

# Masked multi-prediction for multi-aspect anomaly detection

Anonymous authors

Paper under double-blind review

## Abstract

In this paper, we address the anomaly detection problem in the context of heterogeneous normal observations and propose an approach that accounts for this heterogeneity. Although prediction-based methods are common to learn normality, the vast majority of previous work predicts a single outcome, which is generally not sufficient to account for the multiplicity of possible normal observations. To address this issue, we introduce a new masked multi-prediction (MMP) approach that produces multiple likely normal outcomes, and show both theoretically and experimentally that it improves normality learning and leads to a better anomaly detection performance. In addition, we observed that normality can be characterized from multiple aspects, depending on the types of anomalies we would like to detect. Therefore, we propose an adaptation (MMP-AMS) of our approach to cover multiple aspects of normality such as appearance, motion, semantics and location. Since we model each aspect separately, our approach has the advantage of being interpretable and modular, as we can select only a subset of normality aspects. The experiments conducted on several benchmarks show the effectiveness of the proposed approach.

## 1 Introduction

Anomaly detection in the context of videos is crucial for many applications such as video surveillance or autonomous driving for instance. However, it is still an open research problem due to several challenges. The first one is the *scarcity* of anomaly examples and the lack of their corresponding annotations. Indeed, by definition, anomalies are unexpected and usually diverse, therefore, it is infeasible to collect enough representative samples. This makes classical supervised methods ineffective due to the class imbalance issue. Thus, this problem is often considered from the One-Class perspective where a model of normality is learned from normal data only and detects anomalies as outliers. One-class anomaly detection methods can be grouped according to how they model normal data, either *explicitly* or *implicitly*. In the first category, anomalies are detected by measuring their compatibility with a normality model. This category includes probabilistic methods such as diffusion models (Flaborea et al. (2023); Wyatt et al. (2022)) or GANs (Ravanbakhsh et al. (2017); Liu et al. (2018)) which learn a probabilistic model of normality. At inference, samples with low likelihood given the learned density function are considered as anomalies. Another popular family of *explicit* approaches are distance-based methods (Ramachandra et al. (2022); Singh et al. (2023)). Those methods learn an embedding space and the corresponding metric, to ensure that abnormal data are far apart from the normal data. Other approaches perform clustering in some low dimensional space (Ionescu et al. (2019a); Wang et al. (2020); Ionescu et al. (2019b)) to define normality regions, and the anomaly is deduced from the distance to them. Recent approaches such as Park et al. (2020); Gong et al. (2019); Liu et al. (2021); Bergaoui et al. (2022) model the normality via a discrete set of prototypes in the latent space. On the other hand, many existing approaches learn normality patterns *implicitly* via self-supervised pretext tasks that consider different aspects of normality such as appearance and motion (Georgescu et al. (2021a); Barbalau et al. (2022); Wang et al. (2022)). A model is trained on normal data only via those pretext tasks. For a given test sample, the abnormality score can be inferred from the model’s inability to perform the task correctly. There are two main families of *implicit* approaches in the video anomaly detection (VAD) literature: *reconstruction-based* (Gong et al. (2019); Park et al. (2020); Bergaoui et al. (2022); Georgescu et al. (2021b)) and *future prediction-based* (Liu et al. (2018; 2021); Dong et al. (2020); Nguyen & Meunier (2019); Naji et al. (2022); Tang et al. (2020)). Reconstruction-based approaches generally train a neural network which

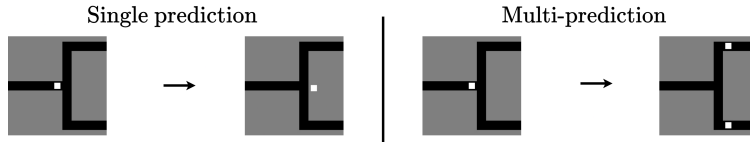


Figure 1: An example illustrating the motivation behind performing multi-prediction instead of a single prediction. We generated a synthetic dataset illustrating a road scenario where a car (in white) is at an intersection, and can perform two actions (turning left or right) which generates two different next states. The left figure shows a single prediction network trained to perform next position prediction, while the right figure shows another one trained to perform multiple predictions using our proposed loss functions.

reproduces the normal training data from a low dimensional space. The fundamental assumption is that the model will not be able to generalize well to anomalies. Differently, future prediction-based approaches train a model to predict a future outcome given the past. The choice between reconstruction-based and prediction-based methods involves trade-offs. While the former reconstruct training data well, they also tend to reconstruct anomalies due to the generalization abilities of neural networks (Gong et al. (2019)). On the contrary, the latter tend to predict anomalies poorly, as the model cannot simply reproduce the input as with reconstruction-based methods. However, these methods have difficulty predicting normal future scenarios because of their diversity. Indeed, many existing future-prediction methods perform single prediction, which is often not enough to characterize all possible future outcomes (Babaeizadeh et al. (2018)). To illustrate this point, let us consider a road scenario (Figure 1), where the task consists in modeling the normal set of car trajectories. If the current state is an intersection, the car can turn left or right. Both cases are normal, but it is impossible to cover both possibilities with a single prediction. In addition, this leads to an abnormal prediction (Figure 1 left). In order to solve this problem, we propose to better model the distribution of possible future scenarios by training a model to perform *multiple predictions* instead of a single one. In order to cover normal possibilities, we use the nearest neighbor loss (Guzmán-rivera et al. (2012); Bhattacharyya et al. (2018); Nguyen et al. (2018)), which is often used in the setting of multiple choice learning, and we introduce a new non-participation loss that ensures a balanced training of all predictors. Our approach offers the advantages of prediction-based methods in terms of poor anomaly prediction, and improves normality learning through multi-prediction. This enables better discrimination between normal and abnormal samples.

Another VAD challenge is that, in order to detect anomalies, it is necessary to determine normality. However, the definition of what is considered normal depends on the context and the application, which also influences the choice of normality aspects to be modeled (e.g. appearance, motion, etc.). While certain aspects of normality are not relevant to detect anomalies for some applications (e.g. the location of a person on a sidewalk when the objective is to detect violence), they tend to be crucial in others (e.g. the location of a person when the objective is to detect jaywalking). Therefore, we model several object-level aspects such as appearance, motion and class-semantics, as well as location-related anomalies. Our approach is interpretable and modular, since it assigns an anomaly score for each aspect. This allows us to adapt our method to applications that require only a subset of the aspects to be modeled while providing information about the anomaly type. In summary, our contributions are as follows: 1) a novel and generic *masked multi-prediction* (MMP) approach for anomaly detection in the context of heterogeneous normal data; 2) an adaptation of MMP to model multiple normality aspects (MMP-AMS); 3) a new *non-participation loss* to better model the multiplicity of normal scenarios; 4) a theoretical analysis and experiments showing the effectiveness of our methodology.

## 2 Related work

**Multi-prediction learning:** also known as multiple choice learning (Lee et al. (2016); Dey et al. (2015); Lee et al. (2017); Guzmán-rivera et al. (2012)) or multiple hypotheses learning (Rupprecht et al. (2017); Nguyen et al. (2018)), is a task where multiple models are learned to produce diverse predictions. During training, samples are assigned to the minimum loss predictor. This technique has been used for tasks involving aleatoric

uncertainty such as future prediction (Bhattacharyya et al. (2018)), human pose estimation (Rupprecht et al. (2017)), image segmentation (Guzmán-rivera et al. (2012); Dey et al. (2015)). In the context of anomaly detection in images (a.k.a novelty detection), Nguyen et al. (2018) proposed a multiple hypothesis auto-encoder which performs multiple reconstructions. Differently from Nguyen et al. (2018), we perform *masked multi-prediction* by combining masking with multi-prediction in order to limit the capacity of the model to recover anomalies. We also propose a new *non-participation loss* that penalizes only the predictors that do not participate in training, and show both theoretically and empirically that it improves the coverage of normal possibilities. In addition, our framework learns multiple aspects of normality, enabling it to detect the corresponding anomaly types in the context of videos. To our knowledge, our work is the first to propose a multiple-prediction approach for VAD.

**Video anomaly detection:** most implicit VAD approaches use self-supervised learning to model normality. The model is trained on a given task using normal data only and it is expected not to generalize well to abnormal samples. Usually, these tasks are designed to characterize a certain aspect of normal data such as appearance or motion, which allows the model to detect the corresponding anomaly types. Hasan et al. (2016) were one of the first to propose a reconstruction-based method by training an auto-encoder to recover handcrafted appearance and motion features. As pointed out by Gong et al. (2019), auto-encoders are able to reconstruct anomalies due to the extrapolation capabilities of deep learning models, which is not suitable for distinguishing between normal and abnormal samples. The reconstruction task can be further constrained using a memory module as proposed in (Gong et al. (2019)). Georgescu et al. (2021b) proposed to train a model via an adversarial objective function, where normal data is well reconstructed and some pseudo-anomalies are explicitly mis-reconstructed. A major self-supervised paradigm to learn normality consists in training a model to perform future frame prediction. Liu et al. (2018) trained a generator using an adversarial objective function to predict a future frame and its optical flow given few past frames. Ravanbakhsh et al. (2017) trained two generators to perform image-to-image translation between RGB and optical flow modalities in order to learn both appearance and motion normality. Liu et al. (2021) introduced a hybrid framework for frame prediction and optical flow reconstruction at the object-level by making use of a pretrained object detector. During inference, the anomaly score is computed based on a sampled future object-level frame. However, one sample may not be representative of the full distribution of future scenarios. In order to cover various modes of this distribution, we propose to train our framework to produce diverse and representative predictions, using the nearest neighbor loss (Guzmán-rivera et al. (2012)) and our novel non-participation loss. Similarly to Liu et al. (2021), we propose to model normality at the object-level to provide better robustness to scene changes and background variety. Recent works introduced other self-supervised tasks for object-level normality learning such as spatio-temporal jigsaw puzzle (Wang et al. (2022)), video event completion (Yu et al. (2020)) or random patch inpainting (Barbalau et al. (2022); Huang et al. (2022)). Barbalau et al. (2022); Georgescu et al. (2021a) proposed to combine multiple tasks such as arrow-of-time prediction, motion irregularity, middle-frame prediction and knowledge distillation to characterize object-level normality. Differently from these previous works, our approach performs multiple predictions instead of a single one for each normality aspect. Furthermore, our method can detect location-dependent anomalies which have not been addressed in the aforementioned methods. To our knowledge, only Doshi & Yilmaz (2020) addressed the modeling of location attributes at the object-level. The authors proposed a non-parametric model of hand-crafted object-level features which included the object position. Differently, we model the distribution of normal bounding boxes separately for each object class which allows to detect class-wise location anomalies.

### 3 Masked multi-prediction for normality modeling: a preliminary study

In this section, we introduce our masked multi-prediction (MMP) approach and motivate it theoretically and experimentally in the context of anomaly detection. We first demonstrate the importance of multi-prediction compared to single prediction. Then, we show the impact of the loss choice to ensure diversity and likelihood of predictions. Finally, we discuss the importance of spatial and temporal masking in improving anomaly detection performance. The proofs of all propositions are provided in the supplementary material. Preliminary experiments to support our analysis have been carried out on the following datasets:

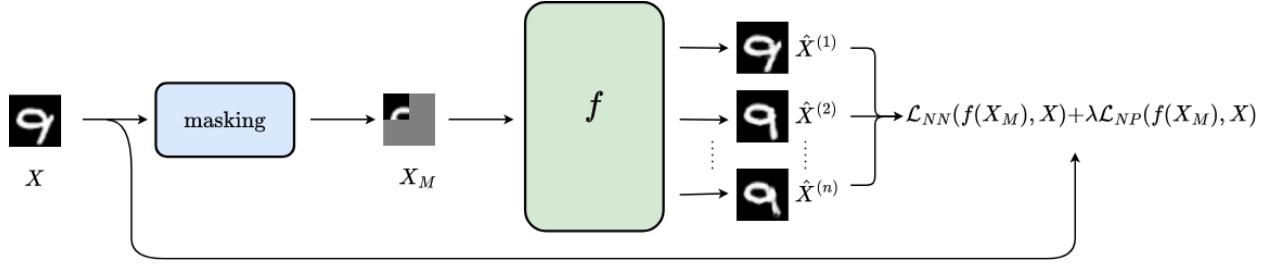


Figure 2: Overview of our methodology: the masked multi-prediction network (MMP) denoted by  $f$  receives a masked sample  $X_M$  and produces  $n$  likely predictions:  $f(X_M) = (\hat{X}^{(k)})_{k \in \llbracket 1, n \rrbracket}$  which are compared to the original sample  $X$ . The model is trained using  $\mathcal{L}_{NN}$  and  $\mathcal{L}_{NP}$ . At inference only  $\mathcal{L}_{NN}$  is used.

**Synthetic roads dataset:** This dataset is used to show the impact of performing multi-prediction instead of single prediction in the context of anomaly detection. Indeed, we generated a synthetic dataset (Figure 1) illustrating a road scenario in which a car is at the intersection of  $m$  roads. In the normal case, the car moves towards one of the roads. More precisely, we assume that the  $x$  axis is oriented from left to right and that the  $y$  axis is oriented from bottom to top. We indicate the position of the car by its 2D coordinates  $(x, y)$ . Our synthetic roads dataset consists of sequences of 2 pairs of coordinates indicating the position of the car at times  $t$  and  $t + 1$ . We will assume that the car is always at position  $p_t = (0, 0)$  at time  $t$  and that it uniformly moves to one of the following positions at time  $t + 1$ :  $\left(p_{t+1}^{(i)} = \left(1, \frac{2 \times (i-1)}{m-1} - 1\right)\right)_{i \in \llbracket 1, m \rrbracket}$ . The parameter  $m \geq 2$  indicates the number of roads the car can take.

**MNIST:** (LeCun et al. (2010)) this dataset is used for showing the importance of masking. It contains handwritten digits from 0 to 9. As this dataset was not designed for the one-class setting, it is adapted by considering a particular class as normal in each case. We use the training-testing split of the dataset and perform training only on samples from the class considered normal. During inference, all test data are used. Test samples from the class on which the model has been trained are considered normal, and other classes are considered abnormal.

### 3.1 Multi-prediction vs. single prediction

In the following analysis, we will focus on the general problem of sample prediction (denoted by  $X$ ) given a masked input  $X_M = X \odot M$ , where  $M$  is the mask applied to  $X$ . The sample prediction problem boils down to learning the conditional probability distribution  $\mathbb{Q} \triangleq \mathbb{P}(X|X_M)$ , which can be done using a model that receives as input a masked sample  $X_M$  and generates candidate samples according to  $\mathbb{P}(X|X_M)$ . However, learning this distribution via a single prediction model does not capture its multi-modality as shown in Figure 1. Indeed, a single prediction model  $g$  is generally trained to predict a sample  $X$  by minimizing:

$$\mathcal{L}_{single}(g(X_M), X) \triangleq \|g(X_M) - X\| \quad (1)$$

However, since  $X$  is not deterministic given  $X_M$ , the network minimizes:

$$\tilde{\mathcal{L}}_{single}(g(X_M)) \triangleq \mathbb{E}_{X \sim \mathbb{Q}}(\|g(X_M) - X\|) \quad (2)$$

This loss can be interpreted as the average distance across all possible samples  $X \sim \mathbb{Q}$ . If mean squared error (MSE) is used, the global minimum of the loss is achieved by a model  $g^*$  that predicts the conditional expectation of  $X$ :

**Proposition 1.** *Let  $g$  be a single prediction model trained using the  $L_2$  norm. The minimum loss is achieved for a model  $g^*$  that predicts the conditional expectation:*

$$g^*(X_M) = \mathbb{E}(X|X_M)$$

This shows that the output of a single prediction model is an average of possible samples ( $X \sim \mathbb{Q}$ ) in the best case, which is sub-optimal. In order to better cover normal possibilities, we propose to train a masked

multi-prediction model (MMP):  $f = (f^{(1)}, f^{(2)}, \dots, f^{(n)}) = (f^{(k)})_{k \in \llbracket 1, n \rrbracket}$  to predict  $n$  likely possibilities using the nearest neighbor objective function (Guzmán-rivera et al. (2012)). An illustration of MMP is provided in Figure 2. The nearest neighbor loss penalizes only the distance to the closest prediction. Formally, the loss between the set of predictions  $f(X_M) = (f^{(k)}(X_M))_{k \in \llbracket 1, n \rrbracket}$  and a sample  $X$  can be written as follows:

$$\mathcal{L}_{NN}(f(X_M), X) \triangleq \min_{k \in \llbracket 1, n \rrbracket} \|f^{(k)}(X_M) - X\| \quad (3)$$

The corresponding expected loss across all possible samples  $X \sim \mathbb{Q}$  is:

$$\bar{\mathcal{L}}_{NN}(f(X_M)) \triangleq \mathbb{E}_{X \sim \mathbb{Q}}(\mathcal{L}_{NN}(f(X_M), X)) \quad (4)$$

Training a MMP model using the previous loss can achieve a better fitting of normal data (i.e a lower training loss) as shown in the following proposition:

**Proposition 2.** *Let  $X$  a sample from  $\mathbb{P}$ ,  $\mathcal{F}$  the space of self-maps of  $[0, 1]^{C \times H \times W}$  and  $f^* \in \arg \min_{f = (f^{(k)})_{k \in \llbracket 1, n \rrbracket} \in \mathcal{F}^n} \bar{\mathcal{L}}_{NN}(f(X_M))$ . The minimum expected loss is lower when using multi-prediction than when using single prediction:*

$$\bar{\mathcal{L}}_{NN}(f^*(X_M)) \leq \bar{\mathcal{L}}_{single}(g^*(X_M))$$

Moreover, a MMP network trained using the nearest neighbor loss provides a better anomaly detection performance than a single prediction model, as shown in section 5.4 and as demonstrated in the case of the synthetic roads dataset. Indeed, Proposition 3 shows that a MMP network trained using the nearest neighbor objective function  $\mathcal{L}_{NN}$  can theoretically achieve a perfect anomaly detection score on the synthetic roads dataset, if it reaches the global minimum of the training loss, whereas a single prediction model cannot, even if the global minimum is reached.

**Proposition 3.** *A MMP model with a number of predictors corresponding to the number of roads  $m$  and reaching the global minimum training loss using  $\mathcal{L}_{NN}$ , achieves an AUC of 100% on the synthetic roads dataset. On the other hand, a single prediction model achieving the global minimum training loss using  $\mathcal{L}_{single}$  cannot achieve an AUC of 100%.*

## 3.2 Impact of loss functions

### 3.2.1 Nearest neighbor loss

The choice of the nearest neighbor objective function is important to ensure diversity of predictions and a better fit to normal data (Guzmán-rivera et al. (2012)). Indeed, if we use instead a naive objective function  $\mathcal{L}_{naive} \triangleq \frac{1}{n} \sum_{k \in \llbracket 1, n \rrbracket} \|f^{(k)}(X_M) - X\|$ . The diversity of predictions is lost, which amounts to making a single prediction, as shown in the following proposition:

**Proposition 4.** *Let  $\bar{\mathcal{L}}_{naive}(f(X_M))$  the expected loss corresponding to  $\mathcal{L}_{naive}(f(X_M), X)$  :*

$$\bar{\mathcal{L}}_{naive}(f(X_M)) = \mathbb{E}_{X \sim \mathbb{Q}} \left( \frac{1}{n} \sum_{k \in \llbracket 1, n \rrbracket} \|f^{(k)}(X_M) - X\| \right) \quad (5)$$

*In case of  $L_2$  norm, the minimum expected loss is achieved by a MMP model:  $f^* = (f^{*(k)})_{k \in \llbracket 1, n \rrbracket}$  such that  $(\forall k \in \llbracket 1, n \rrbracket) : f^{*(k)}(X_M) = \mathbb{E}(X|X_M)$ , which is similar to perform single prediction.*

### 3.2.2 Non-participation loss

In practice,  $\mathbb{P}(X|X_M)$  is often intractable, highly dimensional and we have only access to samples from it (e.g. the future frame observed in a video). Therefore, it is infeasible to train the network via

$\bar{\mathcal{L}}_{NN}(f(X_M))$ , since it involves an expectation over all possible samples  $X$ , which is difficult to compute. Instead, we train the network via  $\mathcal{L}_{NN}(f(X_M), X)$ , where  $X$  is the actual observed sample. This loss encourages the model to produce diverse predictions in order to reduce the distance between the sampled data point  $X$  and its nearest neighbor  $f^{(k^*)}(X_M)$  where  $k^* = \arg \min_k \|f^{(k)}(X_M) - X\|$ . Nevertheless, this objective function only optimizes the best among all predictions ( $f^{(k^*)}(X_M)$ ). Thus, the model receives a sparse signal which may lead to optimizing only predictors which are selected as nearest neighbors during training. Other predictors remain far from the data subspace, as they are never selected and therefore never optimized. In order to overcome this issue, we introduce a novel objective function called *the non-participation loss* in order to optimize these predictors. More specifically, we collect the indices of unoptimized predictors  $\mathcal{U} \subset \llbracket 1, n \rrbracket$  that were not selected as nearest neighbors in the last epoch of training (a small threshold  $\delta$  is used in practice, cf. Section 5.2). Then, we optimize those predictors :  $(f_F^{(p)})_{p \in \mathcal{U}}$  via the non-participation loss that can be written as:

$$\mathcal{L}_{NP}(f(X_M), X) \triangleq \sum_{p \in \mathcal{U}} \|f^{(p)}(X_M) - X\| \quad (6)$$

Thus, we train the network via a weighted combination of the two losses  $\mathcal{L} \triangleq \mathcal{L}_{NN} + \lambda \mathcal{L}_{NP}$ . The following proposition shows that adding the non-participation objective to the nearest neighbor loss can only decrease the prediction error:

**Proposition 5.** *Let  $\mathcal{X}$  the training dataset composed from normal samples. We denote by  $\tilde{f} = (\tilde{f}^{(k)})_{k \in \llbracket 1, n \rrbracket}$  a MMP model optimized via the nearest neighbor objective  $\mathcal{L}_{NN}$  and achieving a loss  $\mathcal{L}_{NN}(\tilde{f}(X_M), X)$  on a sample  $X \in \mathcal{X}$ . Let  $\hat{f} = (\hat{f}^{(k)})_{k \in \llbracket 1, n \rrbracket}$  a MMP model resulting from training the non-optimized predictors of  $\tilde{f}$  via the non-participation loss  $\mathcal{L}_{NP}$ , and  $\mathcal{L}_{NN}(\hat{f}(X_M), X)$  the loss of  $\hat{f}$  on a sample  $X$ . We have  $(\forall X \in \mathcal{X})$ :*

$$\mathcal{L}_{NN}(\hat{f}(X_M), X) \leq \mathcal{L}_{NN}(\tilde{f}(X_M), X)$$

### 3.3 Impact of masking

The role of masking is to ensure that the model does not learn the trivial identity function, which would result in a good reconstruction of both normal and abnormal samples, and therefore a poor discrimination between them. Indeed, masking forces the model to learn the specificities of normal data in order to predict well normal samples while having a poor prediction of abnormal samples. This would result in a better detection of anomalies. It is important to note that masking can be applied either spatially or temporally. For example, patch masking, illustrated in Figure 3, is a clear example of spatial masking. On the other hand, in the case of future frame prediction based on a current frame, masking is implicitly performed at the temporal level. In this case, the sample  $X$  can be considered as a pair of the current and next frames:  $(X_P, X_F)$  and the masked sample  $X_M$  is the current frame  $X_P$ . In order to illustrate the importance of masking, we carried out a comparison between our masked multi-prediction model (MMP) and a reconstruction network, both trained on a single class from MNIST (Figure 3). The training details are provided in the supplementary material. Different masking strategies are presented in Tables 1a,1b. We observe that when masking is used,

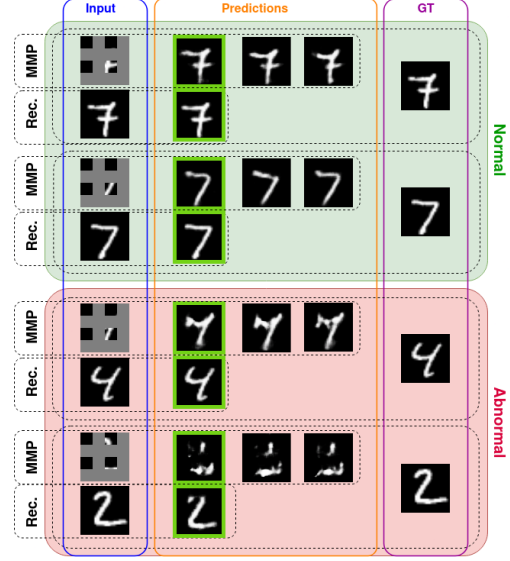


Figure 3: Visualization of two normal samples (class 7) and two abnormal samples (classes 4 and 2) from MNIST. The first column indicates the input of the model, the second column shows the predictions of two models: MMP (ours) which receives a masked input and preforms multiple predictions, and a reconstruction network (Rec.) which receives a non-masked input and preforms a single reconstruction. The nearest predictions are circled in green. The third column indicates the corresponding ground truth.

Table 1: Influence of masking strategies on anomaly detection performance.

(a) Influence of patch size of masks on anomaly detection performance (Mean AUC) on MNIST, using 75% percentage masking.

Exp.	Patch size	Mean AUC
e1	$8 \times 8$	89.0%
e2	$16 \times 16$	<b>94.0%</b>
e3	$32 \times 32$	90.2%

(b) Influence of the percentage of pixels masked on anomaly detection performance (Mean AUC) on MNIST, using a patch size of  $16 \times 16$ .

Exp.	Masking percentage	Mean AUC
e4	0%	85.9%
e5	25%	90.4%
e6	50%	91.3%
e7	75%	<b>94.0%</b>

the recovery of abnormal patterns is more difficult than that of normal patterns, resulting in better anomaly detection performance (e7 vs e4 in Table 1b). Moreover, it can be seen that a high percentage of masking is advantageous for anomaly detection on MNIST. However, a compromise arises when selecting the patch size for masks. Specifically, as the patch size increases, the prediction task becomes more challenging, since the model is required to predict a more global information.

## 4 Multi-aspect normality modeling

In this section, we present our masked multi-prediction framework for modeling appearance, motion and semantics (MMP-AMS), which is an adaptation of MMP to VAD. An illustration of our framework is presented in Figure 4. First we introduce the notations for this section:

We denote the normality aspects as follows:  $F$  for future frame prediction (appearance),  $O$  for optical flow prediction (motion),  $C$  for class prediction (semantics), and  $B_X, B_Y, B_H, B_W$  for the bounding boxes center coordinates, height and width respectively (location). We denote by  $A$  one of these aspects which can belong to the set of appearance, motion and semantics aspects:  $\Gamma = \{F, O, C\}$  or location aspects:  $\Theta = \{B_X, B_Y, B_H, B_W\}$ .  $f_A$  is the network used for modeling an aspect  $A$ , if a network is used for multiple aspects we denote it as  $f_A$ , where  $\mathcal{A}$  is the set of modeled aspects. Let  $X$  be a normal object extracted from a frame at time  $t$ , following the normality distribution  $\mathbb{P}$  and having a bounding box with center coordinates  $(X_{B_X}, X_{B_Y})$ , height  $X_{B_H}$  and width  $X_{B_W}$ . The image of  $X$  at times  $t$  (Present) and  $t + 1$  (Future) are denoted respectively by  $(X_P, X_F)$ . The masked current image of  $X$  is denoted by  $X_{MP} \triangleq X_P \odot M$ , where  $M$  is the mask applied to  $X_P$  and  $\odot$  is the Hadamard product. The one-hot encoding of object classes and the object-level forward optical flow are denoted by  $X_C, X_O$  respectively.

### 4.1 Overview of MMP-AMS

Our framework learns the appearance, motion and semantics aspects of normality, by predicting multiple future frames as well as the corresponding optical flow and class vectors given a masked current frame. In

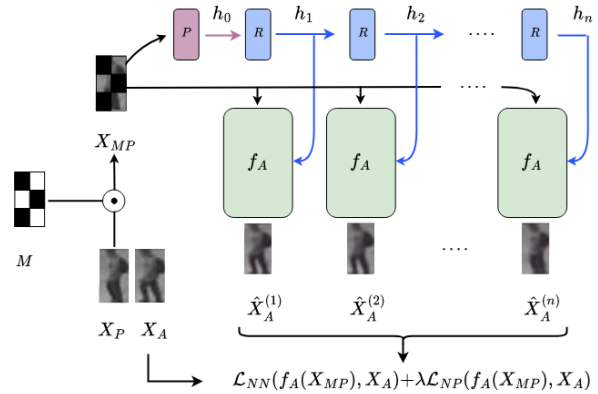


Figure 4: Overview of the proposed masked multi-prediction framework applied to appearance, motion and semantics (MMP-AMS). We show the model predictions for an aspect  $A$ , which can be one of the following aspects: future frame prediction  $F$  (appearance), optical flow prediction  $O$  (motion) or class prediction  $C$  (semantics). We consider  $A = F$  in the illustration. The network receives a masked image at time  $t$  and produces multiple predictions  $(\hat{X}_A^{(k)})_{k \in [1, n]}$  for an aspect  $A$ . The model is trained using  $\mathcal{L}_{NN}$  and  $\mathcal{L}_{NP}$ . At inference only  $\mathcal{L}_{NN}$  is used. Networks of the same color share parameters.

this way, the model learns: appearance features via unmasking, motion features via future prediction and semantic features through class prediction. For each aspect, the model performs multiple predictions in order to take into account the diversity of normal data w.r.t to each aspect. We model the three aspects:  $\Gamma = \{F, O, C\}$  via three MMP networks  $(f_A)_{A \in \Gamma}$  (one for each aspect). Therefore, the full model can be written as:  $f_\Gamma \triangleq (f_A)_{A \in \Gamma} = \left( (f_A^{(k)})_{k \in \llbracket 1, n \rrbracket} \right)_{A \in \Gamma}$ . The model predicts respectively different possible future frames  $(\hat{X}_F^{(k)})_{k \in \llbracket 1, n \rrbracket}$ , optical flows  $(\hat{X}_O^{(k)})_{k \in \llbracket 1, n \rrbracket}$  and classes  $(\hat{X}_C^{(k)})_{k \in \llbracket 1, n \rrbracket}$  given a masked current frame  $X_{MP}$ . The  $k$ -th prediction for an aspect  $A \in \Gamma$  is given by:  $\hat{X}_A^{(k)} \triangleq f_A^{(k)}(X_{MP})$ . The full model  $f_\Gamma$  is trained to minimize the nearest neighbor loss and the non-participation loss for each of the three aspects. Specifically, the following sum is minimized:

$$\mathcal{L}_\Gamma(f_\Gamma(X_{MP}), X) \triangleq \sum_{A \in \Gamma} \mathcal{L}_{NN}(f_A(X_{MP}), X_A) + \lambda \mathcal{L}_{NP}(f_A(X_{MP}), X_A) \quad (7)$$

## 4.2 Network design

In order to allow the number of predictions to be varied without changing the total number of model parameters, we introduce a novel architecture which consists of a recurrent neural network  $R$ , a first state predictor  $P$  as well as three state conditional networks  $(f_A)_{A \in \Gamma} = \left( (f_A(X_{MP}, h_k))_{k \in \llbracket 1, n \rrbracket} \right)_{A \in \Gamma}$  (illustrated in Figure 4). For each aspect  $A$ , we propose to replace the predictors  $(f_A^{(k)})_{k \in \llbracket 1, n \rrbracket}$  by a state conditional network  $f_A$  which receives as input the masked frame  $X_{MP}$  as well as a state  $h_k$  generated using the recurrent neural network  $R$ . Therefore,  $k$ -th predictor  $f_A^{(k)}$  is equivalent to conditioning the network  $f_A$  on the state  $h_k$ . Formally, we have  $(\forall A \in \{F, O, C\}) : \hat{X}_A^{(k)} = f_A^{(k)}(X_{MP}) = f_A(X_{MP}, h_k) = f_A(X_{MP}, R(h_{k-1}))$ . The recurrent architecture ensures that the hidden state  $h_k$  contains information about previously predicted states, and encourages the model to explore new normality regions in order to minimize the nearest neighbor error. The networks architectures and training pseudo-code are provided in the supplementary material.

## 4.3 Location module

In order to detect location related anomalies, we learn the distribution of object positions for each normal class using a simple Gaussian model. More specifically, given a class  $X_C$ , we model the distribution of bounding box centers  $(X_{B_X}, X_{B_Y})$ , height and width  $(X_{B_H}, X_{B_W})$  of objects belonging to  $X_C$ , using 4 Gaussians:  $\mathcal{N}(\alpha_A(X_C), \beta_A(X_C))$  for  $A \in \Theta = \{B_X, B_Y, B_H, B_W\}$ . The mean  $\alpha_A(X_C)$  and the standard deviation  $\beta_A(X_C)$  for a given dimension  $A$  are predicted by a network  $f_\Theta$  using the negative log-likelihood loss:

$$\mathcal{L}_\Theta(f_\Theta(X_C), X) \triangleq \frac{1}{2} \sum_{A \in \Theta} \log(\beta_A(X_C)) + \left( \frac{X_A - \alpha_A(X_C)}{\beta_A(X_C)} \right)^2 \quad (8)$$

## 4.4 Anomaly scoring

At inference time, we perform the same pre-processing steps on a test sample and compute the appearance, motion and semantics anomaly scores  $S_\Gamma(X)$  as well as the location anomaly scores  $S_\Theta(X)$  by summing the  $z$ -scores of the prediction errors across aspects. This allows balancing the contribution of each aspect to the anomaly score. More formally, given an object  $X$  at frame  $\mathcal{F}_t$ :

$$S_\Gamma(X) \triangleq \sum_{A \in \Gamma} w_A \frac{\mathcal{L}_{NN}(f_A(X_{MP}), X_A) - \mu_A}{\sigma_A} \quad (9) \quad S_\Theta(X) \triangleq \sum_{A \in \Theta} w_A \frac{\mathcal{L}_{NLL}(X_A; (\alpha_A, \beta_A)) - \mu_A}{\sigma_A} \quad (10)$$

where  $f_A, \alpha_A, \beta_A$  are the networks and parameters after training.  $\mu_A$  and  $\sigma_A$  are respectively the expectation and the standard deviation of the loss function which is either  $\mathcal{L}_{NN}$  or  $\mathcal{L}_{NLL}$  estimated from normal training



data for a given aspect  $A$ .  $w_A$  is the weight assigned to each aspect, which may vary depending on the application. The final anomaly score is a weighted combination of the appearance, motion, semantics and location anomaly scores:  $S(X) \triangleq S_{\Gamma}(X) + \gamma S_{\Theta}(X)$ . The parameter  $\gamma$  can be chosen depending on whether we aim to detect location anomalies. The frame level score for a frame is the maximum object-level score in the frame.

## 5 Experimental study

### 5.1 Datasets and evaluation metrics

We adopt the following metrics: the Region-Based Detection Criterion (RBDC) and the Track-Based Detection Criterion (TBDC) introduced by Ramachandra & Jones (2020) as an alternative to frame-level AUC widely used in the literature. The latter metric measures the anomaly detection performance at the temporal level only. However, it does not evaluate the capacity of the model to localize anomalies spatially because it does not penalize false positive regions detected in abnormal frames as pointed out by Ramachandra & Jones (2020). We perform experiments on the most commonly used datasets for the one-class and object-centric scenario. **UCSDped2** (Mahadevan et al. (2010)) is a single scene dataset which includes anomalies such as riding a bike and driving a vehicle on a sidewalk. Ramachandra & Jones (2020) provided region-level and track-level annotations for the RBDC and TBDC metrics. **ShanghaiTech** (Luo et al. (2017)) contains scenes of different backgrounds. Anomalies include jumping, running, or stalking on a sidewalk. The region-level and track-level annotations are provided Georgescu et al. (2021b). **CUHK Avenue** (Lu et al. (2013)) is a single scene dataset which consists of videos with abnormal events such as running or walking towards the camera. We use the improved set of annotations proposed by Ramachandra & Jones (2020) which take into account some static anomalies that were not considered in the original annotations.

### 5.2 Implementation details

For a fair comparison to other object-centric approaches, Yolov3 (Redmon & Farhadi (2018)) pretrained on MSCOCO is applied for object detection, using the implementation of MMDetection (Chen et al. (2019)) with an objectness threshold of 0.5 for UCSDped2 since objects have low resolutions and 0.7 for Avenue and ShanghaiTech. The set of objects detected by the used implementation contained very small false positives which are filtered out based on their area (lower than 350 pixels). Optical flow maps are computed using the official implementation of FlowNet2 (Reda et al. (2017)) as in (Liu et al. (2021)). For anomaly scoring, we keep only the optical flow magnitudes since the optical flow orientation maps are not precise enough to be predicted for small displacements. The detected objects as well as the corresponding optical flow maps are resized to 64x64. For the mask  $M$ , we remove 50% of pixels using a grid of 4x4 pixels. Regarding the distance used for anomaly scoring in Section 4, we use the  $L_1$  distance for RGB and optical flow, as well as the cross entropy loss for class probabilities. We train the network for 150 epochs for UCSDped2 and Avenue and for 400 epochs for ShanghaiTech using Adam optimizer with a learning rate of  $10^{-3}$  with a batch size of 640 for the biggest dataset ShanghaiTech and 64 for UCSDped2 and Avenue. We set  $(\gamma, w_I, w_F, w_C, w_{B_X}, w_{B_Y}, w_{B_H}, w_{B_W}, \lambda) = (0, 1, 1, 1, 1, 1, 1, 1, 0.1)$  for UCSDped2 and ShanghaiTech and  $(1, 1, 0.1, 0.1, 1, 1, 1, 1, 0.1)$  for Avenue. We explain the parameters choices in 5.4. For the non-participation loss, we select predictors which have a participation below  $\delta = 5\%$ . Regarding the inference time, our model processes a batch of objects in a frame taken from Avenue in 18ms on a single Nvidia-Titan-X GPU. Therefore, it satisfies the real-time constraints, given real-time object detector and optical flow extractor that can run in parallel. Furthermore, the method requires access to one future frame only to compute the anomaly score which allows online application.

### 5.3 Evaluation results

This section presents the results (Table 2) of our method on the benchmarks UCSDped2, Avenue, ShanghaiTech with respect to recent state-of-the-art object-centric methods. Qualitative results are provided in the supplementary material.

Table 2: Comparison of our approach to state-of-the-art object-centric VAD methods on RBDC and TBDC (%). Best results are in bold, second best are underlined.

Method	UCSDped2		ShanghaiTech		Avenue	
	RBDC	TBDC	RBDC	TBDC	RBDC	TBDC
Ionescu et al. (2019a)	52.8	72.9	20.7	44.5	15.8	27.0
Liu et al. (2021)	-	-	-	-	41.1	<u>86.2</u>
Georgescu et al. (2021b)	69.2	93.2	41.3	78.8	65.1	66.9
Georgescu et al. (2021a)	72.8	91.2	42.8	83.9	57.0	58.3
Bergaoui et al. (2022)	80.1	95.4	51.5	82.2	<u>75.8</u>	70.0
Naji et al. (2022)	77.2	<u>98.5</u>	51.6	84.6	75.3	73.4
Georgescu et al. (2021a) + Ristea et al. (2022)	-	-	40.6	83.5	66.0	64.9
Liu et al. (2018) + Ristea et al. (2022)	-	-	18.5	60.2	20.1	62.3
Liu et al. (2021) + Ristea et al. (2022)	-	-	45.5	84.5	62.3	<b>89.3</b>
Barbalau et al. (2022)	-	-	47.1	<u>85.6</u>	47.8	85.2
Ours (MMP-AMS w/o location module)	<b>84.0</b>	<b>99.0</b>	<b>55.9</b>	<b>85.7</b>	67.4	68.0
Ours (MMP-AMS w/ location module)	<u>80.5</u>	94.2	<u>52.6</u>	83.1	<b>77.7</b>	74.2

**UCSDped2.** On this dataset, MMP-AMS outperforms previous works on RBDC (+3.9p.p), and reaches the state-of-the-art on TBDC. It can be seen that the optical flow prediction is particularly relevant for this dataset (cf. ablation study in Section 5.4) due to the fact that most anomalies have abnormal motion.

**ShanghaiTech.** Unlike other datasets, this one contains multiple scenes for training and testing. Nevertheless it shares a similar normal context across scenes. Since our method is object-centric, it is less sensitive to scene changes. MMP-AMS achieves a significant improvement (+4.3p.p) in terms of RBDC and slightly outperforms other methods on TBDC. The improvements in terms of anomaly localization can be explained by two factors. First, the choice of aspects allows us to detect appearance, semantics and behavior anomalies (cf. ablation study in Section 5.4). Second, we found that multi-prediction is beneficial for this dataset, which is explained by a high scene complexity leading to a multiplicity of normal scenarios.

**Avenue.** This dataset is challenging because it contains several types of anomalies, such as human behavior and unusual objects. In addition, unlike previous datasets, it contains location-dependent anomalies. It is important to note that no method consistently outperforms the others in all metrics. However, MMP-AMS combined with the location module provides a good compromise between RBDC and TBDC. More specifically, our approach achieves the best performance in terms of RBDC (+1.9p.p) and a moderate TBDC. The performance on TBDC can be explained by the fact that our method involves a small temporal context (only two frames), while approaches that surpass ours use a temporal context of at least four consecutive frames. Nevertheless, our method achieves the best TBDC (+0.8 p.p) among methods which use a similar temporal window (Georgescu et al. (2021b); Bergaoui et al. (2022); Naji et al. (2022)). The good performance in RBDC is partly due to simultaneously taking into account the appearance, motion and location aspects which are relevant for this dataset. Moreover, MMP-AMS alone outperforms methods designed for appearance and motion anomalies (Georgescu et al. (2021b;a); Barbalau et al. (2022); Liu et al. (2021); Ristea et al. (2022); Ionescu et al. (2019a)) in terms of RBDC. When combined with the location module, it outperforms all methods under consideration. This shows the complementarity of our two modules for this dataset.

## 5.4 Discussion

### 5.4.1 Impact of multi-prediction

As mentioned in the introduction, there are trade-offs between reconstruction-based methods and future prediction-based methods. While the former reconstruct training data well, they also tend to reconstruct anomalies. On the contrary, the latter predict anomalies poorly, however, they predict less well normal data. Our approach embraces advantages of both families. Indeed, as our model has only access to

a masked current image, it cannot recover anomalies well, and thanks to multi-prediction, it fits normal data better than single-prediction methods, as shown in Proposition 1. In order to empirically verify our claims, we compared MMP-AMS performing only one prediction (v1 in Table 3) with the same framework performing multiple predictions (v2,v3,v4 in Table 3). We observe in Table 3 that a multi-prediction network achieves a lower prediction error than a single prediction network only when  $\mathcal{L}_{NN}$  is used ((v2,v3,v4) vs v1). On the one hand, we observe that the model actually produces similar predictions when trained with  $\mathcal{L}_{naive}$  (v2 in Table 3), which is coherent with Proposition 4. On the other hand, when we train the model using  $\mathcal{L}_{NN}$ , it produces a higher diversity of predictions (v3, v4 in Table 3). Indeed,  $\mathcal{L}_{NN}$  encourages predictor specialization, since it penalizes only the best guess. This observation is consistent with Proposition 1. Nevertheless, optimizing  $\mathcal{L}_{NN}$  alone leads to the non-participation (v3 in Tab. 3) of some predictors in the training, since they are never selected as nearest neighbors. This explains the introduction of the non-participation loss  $\mathcal{L}_{NP}$  (v4 in Table 3) which ensures that all branches get optimized. Empirically, we notice that it allows the model to better fit normal data since it helps to decrease the prediction loss and increases diversity. In terms of anomaly detection performance, we trained the MMP-AMS framework to predict up to 4 predictions on Avenue (Figure 5). We can observe a significant increase in all metrics (RBDC: +4.1p.p, TBDC: +1.8p.p) until 3 predictions. However, in the case of 4 predictions, performance decreases slightly but remains superior to that of a single prediction. This suggests that 3 predictions are enough to model the diversity of normality for this dataset.

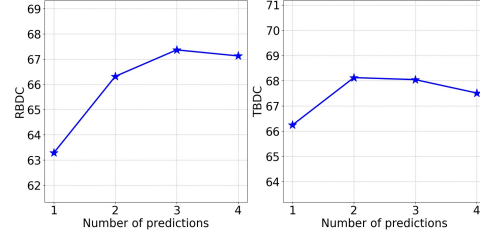


Figure 5: Influence of the number of predictions on MMP-AMS performance (RBDC, TBDC % scores on Avenue).

Table 3: Comparison between different training loss functions using the metrics *participation*, *diversity* and *prediction loss* computed for normal samples from Avenue using the RGB modality. The *participation* is the selection frequency of a predictor. The *diversity* ( $\times 10^2$ ) is the average pixelwise distance between predictions from two different predictors (higher is better). The *prediction loss* is the nearest neighbor loss (lower is better). Green and red respectively indicate best and worst results.

Variant	Losses	Training metrics		
		Participation	Diversity $\uparrow$	Prediction loss $\downarrow$
v1 (1 pred.)	$\mathcal{L}_{naive} = \mathcal{L}_{NN}$	100%	0	0.164
v2 (3 pred.)	$\mathcal{L}_{naive}$	29%, 13%, 58%	0.03	0.165
v3 (3 pred.)	$\mathcal{L}_{NN}$	65%, 35%, 0%	1.4	0.159
v4 (3 pred.)	$\mathcal{L}_{NN} + \lambda \mathcal{L}_{NP}$	35%, 51%, 14%	1.8	0.157

#### 5.4.2 Impact of the choice of normality aspects

In MMP-AMS, we introduced multiple aspects to capture diverse types of normality patterns. This allows the model to jointly learn spatio-temporal fine grained patterns via unmasking and future prediction, as well as the object-level semantics through class prediction. Those aspects are complementary, especially for datasets that contain diverse appearance, motion and semantics anomaly types such as ShanghaiTech. This results in performance improvement when incrementally adding more aspects in the normality modeling (a2, a3, a4 in Table 4). We notice that class and optical flow predictions are less relevant for Avenue dataset. This can be explained by two reasons: 1) most anomalies on this dataset are done by humans for which the class information is not relevant to detect anomalies; 2) optical flow predictions are not enough to characterize complex motion patterns in the scene that would require additional 3D information. Therefore, we give them less weight for anomaly scoring (cf. weights  $w_A$  detailed in Section 5.2). Regarding the masking of the input

Table 4: RBDC and TBDC scores in % obtained by incrementally combining the aspects in MMP-AMS. Best performances are in bold. The normality aspects are denoted as follows:  $F$  for future frame prediction,  $O$  for optical flow prediction,  $C$  for class prediction.

Ablation	Set of aspects	UCSDped2		ShanghaiTech		Avenue	
		RBDC	TBDC	RBDC	TBDC	RBDC	TBDC
a1	$\{F(\text{w/o mask})\}$	61.9	82.1	47.2	82.2	67.5	66.7
a2	$\{F\}$	66.8	85.1	48.4	83.4	68.2	66.7
a3	$\{F, C\}$	76.7	94.9	53.8	83.2	<b>68.5</b>	67.5
a4	$\{F, C, O\}$	<b>84.0</b>	<b>99.0</b>	<b>55.9</b>	<b>85.7</b>	67.4	<b>68.0</b>

(cf. a2 vs. a1 in Table 4), we can observe that it is beneficial for all datasets especially for UCSDped2. This suggests that constraining the prediction task by masking the input makes the prediction even harder for abnormal objects, which leads to a better discrimination between normal and abnormal samples. Concerning the location module, modeling the distribution of bounding boxes (cf. Table 2) significantly improves the results for Avenue dataset. This can be explained by the fact that some anomalies in Avenue are related to the position with respect to the camera, which is not the case for UCSDped2 and ShanghaiTech, that do not include location-dependent anomalies in the definition of what is considered as abnormal. For the sake of consistency with the definition of anomalies in these datasets, the location aspect is not taken into account (cf. parameter  $\gamma = 0$  as detailed in Section 5.2). These results show the importance of defining the aspects of normality that are relevant to each user application, in order to achieve optimal anomaly detection performance. As our approach models these aspects via separate networks, it allows aspects of normality to be weighted according to their relevance to the types of anomalies to be detected. This also provides an end-user explanation of which aspects cause an anomaly score that would trigger an alarm.

### 5.4.3 Limitations and Future work

One downside of our approach is that it depends on supervised object detectors which are usually trained in a closed world manner. However, for some applications, anomaly detection can be aimed at finding objects not seen during training. To address this, it would be interesting to expand our method using open-set object detectors, which are better at adapting to out-of-distribution objects. This could reduce the number of missed anomalies caused by missed detection. Moreover, our method can be improved by adding further normality aspects, which are useful for a given application. For example, it would be interesting to model long-term dependencies such as trajectories. This improvement is valuable for spotting anomalies like loitering, where having a longer context is crucial for an accurate identification.

## 6 Conclusion

In this work, we addressed the problem of modeling a heterogeneous and multi-aspect normality. For this purpose, we proposed a masked multiple prediction approach (MMP) that is adapted to the multiplicity of possible scenarios. We showed both theoretically and experimentally that modeling the distribution of normal data via multiple predictions improves normality learning and anomaly detection performance. We also discussed the importance of determining the relevant aspects of normality for a given application in order to achieve satisfactory performance, and proposed to model several important aspects of normality such as appearance, motion, semantics and localization. As we model each aspect separately, our approach has the advantage of being both interpretable and modular.

## References

Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. Stochastic variational video prediction. 2018. URL <https://openreview.net/forum?id=rk49Mg-CW>.

- Antonio Barbalau, Radu Tudor Ionescu, Mariana-Iuliana Georgescu, Jacob Dueholm, Bharathkumar Ramachandra, Kamal Nasrollahi, Fahad Shahbaz Khan, Thomas B. Moeslund, and Mubarak Shah. Ssmtl++: Revisiting self-supervised multi-task learning for video anomaly detection. 2022.
- Khalil Bergaoui, Yassine Naji, Aleksandr Setkov, Angelique Loesch, Michèle Gouiffès, and Romaric Audigier. Object-centric and memory-guided normality reconstruction for video anomaly detection. pp. 2691–2695, 2022.
- Apratim Bhattacharyya, Bernt Schiele, and Mario Fritz. Accurate and diverse sampling of sequences based on a “best of many” sample objective. pp. 8485–8493, 2018.
- Kai Chen, Jiaqi Wang, and et al. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- Debadeepta Dey, Varun Ramakrishna, Martial Hebert, and J. Andrew Bagnell. Predicting multiple structured visual interpretations. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 2947–2955, 2015. doi: 10.1109/ICCV.2015.337.
- Fei Dong, Yu Zhang, and Xiushan Nie. Dual discriminator generative adversarial network for video anomaly detection. *IEEE Access*, 8, 2020.
- Keval Doshi and Yasin Yilmaz. Continual learning for anomaly detection in surveillance videos. 2020.
- Alessandro Flaborea, Luca Collorone, Guido Maria D’Amely di Melendugno, Stefano D’Arrigo, Bardh Prenkaj, and Fabio Galasso. Multimodal motion conditioned diffusion model for skeleton-based video anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10318–10329, October 2023.
- Mariana-Iuliana Georgescu, Antonio Bărbălău, Radu Tudor Ionescu, Fahad Shahbaz Khan, Marius Claudiu Popescu, and Mubarak Shah. Anomaly detection in video via self-supervised and multi-task learning. *CVPR*, 2021a.
- Mariana Iuliana Georgescu, Radu Ionescu, Fahad Shahbaz Khan, Marius Popescu, and Mubarak Shah. A background-agnostic framework with adversarial training for abnormal event detection in video. *TPAMI*, 2021b.
- 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. 2019.
- Abner Guzmán-rivera, Dhruv Batra, and Pushmeet Kohli. Multiple choice learning: Learning to produce multiple structured outputs. 25, 2012. URL [https://proceedings.neurips.cc/paper\\_files/paper/2012/file/cfbce4c1d7c425baf21d6b6f2babe6be-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/cfbce4c1d7c425baf21d6b6f2babe6be-Paper.pdf).
- M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis. Learning temporal regularity in video sequences. 2016.
- Chaoqin Huang, Qinwei Xu, Yanfeng Wang, Yu Wang, and Ya Zhang. Self-supervised masking for unsupervised anomaly detection and localization. *ArXiv*, abs/2205.06568, 2022.
- Radu Tudor Ionescu, Fahad Shahbaz Khan, Mariana-Iuliana Georgescu, and Ling Shao. Object-centric auto-encoders and dummy anomalies for abnormal event detection in video. 2019a.
- Radu Tudor Ionescu, Sorina Smeureanu, Marius Popescu, and Bogdan Alexe. Detecting abnormal events in video using narrowed normality clusters. 2019b.
- Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.

- Kimin Lee, Changho Hwang, KyoungSoo Park, and Jinwoo Shin. Confident multiple choice learning. *ArXiv*, abs/1706.03475, 2017. URL <https://api.semanticscholar.org/CorpusID:19441737>.
- Stefan Lee, Senthil Purushwalkam, Michael Cogswell, Viresh Ranjan, David J. Crandall, and Dhruv Batra. Stochastic multiple choice learning for training diverse deep ensembles. *ArXiv*, abs/1606.07839, 2016. URL <https://api.semanticscholar.org/CorpusID:655507>.
- W. Liu, D. Lian W. Luo, and S. Gao. Future frame prediction for anomaly detection – a new baseline. 2018.
- Zhian Liu, Yongwei Nie, Chengjiang Long, Qing Zhang, and Guiqing Li. A hybrid video anomaly detection framework via memory-augmented flow reconstruction and flow-guided frame prediction. 2021.
- Cewu Lu, Jianping Shi, and Jiaya Jia. Abnormal event detection at 150 fps in matlab. *ICCV*, 2013.
- Weixin Luo, Wen Liu, and Shenghua Gao. A revisit of sparse coding based anomaly detection in stacked rnn framework. 2017.
- Vijay Mahadevan, Weixin Li, Viral Bhalodia, and Nuno Vasconcelos. Anomaly detection in crowded scenes. 2010.
- Yassine Naji, Aleksandr Setkov, Angélique Loesch, Michèle Gouiffès, and Romaric Audigier. Spatio-temporal predictive tasks for abnormal event detection in videos. pp. 1–8, 2022. doi: 10.1109/AVSS56176.2022.9959669.
- Duc Tam Nguyen, Zhongyu Lou, Michael Klar, and Thomas Brox. Anomaly detection with multiple-hypotheses predictions. 2018. URL <https://api.semanticscholar.org/CorpusID:59316845>.
- Trong-Nguyen Nguyen and Jean Meunier. Anomaly detection in video sequence with appearance-motion correspondence. 2019.
- Hyunjong Park, Jongyoun Noh, and Bumsub Ham. Learning memory-guided normality for anomaly detection. 2020.
- Bharathkumar Ramachandra and Michael Jones. Street scene: A new dataset and evaluation protocol for video anomaly detection. 2020.
- Bharathkumar Ramachandra, Michael J. Jones, and Ranga Raju Vatsavai. A survey of single-scene video anomaly detection. *TPAMI*, 2022.
- Mahdyar Ravanbakhsh, Moin Nabi, Enver Sangineto, Lucio Marcenaro, Carlo Regazzoni, and Nicu Sebe. Abnormal event detection in videos using generative adversarial nets. 2017.
- Fitsum Reda, Robert Pottorff, Jon Barker, and Bryan Catanzaro. flownet2-pytorch: Pytorch implementation of flownet 2.0: Evolution of optical flow estimation with deep networks. <https://github.com/NVIDIA/flownet2-pytorch>, 2017.
- Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- Nicolae-Catalin Ristea, Neelu Madan, Radu Tudor Ionescu, Kamal Nasrollahi, Fahad Shahbaz Khan, Thomas B Moeslund, and Mubarak Shah. Self-supervised predictive convolutional attentive block for anomaly detection. 2022.
- Christian Rupprecht, Iro Laina, Robert DiPietro, Maximilian Baust, Federico Tombari, Nassir Navab, and Gregory D. Hager. Learning in an uncertain world: Representing ambiguity through multiple hypotheses. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 3611–3620, 2017. doi: 10.1109/ICCV.2017.388.
- Ashish Singh, Michael J Jones, and Erik G Learned-Miller. Eval: Explainable video anomaly localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18717–18726, 2023.

- Yao Tang, Lin Zhao, Shanshan Zhang, Chen Gong, Guangyu Li, and Jian Yang. Integrating prediction and reconstruction for anomaly detection. *Pattern Recognit. Lett.*, 2020.
- Guodong Wang, Yunhong Wang, Jie Qin, Dongming Zhang, Xiuguo Bao, and Di Huang. Video anomaly detection by solving decoupled spatio-temporal jigsaw puzzles. pp. 494–511, 2022.
- Ziming Wang, Yuexian Zou, and Zeming Zhang. Cluster attention contrast for video anomaly detection. pp. 2463–2471, 2020. doi: 10.1145/3394171.3413529. URL <https://doi.org/10.1145/3394171.3413529>.
- Julian Wyatt, Adam Leach, Sebastian M. Schmon, and Chris G. Willcocks. Anoddpm: Anomaly detection with denoising diffusion probabilistic models using simplex noise. pp. 650–656, June 2022.
- Guang Yu, Siqi Wang, Zhiping Cai, En Zhu, Chuanfu Xu, Jianping Yin, and Marius Kloft. Cloze test helps: Effective video anomaly detection via learning to complete video events. 2020.