

EXPLORING DIFFUSION MODELS FOR UNSUPERVISED VIDEO ANOMALY DETECTION

Anil Osman Tur^{*◇} Nicola Dall’Asen^{*‡} Cigdem Beyan^{*} Elisa Ricci^{*†}

^{*}Department of Information Engineering and Computer Science, University of Trento, Trento, Italy

[◇]Energy Efficient Embedded Digital Architectures Unit, Fondazione Bruno Kessler, Trento, Italy

[‡] Department of Computer Science, University of Pisa, Pisa, Italy

[†]Deep Visual Learning Research Group, Fondazione Bruno Kessler, Trento, Italy

ABSTRACT

This paper investigates the performance of diffusion models for video anomaly detection (VAD) within the most challenging but also the most operational scenario in which the data annotations are not used. As being sparse, diverse, contextual, and often ambiguous, detecting abnormal events precisely is a very ambitious task. To this end, we rely only on the information-rich spatio-temporal data, and the reconstruction power of the diffusion models such that a high reconstruction error is utilized to decide the abnormality. Experiments performed on two large-scale video anomaly detection datasets demonstrate the consistent improvement of the proposed method over the state-of-the-art generative models while in some cases our method achieves better scores than the more complex models. This is the first study using a diffusion model and examining its parameters’ influence to present guidance for VAD in surveillance scenarios.

Index Terms— Anomaly Detection, unsupervised learning, video understanding, imbalanced data

1. INTRODUCTION

Automated video anomaly detection (VAD) has become an essential task in the computer vision community as a consequence of the exponential increase in the number of videos being captured. VAD is relevant to several applications in intelligent surveillance, and behavior understanding [1, 2, 3, 4, 5, 6], to name a few. Anomaly is commonly defined as a rare or unexpected or unusual entity, that diverges significantly from normality, which is defined as expected and common. Despite being sparse and diverse, the abnormal events are also very contextual, and often ambiguous, thus they challenge the performance of the VAD models [7].

Data labeling is already a notoriously expensive and time-consuming task and considering the aforementioned characteristics of the abnormal events, it is almost infeasible to collect all possible anomaly samples to perform *fully-supervised* learning methods. Therefore, a typical approach in VAD, is to train a *one-class* classifier that learns from the *normal* training data [8, 9, 10]. However, the data collection problem occurring for fully-supervised learning almost remains for the one-class classifier, since it is unfeasible to have access to every variety of normal training data, given the dynamic nature of real-world applications and the wide range of

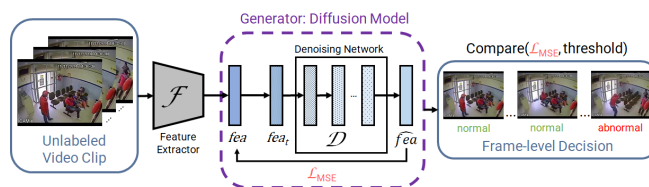


Fig. 1: The proposed method takes a batch of unlabeled video clips as the input and learns to determine whether each frame is anomalous or exhibits normal behavior. Our model is a generative one, which leverages the reconstruction capability of diffusion models for unsupervised VAD. Mean Squared Error (MSE) distribution over a batch is used in conjunction with a data-driven threshold to decide which frames are anomalous. For the purpose of clarity, we show only one video clip.

normal classes [4, 5]. In a one-class classifier setting, it is highly possible that an unseen normal event can be misclassified as abnormal since its representation is remarkably different from the representations learned from normal training data.

The data availability problem led some researchers to define the weakly supervised VAD, which does not rely on fine-grained per-frame annotations but yields the video-level labels [11, 12]. In detail, in fully-supervised VAD each *individual frame* has an annotation as normal or abnormal. Instead, in weakly supervised VAD, a *video* is labeled as anomalous even if only one frame of it is anomalous, and labeled as normal when all frames of it are normal. Even though performing such annotations seems relatively cheaper, it is important to notice that, in the weakly supervised setting, (a) labeling a video as normal still requires inspection of whole frames (similar to the fully-supervised setting), and (b) such methodologies often fail to localize the abnormal portion of the video, which can be impractical, e.g., when the video footage is long.

Recently, Zaheer et al. [13] defined *unsupervised* VAD, which takes *unlabelled* videos as the input and learns to make the decision of anomaly or normality for each frame. Such a fashion is undoubtedly more challenging compared to fully, weakly, and one-class counterparts, but it literally brings in the advantage of not requiring data annotations at all. It is worth differentiating the definition of unsupervised VAD [13] from one-class VAD since the latter is being referred to as unsupervised in some studies [14, 15, 10, 16, 17, 18]. In the case of one-class VAD, the training data distribution represents only the normality, meaning that there still exists a notion of labeling. Whereas the implementation of *unsupervised* VAD [13] does not make any assumption regarding the distribution of the training data, and never uses the labels for model training, instead, it

relies only on the spatiotemporal features of the data.

In this study, we perform *unsupervised VAD* by leveraging the information-rich unlabelled videos. To do so, we only depend on the reconstruction capability of the diffusion models [19] (see Fig. 1 for the proposed method’s description). This is the first attempt that the effectiveness of the diffusion models is being investigated for VAD in surveillance scenarios. The aim of this work is to present an exploratory study: (a) to understand whether diffusion models can be effectively used for unsupervised VAD, and (b) to discover the behavior of the diffusion model [19] in terms of several parameters of it for VAD. Experimental analysis performed on two large-scale datasets: UCF-Crime [2] and ShanghaiTech [3], demonstrate that the proposed approach always performs better than the state-of-the-art (SOTA) generative model of VAD. Moreover, in some cases, the proposed method is able to surpass more complex SOTA methods [13, 20]. The code of our method and the SOTA [13] is publicly available [HERE](#).

2. RELATED WORK

Anomaly detection is a widely studied topic that regards several tasks such as medical diagnosis, fault detection, animal behavior understanding, and fraud detection. Interested readers can refer to a recent survey: [4]. Below, our review focuses on VAD in *surveillance scenarios*. We also present the definition and notations of diffusion models and state the methodology we follow for VAD.

Video Anomaly Detection in Surveillance Scenarios. VAD has been typically solved as an outlier detection task (i.e., one-class classifier), in which a model is learned from the normal training data (requiring data annotations), and during testing, an abnormality is detected with the approaches such as distance-based [21], reconstruction-based [8] or probability-based [22]. Such approaches might result in an ineffective classifier since they exclude the abnormal classes during training. This might occur particularly when a sufficient amount of data representing each variety of the normal class cannot be used in training. An alternative approach is using *unlabelled* training data without assuming any normalcy [13], referred to as (*fully*) *unsupervised VAD*. Unlike one-class classifiers, unsupervised VAD does not require data labeling and can potentially generalize well by not excluding the abnormal data from training. Zaheer et al. [13] proposed a Generative Cooperative Learning composed of a generator and a discriminator mutually being trained together with the negative learning paradigm. The generator which is an autoencoder reconstructs the normal and abnormal representations while using the negative learning approach to help the discriminator to estimate the probability of an instance being abnormal with a data-driven threshold. That approach [13] conforms that anomalies are less frequent than normal events and events are often temporally consistent. In this study, we follow the unsupervised VAD definition in [13]. Unlike [13], our method relies only on a generative architecture, which is based on a diffusion model. We, first time in this study, investigate the effectiveness of the diffusion models for VAD in surveillance scenarios, by reporting how individual parameters affect the model performance, and by comparing them with the SOTA.

Diffusion Models. Diffusion models (DMs) [23, 24] are a type of generative model that gains the ability to generate diverse samples by corrupting training samples with noise and learning to reverse the process. These models have achieved SOTA performance in tasks such as text-to-image synthesis [25], semantic editing [26], and audio synthesis [27]. They have also been used in representation learning for discriminative tasks like object detection [28], image segmentation [29], and disease detection [30]. This study is the first

attempt to apply DMs for video anomaly detection.

DMs are formulated as a progressive addition of Gaussian noise of standard deviation σ to an input data point x_T sampled from a distribution $p_{data}(x)$ with standard deviation σ_{data} . The noised distribution $p(x, \sigma)$, for $\sigma \gg \sigma_{data}$, becomes isotropic Gaussian and allows to sample a point $x_0 \sim \mathcal{N}(0, \sigma_{max}\mathbf{I})$. This point is gradually denoised with noise levels $\sigma_0 = \sigma_{max} > \sigma_{T-1} > \dots > \sigma_1 > \sigma_T = 0$ into new samples distributed according to the dataset distribution. DMs are trained with Denoising Score Matching [31], where a denoiser function $D_\theta(x; \sigma)$ minimizes the expected L_2 denoising error for samples drawn from p_{data} for every σ :

$$\mathbb{E}_{x \sim p_{data}} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma \mathbf{I})} \|D_\theta(x + \epsilon; \sigma) - x\|_2^2, \quad (1)$$

and the score functions used in the reverse process become:

$$\nabla \log p(x; \sigma) = (D_\theta(x; \sigma) - x)s/\sigma^2. \quad (2)$$

In this paper, we adapt the diffusion model of [19], whose details are described in the next section.

3. METHOD

Given a video clip, we first extract features from a 3D-CNN (F) both in training and testing. These features are supplied to the generator, which is a diffusion model, to reconstruct them without using the labels. We follow the diffusion model variant proposed in [19] and refer to it as *k-diffusion*. It disentangles the design choices of previous diffusion models and provides a framework where each component can be adjusted separately, as shown in Table 1. In particular, Karras et al. [19] exposes the issue of expecting the network D_θ to work well in high noise regimes, i.e. when σ_i is high. To solve this, *k-diffusion* proposes a σ -dependent skip connection, allowing the network to perform x_0 or ϵ -prediction, or something in between based on the noise magnitude. The denoising network D_θ , therefore, is formulated as follows:

$$D_\theta(x; \sigma) = c_{skip}(\sigma) x + c_{out}(\sigma) F_\theta(c_{in}(\sigma) x; c_{noise}(\sigma)), \quad (3)$$

where F_θ becomes the effective network to train, c_{skip} modulates the skip connection, $c_{in}(\cdot)$ and $c_{out}(\cdot)$ scale input and output magnitudes, and $c_{noise}(\cdot)$ scales σ to become suitable as input for F_θ .

Several hyperparameters control the diffusion process in *k-diffusion*, and we extensively explore the role of training noise – distributed according to a log-normal distribution with parameters (P_{mean}, P_{std}) – and sampling noise with boundary values of σ_{min} and σ_{max} . These distributions are crucial choices depending on the task and on the dataset [32]. Given we use diffusion models on an

Table 1: The design choice of *k-diffusion*. T is the Number of Function Evaluations (NFEs) executed during sampling. The corresponding sequence of time steps is $\{t_0, t_1, \dots, t_T\}$, where $t_T = 0$. F_θ represents the raw neural network.

Sampling	
ODE solver	LMS
Time steps	$(\sigma_{max}^{\frac{1}{p}} + \frac{i}{T-1}(\sigma_{min}^{\frac{1}{p}} - \sigma_{max}^{\frac{1}{p}}))^p$
Network and preconditioning	
Architecture of F_θ	Any, MLP in our case
Skip scaling $c_{skip}(\sigma)$	$\sigma_{data}^2 / (\sigma^2 + \sigma_{data}^2)$
Output scaling $c_{out}(\sigma)$	$\sigma \cdot \sigma_{data} / \sqrt{\sigma_{data}^2 + \sigma^2}$
Input scaling $c_{in}(\sigma)$	$1 / \sqrt{\sigma^2 + \sigma_{data}^2}$
Noise cond. $c_{noise}(\sigma)$	$\frac{1}{4} \ln(\sigma)$
Training	
Noise distribution	$\ln(\sigma) \sim \mathcal{N}(P_{mean}, P_{std}^2)$
Loss weighting	$(\sigma^2 + \sigma_{data}^2) / (\sigma \cdot \sigma_{data})^2$

Algorithm 1 Anomaly detection with denoising diffusion

Require: Batch of video clips x , feature extractor network F , denoising network D , denoising step t , threshold sensitivity k

- 1: $\text{fea} = F(x)$ # Feature extraction with the backbone F
- 2: $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ # Noise sampling for k -diffusion
- 3: $\text{fea}_t = \text{fea} + \epsilon * \sigma_t$ # Diffusion input corruption
- 4: $\widehat{\text{fea}} = \text{sampling}(D(\text{fea}_t, \sigma_t))$ # Reconstruction of a feature vector with k -diffusion algorithm
- 5: $L_b = \text{MSE}(\text{fea}, \widehat{\text{fea}})$ # Reconstruction loss computation
- 6: $L_{th} = \mu_p + k\sigma_p$ # Data-driven threshold L_{th}
- 7: If $L_{th} < L_b$ then abnormal, otherwise normal # L_b is the loss in a batch
- 8: **Return** Normal or Abnormal

unprecedented task and on new datasets we do not rely on parameters from literature, instead, we perform an extensive study of the correlation between noise and performance on the task in Sec. 4.2.

The reverse process of a DM does not need to start from noise with variance σ_{max}^2 but it can place at any arbitrary step $t \in (0, T)$, with $\sigma_{max}^2 = \sigma_0^2$ as shown in [26]. Given a real data point x , we can sample $x_t \sim \mathcal{N}(x, \sigma_t \mathbf{I})$ and then apply the reverse process to x_T . This allows for retaining part of the information of the original data point – the low-frequency component – and removing the high-frequency component. We exploit this property to remove the components associated with abnormal parts of the clip by adding Gaussian noise. Then, we measure the goodness of reconstruction using mean squared error (MSE), meaning that a high reconstruction error might indicate the presence of abnormal activity. The choice of the starting point t for this procedure is a crucial hyperparameter of the method, as it controls the realism-faithfulness tradeoff as described in [26]. Refer to Sec. 4.2 presenting a study to understand the influence of this tradeoff on VAD.

We adopt [13] to decide whether a video frame is anomalous. In detail, the decision for a single video frame is made by keeping the distribution of the reconstruction loss (MSE) of each instance over a batch. The feature vectors resulting in higher loss refer to anomalous and smaller loss refers to normal while this decision is made through a data-driven threshold (L_{th}), defined as $L_{th} = \mu_p + k\sigma_p$ where k is a constant, μ_p and σ_p are the mean and standard deviations of the MSE loss for each batch. The anomaly detection phase is given in Algorithm 1.

4. EXPERIMENTAL ANALYSIS AND RESULTS

As the **evaluation metric**, we use the area under the Receiver Operating Characteristic (ROC) curve (AUC), which is computed based on frame-level annotations of the test videos of the datasets, in line with the prior arts. To evaluate and compare the performance of the proposed method, experiments are conducted on two large-scale unconstrained **datasets**: UCF-Crime [2] and ShanghaiTech [3]. The UCF-Crime dataset [2] is collected from various CCTV cameras, having different field-of-views. It is composed of in total 128 hours of videos with the annotation of 13 different real-world anomalous events such as road accidents, stealing, and explosion. We use the standard training (810 abnormal and 800 normal videos, without using the labels) and testing (130 abnormal and 150 normal videos) splits of the dataset to provide fair comparisons with SOTA. ShanghaiTech dataset [3] is captured in 13 different camera angles with complex lighting conditions. We use the training split that contains 63 abnormal and 174 normal videos and the testing split composed of 44 abnormal and 154 normal videos curated in line with SOTA.

We use 3D-ResNext101 and 3D-ResNet18 as feature extractors F due to their popularity in VAD [4, 5, 13]. 3D-ResNext101 has a dimensionality of 2048, and 3D-ResNet18 has 512 dimensions. The denoising network D is an MLP with an encoder-decoder structure.

The encoder is composed of 3 layers of size $\{1024, 512, 256\}$ while the decoder has the structure of $\{256, 512, 1024\}$. The learning rate scheduler and EMA of the model are taken as the default values of k -diffusion, with an initial learning rate of 2×10^{-4} and InverseLR scheduling; weight decay is set at 1×10^{-4} . Segment size for feature extraction is set to 16 non-overlapping frames and the training is performed up to 50 epochs with a batch size of 8192 in line with [13]. The timestep σ_t is transformed via Fourier embedding and integrated into the network through FiLM layers [34] both in the encoder and the decoder segments of the network. The hyperparameters (e.g., P_{mean}, P_{std}, t) used to realize k -diffusion are given in Sec. 4.2.

4.1. Comparisons with State-Of-The-Art (SOTA)

The performance of the proposed method is compared with the SOTA [20, 13] in Table 2. Kim et al. [20] proposed a one-class VAD method, which was then adopted to perform *unsupervised* VAD in [13]. In our comparisons, we use the unsupervised version of [20]. The proposed approach surpasses [20] within a large margin of: 10.91-12.41% AUC. The comparisons between the proposed method and the autoencoder of [13] demonstrate that as a generative model, the proposed method is more favorable than [13] by better performing VAD within a margin of: 6.15-14.44% AUC. When the features extracted from 3D-ResNext101 are used, the full model of [13] achieves better results than the proposed method. This is rather not surprising given that the full model of [13] is more complex than a generative model (i.e., autoencoder or diffusion model) since it additionally includes a discriminator and a negative learning component. Importantly, when 3D-ResNet18 is used as the backbone, the proposed method exceeds the full model of [13] within a large margin of: 4.9-8.36% AUC. Such results confirm the remarkable effectiveness of k -diffusion to perform VAD.

4.2. Diffusion Model Analysis

Below, the effect of different hyperparameters of k -diffusion model and a comparative study regarding timestep embeddings are given.

Noise. The training and sampling noise distributions are not independent in k -diffusion model, and we computed the relation between (P_{mean}, P_{std}) and $(\sigma_{min}, \sigma_{max})$ to be governed by the following formula: $\sigma_{max}, \sigma_{min} = e^{P_{mean} \pm 5P_{std}}$. This allows us to restrict our search to two parameters instead of four. We also extracted the formula using the default parameters of k -diffusion: $P_{mean} = -1.2, P_{std} = 1.2, \sigma_{min} = 0.02$ and $\sigma_{max} = 80$. The corresponding results are given in Fig. 2 when the k of L_{th} is taken as 1 for ShanghaiTech dataset [3] with 3D-ResNet18. One can observe that, in general, a smaller value of P_{mean} leads to higher results. This shows that we perform diffusion in a latent space that is well-behaved, and therefore a lower amount of noise is needed to reach an isotropic Gaussian distribution.

Starting point of the reverse process. Similar to SDEdit [26] and their realism-faithfulness tradeoff, we explore the effect of different t as the starting point of the reverse process. Recall that $\sigma_t > \sigma_{t+1}$ means that a t close to zero indicates a noised x_t closer to isotropic Gaussian, instead for t closer to T means, the features used are closer to the original data distribution. We target to find the best value of t such that sufficient information about the structure of the clip is retained while the information about the possible anomaly is destroyed. In this way, one can obtain a higher reconstruction error, leading to deciding the associated video frame as anomalous. The corresponding results are given in Fig. 2 when k of L_{th} is 1. $t = \text{best}$ refers to the best results obtained from $t = 0$ to $t = 9$, given a fixed P_{mean}, P_{std} combination. For ShanghaiTech with 3D-ResNet18 backbone, the majority of the time the starting point $t = 4$

Table 2: Performance comparisons with the SOTA on (a) UCF-Crime [2] and (b) ShanghaiTech [3] datasets. The best results are in bold. The second best results are underlined. The full model of [13] includes generator, negative learning, and discriminator. * indicates our implementation since the corresponding code is not publicly available. The results indicated with \diamond were taken from [13].

Method	Feature	AUC (%)
Kim et al. [20] \diamond	3D-ResNext 101	52.00
Autoencoder [13]	3D-ResNext 101	56.32
Autoencoder [13]*	3D-ResNext 101	56.27
Full model [13]	3D-ResNext 101	68.17
Full model [13]*	3D-ResNext 101	58.30
Proposed w/ [33]	3D-ResNext101	59.42
Proposed	3D-ResNext 101	<u>62.91</u>
Autoencoder [13]*	3D-ResNet18	49.78
Full model [13]*	3D-ResNet18	56.86
Proposed w/ [33]	3D-ResNet18	<u>60.52</u>
Proposed	3D-ResNet18	65.22

(a) Results on UCF-Crime [2]

Method	Feature	AUC (%)
Kim et al. [20] \diamond	3D-ResNext 101	56.47
Autoencoder [13]	3D-ResNext 101	62.73
Autoencoder [13]*	3D-ResNext 101	62.05
Full model [13]	3D-ResNext 101	72.41
Full model [13]*	3D-ResNext 101	65.62
Proposed w/ [33]	3D-ResNext101	62.41
Proposed	3D-ResNext 101	<u>68.88</u>
Autoencoder [13]*	3D-ResNet18	69.02
Full model [13]*	3D-ResNet18	71.20
Proposed w/ [33]	3D-ResNet18	<u>74.23</u>
Proposed	3D-ResNet18	76.10

(b) Results on ShanghaiTech [3]

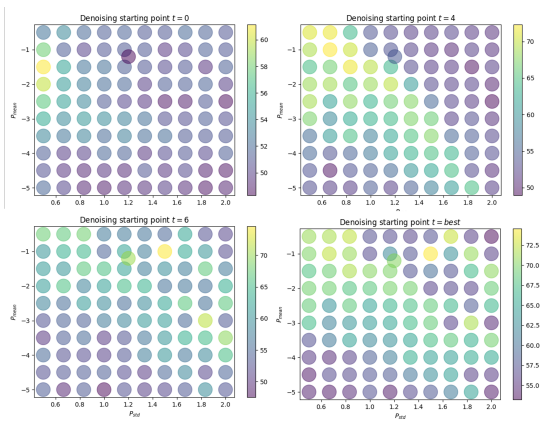


Fig. 2: The effect of noise and the starting point of the reverse process while performing k -diffusion for VAD. The results (AUC %) belong to ShanghaiTech dataset [3] with 3D-ResNet18 backbone.

results in the best performances. The best of all results was observed when $t = 6$. For all other datasets, and backbone combinations, the best results were obtained with $t = 9$. Overall, increasing the t value for a fixed combination of P_{mean} , P_{std} improves the VAD results.

Threshold L_{th} . Given the abnormality threshold $L_{th} = \mu_p + k \sigma_p$, the effect of k was investigated by setting its value to: 0.1, 0.3, 0.5, 0.7, and 1. For 3D-ResNext101, the best results correspond to $k=0.5$ for both ShanghaiTech and UCF Crime. For 3D-ResNet18, $k=0.7$ and 0.1 for ShanghaiTech and UCF Crime, respectively, result in the best scores. The difference between the highest and lowest performances upon changing the value of k is up to 3% AUC when the values of all other hyperparameters are kept the same.

Timestep embeddings. As mentioned before, our method includes transforming the timestep σ_t via Fourier embedding and integrating it into the network through FiLM layers [34]. We also adopted the implementation of [33], which concatenates the timestep embeddings together with its sinus and cosinus values (shown as Proposed w/ [33] in Table 2). The results confirm the better performance of our proposal *wrt* to adapting [33] for all cases while both surpassing the SOTA with the 3D-ResNet18 features.

4.3. Qualitative Results

Fig. 3 shows the anomaly scores produced by our approach for example video clips. As seen, independent of the type of anomaly, the anomaly scores increase immediately when ground-truth anomalies start and decrease right after the ground-truth anomalies finish, showing that the proposed method is favorable for VAD.

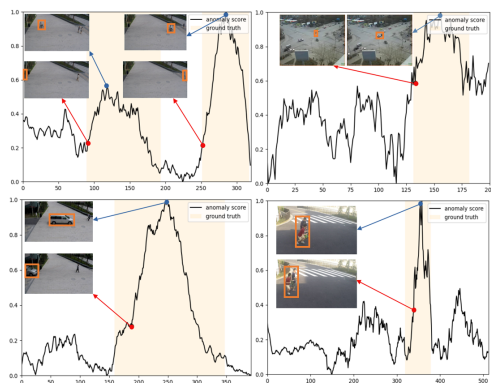


Fig. 3: Predicted frame-level anomaly scores (black), predicted starting frames of an anomaly interval (red arrow), predicted highest anomaly score in an interval (blue arrow), ground-truth anomalies (yellow shadows), spatial reasoning of the anomaly (orange boxes).

5. CONCLUSION

Unsupervised video anomaly detection (VAD) presents the advantage of not requiring data annotation for learning. This solves the problems posed by the heterogeneity of normal and anomalous instances and the scarcity of anomalous data. This paper is the first attempt to investigate the capability of diffusion models for VAD in video surveillance in which we have specifically investigated the use of high reconstruction error as an indicator of abnormality. The experiments performed on popular benchmarks show that the proposed model achieves better performance compared to SOTA generative model: autoencoders independent of the feature extractor used. Our model, although relying only on the reconstruction of the spatial-temporal data, is able, in some cases, to surpass the performance of more complex methods, e.g. the ones performing collaborative learning of generative and discriminative networks. We have also presented a guideline on how the diffusion models (particularly the k -diffusion [19] formulation) should be utilized in terms of its several parameters for VAD. The future work includes investigating the generalization ability of our method in cross-dataset settings.

6. ACKNOWLEDGMENT

We acknowledge the support of the MUR PNRR project FAIR - Future AI Research (PE00000013) funded by the NextGenerationEU. E.R. is partially supported by the PRECRISIS, funded by the EU Internal Security Fund (ISFP-2022-TFI-AG-PROTECT-02-101100539). The work was carried out in the Vision and Learning joint laboratory of FBK and UNITN.

7. REFERENCES

- [1] Sabah Abdulazeez Jebur, Khalid A Hussein, Haider Kadhim Hoomod, Laith Alzubaidi, and José Santamaría, “Review on deep learning approaches for anomaly event detection in video surveillance,” *Electronics*, vol. 12, no. 1, pp. 29, 2022.
- [2] W. Sultani, C. Chen, and M. Shah, “Real-world anomaly detection in surveillance videos,” in *CVPR*, 2018, pp. 6479–6488.
- [3] W. Liu, D. Lian W. Luo, and S. Gao, “Future frame prediction for anomaly detection – a new baseline,” in *Proc. CVPR*, 2018.
- [4] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly detection: A survey,” *ACM CSUR*, vol. 41, no. 3, pp. 1–58, 2009.
- [5] Bahram Mohammadi, Mahmood Fathy, and Mohammad Sabokrou, “Image/video deep anomaly detection: A survey,” *arXiv preprint arXiv:2103.01739*, 2021.
- [6] Cigdem Beyan and Robert B Fisher, “Detecting abnormal fish trajectories using clustered and labeled data,” in *ICIP*, 2013, pp. 1476–1480.
- [7] Jing Ren, Feng Xia, Yemeng Liu, and Ivan Lee, “Deep video anomaly detection: Opportunities and challenges,” in *ICDM workshops*, 2021, pp. 959–966.
- [8] Mahdyar Ravanbakhsh, Moin Nabi, Enver Sangineto, Lucio Marcenaro, Carlo Regazzoni, and Nicu Sebe, “Abnormal event detection in videos using generative adversarial nets,” in *ICIP*, 2017, pp. 1577–1581.
- [9] M. Sabokrou, M. Fayyaz, M. Fathy, and R. Klette, “Deep-cascade: Cascading 3d deep neural networks for fast anomaly detection and localization in crowded scenes,” *IEEE TIP*, vol. 26, no. 4, pp. 1992–2004, 2017.
- [10] Muhammad Zaigham Zaheer, Jin-ha Lee, Marcella Astrid, and Seung-Ik Lee, “Old is gold: Redefining the adversarially learned one-class classifier training paradigm,” in *CVPR*, 2020, pp. 14183–14193.
- [11] Snehashis Majhi, Srijan Das, and François Brémond, “Dam: Dissimilarity attention module for weakly-supervised video anomaly detection,” in *AVSS*, 2021, pp. 1–8.
- [12] Yu Tian, Guansong Pang, Yuanhong Chen, Rajvinder Singh, Johan W Verjans, and Gustavo Carneiro, “Weakly-supervised video anomaly detection with robust temporal feature magnitude learning,” in *ICCV*, 2021, pp. 4975–4986.
- [13] M Zaigham Zaheer, Arif Mahmood, M Haris Khan, Mattia Segu, Fisher Yu, and Seung-Ik Lee, “Generative cooperative learning for unsupervised video anomaly detection,” in *CVPR*, 2022, pp. 14744–14754.
- [14] Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton van den Hengel, “Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection,” in *CVPR*, 2019, pp. 1705–1714.
- [15] Hyunjong Park, Jongyoun Noh, and Bumsub Ham, “Learning memory-guided normality for anomaly detection,” in *CVPR*, 2020, pp. 14372–14381.
- [16] Muhammad Zaigham Zaheer, Jin-Ha Lee, Arif Mahmood, Marcella Astrid, and Seung-Ik Lee, “Stabilizing adversarially learned one-class novelty detection using pseudo anomalies,” *IEEE TIP*, vol. 31, pp. 5963–5975, 2022.
- [17] Pascal Schneider, Jason Rambach, Bruno Mirbach, and Didier Stricker, “Unsupervised anomaly detection from time-of-flight depth images,” in *CVPR workshops*, June 2022, pp. 231–240.
- [18] ZhenJiang Li, Wenbo Yang, Guangli Wu, and Liping Liu, “Unsupervised video anomaly detection based on sparse reconstruction,” in *BDCPS*. Springer, 2021, pp. 994–1001.
- [19] T. Karras, M. Aittala, Timo Aila, and Samuli Laine, “Elucidating the design space of diffusion-based generative models,” in *NeurIPS*, 2022.
- [20] J.-H. Kim, D.-H. Kim, S. Yi, and T. Lee, “Semi-orthogonal embedding for efficient unsupervised anomaly segmentation,” *arXiv preprint:2105.14737*, 2021.
- [21] Bharathkumar Ramachandra, Michael Jones, and Ranga Vasavai, “Learning a distance function with a siamese network to localize anomalies in videos,” in *WACV*, 2020, pp. 2598–2607.
- [22] Ryota Hinami, Tao Mei, and Shin’ichi Satoh, “Joint detection and recounting of abnormal events by learning deep generic knowledge,” in *ICCV*, 2017, pp. 3619–3627.
- [23] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics,” in *ICML*, 2015, pp. 2256–2265.
- [24] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *NeurIPS*, vol. 33, pp. 6840–6851, 2020.
- [25] Chitwan Saharia, William Chan, and Saurabh et al. Saxena, “Photorealistic text-to-image diffusion models with deep language understanding,” *arXiv preprint arXiv:2205.11487*, 2022.
- [26] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon, “Sdedit: Guided image synthesis and editing with stochastic differential equations,” in *ICLR*, 2021.
- [27] Curtis Hawthorne, Ian Simon, and Adam et al. Roberts, “Multi-instrument music synthesis with spectrogram diffusion,” *arXiv preprint arXiv:2206.05408*, 2022.
- [28] S. Chen, P. Sun, Y. Song, and P. Luo, “Diffusiondet: Diffusion model for object detection,” *arXiv preprint:2211.09788*, 2022.
- [29] Zhangxuan Gu, Haoxing Chen, Zhuoer Xu, and Lan et al., “Diffusioninst: Diffusion model for instance segmentation,” *arXiv preprint:2212.02773*, 2022.
- [30] Julia Wolleb, Florentin Bieder, Robin Sandkühler, and Philippe C Cattin, “Diffusion models for medical anomaly detection,” in *MICCAI*. Springer, 2022, pp. 35–45.
- [31] Aapo Hyvärinen and Peter Dayan, “Estimation of non-normalized statistical models by score matching,” *JMLR*, vol. 6, no. 4, 2005.
- [32] Ting Chen, “On the importance of noise scheduling for diffusion models,” *arXiv preprint arXiv:2301.10972*, 2023.
- [33] Shitong Luo and Wei Hu, “Diffusion probabilistic models for 3d point cloud generation,” in *Proc. CVPR*, 2021, pp. 2837–2845.
- [34] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville, “Film: Visual reasoning with a general conditioning layer,” in *Proc. AAAI*, 2018, vol. 32.