frames
016 relevant to downstream anomaly detection tasks. We also
017 introduce multi-level prompt-guided feature learning integrat-
018 ing entity-level and scene-level prompts to mimic hu-
019 man cognitive processes more closely in anomaly detec-
020 tion task. Extensive experiments on the UCF-Crime and
021 XD-Violence benchmarks demonstrate that MLPVAD sur-
022 passes current state-of-the-art methods, achieving an AUC
023 of 88.84% and an AP of 85.37%, respectively.

024 1. Introduction

025 Weakly supervised video anomaly detection (WSVAD) has
026 recently gained significant attention in the fields of video
027 understanding and computer vision. The primary objective
028 of WSVAD is to utilize video-level annotations to detect
029 abnormal activities within video frames, with applications
030 spanning surveillance systems [2], traffic monitoring [1],
031 and beyond. Given the scarcity of frame-level annotations,
032 video-level supervision has become prevalent in WSVAD
033 research.

034 Existing approaches typically employ backbone net-
035 works such as I3D/C3D [5] or transformer-based models

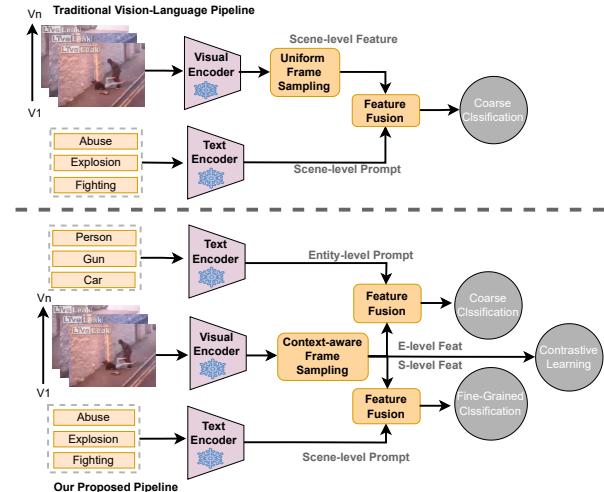


Figure 1. Traditional multimodal VAD vs. our proposed MLP-VAD paradigm.

like ViT [10] to extract frame-level representations. These features are processed with temporal modeling within a Multiple Instance Learning (MIL) framework [11] to predict normality or abnormality scores. Although effective on several datasets, these approaches mainly rely on single-domain video data, which limits their generalizability in complex scenarios where distinctions between normal and abnormal events are very subtle [7, 13]. The vision-language models have emerged as a promising direction for video anomaly detection [4, 6, 26]. By leveraging the CLIP model [19], these methods obtain rich, generalized visual representations aligned with textual descriptions via contrastive learning on large-scale image-text datasets. Incorporating CLIP as a feature extractor enhances contextual understanding of challenging scenes and significantly improves performance compared to unimodal approaches. Despite these advances, several challenges remain:

First, most existing methods are designed for short to medium-length videos. Processing long-duration videos

036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

055 (e.g., over 30 minutes at 30 fps, containing upwards
 056 of 60,000 frames) necessitates frame sampling strategies,
 057 commonly uniform or stepwise, to reduce input size to
 058 manageable lengths (e.g., 128–256 frames). Note that the
 059 coarse sampling induces temporal noise, disregards impor-
 060 tant fine-grained frame information, and risks excluding
 061 critical frames, thereby impairing overall performance.

062 **Second**, current multimodal WSVAD methods predom-
 063 inantly utilize coarse class-level textual prompts [26] at-
 064 tended over aggregated video-level features, overlooking
 065 fine-grained entity-level information and frame-wise inter-
 066 actions. This coarse granularity limits their ability to capture
 067 subtle anomaly cues in the frames.

068 **Third**, effective discrimination between normal and ab-
 069 normal scenes with subtle visual differences remains under-
 070 explored, yet is essential for realistic applications.

071 To address these challenges, we propose a novel
 072 paradigm for VAD that emulates human cognitive pro-
 073 cesses. Our model, **MLPVAD**, illustrated in Figure 1, in-
 074 incorporates the following key components:

- 075 • **Context-Aware Frame Sampling (CAFS):** To effi-
 076 ciently process long videos, CAFS selects a compact
 077 subset of informative frames by leveraging the CLIP-
 078 based pre-trained model *BLIP2* [15]. Through frame-
 079 wise image-text matching, CAFS filters redundant frames
 080 while retaining those most representative of anomalies.
 081 Here, text prompts are textual summaries comprising
 082 class-specific entities and actions.

- 083 • **Multi-level Text Prompt Learning:** To accurately un-
 084 derstand a scene, it is necessary to focus on each object in
 085 the frame as well as their interactions over time. Entity-
 086 level feature learning emphasizes class entities in frames
 087 (e.g., *Fight*: person, fist, kick, face; *RoadAccident*: car,
 088 truck, collision, fire). This enables fine-grained frame-
 089 level attention through entity prompts, concentrating on
 090 object relations within the local context. Scene-guided
 091 learning further addresses entity interactions over time
 092 by applying attention over the summarized video features
 093 and scene prompts. This hierarchical approach improves
 094 robustness in complex, long-duration scenarios.

- 095 • **Contrastive Representation Learning:** To improve dis-
 096 crimination between subtly different classes, a contrastive
 097 learning branch leverages multiple positive and negative
 098 samples, encouraging anomaly representations to align
 099 closely with queries while separating normal representa-
 100 tions, thereby complementing classification objectives.

101 By closely modeling human visual reasoning and ex-
 102 ploiting advanced vision-language representations, **MLP-**
VAD effectively addresses the challenges above, demon-
 103 strating superior performance on large-scale benchmarks.

104 The main contributions of this work are:

- 105 1. A novel VAD framework that emulates human cogni-
 106 tive processes, featuring a context-aware frame sampling

(CAFS) module optimized for long-duration videos by
 107 leveraging large language models (LLMs).

- 108 2. A *Multi-level Text Prompt Guided Network* enhancing
 109 contextual understanding at both entity and scene lev-
 110 els, supported by contrastive learning for discriminative
 111 coarse- and fine-grained classification.
- 112 3. Empirical validation of the **MLPVAD** model on two
 113 large-scale benchmark datasets, achieving state-of-the-
 114 art performance with an AUC of 88.84% and an AP of
 115 85.37% on UCF-Crime and XD-Violence, respectively.

2. Related Works

2.1. Weakly Supervised Video Anomaly Detection:

Leveraging video-level labels, weakly supervised methods have made significant improvements over the past few years in the video anomaly detection task. **WSVAD** was first introduced by Sultani et al. [20] using an MIL approach for surveillance video and outperformed previous unsupervised models. Since then, numerous advancements have been made for WSVAD, improving the temporal network and better feature learning with transformer and CLIP-based pretrained models. To better capture the temporal dynamics across video frames, Tian et al. [21] used a self-attention network. Zhong et al. reformulated the WSVAD problem as a supervised learning task using a graph convolutional network (GCN) method, generating noisy labels from an off-the-shelf video [30]. Feng et al. attempted to remove noisy labels with a multiple instance pseudo-label generator that produces more reliable pseudo-labels for fine-tuning a task-specific feature encoder [11]. More recently, the CLIP model has also attracted great attention in the VAD community [26]. Based on the visual features of the CLIP model, Lv et al [17] proposed a new MIL framework: Unbiased MIL (UMIL), to learn unbiased anomaly features that improve WSVAD performance. Furthermore, Wu et al. utilize a learnable text prompt that leverages the frozen CLIP model, performing both coarse-grained and fine-grained classification at the video level and class level [26]. Most of these weakly supervised methods perform classification using video labels and do not focus on the entity that plays a crucial role in the anomaly activities. This performance is suboptimal and does not reflect real-world understanding of challenging anomaly activities.

2.2. Large Vision-Language Models:

Multimodal training has gained popularity due to its effectiveness in extracting rich, context-aware features from multimodal data. The application of multimodal training began with the most straightforward tasks, such as image classification, and expanded to more challenging tasks, including video understanding and VAD. One of the first works to include visual and textual information for the VAD task was

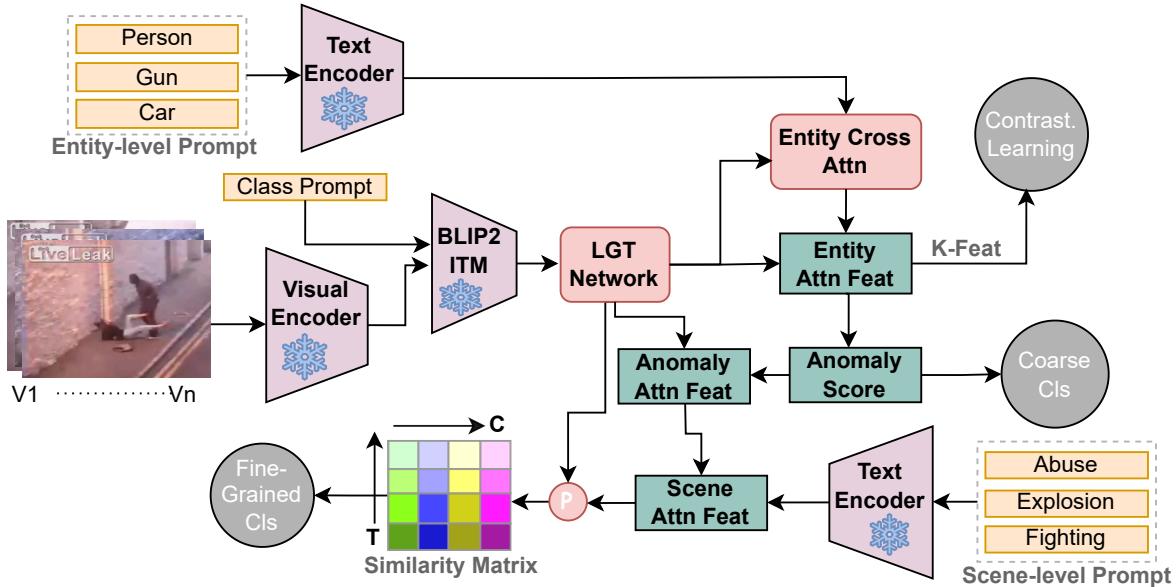


Figure 2. Proposed MLPVAD Architecture with each major component and objective functions.

introduced by Chen et al. [6]. In this work, the authors generated text captions from video frames using the Swin-BERT pretrained model and later performed classification on multimodal features. Later on, the CLIP model demonstrated significant progress in capturing text-vision relations and has been effective for several downstream tasks, such as classification [28], visual question-answering (VQA) [15], semantic segmentation [16], and others. The CLIP model has recently been effectively extended from static image understanding to the spatiotemporal domain of video, enabling joint modeling of visual and textual information over time. VideoCLIP [27] is proposed to align video and textual representations by contrasting temporally overlapping video-text pairs with hard negatives mining. Azad et al. [3] employed hierarchical context learning using multi-scale text prompts for better scene understanding. MA-LMM [12] also leverages the CLIP model for the video-understanding task and incorporates the memory-bank technique to retain long-term temporal information. Next, Zanella et al. [29] and Wu et al. [26] use learnable textual features from CLIP to enhance the representation of the overall features, followed by an MIL-based video anomaly classifier. Existing approaches rarely leverage CLIP models to sample a large number of video frames based on the CLIP text features. Instead, most prior work incorporates language information only at a coarse level for representation learning, primarily focusing on improving classification performance. There is a pressing need to investigate the power of the CLIP model and fully leverage its capabilities in language representation, as well as its effectiveness in the VAD pipeline.

3. Methodology

In this section, we first present the problem definition of the WSVAD task. Next, we briefly introduce the overall architecture of the proposed MLPVAD method. Next, we subsequently elaborate on each of the modules of the MLPVAD method and its role in solving different challenges discussed in earlier sections.

3.1. Problem Definition and Overall Architecture

The goal of this work is to perform the video anomaly detection task in a weakly supervised manner using video-level labels. For example, given a video v , if all the frames of the particular video are normal, the video is annotated as a normal video ($y = 0$). On the other hand, if even a few frames from video v contain abnormal frames, the complete video is annotated as an abnormal video ($y = 1$). The task is to perform binary classification and predict the \hat{y} for a given video v , with annotation $y = 0/1$. The recent successful approaches for the WSVAD task mostly used the MIL approach on frames to predict abnormality. In the MIL approach, previous works use CNN methods such as I3D [5] and C3D [22] to extract consecutive features, which are grouped to create video snippets for the binary classification task. This research focuses equally on the visual and textual information and leverages the power of Large-Language Models (LLMs) to accompany the visual pipeline. For feature extraction, we utilize the encoder of CLIP as the backbone due to its significant performance over the past few years across a wide range of downstream tasks such as VQA [18], VAD [26], Video Understanding

158 188

159 189

160 190

161 191

162 192

163 193

164 194

165 195

166 196

167 197

168 198

169 199

170 200

171 201

172 202

173 203

174 204

175 205

176 206

177 207

178 208

179 209

180 210

181 211

182 212

183 213

184 214

185 215

186 216

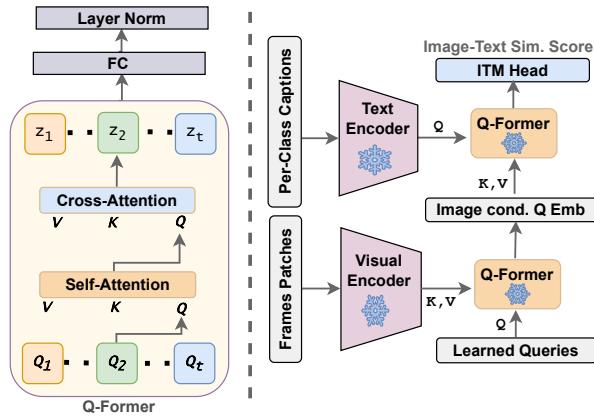


Figure 3. Context-Aware frame sampling (CAFS) with Q-Former module.

[12]. We have three objective functions (See Figure 2): 1) Coarse Classification on entity-attended visual features, 2) Randomized Contrastive Learning on top-k/bottom-k features, and 3) Fine-grained classification from the similarity matrix generated from scene-attended visual features.

3.2. Context Aware Frame Sampling (CAFS)

The current video understanding and VAD methods are prone to information loss and sometimes noisy due to the selection of irrelevant frames through uniform or stepwise frame selection. For the first time, we propose a CAFS method as illustrated in Figure 3, to select frames related to the downstream task. For a video v , we pass it to the CLIP visual encoder (ViT16/ViT14) and extract the visual feature $X \in B \times N \times d_v$. Here, B is the batch size, N is the number of frames (CLS embedding), and d_v is the visual feature dimension. The text-prompt is derived by randomly generating dense captions from several videos per class using VideoLLaMA2 [8]. The robustness of modern LLM-based video models for the caption generation task across different domains is demonstrated in several recent benchmark works [8, 23, 24]. Our experiment also found that the captions generated from VideoLLaMA2 are satisfactorily accurate to guide the frame selection with several strong actions/nouns. We present a few captions generated by VideoLLaMA2 in Figure 4 to prove the caption confidence. We randomly pick 30 videos per class from train set, create sub-samples of the videos, and run caption generations on them. The highest occurring entities/nouns form a final text summarization prompt and represent the overall actions/activities per class. We perform sub-sampling on videos and majority voting on generated entities to make the final prompt stable and not affected by the model’s hallucination. The final prompt is passed to the BLIP2-pretrained BERT text-encoder [9] and generates text embedding $T \in B \times \text{max_token} \times d_t$ for the further frame matching task.

$$\mathbf{Q}_{\text{self}} = \text{softmax} \left(\frac{\mathbf{Q}_l \mathbf{W}_Q (\mathbf{Q}_l \mathbf{W}_K)^{\top}}{\sqrt{d_k}} \right) (\mathbf{Q}_l \mathbf{W}_V) \quad 252$$

$$\mathbf{Q}'_l = \text{LayerNorm}(\mathbf{Q}_l + \mathbf{Q}_{\text{self}}) \quad 253$$

$$\mathbf{Q}_{\text{cross}} = \text{softmax} \left(\frac{\mathbf{Q}'_l \mathbf{W}'_Q (\mathbf{V} \mathbf{W}'_K)^{\top}}{\sqrt{d_k}} \right) (\mathbf{V} \mathbf{W}'_V) \quad 254$$

$$\mathbf{Q}_{l+1} = \text{LayerNorm}(\mathbf{Q}'_l + \mathbf{Q}_{\text{cross}}) \quad 255$$

The CAFS module builds upon the idea of Q-Former as illustrated in Figure 3(Left). The work process of Q-Former starts with M learnable query tokens $\mathbf{Q}_0 \in \mathbb{R}^{M \times D_q}$, where D_q is the query token dimension. Next, the self-attention query to query is defined in Eq. 1 and the cross-attention query-to-key tokens are defined in Eq. 2. We repeat the work process of Eq. 1 to 2 for L layers. Figure 3 shows how the CAFS is designed using two stages of the Q-Former block. In the first stage, we use pre-trained learned queries as the \mathbf{Q}_0 , whereas key (K) and values (V) come from the Visual Encoder output X . The production of the first stage is Visual attention query embeddings (VA_{emb}). This VA_{emb} is later used as the key (k) and value (v), whereas the Per-class text prompt embedding (T) is used as the Query embeddings. The output of the second-stage Q-former is derived from the last-layer cross-attention (A_{cross}), which represents the relevance between the class prompt and each frame in the particular video. Finally, the cross-attention output (A_{cross}) is passed into the ITM Head as provided in Eq. 3.

$$\text{ITM}_{\text{score}} = \text{Sort}(\text{Sigmoid}(\text{FC}(\mathbf{A}_{\text{cross}}))) \in \mathbb{R}^{T \times 2} \quad 255$$

From the set of T temporally sorted frames, the algorithm selects the Top- K frames with the highest matching scores for downstream anomaly detection. To enhance computational efficiency, the indices of these Top- K frames are cached after the first training epoch, thereby avoiding redundant recomputation in subsequent epochs. This module is task-agnostic and can be seamlessly integrated as an off-the-shelf component within a wide range of video understanding tasks.

3.3. Local-Global Temporal (LGT) Network

Preserving temporal information is crucial for the VAD task. There are several ways to model temporal information, such as adding positional embedding, Multi-scale Temporal Network with different strides of 1-D convolution [6], and Graph Convolutional Network (GCN) for learning the long-term and short-term temporal relationships [30]. Due to the recent effectiveness of the GCN for temporal modeling [26, 30], we leverage a similar implementation from a recent work on WSVAD [26]. The LGT Network in Figure 2 utilizes a self-attention module masked within a local

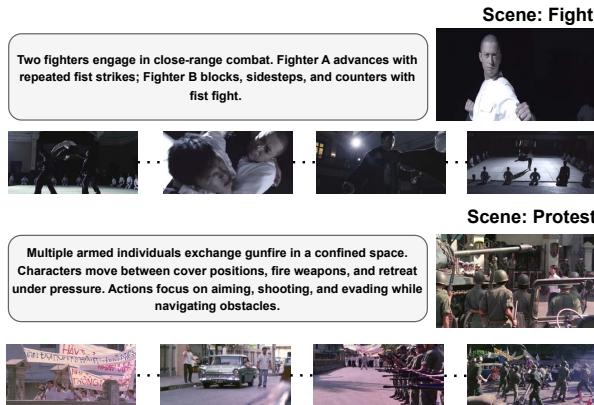


Figure 4. Generated caption using VideoLLaMA2 for different scene types.

296 window for **local temporal dependencies**. Next, we split
 297 the Top-k frame-level features into several uniform-length
 298 segments over the temporal dimension to achieve the lo-
 299 cal self-attention. The goal is to learn the temporal context
 300 from nearby features and establish local-context relation-
 301 ships across the video frames, denoted as \mathbf{X}_l . Windowed
 302 Attention Mask only attends to a fixed window or block,
 303 reducing quadratic complexity. To capture the **global tem-**
 304 **poral dependency**, we utilize the GCN module following
 305 the idea introduced by [26]. Wu et al. address the long-
 306 term temporal relationship based on feature similarity and
 307 relative distance. The following formula summarizes the
 308 concept of global temporal learning through the similarity
 309 and distance matrices.

$$\begin{aligned} \text{sim}_{\text{adj}} &= \text{GELU}(\text{Softmax}(H_{\text{sim}}(\mathbf{X}_l))) \\ \text{dist}_{\text{adj}} &= \text{GELU}(\text{Softmax}(H_{\text{dist}}(\mathbf{X}_l))) \\ \mathbf{X}_g &= \text{sim}_{\text{adj}} \parallel \text{dist}_{\text{adj}} \end{aligned} \quad (4)$$

311 The sim_{adj} is the cosine similarity function. The dist_{adj}
 312 is the proximity adjacency matrix defined as $-|i - j|/\sigma$,
 313 where i and j are frame indices and σ is a hyperparameter
 314 to control the influence of the distance relation.

315 3.4. Multi-Level Prompt Network

316 To improve VAD performance, we propose a novel
 317 paradigm that mimics the human cognitive process through
 318 hierarchical focus and multi-level text prompt attention. As
 319 shown in Figure 2, the pipeline consists of three branches
 320 overall. The top and bottom branches deal with multi-level
 321 text prompt learning. On the other hand, the middle branch
 322 provides the temporal visual information to accompany the
 323 text information.

324 **Entity-level Prompt Attention** focuses on important enti-
 325 ties that might play a role in a particular class of anomaly.
 326 For example: "Arrest" anomaly class focuses on enti-

ties/nouns such as: ["police", "handcuffs", "officer", "per-
 327 son", "uniform", "vehicle", "badge"]. The number of en-
 328 tities per class is denoted as N_c . We utilize a pre-trained
 329 CLIP text encoder (ViT-B/16) to generate entity embed-
 330 dings (T_e). Next, we perform entity cross-attention with
 331 temporal visual features. For cross attention, we use visual
 332 feature as Query ($\mathbf{Q} = \mathbf{X}_g$), and per-class entity prompts
 333 as key and value ($\mathbf{K} = \mathbf{V} = \mathbf{T}_e^c$). Eq. 5 outlines the cross-
 334 attention formula with a learnable class-entity weight W_c to
 335 focus more on important entities.
 336

$$\mathbf{A}_o = \text{Norm} \left(\text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \cdot \mathbf{W}_c \right) \quad (5)$$

337 Next, we normalize the weighted attention by dividing it by
 338 the total attention and add a very small bias (ϵ) to avoid in-
 339 valid operations and vanishing gradient descent. Finally, we
 340 compute the attention output by multiplying the weighted
 341 attention by the value matrix (V). The **entity attn. features**
 342 becomes $\mathbf{A}_o \in \mathbb{R}^{B \times K \times d_v}$.

343 **Scene-level Prompt Attention** Focuses on overall activi-
 344 ties in the video over time using a summarized video feature
 345 and a scene text prompt. Next, an anomaly-attended video
 346 feature is created by the dot product between the anomaly
 347 scores and each frame feature in the video, with the follow-
 348 ing formula:
 349

$$\mathbf{A} = \left(\text{FC}(\mathbf{A}_o)^\top \cdot \mathbf{X}_g \right) \in \mathbb{R}^{B \times 1 \times d_v} \quad (6)$$

350 Next, we introduce learnable scene prompts to guide
 351 the scene-level feature learning [26, 31]. The initial scene
 352 prompts are anomaly class labels, such as "Fighting," "Ex-
 353 plosion," and "Abuse," etc. Next, we tokenize the class label
 354 and place the CLIP tokens between the learnable prefix and
 355 the postfix tokens. The complete scene prompt structure fol-
 356 lows: [CLS] $\underbrace{[P_1], \dots, [P_m]}_{\text{Prompt Prefix Tokens}}$ [T] $\underbrace{[P_{m+1}], \dots, [P_n]}_{\text{Prompt Postfix Tokens}}$ [EOS]. To
 357 obtain the scene embedding (T_s), we utilize the same pre-
 358 trained CLIP text encoder (ViT-B/16) that we employed
 359 during entity embedding. Next, we aggregate the scene
 360 embedding with anomaly attn. feature, and feed to a feed-
 361 forward-network (FFN), with operations in Eq. 7.
 362

$$\mathbf{S} = (\text{FFN}(A + T_s) + T_s) \quad (7)$$

363 Finally, we perform matrix multiplication on the normal-
 364 ized frame feature and textual scene attention to derive a
 365 similarity map between each frame and anomaly classes.
 366 The similarity matrix is calculated in Eq. 8.

$$\text{Sim}_{\text{mat}} = \left(\text{Norm}(\mathbf{X}_g) \text{Norm}(\mathbf{S})^\top \right) \in \mathbb{R}^{B \times K \times C} \quad (8)$$

367 3.5. Objective Functions

368 In our MLPVAD pipeline, we use three objective functions
 369 to optimize the model: (1) Coarse Classification, (2) Con-

372 trastive Learning, and (3) Fine-grained Classification. The
 373 functions are labeled with the circle shape in Figure 2.
374 Coarse Classification Determines whether a video is normal/abnormal.
 375 To perform MIL, we select top-K confidence frames from the anomaly video and bottom-K confidence
 376 frames from the normal video. Next, we perform mean ag-
 377 gregation and derive a single confidence value as the video-
 378 level confidence and feed it to the binary cross-entropy loss
 379 function. Eq. 9 outlines the loss calculation formula.
 380

381 $\mathcal{L}_{bce} = (\text{Sigmoid}(\text{FC}(A_o))) \in \mathbb{R}^{B \times 1}$ (9)

382 **Contrastive Learning** For challenging datasets, the differ-
 383 ence between normal and abnormal cases can be very subtle.
 384 We require efficient feature learning techniques to clas-
 385 sify subtle cases better. We propose using randomized
 386 contrastive learning with multiple positive and negative keys to
 387 achieve robust separation in feature space.

388 $\mathcal{L}_{\text{cons}} = \text{Mean} \left(-\log \frac{\sum_{i=1}^P \exp(\text{sim}(Q, K_i^+)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(Q, K_j^-)/\tau)} \right)$ (10)

389 The P and N are the number of positive and negative keys.
 390 The Query and Keys come from top-k/bottom-k confidence
 391 frames. Using contrastive learning, we can cluster all ab-
 392 normal frames together and separate normal frames from
 393 the abnormal frames in the feature space.

394 **Fine-Grained Classification** The imbalanced dataset
 395 makes accurate and unbiased VAD detection challenging.
 396 Also, actual anomaly frames may span a very short amount
 397 of time in the video. For multi-class anomaly classification,
 398 some videos might fall under more than one anomaly class.
 399 To address these challenges, we introduce fine-grained clas-
 400 sification based on a scene-guided frame-class alignment
 401 matrix. Similar to the coarse classification, we rely on the
 402 MIL approach, considering only the top-K candidates. The
 403 steps to calculate fine-grained loss are:

404
$$\mathbf{p}_i = \text{Softmax} \left(\frac{1}{k_i} \sum_{j=1}^{k_i} \text{TopK}(\mathbf{L}_i) \right) \in \mathbb{R}^C$$

 405
$$\mathcal{L}_{\text{ce}} = \text{Mean} \left(-\sum_{c=1}^C y_i^{(c)} \log p_i^{(c)} \right)$$
 (11)

406 The final objective function for the MLPVAD model is:

407
$$\mathcal{L}_{\text{total}} = \mathcal{L}_{bce} + \mathcal{L}_{\text{cons}} + \mathcal{L}_{\text{ce}}$$
 (12)

4. Experiments

408 **Datasets** are (1) UCF-Crime and (2) XD-Violence bench-
 409 marks. Both of the datasets have video-level labels (weakly
 410 labeled) for the training set and frame-level labels for the

Methods	UCF (AUC)	XD (AP)
Sultani et al. (2018)	84.06	74.78
Wu et al. (2020)	84.40	78.66
RTFM (2021)	85.22	78.27
AVVD (2022)	82.38	77.50
TEVAD (2023)	84.85	79.52
MGFN (2023)	86.67	80.11
VadCLIP (2024)	87.65	83.88
TPWNG (2024)	87.79	83.68
MLPVAD	88.84	85.37

Table 1. Coarse-grained comparison using AUC and AP on UCF-Crime and XD-Violence dataset.

411 testing set. The UCF-Crime CCTV surveillance videos
 412 cover 13 anomaly event categories, whereas XD-Violence
 413 covers six anomaly event categories from movie clips.

414 **Evaluation Metrics** Following earlier benchmarks, Av-
 415 erage Precision (AP) and Area Under the Curve (AUC)
 416 are used as classification metrics [6, 26]. Particularly for
 417 coarse-grained classification, we calculate frame-level AUC
 418 and frame-level AP for the dataset. For fine-grained classifi-
 419 cation, we calculate the mean Average Precision (mAP) un-
 420 der different intersection over union (IoU) thresholds rang-
 421 ing from 0.1 to 0.5.

422 **Implementation** is adopted from LAVIS [14] and VadCLIP
 423 [26] github repositories. For feature extraction, we util-
 424 ize the pre-trained CLIP (ViT-B/16 or ViT-B/14). Our im-
 425 plementation seamlessly supports both versions of the ViT
 426 encoder. The BLIP-2 Q-Former pretrained model is also
 427 adopted from LAVIS. The VadCLIP implementation serves
 428 as inspiration for the LGT Network implementation. In the
 429 CAFS module, the value of K is 256. The batch size is set to
 430 64, and the text and visual-feature dimensions are 768 and
 431 512, respectively. For contrastive learning, we select k = 8
 432 out of K (256) frames. For model training, a single NVIDIA
 433 RTX A6000 GPU is used with the PyTorch framework. For
 434 both datasets, the learning rate and # of epochs were set to
 435 2×10^{-4} and 20, respectively.

4.1. Comparison with State-of-the-art methods

436 In this section, we compare the proposed MLPVAD model
 437 with several of the latest Weakly-supervised models. To
 438 perform a fair comparison, we use the same visual and tex-
 439 tual encoders (ViT-B/16) and the same feature dimension
 440 for all SOTA methods. We also maintain the batch size and
 441 the number of epochs constant throughout the evaluation.
 442 Most existing SOTA methods downsize the video input to a
 443 fixed k=256 frames for training. So, we also select the Top
 444 k = 256 frames from the CAFS modules to match the per-
 445 formance comparison with other SOTA methods. Moving
 446 forward, we use K=256 for all experiments, except in the

mAP@IOU %	0.1	0.2	0.3	0.4	0.5
Method ↓	Dataset: UCF-Crime				
Sultani et al.	5.72	4.41	2.69	0.02	0.01
AVVD	10.27	7.01	6.25	3.42	3.29
VadCLIP	11.72	7.83	6.40	4.53	2.93
MLPVAD	13.21	9.74	7.78	6.20	4.11
Method ↓	Dataset: XD-Violence				
Sultani et al.	22.72	15.57	9.98	6.20	3.78
AVVD	30.51	25.75	20.18	14.83	9.79
VadCLIP	37.03	30.84	23.38	17.90	14.31
MLPVAD	38.55	32.63	25.82	19.39	15.60

Table 2. Fine-grained comparisons on XD-Violence and UCF-Crime.

CAFS	Vis-P	Scn.-P	Ent.-P	C-Lrn	AP(%)
Baseline	✓	✓			84.53
✓	✓	✓			86.25
✓	✓		✓		85.33
	✓	✓	✓		86.56
✓	✓	✓	✓		88.20
✓	✓	✓	✓	✓	88.84

Table 3. Ablation study on UCF-Crime dataset for different MLPVAD modules.

ablation study, where different values of K are examined.

Coarse-grained WSVAD Performance Table 1 presents the coarse-grained performance comparison with several SOTA methods. The coarse classification based on visual data performs worse than the MLPVAD method [21, 25, 26]. Earlier works, such as [21, 25], showed decent AUC gains of 85.22%, and 84.57%, respectively, on UCF-Crime. Whereas our MLPVAD method achieved an AUC of 88.84%, outperforming its closest competitors, such as VadCLIP (87.65%) and TPWNG (87.79%).

MLPVAD also outperforms the state of the art for the XD-Violence dataset. Most previous works, such as [6], [21], and [7], achieve AP close to 80%. VadCLIP and TPWNG achieved a noticeable improvement by utilizing a more effective temporal network and feature learning, resulting in an AP of 83.88%. Our efficient frame selection technique and entity-guided approach outperformed existing SOTA methods with a margin of $+\Delta 1.5\%$ and gained an AP of 85.37%.

Fine-grained WSVAD Performance. In Table 2, the fine-grained classification is evaluated using mAP at different IOU thresholds. We compare the fine-grained MLPVAD performance with other SOTA work that introduced fine-grained analysis. Fine-grained classification was first pro-

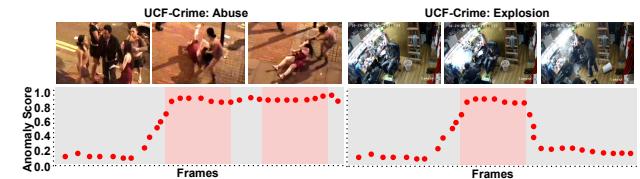


Figure 5. Frame-level anomaly detection using the MLPVAD model for the UCF-Crime dataset.

posed by AVVD, later fine-tuned by Sultani et al. and VadCLIP with CLIP features. Our work modifies the fine-grained pipeline, integrating scene attention and task-centric anomaly attention features guided by entity prompts. Compared to AVVD and VadCLIP, MLPVAD achieves a performance improvement of 8% and 1.5% on the UCF-Crime dataset (IOU = 0.1). MLPVAD also showed a decent performance gain of 13.21% mAP (IOU=0.1) in challenging the XD-Violence dataset.

4.2. Ablation Study

In this section, we try to answer several questions regarding the effectiveness of the proposed modules in the MLPVAD pipeline. The ablation study begins with the baseline components in the pipeline (See Table 3) for the UCF-Crime dataset. We start with textual and visual data in the baseline, as our MLPVAD focuses more on utilizing multimodal features for improved performance. The baseline model achieves an AP of 84.53%. Next, we incorporate the proposed CAFS module for task-specific frame selection, and we observed a significant gain of $+\Delta 1.75\%$. It proves that careful frame selection is very crucial for downstream tasks. Next, we turn off scene prompt learning and use only the entity-visual prompt with the CAFS module. We found that entity learning alone outperforms scene learning, achieving an AP of 85.33%. In the next two experiments, we use visual, entity, and scene prompts with/without the CAFS module. The pipeline with CAFS achieved a $+\Delta 1.64\%$ improvement compared to the one without it, demonstrating the effectiveness of the CAFS module. Finally, we introduce randomized contrastive learning in the pipeline and optimize the model for better feature separation. Using multiple positive and negative keys for contrastive learning, MLPVAD achieved an AP of **88.84%** on the benchmark dataset.

Effectiveness based on video duration The duration of the video is a very important factor to consider for any video-related downstream task. In Table 4, we present a performance comparison between MLPVAD and other SOTA methods, dividing the test videos into short ($< 15mins$) and long ($\geq 15mins$) categories. We pick *15 minutes* as the threshold for a balanced evaluation. The UCF-Crime dataset contains 1,795 short and 112 long videos. It is ev-

Methods	Short Videos (< 15mins)	Long Videos (>= 15mins)
TEVAD	83.76	77.35
MGFN	89.47	84.75
VadCLIP	90.22	85.30
TPWNG	90.08	84.90
MLPVAD	91.48	87.48

Table 4. Duration-based performance comparison using AUC on UCF-Crime dataset.

Frame Select.	UCF AUC(%)	XD AP(%)
Stepwise	85.78	82.10
Mean-Pool	87.30	83.88
CAFS (K=128)	88.25	84.68
CAFS (K=256)	88.84	85.37
CAFS (K=320)	89.06	85.60

Table 5. Effectiveness of different frame selection methods and performance analysis for CAFS with different values of K.

514 ident from Table 4 that methods perform better for short
 515 videos than long videos. The frame selection and tempo-
 516 ral modeling are more effective when the duration is short.
 517 The MLPVAD method still outperforms other methods by
 518 at least a minimum margin of $+ \Delta 1.3\%$ for the UCF-Crime
 519 dataset. For long videos, our MLPVAD demonstrates its su-
 520 periority through efficient frame selection and a multi-level
 521 prompt pipeline. It achieves 87.48% AUC on UCF-Crime
 522 and beats the closest competitors, such as VadCLIP and TP-
 523 WNG, by a large margin of $\Delta 2\%$.

524 **Ablation on different frame selection** Table 5 illustrates
 525 the MLPVAD evaluation results when trained using differ-
 526 ent frame selection methods. The performance drops by
 527 3% when we use traditional interval-based frame selection.
 528 The mean-pool method shows a decent performance gain
 529 of $\sim 2\%$ for both experimental datasets. The technique
 530 still suffers when the video length increases (See Table 4),
 531 as more noise is introduced in frames during the mean-
 532 pooling operations. We show the CAFS effectiveness for
 533 three different values of K. CAFS performs better than most
 534 of the SOTA methods, even with half the amount of tempo-
 535 ral frames. With K=256, the CAFS stands out as the super-
 536 ior method, achieving $\sim 3\%$ gain in each dataset. Further
 537 increasing K to 320 improves performance by $+ \Delta 0.3\%$;
 538 however, it also requires additional GPU memory and com-
 539 putational power.

540 **Qualitative Results** Figure 5 and Figure 6 illustrate the
 541 frame-level detection result for a few videos with the
 542 anomaly score. The red dots in the Figure denote the
 543 anomaly scores for a particular frame. The anomaly scores

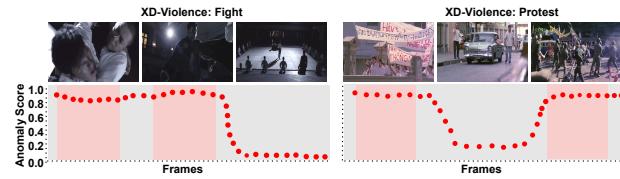


Figure 6. Frame-level anomaly detection using the MLPVAD model for the XD-Violence dataset.

exhibit a sharp increase during the presence of anomalous
 544 events, followed by a rapid decline once these events
 545 conclude. In contrast, normal events consistently yield lower
 546 anomaly scores over time. This behavior indicates that
 547 our approach is highly sensitive to abnormal activities, ef-
 548 fectively identifying them promptly while preserving low
 549 anomaly predictions during normal scenarios.

5. Conclusion and Future Work

In this work, we propose a new paradigm, named MLP-
 552 VAD, for the WSVAD task, focusing on long-duration
 553 videos. We propose an elegant technique for selecting a
 554 limited number of task-specific frames from a larger pool.
 555 We utilize a frozen Image-Text Matching CLIP Network
 556 to perform text-guided frame selection. Next, we pro-
 557 pose a multi-level prompt learning architecture that follows
 558 human-like cognitive processes. We leverage the power
 559 of the latest LLMs and CLIP models to guide the learn-
 560 ing at the entity and scene levels. Finally, we introduce
 561 MIL learning with three objective functions: coarse and
 562 fine-grained classification, and contrastive learning. We
 563 empirically verify the effectiveness of MLPVAD modules
 564 through SOTA comparison and extensive ablations on WS-
 565 VAD benchmarks. The ablation study also verifies the su-
 566 periority of MLPVAD for long-duration VAD. In the future,
 567 we plan to improve the VAD performance for cases where
 568 the abnormality is very subtle and exhibits high interclass
 569 similarity. We also plan to extend the performance anal-
 570 ysis for different LLM models for caption generation and
 571 verify the LLMs' hallucination for challenging and subtle
 572 scenes.

574 References

- [1] Armstrong Aboah. A vision-based system for traffic anomaly detection using deep learning and decision trees. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4207–4212, 2021. 1
- [2] Jungmo Ahn, JaeYeon Park, Sung Sik Lee, Kyu-Hyuk Lee, Heesung Do, and JeongGil Ko. Safefac: Video-based smart safety monitoring for preventing industrial work accidents. *Expert Systems with Applications*, 215:119397, 2023. 1
- [3] Shehreen Azad, Vibhav Vineet, and Yogesh Singh Rawat. Hierarq: Task-aware hierarchical q-former for enhanced video understanding. In *Proceedings of the Computer Vi-*

- 586 sion and Pattern Recognition Conference, pages 8545–8556,
 587 2025. 3
- 588 [4] Debojoyoti Biswas and Jelena Tesic. Mmvad: A vision-
 589 language model for cross-domain video anomaly detection
 590 with contrastive learning and scale-adaptive frame segmenta-
 591 tion. *Expert Systems with Applications*, page 127857, 2025.
 592 1
- 593 [5] Joao Carreira and Andrew Zisserman. Quo vadis, action
 594 recognition? a new model and the kinetics dataset. In *pro-
 595 ceedings of the IEEE Conference on Computer Vision and
 596 Pattern Recognition*, pages 6299–6308, 2017. 1, 3
- 597 [6] Weiling Chen, Keng Teck Ma, Zi Jian Yew, Minhoe Hur, and
 598 David Aik-Aun Khoo. Tevad: Improved video anomaly de-
 599 tection with captions. In *Proceedings of the IEEE/CVF Con-
 600 ference on Computer Vision and Pattern Recognition*, pages
 601 5548–5558, 2023. 1, 3, 4, 6, 7
- 602 [7] Yingxian Chen, Zhengzhe Liu, Baoheng Zhang, Wilton
 603 Fok, Xiaojuan Qi, and Yik-Chung Wu. Mgfn: Magnitude-
 604 contrastive glance-and-focus network for weakly-supervised
 605 video anomaly detection. In *Proceedings of the AAAI Con-
 606 ference on Artificial Intelligence*, pages 387–395, 2023. 1,
 607 7
- 608 [8] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin
 609 Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang
 610 Luo, Deli Zhao, and Lidong Bing. Videollama 2: Advancing
 611 spatial-temporal modeling and audio understanding in video-
 612 llms. *arXiv preprint arXiv:2406.07476*, 2024. 4
- 613 [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina
 614 Toutanova. Bert: Pre-training of deep bidirectional trans-
 615 formers for language understanding. In *Proceedings of the
 616 2019 conference of the North American chapter of the asso-
 617 ciation for computational linguistics: human language tech-
 618 nologies, volume 1 (long and short papers)*, pages 4171–
 619 4186, 2019. 4
- 620 [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
 621 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
 622 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-
 623 vain Gelly, et al. An image is worth 16x16 words: Trans-
 624 formers for image recognition at scale. *arXiv preprint
 625 arXiv:2010.11929*, 2020. 1
- 626 [11] Jia-Chang Feng, Fa-Ting Hong, and Wei-Shi Zheng. Mist:
 627 Multiple instance self-training framework for video anomaly
 628 detection. In *Proceedings of the IEEE/CVF conference
 629 on computer vision and pattern recognition*, pages 14009–
 630 14018, 2021. 1, 2
- 631 [12] Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xue-
 632 fei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam
 633 Lim. Ma-lmm: Memory-augmented large multimodal model
 634 for long-term video understanding. In *Proceedings of the
 635 IEEE/CVF Conference on Computer Vision and Pattern
 636 Recognition*, pages 13504–13514, 2024. 3, 4
- 637 [13] Noussaiba Jaafar and Zied Lachiri. Multimodal fusion meth-
 638 ods with deep neural networks and meta-information for ag-
 639 gression detection in surveillance. *Expert Systems with Ap-
 640 plications*, 211:118523, 2023. 1
- 641 [14] Dongxu Li, Junnan Li, Hung Le, Guangsen Wang, Silvio
 642 Savarese, and Steven C.H. Hoi. LAVIS: A one-stop library
 for language-vision intelligence. In *Proceedings of the 61st
 643 Annual Meeting of the Association for Computational Lin-
 644 guistics (Volume 3: System Demonstrations)*, pages 31–41,
 645 Toronto, Canada, 2023. Association for Computational Lin-
 646 guistics. 6
- 647 [15] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.
 Blip-2: Bootstrapping language-image pre-training with
 648 frozen image encoders and large language models. In *In-
 649 ternational conference on machine learning*, pages 19730–
 650 19742. PMLR, 2023. 2, 3
- 651 [16] Yuqi Lin, Minghao Chen, Wenxiao Wang, Boxi Wu, Ke
 652 Li, Binbin Lin, Haifeng Liu, and Xiaofei He. Clip is also
 653 an efficient segmenter: A text-driven approach for weakly
 654 supervised semantic segmentation. In *Proceedings of the
 655 IEEE/CVF Conference on Computer Vision and Pattern
 656 Recognition*, pages 15305–15314, 2023. 3
- 657 [17] Hui Lv, Zhongqi Yue, Qianru Sun, Bin Luo, Zhen Cui, and
 658 Hanwang Zhang. Unbiased multiple instance learning for
 659 weakly supervised video anomaly detection. In *Proceedings
 660 of the IEEE/CVF conference on computer vision and pattern
 661 recognition*, pages 8022–8031, 2023. 2
- 662 [18] Övgü Özdemir and Erdem Akagündüz. Enhancing visual
 663 question answering through question-driven image captions
 664 as prompts. In *Proceedings of the IEEE/CVF Conference
 665 on Computer Vision and Pattern Recognition*, pages 1562–
 666 1571, 2024. 3
- 667 [19] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
 668 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
 669 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
 670 transferable visual models from natural language supervi-
 671 sion. In *International conference on machine learning*, pages
 672 8748–8763. PMLR, 2021. 1
- 673 [20] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world
 674 anomaly detection in surveillance videos. In *Proceedings of
 675 the IEEE conference on computer vision and pattern recog-
 676 nition*, pages 6479–6488, 2018. 2
- 677 [21] Yu Tian, Guansong Pang, Yuanhong Chen, Rajvinder Singh,
 678 Johan W Verjans, and Gustavo Carneiro. Weakly-supervised
 679 video anomaly detection with robust temporal feature magni-
 680 tude learning. In *Proceedings of the IEEE/CVF international
 681 conference on computer vision*, pages 4975–4986, 2021. 2,
 682 7
- 683 [22] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani,
 684 and Manohar Paluri. Learning spatiotemporal features with
 685 3d convolutional networks. In *Proceedings of the IEEE inter-
 686 national conference on computer vision*, pages 4489–4497,
 687 2015. 3
- 688 [23] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,
 689 Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin
 690 Ge, et al. Qwen2-vl: Enhancing vision-language model's
 691 perception of the world at any resolution. *arXiv preprint
 692 arXiv:2409.12191*, 2024. 4
- 693 [24] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He,
 694 Guo Chen, Baoqi Pei, Rongkun Zheng, Zun Wang, Yansong
 695 Shi, et al. Internvideo2: Scaling foundation models for mul-
 696 timodal video understanding. In *European Conference on
 697 Computer Vision*, pages 396–416. Springer, 2024. 4

- 700 [25] Peng Wu, Jing Liu, Yujia Shi, Yujia Sun, Fangtao Shao,
 701 Zhaoyang Wu, and Zhiwei Yang. Not only look, but also
 702 listen: Learning multimodal violence detection under weak
 703 supervision. In *Computer Vision–ECCV 2020: 16th Euro-*
 704 *pean Conference, Glasgow, UK, August 23–28, 2020, Pro-*
 705 *ceedings, Part XXX 16*, pages 322–339. Springer, 2020. 7
- 706 [26] Peng Wu, Xuerong Zhou, Guansong Pang, Lingru Zhou,
 707 Qingsen Yan, Peng Wang, and Yanning Zhang. Vadclip:
 708 Adapting vision-language models for weakly supervised
 709 video anomaly detection. In *Proceedings of the AAAI Con-*
 710 *ference on Artificial Intelligence*, pages 6074–6082, 2024. 1,
 711 2, 3, 4, 5, 6, 7
- 712 [27] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko,
 713 Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and
 714 Christoph Feichtenhofer. Videoclip: Contrastive pre-training
 715 for zero-shot video-text understanding. *arXiv preprint*
 716 *arXiv:2109.14084*, 2021. 3
- 717 [28] Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel
 718 Jin, Chris Callison-Burch, and Mark Yatskar. Language
 719 in a bottle: Language model guided concept bottlenecks
 720 for interpretable image classification. In *Proceedings of*
 721 *the IEEE/CVF Conference on Computer Vision and Pattern*
 722 *Recognition*, pages 19187–19197, 2023. 3
- 723 [29] Luca Zanella, Benedetta Liberatori, Willi Menapace, Fabio
 724 Poiesi, Yiming Wang, and Elisa Ricci. Delving into clip la-
 725 tent space for video anomaly recognition. *Computer Vision*
 726 *and Image Understanding*, 249:104163, 2024. 3
- 727 [30] Jia-Xing Zhong, Nannan Li, Weijie Kong, Shan Liu,
 728 Thomas H Li, and Ge Li. Graph convolutional label noise
 729 cleaner: Train a plug-and-play action classifier for anomaly
 730 detection. In *Proceedings of the IEEE/CVF conference on*
 731 *computer vision and pattern recognition*, pages 1237–1246,
 732 2019. 2, 4
- 733 [31] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei
 734 Liu. Learning to prompt for vision-language models. *In-*
 735 *ternational Journal of Computer Vision*, 130(9):2337–2348,
 736 2022. 5