# Towards Adaptive Human-centric Video Anomaly Detection: A Comprehensive Framework and A New Benchmark

Armin Danesh Pazho*[†1], Shanle Yao*[1], Ghazal Alinezhad Noghre*[1], Babak Rahimi Ardabili[2], Vinit Katariya[1], Hamed Tabkhi[1]

*Abstract*—Human-centric Video Anomaly Detection (VAD) aims to identify human behaviors that deviate from normal. At its core, human-centric VAD faces substantial challenges, such as the complexity of diverse human behaviors, the rarity of anomalies, and ethical constraints. These challenges limit access to high-quality datasets and highlight the need for a dataset and framework supporting continual learning. Moving towards adaptive human-centric VAD, we introduce the HuVAD (Human-centric privacy-enhanced Video Anomaly Detection) dataset and a novel Unsupervised Continual Anomaly Learning (UCAL) framework. UCAL enables incremental learning, allowing models to adapt over time, bridging traditional training and real-world deployment. HuVAD prioritizes privacy by providing de-identified annotations and includes seven indoor/outdoor scenes, offering over $5\times$ more pose-annotated frames than previous datasets. Our standard and continual benchmarks, utilize a comprehensive set of metrics, demonstrating that UCAL-enhanced models achieve superior performance in $82.14\%$ of cases, setting a new state-of-the-art (SOTA). The dataset can be accessed at https://github.com/TeCSAR-UNCC/HuVAD.

*Index Terms*—Video Processing, Anomaly Detection, Dataset, Adaptive Learning

## I. Introduction

Human-centric Video Anomaly Detection (VAD) refers to identifying events or patterns in human behavior that deviate from the expected behavior. Advancing human-centric VAD technologies encounter several substantial challenges emerging from its context-specific and open-set nature, where anomalies vary across different environments and new unseen events regularly happen, making it difficult to generalize VAD algorithms [1]–[3]. Furthermore, anomalous events are, by definition, rare, creating a scarcity of positive samples in datasets, making it challenging for models to learn solely from labeled examples, pushing the field towards unsupervised approaches [4], [5].

To adjust models to tackle such challenges, and adapt to real-world scenarios, research is moving towards continual learning [4], [6]–[9]. Especially, context-specificity and open-set nature of human-centric VAD underscore a critical need for a new framework and dataset to address the data limitation and support continual learning [10], [11]. Such a framework would allow VAD systems to adapt to new patterns over time, thereby improving detection accuracy in dynamic, real-world

* Equal contribution.
† Corresponding author (adaneshp@charlotte.edu).
[1] Department of Electrical and Computer Engineering, University of North Carolina at Charlotte ([syao, galinezh, vkatariy, htabkhiv]@charlotte.edu
[2] Public Policy Program, University of North Carolina at Charlotte (brahimia@charlotte.edu)

environments. As a result, addressing these challenges requires a holistic approach to data curation and framework design.

To support extensive data access and facilitate continual learning, we introduce HuVAD, the largest continuously recorded VAD dataset (see samples in fig. 1). Effective continual learning requires datasets that mirror real-world surveillance, capturing diverse environments with ongoing recording and ample samples. HuVAD provides over 5x the training frames and 4x the testing frames compared to previous datasets, supporting unsupervised learning and continual model adaptation [12]–[14]. It maintains real-world fidelity with continuous recordings of commuters' real-world traffic data across seven scenes, including indoor/outdoor spaces. While 6 cameras capture typical environments within a community space with, a specialized Context-Specific Camera (CSC) observes law enforcement training. This unique setup allows HuVAD to evaluate model performance in scenarios where normal behaviors differ from typical public environments. To address potential privacy concerns and biases against minority groups [15]–[18], HuVAD employs anonymization, publishing only de-identified human annotations, such as bounding boxes, tracking IDs, and poses [1], [2], [19].

With HuVAD providing an ideal environment for continual learning, we introduce a novel Unsupervised Continual Anomaly Learning (UCAL) framework. UCAL tackles VAD's context-specificity and open-set nature by incrementally adapting the model for each environment in the HuVAD dataset, allowing it to capture and evolve with dynamic patterns. The framework simulates real-world streaming data by organizing data per camera and randomly injecting anomalies into the training stream. Structured in multiple incremental training steps, UCAL leverages a pre-trained model to refine its performance progressively, with continual evaluation against the HuVAD test set.

For a thorough benchmark, we statistically compare HuVAD with peer datasets and evaluate state-of-the-art (SOTA) pose-based VAD algorithms. Unlike prior works that rely on a single metric [19]–[25], we evaluate these algorithms using a comprehensive set of metrics, including Area Under the Receiver Operating Characteristic Curve (AUC-ROC), Area Under the Precision-Recall Curve (AUC-PR), Equal Error Rate (EER), and the 10% Error Rate (10ER), a metric brought to human-centric VAD for the first time to reflect the unequal costs of false positives and negatives in real-world scenarios. Recognizing the importance of continual learning, we conduct a tailored benchmark of the UCAL framework on the HuVAD
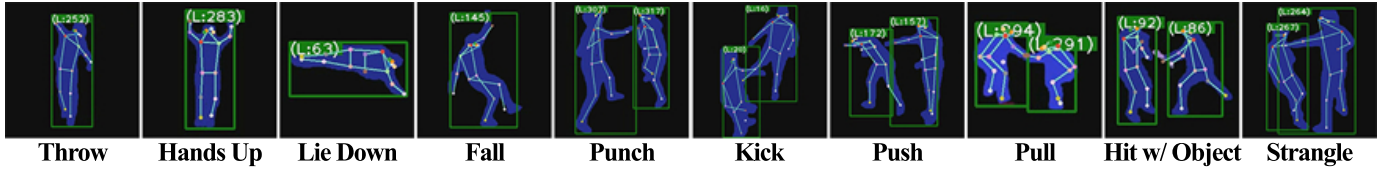
Fig. 1: Sample of anomalies and their annotations in the new proposed benchmark: HuVAD dataset. Cropped for visualization purposes. Segmentation is solely used for demonstration purposes.

dataset. Results demonstrate that models enhanced by the UCAL framework achieve superior performance in 82.14% of cases, setting a new SOTA in human-centric VAD.

The contributions of this paper are:

- HuVAD, the largest continuously recorded dataset in real-world community spaces as well as a especial context-specific law enforcement training environment, emphasizing the impact of context awareness and social interaction on VAD, with comprehensive de-identified human annotations to ensure privacy.
- UCAL, a novel Unsupervised Continual Anomaly Learning framework, enabling adaptive learning for VAD by allowing models to evolve continuously with new data.
- Providing detailed statistical comparisons, conventional algorithmic evaluations, and offering the first comprehensive benchmark for continual learning in human-centric VAD, showcasing SOTA results of UCAL-enhanced models in 82.14% of cases.

## II. RELATED WORKS

### A. Video Anomaly Detection Dataset

In this paper, we focus on unsupervised VAD, which differs fundamentally from weakly-supervised tasks [26] addressed by datasets such as UCF-Crime [27] and XD-Violence [28]. Additionally, UBnormal [29], composed entirely of synthetic videos, may not reliably reflect real-world environments, limiting its relevance for unsupervised detection. Datasets such as UCF-Crime, XD-Violence, MSAD [30], and CUVA [31], compiled from varied sources rather than continuous recordings, reflect a different task formulation and are thus not comparable to our approach [32].

A few recent datasets feature anomalies for both vehicles and pedestrians. Street Scene [32] dataset is distinguished by its near bird's eye view perspective of a street, and includes non-human anomalies such as illegally parked cars. The NOLA dataset [8] with over 1.4 million frames is captured from a single moving camera. In addition to human-centric anomalies, it includes anomalies such as a vehicle moving in the wrong direction.

Several early datasets started the field of VAD such as the Subway dataset [33], the UCSD Pedestrian dataset [34], and the CUHK Avenue dataset [35]. While pivotal for VAD's early development, these datasets are now considered relatively small and feature a limited variety of scenes.

More recent datasets have emerged to advance the field of VAD. The ShanghaiTech Campus (SHT) dataset [36] stands as a primary benchmark within the realm of VAD, particularly for human-centric approaches. It includes anomalies such as

chases and fights. A shortcoming of the dataset is that the frames per camera are limited, posing challenges for continual learning. The IITB dataset [37] captures human activities within a corridor, recorded using a single fish-eye camera. The CHAD dataset [5] is a large-scale VAD dataset recorded within a parking lot setting. CHAD features 22 anomaly classes and is captured using four high-resolution cameras at 30 FPS. The NWPUC dataset [26] was developed to introduce greater diversity in a campus setting, featuring 28 classes of anomalies across 43 distinct scenes. However, the dataset's per-scene volume remains relatively limited, which may constrain its comprehensiveness for certain large-scale anomaly detection tasks such as continual learning.

### B. Continual Learning

Continual learning has been applied to anomaly detection across various fields, addressing challenges such as data drift and evolving data streams. Bugarin et al. [38] introduce a benchmark for anomaly detection in industrial images, while Chavan et al. [39] propose a framework to handle data drift in industrial settings. Mozaffari et al. [40] develop an approach for high-dimensional data streams in cybersecurity, and Mazarbhuiya and Shenify [41] focus on clustering for real-time anomaly detection in IoT. Together, these works showcase continual learning's potential to enhance anomaly detection in complex, evolving environments.

Initial efforts to integrate continual learning into VAD are led by Doshi and Yilmaz, who developed several frameworks enhancing VAD through continual learning. Their 2020 work [6] combines transfer learning with k-nearest neighbor (kNN) techniques for low-complexity, continual anomaly detection. In 2022 [8] they focus on pedestrian and vehicle VAD with continual and few-shot learning to improve detection delay and alarm precision. Their 2021 MONAD framework [7] integrates GAN-based prediction with statistical guarantees on false alarms, ensuring robustness. Although these works contribute significantly to adaptive VAD, progress can accelerate further with the establishment of a unified and comprehensive benchmark, particularly in continual human-centric VAD.

## III. PRIVACY AND ETHICAL CONSIDERATIONS

Integrating computer vision algorithms into various sectors of society has underscored the importance of responsibly developing these technologies [42]. VAD and Smart Video Surveillance (SVS) touch upon critical ethical concerns such as bias, discrimination, and privacy violations [43]. These concerns have practical implications for the performance and

Fig. 2: The camera views excluding people. The ratio has been adjusted to fit the manuscript.
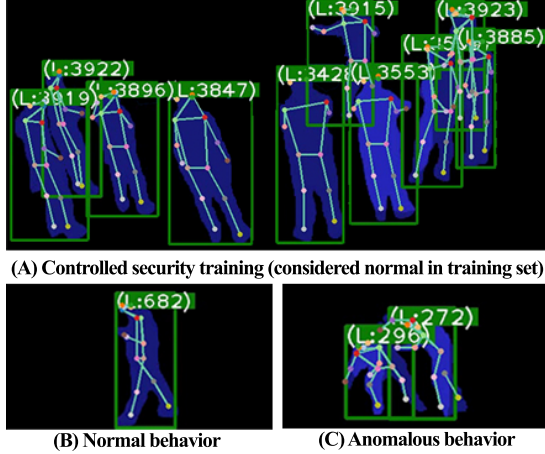


Fig. 3: Samples from CSC Camera. Segmentation is solely used for demonstration purposes.

fairness of computer vision systems in the real world. Bias in VAD algorithms can lead to discriminatory practices, where specific demographics may be unfairly targeted or misrepresented [44], [45]. Privacy concerns are equally significant, as the pervasive monitoring and analysis of individuals without safeguards can lead to unwarranted privacy invasions, raising ethical and legal issues [16].

Addressing these privacy and ethical concerns necessitates a multifaceted approach. Such considerations are particularly crucial during the dataset creation [46] and algorithm development [16]. Adopting more abstract representation techniques, such as only focusing on human pose information, has been proposed as an intermediate solution to address some privacy, biases, and ethical concerns if not fully mitigating them [16]. Our interaction with communities also proves that the public and stakeholders are more responsive and participatory to dataset collection when focusing solely on human pose information rather than actual pixels. However, transitioning to a more abstract approach (e.g. human pose rather than pixel-level data) raises questions regarding the potential compromise in model accuracy. Recent advancements challenge this notion. On the SHT dataset [36], pixel-based methods like SSMTL++v2 [47] and Jigsaw-VAD [48] achieve AUC-ROC scores of 83.80 and 84.30, while pose-based approaches like MoPRL [20] and STG-NF [19] attain AUC-ROC scores of 83.35 and 85.90.

Our objective is to initiate the trend of de-identifying PII at the data collection phase instead of the algorithmic phase. This methodology not only addresses legal challenges but also promotes expedited data gathering and dissemination. Such an approach is anticipated to facilitate research and technological advancements more finely attuned to the complexities of human behavior and societal requirements.

## IV. DATA COLLECTION AND ANNOTATION

The HuVAD dataset was captured using seven Closed-Circuit Television (CCTV) cameras at a resolution of $1280 \times 720$ pixels, recording over five days from 6:00 AM to 6:00 PM. It includes varied scenes: three outdoor parking areas, three hallways, and a building entrance/exit (see fig. 2). Six cameras (C0 to C5) cover normal environments, while a context-specific camera (CSC) focuses on an outdoor area for security and law enforcement training (fig. 3). HuVAD includes anomalies such as throwing, hands up, lying down, falling, punching, kicking, pushing, and strangling (see fig. 1), enhancing real-world applicability by focusing on behaviors relevant to public safety while excluding less relevant actions such as jumping.

### A. Annotation Methodology

**Anomaly Annotation:** The HuVAD dataset has been meticulously annotated for anomalies to ensure high precision, with each frame reviewed by at least three trained annotators and random testing for label consistency. Anomalous behaviors from section IV are marked abnormal on cameras C0 to C5, while routine activities are labeled as normal. For the CSC camera, if these actions occur during a controlled training session, they are labeled as normal; otherwise, they are labeled as anomalies. The dataset includes frame-level anomaly labels and detailed spatial annotations through region masks that localize anomalies, enabling granular analysis of events. Following the standard approach in unsupervised and self-supervised video anomaly detection benchmarks such as SHT [36], IITB [37], CHAD [5], and NWPUC [26], we focus solely on the normal vs. anomalous distinction without class-specific annotations.

**Peson Annotations:** For realistic anomaly detection, HuVAD simulates real-world reliance on algorithmically extracted person data. Using a semi-automated approach, human annotations are first generated by models, followed by manual human verification to ensure accuracy.

The HuVAD dataset provides comprehensive de-identified annotations [49], including bounding boxes, person IDs, and pose annotations. Bounding boxes, generated with YOLOv8 [50], mark detected individuals' positions. Person IDs, extracted with ByteTrack [51], ensure temporal consistency and identity preservation across frames. Human poses, generated with HRNet [52] in the COCO17 format [53], capture body keypoints, with linear interpolation and a 15-frame smoothing window for quality assurance, and interpolated poses flagged for optional exclusion.
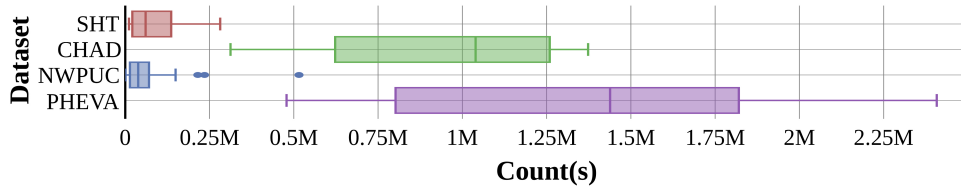
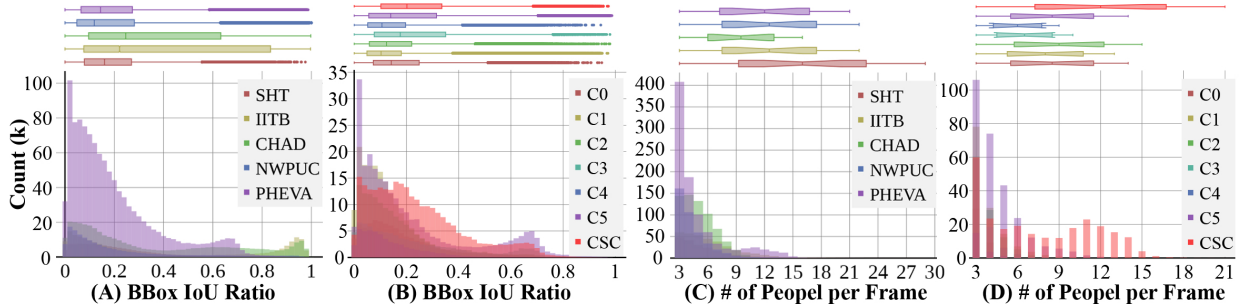Fig. 4: Pose counts per camera across major multi-camera datasets (SHT [36], CHAD [5], and NWPUC [26]).



Fig. 5: IoU and crowd density across key datasets (SHT [36], IITB [37], CHAD [5], and NWPUC [26]) and HuVAD's camera views.

| Dataset | Conference | Total | Train | Test | Test Normal | Test Anomaly | Camera(s) |
|---------|-----------|-------|-------|------|-------------|--------------|-----------|
| SHT | CVPR 18 | 295,495 | 257,650 | 37,845 | 21,141 | 16,704 | 13 |
| IITB | WACV 20 | 459,341 | 279,880 | 179,461 | 71,316 | 108,145 | 1 |
| CHAD | SCIA 23 | 922,034 | 802,167 | 119,867 | 60,969 | 58,898 | 4 |
| NWPUC | CVPR 23 | 1,000,129 | 715,901 | 284,228 | 235,957 | 48,271 | 43 |
| **HuVAD (Ours)** | - | 5,196,675 | 4,467,271 | 729,404 | 517,286 | 212,118 | 7 |

TABLE I: Number of frames with at least one human in major datasets (SHT [36], IITB [37], CHAD [5], and NWPUC [26]).

## B. Data Statistics

HuVAD is the largest continuous VAD dataset to date, with over five million frames (~46 hours) of annotated pose sequences. As shown in fig. 4, HuVAD's per-camera pose counts significantly exceed those of other datasets, providing a robust foundation for developing advanced learning methods crucial for real-world anomaly detection. This extensive data volume per environment directly supports Unsupervised Continual Anomaly Detection (UCAL), discussed in section VII.

In fig. 5.A and fig. 5.B, the distribution of the highest Intersection over Union (IoU) ratio per frame highlights HuVAD's occlusion characteristics, offering a challenging environment comparable to NWPUC [26] and SHT [36]. Per camera, stats show a high median level of occlusion, particularly in the CSC camera. This diversity in occlusion levels ensures a challenging benchmark for models. fig. 5.C and fig. 5.D illustrate that HuVAD provides a wide range of crowd densities, from 3 to 15 individuals per frame, with higher densities in CSC due to capturing large-group security exercises. While CHAD [5] and SHT [36] show variable crowdedness, HuVAD's density distribution is more comprehensive, offering a broad benchmark for crowd-based anomaly detection tasks.

## V. METRICS

AUC-ROC, AUC-PR, EER, and 10ER each provide unique insights and limitations, making their combined use essential for a comprehensive assessment. **AUC-ROC** measures class distinction ability by plotting the True Positive Rate (TPR) against False Positive Rate (FPR) across thresholds; however, it lacks consideration of the False Negative Rate (FNR), is sensitive to data imbalance, and may obscure key trade-offs [55]–[57]. **AUC-PR** calculates the area under the Precision-Recall curve, which better handles imbalanced datasets but falls short in analyzing negative predictions and overall error distribution [57], [58]. **EER** identifies the threshold where FPR and FNR are equal, offering a balance of sensitivity (detecting anomalies) and specificity (recognizing normals), valuable in real-world deployments where the costs of false positives and negatives are balanced [27], [59]. Inspired by the False Match Rate 100 (FMR100) metric [60], [61], this paper brings **10ER** to VAD, measuring the FPR at a fixed 10% FNR—a threshold regarded as acceptable for VAD [62]. 10ER provides a practical, threshold-based perspective and when combined with other metrics, offers a more robust evaluation of VAD performance.

## VI. HUVAD STANDARD BENCHMARKS (HUVAD-S)

To align with the conventional anomaly detection training paradigm, we introduce HuVAD-S, constructed similarly to existing datasets [5], [26], [36], [37]. This subset serves as a foundation for the primary goal of this paper: advancing

| Model | Conference | AUC-ROC | AUC-PR | EER | 10ER | AUC-ROC | AUC-PR | EER | 10ER |
|---|---|---|---|---|---|---|---|---|---|
| | | C0 | | | | C1 | | | |
| MPED-RNN | CVPR 19 | **79.57** | 46.76 | **0.26** | **0.37** | 83.57 | 53.62 | **0.22** | **0.39** |
| GEPC | CVPR 20 | 59.07 | 28.00 | 0.44 | 0.71 | 56.27 | 23.20 | 0.45 | 0.78 |
| STG-NF | ICCV 23 | 58.96 | **83.45** | 0.46 | 0.84 | 47.82 | **78.31** | 0.52 | 0.89 |
| TSGAD | WACV 24 | 64.18 | 31.93 | 0.40 | 0.66 | 68.88 | 39.02 | 0.35 | 0.72 |
| | | C2 | | | | C3 | | | |
| MPED-RNN | CVPR 19 | **73.36** | 47.66 | **0.32** | **0.57** | 83.62 | 63.87 | **0.23** | **0.42** |
| GEPC | CVPR 20 | 55.09 | 30.13 | 0.45 | 0.79 | 52.40 | 27.09 | 0.50 | 0.77 |
| STG-NF | ICCV 23 | 51.06 | **74.20** | 0.49 | 0.94 | 49.15 | **74.10** | 0.50 | 0.90 |
| TSGAD | WACV 24 | 62.81 | 35.53 | 0.39 | 0.73 | 54.64 | 24.25 | 0.48 | 0.76 |
| | | C4 | | | | C5 | | | |
| MPED-RNN | CVPR 19 | 74.59 | 44.55 | 0.29 | **0.31** | 79.23 | 49.09 | **0.27** | **0.41** |
| GEPC | CVPR 20 | 72.44 | 30.71 | 0.29 | 0.99 | 66.05 | 33.04 | 0.36 | 0.56 |
| STG-NF | ICCV 23 | 72.24 | **93.95** | **0.28** | 0.71 | 56.48 | **82.36** | 0.46 | 0.96 |
| TSGAD | WACV 24 | **75.11** | 37.13 | **0.28** | 1.00 | 67.15 | 35.20 | 0.37 | 0.57 |
| | | CSC | | | | Combined | | | |
| MPED-RNN | CVPR 19 | 56.85 | 37.13 | 0.43 | **0.69** | 76.05 | 42.83 | **0.28** | **0.49** |
| GEPC | CVPR 20 | 58.32 | 41.09 | **0.42** | 0.86 | 62.25 | 28.62 | 0.41 | 0.67 |
| STG-NF | ICCV 23 | 53.60 | **66.28** | 0.47 | 0.92 | 57.57 | **83.77** | 0.46 | 0.90 |
| TSGAD | WACV 24 | **58.91** | 43.28 | 0.43 | 0.86 | 68.00 | 34.61 | 0.36 | 0.64 |

TABLE II: Benchmark of available SOTA pose-based models (MPED-RNN [54], GEPC [25], STG-NF [19], and TSGAD [2]) on HuVAD-S.
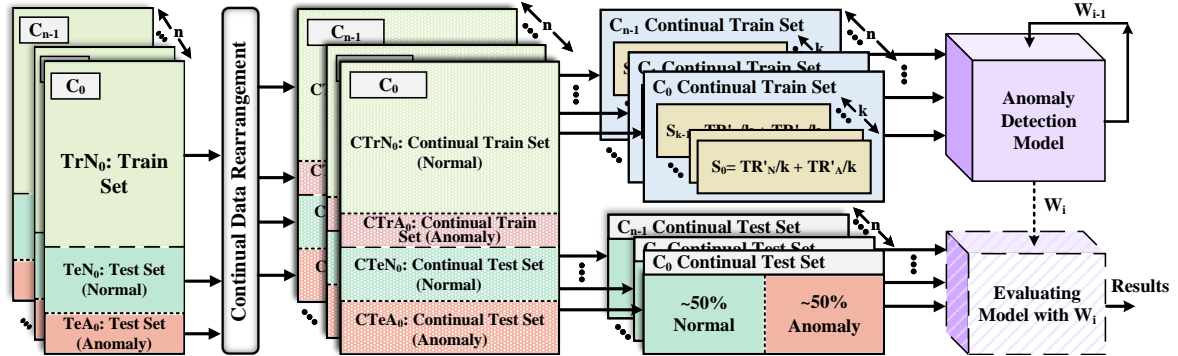


Fig. 6: The proposed UCAL framework begins with per-camera data from HuVAD-S, using the Continual Data Rearrangement module to simulate a real-world data stream. The data is then divided into k slices, and the model is incrementally trained on each slice, with performance evaluated on the continual test set at each step. This process is repeated across all cameras.

continual learning. This approach not only enables comparative analysis with traditional datasets but also allows for the application of the proposed Unsupervised Continual Anomaly Learning (UCAL) method to previous datasets, thereby broadening the impact of this new paradigm.

Following the principles of unsupervised anomaly detection, the training set is curated to contain only normal behavior. The test set, by contrast, includes both normal and anomalous behaviors in a 70/30 ratio to facilitate robust evaluation across varied conditions. As outlined in table I, HuVAD provides over five times more normal training frames than the next largest dataset, alongside a comprehensive test set exceeding 700K frames (approximately 6.5 hours), offering four times the scale of comparable datasets.

### A. Standard Results

To demonstrate the effectiveness of the proposed dataset, we leverage SOTA pose-based anomaly detection models with an available code repository. Specifically, we employ MPED-RNN [54], GEPC [25], STG-NF [19] and TSGAD [2] as discussed in section II. For the TSGAD model, we opted to use only the pose branch, aligning with this study's primary focus on pose-based anomaly detection. As for all chosen models, hyper-parameters follow the original study, and detailed training parameters are available in the supplementary materials for replication and further investigation.

table II illustrates a comprehensive analysis of models benchmarked on the HuVAD dataset. MPED-RNN [54] consistently achieves the highest overall performance, both across combined and individual cameras, with TSGAD [2] consis-
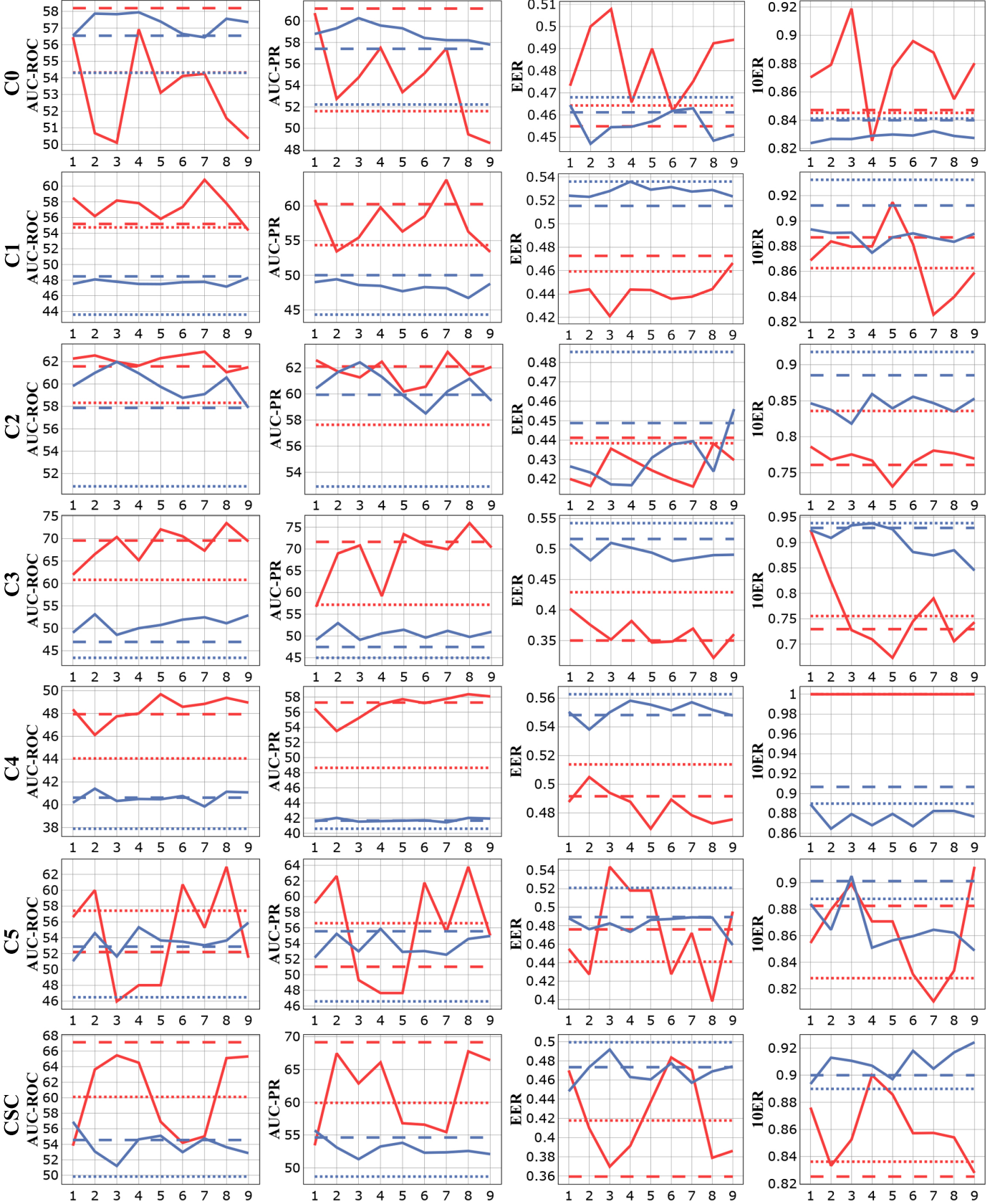
Fig. 7: Continual learning benchmark using TSGAD (red) [2] and STG-NF (blue) [19]. The dotted line is the baseline (trained on SHT's training set [36] and tested on HuVAD-C's test set), the dashed line is standard training, and the solid line is continual learning.

| Model | AUC-ROC | AUC-PR | EER | 10ER | AUC-ROC | AUC-PR | EER | 10ER |
|---|---|---|---|---|---|---|---|---|
| | C0 | | | | C1 | | | |
| TSGAD | **58.19** | **61.14** | **0.45** | 0.84 | 55.17 | 60.25 | 0.47 | 0.88 |
| TSGAD+UCAL(Ours) | 56.45 | 60.72 | 0.46 | **0.82** | **60.81** | **63.77** | **0.42** | **0.82** |
| STG-NF | 56.54 | 57.40 | 0.46 | 0.84 | **48.47** | **50.01** | **0.51** | 0.91 |
| STG-NF+UCAL(Ours) | **57.94** | **60.25** | **0.44** | **0.82** | 48.30 | 49.42 | 0.52 | **0.87** |
| | C2 | | | | C3 | | | |
| TSGAD | 61.57 | 62.10 | 0.44 | 0.76 | 69.57 | 71.62 | 0.35 | 0.72 |
| TSGAD+UCAL(Ours) | **62.88** | **63.22** | **0.41** | **0.73** | **73.47** | **75.93** | **0.32** | **0.70** |
| STG-NF | 57.86 | 59.94 | 0.44 | 0.88 | 46.96 | 47.45 | 0.51 | 0.92 |
| STG-NF+UCAL(Ours) | **61.97** | **62.42** | **0.41** | **0.81** | **53.11** | **52.93** | **0.48** | **0.84** |
| | C4 | | | | C5 | | | |
| TSGAD | 47.93 | 57.27 | 0.49 | 1.00 | 52.19 | 51.01 | 0.47 | 0.88 |
| TSGAD+UCAL(Ours) | **49.68** | **58.37** | **0.46** | 1.00 | **62.95** | **63.81** | **0.39** | **0.81** |
| STG-NF | 40.62 | 41.67 | 0.54 | 0.90 | 52.87 | **55.56** | 0.48 | 0.90 |
| STG-NF+UCAL(Ours) | **41.40** | **42.03** | **0.53** | **0.80** | **55.89** | 55.23 | **0.47** | **0.84** |
| | CSC | | | | | | | |
| | AUC-ROC | | AUC-PR | | EER | | 10ER | |
| TSGAD | **67.15** | | **69.14** | | **0.35** | | 0.82 | |
| TSGAD+UCAL(Ours) | 65.46 | | 67.75 | | 0.36 | | 0.82 | |
| STG-NF | 54.55 | | 54.63 | | 0.47 | | 0.90 | |
| STG-NF+UCAL(Ours) | **56.90** | | **55.71** | | **0.44** | | **0.89** | |

TABLE III: Continual learning benchmarks using TSGAD [2] and STG-NF [19]. Results for TSGAD and STG-NF are based on conventional training and evaluation on the HuVAD-C. TSGAD+UCAL(Ours) and STG-NF+UCAL(Ours) show the best results of UCAL.

---

**Algorithm 1** UCAL Framework

**Input** : $TeA_j, TeN_j, TrN_j$: Test Anomaly, Test Normal, and Train Normal sets from HuVAD-S for each of $n$ cameras
$k$: Number of data slices per camera
$\theta_{\text{pretrained}}$: Pretrained model weights

**Output:** $CTeA_j, CTeN_j, CTrN_j$: Updated Test Anomaly, Test Normal, and Train Normal sets for each camera in HuVAD-C

**for** $j = 1$ **to** $n$ **cameras do**
  **Step 1: Continual Data Rearrangement**
    $CTrA_j \leftarrow \frac{1}{100} \times \mathcal{N}(TrN_j)$, $CTeA_j \leftarrow TeA_j - CTrA_j$
    $CTeN_j \leftarrow \mathcal{N}(CTeA_j)$, $CTrN_j \leftarrow TrN_j + (TeN_j - CTeN_j) + CTrA_j$
  **Step 2: Split Training Set**
    Divide $(CTrN_j, CTrA_j)$ into $k$ slices: $S_{j1}, S_{j2}, \ldots, S_{jk}$
  **Step 3: Initialize and Train Model**
    $\theta_j \leftarrow \theta_{\text{pretrained}}$
    **for** $i = 1$ **to** $k$ **do**
      Load and train on $S_{ji}$, updating $\theta_j$
    **end**
**end**
**return** $\theta_1, \theta_2, \ldots, \theta_n$

---

tently ranking second. STG-NF [19] shows irregular performance, with the highest EER, lowest AUC-ROC, and highest AUC-PR across all experiments.

table II also reveals that the CSC camera offering a parking lot view (see fig. 2) is the most challenging for algorithms. As noted in section IV, this camera captures context-specific behaviors such as security training, making it harder to distinguish anomalies in the test set. fig. 5 also reveals that CSC has the highest crowd density and occlusion, further complicating its challenge for VAD models. These characteristics further highlight the need for solutions such as continual anomaly learning.

While EER indicates the lowest equal FNR and FPR possible, real-world VAD often prioritizes minimizing FNR. Therefore, a model with a lower 10ER is preferred [62]. table II shows that a lower EER does not always align with a lower 10ER. Comparisons on C4 and CSC highlight that models with similar EERs can have different 10ERs, underscoring the importance of 10ER in model selection.

## VII. Unsupervised Continual Anomaly Learning

As discussed in section I, the context-specific and open-set nature of unsupervised anomaly detection necessitates learning of the normal distribution within each unique environment, limiting model generalizability. Consequently, models trained on specific datasets often lose effectiveness in real-world applications. To address this, we propose an Unsupervised Continual Anomaly Learning framework for each camera available in HuVAD that incrementally evolves the anomaly model, offering a more adaptive solution.

### A. Methodology

fig. 6 illustrates the UCAL framework, which configures HuVAD-S to simulate real-world streaming data and performs continual anomaly learning. It begins by separating data by camera, with the Continual Data Rearrangement module applied independently to each camera's dataset. Since anomalies are rare in the real world, Continual Data Rearrangement starts with HuVAD-S (section VI) and injects anomalies randomly into the training stream, maintaining an anomaly ratio below 1%, as specified in algorithm 1. The test set is balanced to achieve a 1:1 ratio of normal to anomalous frames, enhancing metrics such as AUC-ROC and AUC-PR. Normal frames removed from the test set are added to the training stream as detailed in algorithm 1. The new training and test sets form the HuVAD Continual benchmark (HuVAD-C). More

details on the test and training sets for each camera are in the supplementary material.

The HuVAD-C training set is then divided into $k$ slices for $k$ training steps. As depicted in algorithm 1, the UCAL training process begins with a pre-trained model from an external dataset and incrementally trains on the HuVAD-C training set. Model performance is tested on the HuVAD-C test set at each step. This structured, incremental approach allows for gradual model adaptation, effectively addressing the complexities of continual learning in VAD (see fig. 6).

### B. UCAL Results

We evaluate the two most recent VAD models, STG-NF [19] and TSGAD (pose branch) [2], using this framework. To simulate distribution shifts, we pre-train models on the SHT dataset [36] as the origin data domain. For our experiments, we chose the number of slices (k) equal to 9 and trained the models for 9 steps with the learning rate of $5e - 3$. More detailed parameters are provided in the supplementary material to ensure reproducibility.

The results of continual learning benchmarking are illustrated in fig. 7. The baseline, represented by the dotted lines, involves models trained on SHT [36] and tested on the HuVAD-C test set without any fine-tuning for HuVAD-C. The dashed lines indicate conventional training and testing on HuVAD-C train and test sets. Our findings reveal that continual training outperforms the baseline in 98.21% of cases, demonstrating significant performance improvement. Furthermore, as depicted in table III a comparison with conventional training shows that UCAL yields better results in 82.14% of cases, highlighting its advantages over standard training. The HuVAD dataset has enabled benchmarking of continual learning, further validating these improvements.

## VIII. CONCLUSION

This study presents HuVAD and UCAL, introducing a novel dataset and continual learning framework as a new paradigm for enhancing the effectiveness of VAD. By enabling continuous model adaptation, this paper addresses the limitations of conventional static benchmarks and paves the way for VAD systems that remain robust in dynamic, real-world environments. HuVAD also prioritizes privacy with comprehensive annotations across varied scenes. Our evaluation with both standard and continual benchmark, highlights the potential to drive future research in anomaly detection through adaptive, context-aware approaches.

### ACKNOWLEDGEMENT

## REFERENCES

[1] G. A. Noghre, A. D. Pazho, V. Katariya, and H. Tabkhi, "Understanding the challenges and opportunities of pose-based anomaly detection," *arXiv preprint arXiv:2303.05463*, 2023. 1

[2] G. A. Noghre, A. D. Pazho, and H. Tabkhi, "An exploratory study on human-centric video anomaly detection through variational autoencoders and trajectory prediction," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 995–1004. 1, 5, 6, 7, 8

[3] A. D. Pazho, G. A. Noghre, A. A. Purkayastha, J. Vempati, O. Martin, and H. Tabkhi, "A survey of graph-based deep learning for anomaly detection in distributed systems," *IEEE Transactions on Knowledge and Data Engineering*, vol. 36, no. 1, pp. 1–20, 2023. 1

[4] K. Doshi and Y. Yilmaz, "Towards interpretable video anomaly detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 2655–2664. 1

[5] A. Danesh Pazho, G. Alinezhad Noghre, B. Rahimi Ardabili, C. Neff, and H. Tabkhi, "Chad: Charlotte anomaly dataset," in *Scandinavian Conference on Image Analysis*. Springer, 2023, pp. 50–66. 1, 2, 3, 4

[6] K. Doshi and Y. Yilmaz, "Continual learning for anomaly detection in surveillance videos," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 254–255. 1, 2

[7] ——, "Online anomaly detection in surveillance videos with asymptotic bound on false alarm rate," *Pattern Recognition*, vol. 114, p. 107865, 2021. 1, 2

[8] ——, "Rethinking video anomaly detection-a continual learning approach," in *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2022, pp. 3961–3970. 1, 2

[9] F. Zhu, Z. Cheng, X.-Y. Zhang, C.-L. Liu, and Z. Zhang, "Rcl: Reliable continual learning for unified failure detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 12 140–12 150. 1

[10] Z. Yang and R. J. Radke, "Context-aware video anomaly detection in long-term datasets," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 4002–4011. 1

[11] K. Faber, R. Corizzo, B. Sniezynski, and N. Japkowicz, "Lifelong continual learning for anomaly detection: New challenges, perspectives, and insights," *IEEE Access*, vol. 12, pp. 41 364–41 380, 2024. 1

[12] S. Yao, G. A. Noghre, A. D. Pazho, and H. Tabkhi, "Evaluating the effectiveness of video anomaly detection in the wild: Online learning and inference for real-world deployment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2024, pp. 4832–4841. 1

[13] G. Saha and K. Roy, "Continual learning with scaled gradient projection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 8, 2023, pp. 9677–9685. 1

[14] Y.-C. Hsu, Y.-C. Liu, A. Ramasamy, and Z. Kira, "Re-evaluating continual learning scenarios: A categorization and case for strong baselines," *arXiv preprint arXiv:1810.12488*, 2018. 1

[15] T. Saheb, "Ethically contentious aspects of artificial intelligence surveillance: a social science perspective," *AI and Ethics*, vol. 3, no. 2, pp. 369–379, 2023. 1

[16] B. R. Ardabili, A. D. Pazho, G. A. Noghre, C. Neff, S. D. Bhaskararayuni, A. Ravindran, S. Reid, and H. Tabkhi, "Understanding policy and technical aspects of ai-enabled smart video surveillance to address public safety," *Computational Urban Science*, vol. 3, no. 1, p. 21, 2023. 1, 3

[17] I. Joshi, M. Grimmer, C. Rathgeb, C. Busch, F. Bremond, and A. Dantcheva, "Synthetic data in human analysis: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1

[18] A. Kunchala, M. Bouroche, and B. Schoen-Phelan, "Towards a framework for privacy-preserving pedestrian analysis," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 4370–4380. 1

[19] O. Hirschorn and S. Avidan, "Normalizing flows for human pose anomaly detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 13 545–13 554. 1, 3, 5, 6, 7, 8

[20] S. Yu, Z. Zhao, H. Fang, A. Deng, H. Su, D. Wang, W. Gan, C. Lu, and W. Wu, "Regularity learning via explicit distribution modeling for skeletal video anomaly detection," *IEEE Transactions on Circuits and Systems for Video Technology*, 2023. 1, 3

[21] C. Huang, Y. Liu, Z. Zhang, C. Liu, J. Wen, Y. Xu, and Y. Wang, "Hierarchical graph embedded pose regularity learning via spatio-temporal transformer for abnormal behavior detection," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 307–315. 1

[22] X. Zeng, Y. Jiang, W. Ding, H. Li, Y. Hao, and Z. Qiu, "A hierarchical spatio-temporal graph convolutional neural network for anomaly detec-

tion in videos," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021. 1

[23] X. Chen, S. Kan, F. Zhang, Y. Cen, L. Zhang, and D. Zhang, "Multiscale spatial temporal attention graph convolution network for skeleton-based anomaly behavior detection," *Journal of Visual Communication and Image Representation*, vol. 90, p. 103707, 2023. 1

[24] Y. Jain, A. K. Sharma, R. Velmurugan, and B. Banerjee, "Posecvae: Anomalous human activity detection," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 2927–2934. 1

[25] A. Markovitz, G. Sharir, I. Friedman, L. Zelnik-Manor, and S. Avidan, "Graph embedded pose clustering for anomaly detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 539–10 547. 1, 5

[26] C. Cao, Y. Lu, P. Wang, and Y. Zhang, "A new comprehensive benchmark for semi-supervised video anomaly detection and anticipation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 20 392–20 401. 2, 3, 4

[27] W. Sultani, C. Chen, and M. Shah, "Real-world anomaly detection in surveillance videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6479–6488. 2, 4

[28] P. Wu, J. Liu, Y. Shi, Y. Sun, F. Shao, Z. Wu, and Z. 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*. Springer, 2020, pp. 322–339. 2

[29] A. Acsintoae, A. Florescu, M.-I. Georgescu, T. Mare, P. Sumedrea, R. T. Ionescu, F. S. Khan, and M. Shah, "Ubnormal: New benchmark for supervised open-set video anomaly detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 20 143–20 153. 2

[30] L. Zhu, L. Wang, A. Raj, T. Gedeon, and C. Chen, "Advancing video anomaly detection: A concise review and a new dataset," 2024. [Online]. Available: https://arxiv.org/abs/2402.04857 2

[31] H. Du, S. Zhang, B. Xie, G. Nan, J. Zhang, J. Xu, H. Liu, S. Leng, J. Liu, H. Fan, D. Huang, J. Feng, L. Chen, C. Zhang, X. Li, H. Zhang, J. Chen, Q. Cui, and X. Tao, "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 (CVPR)*, June 2024, pp. 18 793–18 803. 2

[32] B. Ramachandra and M. Jones, "Street scene: A new dataset and evaluation protocol for video anomaly detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 2569–2578. 2

[33] A. Adam, E. Rivlin, I. Shimshoni, and D. Reinitz, "Robust real-time unusual event detection using multiple fixed-location monitors," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 3, pp. 555–560, 2008. 2

[34] V. Mahadevan, W. Li, V. Bhalodia, and N. Vasconcelos, "Anomaly detection in crowded scenes," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2010, pp. 1975–1981. 2

[35] C. Lu, J. Shi, and J. Jia, "Abnormal event detection at 150 fps in matlab," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 2720–2727. 2

[36] W. Liu, W. Luo, D. Lian, and S. Gao, "Future frame prediction for anomaly detection–a new baseline," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6536–6545. 2, 3, 4, 6, 8

[37] R. Rodrigues, N. Bhargava, R. Velmurugan, and S. Chaudhuri, "Multi-timescale trajectory prediction for abnormal human activity detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 2626–2634. 2, 3, 4

[38] N. Bugarin, J. Bugaric, M. Barusco, D. D. Pezze, and G. A. Susto, "Unveiling the anomalies in an ever-changing world: A benchmark for pixel-level anomaly detection in continual learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 4065–4074. 2

[39] V. Chavan, P. Koch, M. Schlüter, C. Briese, and J. Krüger, "Active data collection and management for real-world continual learning via pretrained oracle," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 4085–4096. 2

[40] M. Mozaffari, K. Doshi, and Y. Yilmaz, "Self-supervised learning for online anomaly detection in high-dimensional data streams," *Electronics*, vol. 12, no. 9, p. 1971, 2023. 2

[41] F. A. Mazarbhuiya and M. Shenify, "Real-time anomaly detection with subspace periodic clustering approach," *Applied Sciences*, vol. 13, no. 13, p. 7382, 2023. 2

[42] K. Ahmad, M. Maabreh, M. Ghaly, K. Khan, J. Qadir, and A. Al-Fuqaha, "Developing future human-centered smart cities: Critical analysis of smart city security, data management, and ethical challenges," *Computer Science Review*, vol. 43, p. 100452, 2022. 2

[43] B. R. Ardabili, A. D. Pazho, G. A. Noghre, C. Neff, A. Ravindran, and H. Tabkhi, "Understanding ethics, privacy, and regulations in smart video surveillance for public safety," *arXiv preprint arXiv:2212.12936*, 2022. 2

[44] F. A. Raso, H. Hilligoss, V. Krishnamurthy, C. Bavitz, and L. Kim, "Artificial intelligence & human rights: Opportunities & risks," *Berkman Klein Center Research Publication*, no. 2018-6, 2018. 3

[45] M. Noriega, "The application of artificial intelligence in police interrogations: An analysis addressing the proposed effect ai has on racial and gender bias, cooperation, and false confessions," *Futures*, vol. 117, p. 102510, 2020. 3

[46] S. E. Whang, Y. Roh, H. Song, and J.-G. Lee, "Data collection and quality challenges in deep learning: A data-centric ai perspective," *The VLDB Journal*, vol. 32, no. 4, pp. 791–813, 2023. 3

[47] A. Barbalau, R. T. Ionescu, M.-I. Georgescu, J. Dueholm, B. Ramachandra, K. Nasrollahi, F. S. Khan, T. B. Moeslund, and M. Shah, "Ssmtl++: Revisiting self-supervised multi-task learning for video anomaly detection," *Computer Vision and Image Understanding*, vol. 229, p. 103656, 2023. 3

[48] G. Wang, Y. Wang, J. Qin, D. Zhang, X. Bao, and D. Huang, "Video anomaly detection by solving decoupled spatio-temporal jigsaw puzzles," in *European Conference on Computer Vision*. Springer, 2022, pp. 494–511. 3

[49] A. D. Pazho, C. Neff, G. A. Noghre, B. R. Ardabili, S. Yao, M. Baharani, and H. Tabkhi, "Ancilia: Scalable intelligent video surveillance for the artificial intelligence of things," *IEEE Internet of Things Journal*, 2023. 3

[50] G. Jocher, A. Chaurasia, and J. Qiu, "Ultralytics yolov8," 2023. [Online]. Available: https://github.com/ultralytics/ultralytics 3

[51] Y. Zhang, P. Sun, Y. Jiang, D. Yu, F. Weng, Z. Yuan, P. Luo, W. Liu, and X. Wang, "Bytetrack: Multi-object tracking by associating every detection box," 2022. 3

[52] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5693–5703. 3

[53] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755. 3

[54] R. Morais, V. Le, T. Tran, B. Saha, M. Mansour, and S. Venkatesh, "Learning regularity in skeleton trajectories for anomaly detection in videos," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 11 996–12 004. 5

[55] J. Davis and M. Goadrich, "The relationship between precision-recall and roc curves," in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 233–240. 4

[56] A. Fernández, S. García, M. Galar, R. C. Prati, B. Krawczyk, and F. Herrera, *Learning from imbalanced data sets*. Springer, 2018, vol. 10, no. 2018. 4

[57] H. He and Y. Ma, "Imbalanced learning: foundations, algorithms, and applications," 2013. 4

[58] T. Saito and M. Rehmsmeier, "The precision-recall plot is more informative than the roc plot when evaluating binary classifiers on imbalanced datasets," *PloS one*, vol. 10, no. 3, p. e0118432, 2015. 4

[59] W. Li, V. Mahadevan, and N. Vasconcelos, "Anomaly detection and localization in crowded scenes," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 1, pp. 18–32, 2013. 4

[60] D. Maio, D. Maltoni, R. Cappelli, J. L. Wayman, and A. K. Jain, "Fvc2002: Second fingerprint verification competition," in *2002 International conference on pattern recognition*, vol. 3. IEEE, 2002, pp. 811–814. 4

[61] P. C. Neto, F. Boutros, J. R. Pinto, N. Damer, A. F. Sequeira, J. S. Cardoso, M. Bengherabi, A. Bousnat, S. Boucheta, N. Hebbadj *et al.*, "Ocfr 2022: Competition on occluded face recognition from synthetically generated structure-aware occlusions," in *2022 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 2022, pp. 1–9. 4

[62] B. R. Ardabili, A. D. Pazho, G. A. Noghre, V. Katariya, G. Hull, S. Reid, and H. Tabkhi, "Exploring public's perception of safety and video surveillance technology: A survey approach," *Technology in Society*, vol. 78, p. 102641, 2024. 4, 7