# Latent Space Autoregression for Novelty Detection

Davide Abati    Angelo Porrello    Simone Calderara    Rita Cucchiara

University of Modena and Reggio Emilia

{name.surname}@unimore.it

## Abstract

*Novelty detection is commonly referred to as the discrimination of observations that do not conform to a learned model of regularity. Despite its importance in different application settings, designing a novelty detector is utterly complex due to the unpredictable nature of novelties and its inaccessibility during the training procedure, factors which expose the unsupervised nature of the problem. In our proposal, we design a general framework where we equip a deep autoencoder with a parametric density estimator that learns the probability distribution underlying its latent representations through an autoregressive procedure. We show that a maximum likelihood objective, optimized in conjunction with the reconstruction of normal samples, effectively acts as a regularizer for the task at hand, by minimizing the differential entropy of the distribution spanned by latent vectors. In addition to providing a very general formulation, extensive experiments of our model on publicly available datasets deliver on-par or superior performances if compared to state-of-the-art methods in one-class and video anomaly detection settings. Differently from prior works, our proposal does not make any assumption about the nature of the novelties, making our work readily applicable to diverse contexts.*

## 1. Introduction

Novelty detection is defined as the identification of samples which exhibit significantly different traits with respect to an underlying model of regularity, built from a collection of normal samples. The awareness of an autonomous system to recognize unknown events enables applications in several domains, ranging from video surveillance [7, 11], to defect detection [22] to medical imaging [38]. Moreover, the surprise inducted by unseen events is emerging as a crucial aspect in reinforcement learning settings, as an enabling factor in curiosity-driven exploration [34].

However, in this setting, the definition and labeling of novel examples are not possible. Accordingly, the literature agrees on approximating the ideal shape of the boundary separating normal and novel samples by modeling the intrinsic characteristics of the former. Therefore, prior works tackle such problem by following principles derived from the unsupervised learning paradigm [9, 37, 11, 26, 30]. Due to the lack of a supervision signal, the process of feature extraction and the rule for their normality assessment can only be guided by a proxy objective, assuming the latter will define an appropriate boundary for the application at hand.

According to cognitive psychology [4], novelty can be expressed either in terms of capabilities to *remember* an event or as a degree of *surprisal* [42] aroused by its observation. The latter is mathematically modeled in terms of low probability to occur under an expected model, or by lowering a variational free energy [16]. In this framework, prior models take advantage of either parametric [49] or non-parametric [14] density estimators. Differently, remembering an event implies the adoption of a memory represented either by a dictionary of normal prototypes - as in sparse coding approaches [9] - or by a low dimensional representation of the input space, as in the self-organizing maps [20] or, more recently, in deep autoencoders. Thus, in novelty detection, the remembering capability for a given sample is evaluated either by measuring reconstruction errors [11, 26] or by performing discriminative in-distribution tests [37].

Our proposal contributes to the field by merging remembering and surprisal aspects into a unique framework: we design a generative unsupervised model (i.e., an autoencoder, represented in Fig. 1i) that exploits end-to-end training in order to maximize remembering effectiveness for normal samples whilst minimizing the surprisal of their latent representation. This latter point is enabled by the maximization of the likelihood of latent representations through an autoregressive density estimator, which is performed in conjunction with the reconstruction error minimization. We show that, by optimizing both terms jointly, the model implicitly seeks for minimum entropy representations maintaining its remembering/reconstructive power. While entropy minimization approaches have been adopted in deep neural compression [3], to our knowledge this is the first proposal

tailored for novelty detection. In memory terms, our procedure resembles the concept of prototyping the normality using as few templates as possible. Moreover, evaluating the output of the estimator enables the assessment of the surprisal aroused by a given sample.

## 2. Related work

**Reconstruction-based methods.** On the one hand, many works lean toward learning a parametric projection and reconstruction of normal data, assuming outliers will yield higher residuals. Traditional sparse-coding algorithms [48, 9, 27] adhere to such framework, and represent normal patterns as a linear combination of a few basis components, under the hypotheses that novel examples would exhibit a non-sparse representation in the learned subspace. In recent works, the projection step is typically drawn from deep autoencoders [11]. In [30] the authors recover sparse coding principles by imposing a sparsity regularization over the learned representations, while a recurrent neural network enforces their smoothness along the time dimension. In [37], instead, the authors take advantage of an adversarial framework in which a discriminator network is employed as the actual novelty detector, spotting anomalies by performing a discrete in-distribution test. Oppositely, future frame prediction [26] maximizes the expectation of the next frame exploiting its knowledge of the past ones; at test time, observed deviations against the predicted content advise for abnormality. Differently from the above-mentioned works, our proposal relies on modeling the prior distribution of latent representations. This choice is coherent with recent works from the density estimation community [41, 6]. However, to the best of our knowledge, our work is the first advocating for the importance of such a design choice for novelty detection.

**Probabilistic methods.** A complementary line of research investigates different strategies to approximate the density function of normal appearance and motion features. The primary issue raising in this field concerns how to estimate such densities in a high-dimensional and complex feature space. In this respect, prior works involve hand-crafted features such as optical flow or trajectory analysis and, on top of that, employ both non-parametric [1] and parametric [5, 31, 25] estimators, as well as graphical modeling [17, 23]. Modern approaches rely on deep representations (e.g., captured by autoencoders), as in Gaussian classifiers [36] and Gaussian Mixtures [49]. In [14] the authors involve a Kernel Density Estimator (KDE) modeling activations from an auxiliary object detection network. A recent research trend considers training Generative Adversarial Networks (GANs) on normal samples. However, as such models approximate an implicit density function, they can be queried for new samples

but not for likelihood values. Therefore, GAN-based models employ different heuristics for the evaluation of novelty. For instance, in [38] a guided latent space search is exploited to infer it, whereas [35] directly queries the discriminator for a normality score.

## 3. Proposed model

Maximizing the probability of latent representations is analogous to lowering the surprisal of the model for a normal configuration, defined as the negative log-density of a latent variable instance [42]. Conversely, remembering capabilities can be evaluated by the reconstruction accuracy of a given sample under its latent representation.

We model the aforementioned aspects in a latent variable model setting, where the density function of training samples $p(\mathbf{x})$ is modeled through an auxiliary random variable $\mathbf{z}$, describing the set of causal factors underlying all observations. By factorizing

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}, \qquad (1)$$

where $p(\mathbf{x}|\mathbf{z})$ is the conditional likelihood of the observation given a latent representation $\mathbf{z}$ with prior distribution $p(\mathbf{z})$, we can explicit both the memory and surprisal contribution to novelty. We approximate the marginalization by means of an inference model responsible for the identification of latent space vector for which the contribution of $p(\mathbf{x}|\mathbf{z})$ is maximal. Formally, we employ a deep autoencoder, in which the reconstruction error plays the role of the negative logarithm of $p(\mathbf{x}|\mathbf{z})$, under the hypothesis that $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\tilde{\mathbf{x}}, I)$ where $\tilde{\mathbf{x}}$ denotes the output reconstruction. Additionally, surprisal is injected in the process by equipping the autoencoder with an auxiliary deep parametric estimator learning the prior distribution $p(\mathbf{z})$ of latent vectors, and training it by means of Maximum Likelihood Estimation (MLE). Our architecture is therefore composed of three building blocks (Fig. 1i): an encoder $f(\mathbf{x}; \theta_f)$, a decoder $g(\mathbf{z}; \theta_g)$ and a probabilistic model $h(\mathbf{z}; \theta_h)$:

$$f(\mathbf{x}; \theta_f) : \mathbb{R}^m \to \mathbb{R}^d, \qquad g(\mathbf{z}; \theta_g) : \mathbb{R}^d \to \mathbb{R}^m,$$
$$h(\mathbf{z}; \theta_h) : \mathbb{R}^d \to [0, 1]. \qquad (2)$$

The encoder processes input $\mathbf{x}$ and maps it into a compressed representation $\mathbf{z} = f(\mathbf{x}; \theta_f)$, whereas the decoder provides a reconstructed version of the input $\tilde{\mathbf{x}} = g(\mathbf{z}; \theta_g)$. The probabilistic model $h(\mathbf{z}; \theta_h)$ estimates the density in $\mathbf{z}$ via an autoregressive process, allowing to avoid the adoption of a specific family of distributions (i.e., Gaussian), potentially unrewarding for the task at hand. On this latter point, please refer to supplementary materials for comparison w.r.t. variational autoencoders [19].

With such modules, at test time, we can assess the two sources of novelty: elements whose observation is poorly
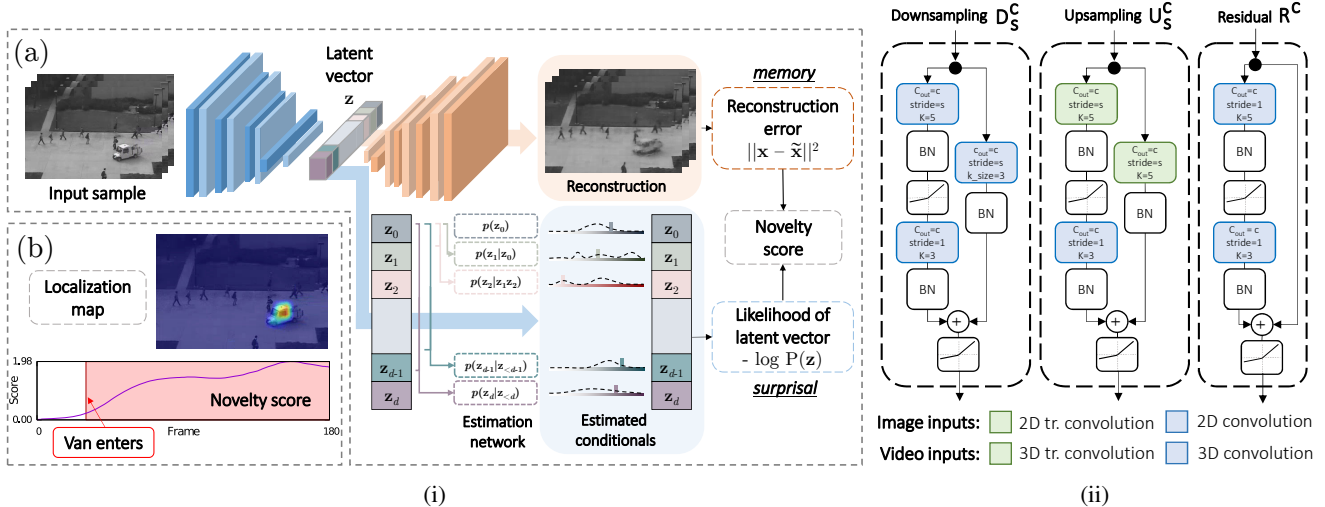
Figure 1: (i) The proposed novelty detection framework. The overall architecture, depicted in (a), consists of a deep autoencoder and an autoregressive estimation network operating on its latent space. The joint minimization of their respective objective leads to a measure of novelty - (b) - obtained by assessing the remembrance of the model when looking to a new sample, combined with its surprise aroused by causal factors. (ii) Building blocks employed in the autoencoder's architecture.

explained by the causal factors inducted by normal samples (i.e., high reconstruction error); elements exhibiting good reconstructions whilst showing surprising underlying representations under the learned prior.

**Autoregressive density estimation.** Autoregressive models provide a general formulation for tasks involving sequential predictions, in which each output depends on previous observations [28, 32]. We adopt such a technique to factorize a joint distribution, thus avoiding to define its landscape a priori [24, 43]. Formally, $p(\mathbf{z})$ is factorized as

$$p(\mathbf{z}) = \prod_{i=1}^{d} p(z_i|\mathbf{z}_{<i}), \qquad (3)$$

so that estimating $p(\mathbf{z})$ reduces to the estimation of each single Conditional Probability Density (CPD) expressed as $p(z_i|\mathbf{z}_{<i})$, where the symbol $<$ implies an order over random variables. Some prior models obey handcrafted orderings [46, 45], whereas others rely on order agnostic training [44, 10]. Nevertheless, it is still not clear how to estimate the proper order for a given set of variables. In our model, this issue is directly tackled by the optimization. Indeed, since we perform autoregression on learned latent representations, the MLE objective encourages the autoencoder to impose over them a pre-defined causal structure. Empirical evidence of this phenomenon is given in the supplementary material.

From a technical perspective, the estimator $h(\mathbf{z}; \theta_h)$ outputs parameters for $d$ distributions $p(z_i|\mathbf{z}_{<i})$. In our implementation, each CPD is modeled as a multinomial over B=100 quantization bins. To ensure a conditional estimate of each

underlying density, we design proper layers guaranteeing that the CPD of each symbol $z_i$ is computed from inputs $\{z_1, \ldots, z_{i-1}\}$ only.

**Objective and connection with differential entropy.** The three components $f$, $g$ and $h$ are jointly trained to minimize $\mathcal{L} \equiv \mathcal{L}(\theta_f, \theta_g, \theta_h)$ as follows:

$$\mathcal{L} = \mathcal{L}_{\text{REC}}(\theta_f, \theta_g) + \lambda \mathcal{L}_{\text{LLK}}(\theta_f, \theta_h)$$
$$= \mathbb{E}_{\mathbf{x}} \left[ \underbrace{||\mathbf{x} - \tilde{\mathbf{x}}||^2}_{\text{reconstruction term}} - \lambda \underbrace{\log(h(\mathbf{z}; \theta_h))}_{\text{log-likelihood term}} \right], \qquad (4)$$

where $\lambda$ is a hyper-parameter controlling the weight of the $\mathcal{L}_{\text{LLK}}$ term. It is worth noting that it is possible to express the log-likelihood term as

$$\mathbb{E}_{\mathbf{z} \sim p^*(\mathbf{z}; \theta_f)} \left[ -\log h(\mathbf{z}; \theta_h) \right]$$
$$= \mathbb{E}_{\mathbf{z} \sim p^*(\mathbf{z}; \theta_f)} \left[ -\log h(\mathbf{z}; \theta_h) + \log p^*(\mathbf{z}; \theta_f) - \log p^*(\mathbf{z}; \theta_f) \right]$$
$$= D_{\text{KL}}(p^*(\mathbf{z}; \theta_f) \parallel h(\mathbf{z}; \theta_h)) + \mathbb{H}[p^*(\mathbf{z}; \theta_f)], \qquad (5)$$

where $p^*(\mathbf{z}; \theta_f)$ denotes the true distribution of the codes produced by the encoder, and is therefore parametrized by $\theta_f$. This reformulation of the MLE objective yields meaningful insights about the entities involved in the optimization. On the one hand, the Kullback-Leibler divergence ensures that the information gap between our parametric model $h$ and the true distribution $p^*$ is small. On the other hand, this framework leads to the minimization of the differential entropy of the distribution underlying the codes produced by the encoder $f$. Such constraint constitutes a crucial point when learning normality. Intuitively, if we think
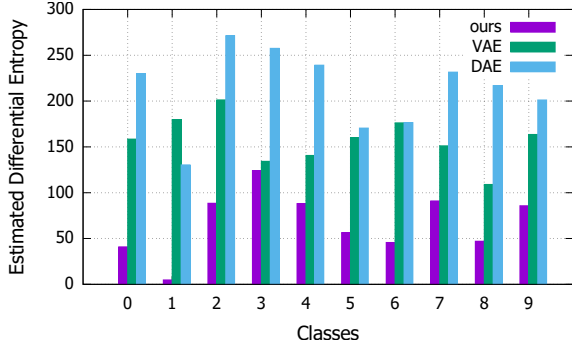
Figure 2: Estimated differential entropies delivered on each MNIST class in the presence of different regularization strategies: our, divergence w.r.t a Gaussian prior (VAE) and input perturbation (DAE). For each class, the estimate is computed on the training samples' hidden representations, whose distribution are fit utilizing a Gaussian KDE in a 3D-space. All models being equal, ours exhibits lower entropies on all classes.

about the encoder as a source emitting symbols (namely, the latent representations), its desired behavior, when modeling normal aspects in the data, should converge to a 'boring' process characterized by an intrinsic low entropy, since surprising and novel events are unlikely to arise during the training phase. Accordingly, among all the possible settings of the hidden representations, the objective begs the encoder to exhibit a low differential entropy, leading to the extraction of features that are easily predictable, therefore common and recurrent within the training set. This kind of features is indeed the most useful to distinguish novel samples from the normal ones, making our proposal a suitable regularizer in the anomaly detection setting.

We report empirical evidence of the decreasing differential entropy in Fig. 2, that compares the behavior of the same model under different regularization strategies.

### 3.1. Architectural Components

**Autoencoder blocks.** Encoder and decoder are respectively composed by downsampling and upsampling residual blocks depicted in Fig. 1ii. The encoder ends with fully connected (FC) layers. When dealing with video inputs, we employ *causal* 3D convolutions [2] within the encoder (i.e., only accessing information from previous time-steps). Moreover, at the end of the encoder, we employ a temporally-shared full connection (TFC, namely a linear projection sharing parameters across the time axis on the input feature maps) resulting in a temporal series of feature vectors. This way, the encoding procedure does not shuffle information across time-steps, ensuring temporal ordering.

**Autoregressive layers.** To guarantee the autoregressive nature of each output CPD, we need to ensure proper connectivity patterns in each layer of the estimator $h$. Moreover, since latent representations exhibit different shapes depending on the input nature (image or video), we propose two different solutions.

When dealing with images, the encoder provides feature vectors with dimensionality $d$. The autoregressive estimator is composed by stacking multiple Masked Fully Connections (MFC, Fig. 3-(a)). Formally, it computes output feature map $\mathbf{o} \in \mathbb{R}^{d \times co}$ (where $co$ is the number of output channels) given the input $\mathbf{h} \in \mathbb{R}^{d \times ci}$ (assuming $ci = 1$ at the input layer). The connection between the input element $\mathbf{h}_i^k$ in position $i$, channel $k$ and the output element $\mathbf{o}_j^l$ is parametrized by

$$
\begin{cases}
w_{i,j}^{k,l} & \text{if } i < j \\
\begin{cases} w_{i,j}^{k,l} & \text{if type = B} \\ 0 & \text{if type = A} \end{cases} & \text{if } i = j \\
0 & \text{if } i > j.
\end{cases}
\tag{6}
$$

Type A forces a strict dependence on previous elements (and is employed only as the first estimator layer), whereas type B masks only succeeding elements. Assuming each CPD modeled as a multinomial, the output of the last autoregressive layer (in $\mathbb{R}^{d \times B}$) provides probability estimates for the $B$ bins that compose the space quantization.

On the other hand, the compressed representation of video clips has dimensionality $t \times d$, being $t$ the number of temporal time-steps and $d$ the length of the code. Accordingly, the estimation network is designed to capture two-dimensional patterns within observed elements of the code. However, naively plugging 2D convolutional layers would assume translation invariance on both axes of the input map, whereas, due to the way the compressed representation is built, this assumption is only correct along the temporal axis. To cope with this, we apply $d$ different convolutional kernels along the code axis, allowing the observation of the whole feature vector in the previous time-step as well as a portion of the current one. Every convolution is free to stride along the time axis and captures temporal patterns. In such operation, named Masked Stacked Convolution (MSC, Fig. 3-(b)), the $i$-th convolution is equipped with a kernel $\mathbf{w}^{(i)} \in \mathbb{R}^{3 \times d}$ kernel, that gets multiplied by the binary mask $\mathbf{M}^{(i)}$, defined as

$$
m_{j,k}^{(i)} \in \mathbf{M}^{(i)} =
\begin{cases}
1 & \text{if } j = 0 \\
1 & \text{if } j = 1 \text{ and } k < i \text{ and type=A} \\
1 & \text{if } j = 1 \text{ and } k \leq i \text{ and type=B} \\
0 & \text{otherwise,}
\end{cases}
\tag{7}
$$

where $j$ indexes the temporal axis and $k$ the code axis. Every single convolution yields a column vector, as a result of its stride along time. The set of column vectors resulting
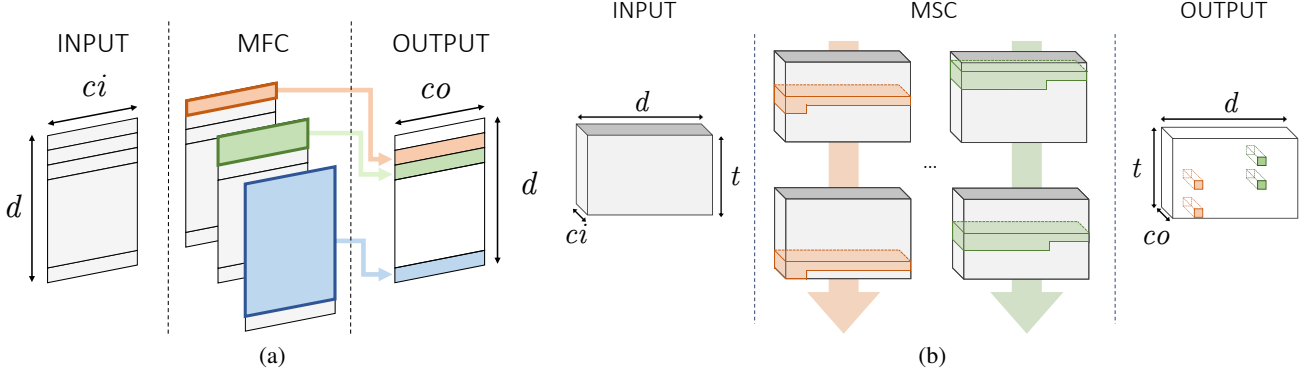
Figure 3: Proposed autoregressive layers, namely the Masked Fully Connection (a, Eq. 6) and the Masked Stacked Convolution (b, Eq. 7). For both layers, we represent type A structure. Different kernel colors represent different parametrizations.

from the application of the $d$ convolutions to the input tensor $\mathbf{h} \in \mathbb{R}^{t \times d \times ci}$ are horizontally stacked to build the output tensor $\mathbf{o} \in \mathbb{R}^{t \times d \times co}$, as follows:

$$\mathbf{o} = \Big\|_{i=1}^{d} [(\mathbf{M}^{(i)} \odot \mathbf{w}^{(i)}) * \mathbf{h}], \qquad (8)$$

where $\|$ represents the horizontal concatenation operation.

## 4. Experiments[1]

We test our solution in three different settings: images, videos, and cognitive data. In all experiments the novelty assessment on the $i$-th example is carried out by summing the reconstruction term ($REC_i$) and the log-likelihood term ($LLK_i$) in Eq. 4 in a single novelty score $NS_i$:

$$NS_i = norm_S(REC_i) + norm_S(LLK_i). \qquad (9)$$

Individual scores are normalized using a reference set of examples $S$ (different for every experiment),

$$norm_S(L_i) = \frac{L_i - \max_{j \in S} L_j}{\max_{j \in S} L_j - \min_{j \in S} L_j}. \qquad (10)$$

Further implementation details and architectural hyperparameters are in the supplementary material.

### 4.1. One-class novelty detection on images

To assess the model's performances in one class settings, we train it on each class of either MNIST or CIFAR-10 separately. In the test phase, we present the corresponding test set, which is composed of 10000 examples of all classes, and expect our model to assign a lower novelty score to images sharing the label with training samples. We use standard train/test splits, and isolate 10% of training samples for

---

[1]Code to reproduce results in this section is released at `https://github.com/aimagelab/novelty-detection`.

validation purposes, and employ it as the normalization set ($S$ in Eq. 9) for the computation of the novelty score. As for the baselines, we consider the following:

- standard methods such as OC-SVM [39] and Kernel Density Estimator (KDE), employed out of features extracted by PCA-whitening;
- a denoising autoencoder (DAE) sharing the same architecture as our proposal, but defective of the density estimation module. The reconstruction error is employed as a measure of normality vs. novelty;
- a variational autoencoder (VAE) [19], also sharing the same architecture as our model, in which the Evidence Lower Bound (ELBO) is employed as the score;
- Pix-CNN [45], modeling the density by applying autoregression directly in the image space;
- the GAN-based approach illustrated in [38].

We report the comparison in Tab. 1 in which performances are measured by the Area Under Receiver Operating Characteristic (AUROC), which is the standard metric for the task. As the table shows, our proposal outperforms all baselines in both settings.

Considering MNIST, most methods perform favorably. Notably, Pix-CNN fails in modeling distributions for all digits but one, possibly due to the complexity of modeling densities directly on pixel space and following a fixed autoregression order. Such poor test performances are registered despite good quality samples that we observed during training: indeed, the weak correlation between sample quality and test log-likelihood of the model has been motivated in [40]. Surprisingly, OC-SVM outperforms most deep learning based models in this setting.

On the contrary, CIFAR10 represents a much more significant challenge, as testified by the low performances of most models, possibly due to the poor image resolution and visual clutter between classes. Specifically, we observe

| | MNIST | | | | | | | CIFAR10 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | OC SVM | KDE | DAE | VAE | Pix CNN | GAN | ours | OC SVM | KDE | DAE | VAE | Pix CNN | GAN | ours |
| 0 | 0.988 | 0.885 | 0.991 | 0.998 | 0.531 | 0.926 | 0.993 | 0.630 | 0.658 | 0.718 | 0.688 | 0.788 | 0.708 | 0.735 |
| 1 | 0.999 | 0.996 | 0.999 | 0.999 | 0.995 | 0.995 | 0.999 | 0.440 | 0.520 | 0.401 | 0.403 | 0.428 | 0.458 | 0.580 |
| 2 | 0.902 | 0.710 | 0.891 | 0.962 | 0.476 | 0.805 | 0.959 | 0.649 | 0.657 | 0.685 | 0.679 | 0.617 | 0.664 | 0.690 |
| 3 | 0.950 | 0.693 | 0.935 | 0.947 | 0.517 | 0.818 | 0.966 | 0.487 | 0.497 | 0.556 | 0.528 | 0.574 | 0.510 | 0.542 |
| 4 | 0.955 | 0.844 | 0.921 | 0.965 | 0.739 | 0.823 | 0.956 | 0.735 | 0.727 | 0.740 | 0.748 | 0.511 | 0.722 | 0.761 |
| 5 | 0.968 | 0.776 | 0.937 | 0.963 | 0.542 | 0.803 | 0.964 | 0.500 | 0.496 | 0.547 | 0.519 | 0.571 | 0.505 | 0.546 |
| 6 | 0.978 | 0.861 | 0.981 | 0.995 | 0.592 | 0.890 | 0.994 | 0.725 | 0.758 | 0.642 | 0.695 | 0.422 | 0.707 | 0.751 |
| 7 | 0.965 | 0.884 | 0.964 | 0.974 | 0.789 | 0.898 | 0.980 | 0.533 | 0.564 | 0.497 | 0.500 | 0.454 | 0.471 | 0.535 |
| 8 | 0.853 | 0.669 | 0.841 | 0.905 | 0.340 | 0.817 | 0.953 | 0.649 | 0.680 | 0.724 | 0.700 | 0.715 | 0.713 | 0.717 |
| 9 | 0.955 | 0.825 | 0.960 | 0.978 | 0.662 | 0.887 | 0.981 | 0.508 | 0.540 | 0.389 | 0.398 | 0.426 | 0.458 | 0.548 |
| avg | 0.951 | 0.814 | 0.942 | 0.969 | 0.618 | 0.866 | **0.975** | 0.586 | 0.610 | 0.590 | 0.586 | 0.551 | 0.592 | **0.641** |

Table 1: AUROC results for novelty detection on MNIST and CIFAR10. Each row represents a different class on which baselines and our model are trained.

that our proposal is the only model outperforming a simple KDE baseline; however, this finding should be put into perspective by considering the nature of non-parametric estimators. Indeed, non-parametric models are allowed to access the whole training set for the evaluation of each sample. Consequently, despite they benefit large sample sets in terms of density modeling, they lead into an unfeasible inference as the dataset grows in size.

The possible reasons behind the difference in performance w.r.t. DAE are twofold. Firstly, DAE can recognize novel samples solely based on the reconstruction error, hence relying on its memorization capabilities, whereas our proposal also considers the likelihood of their representations under the learned prior, thus exploiting surprisal as well. Secondly, by minimizing the differential entropy of the latent distribution, our proposal increases the discriminative capability of the reconstruction. Intuitively, this last statement can be motivated observing that novelty samples are forced to reside in a high probability region of the latent space, the latter bounded to solely capture unsurprising factors of variation arising from the training set. On the other hand, the gap w.r.t. VAE suggests that, for the task at hand, a more flexible autoregressive prior should be pre-

ferred over the isotropic multivariate Gaussian. On this last point, VAE seeks representations whose average surprisal converges to a fixed and expected value (i.e., the differential entropy of its prior), whereas our solution minimizes such quantity within its MLE objective. This flexibility allows modulating the richness of the latent representation vs. the reconstructing capability of the model. On the contrary, in VAEs, the fixed prior acts as a blind regularizer, potentially leading to over-smooth representations; this aspect is also appreciable when sampling from the model as shown in the supplementary material.

Fig. 4 reports an ablation study questioning the loss functions aggregation presented in Eq. 9. The figure illustrates ROC curves under three different novelty scores: i) the log-likelihood term, ii) the reconstruction term, and iii) the proposed scheme that accounts for both. As highlighted in the picture, accounting for both memorization and surprisal aspects is advantageous in each dataset. Please refer to the supplementary material for additional evidence.

### 4.2. Video anomaly detection

In video surveillance contexts, novelty is often considered in terms of abnormal human behavior. Thus, we evaluate our proposal against state-of-the-art anomaly detection models. For this purpose, we considered two standard benchmarks in literature, namely UCSD Ped2 [8] and ShanghaiTech [30]. Despite the differences in the number of videos and their resolution, they both contain anomalies that typically arise in surveillance scenarios (e.g., vehicles in pedestrian walkways, pick-pocketing, brawling). For UCSD Ped, we preprocessed input clips of 16 frames to extract smaller patches (we refer to supplementary materials for details) and perturbed such inputs with random Gaussian noise with $\sigma = 0.025$. We compute the novelty score of each input clip as the mean novelty score among all patches. Concerning ShanghaiTech, we removed the dependency on
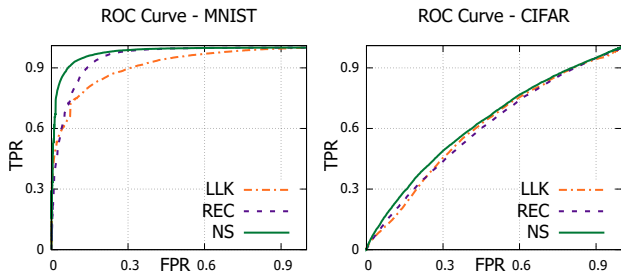


Figure 4: ROC curves delivered by different scoring strategies on MNIST and CIFAR-10 test sets. Each curve is an interpolation over the ten classes.
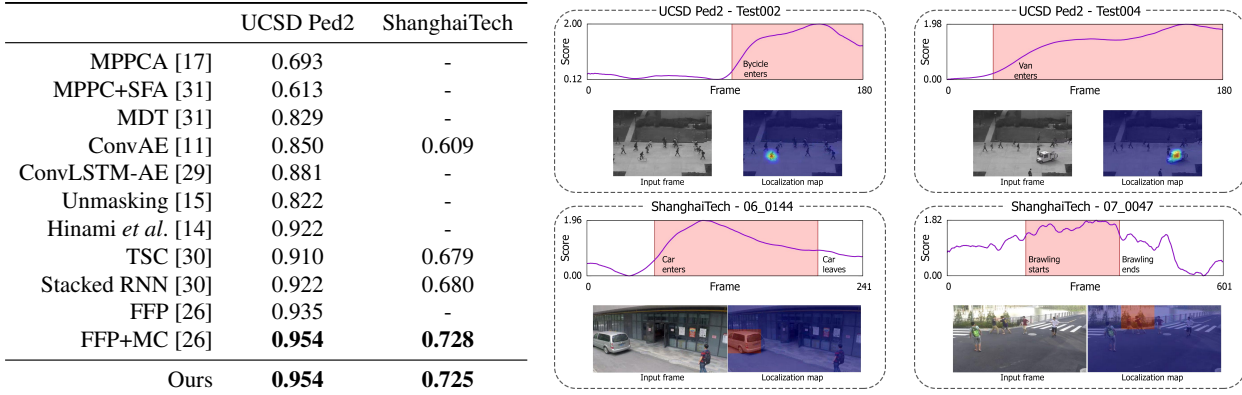
| | UCSD Ped2 | ShanghaiTech |
|---|---|---|
| MPPCA [17] | 0.693 | - |
| MPPC+SFA [31] | 0.613 | - |
| MDT [31] | 0.829 | - |
| ConvAE [11] | 0.850 | 0.609 |
| ConvLSTM-AE [29] | 0.881 | - |
| Unmasking [15] | 0.822 | - |
| Hinami *et al.* [14] | 0.922 | - |
| TSC [30] | 0.910 | 0.679 |
| Stacked RNN [30] | 0.922 | 0.680 |
| FFP [26] | 0.935 | - |
| FFP+MC [26] | **0.954** | **0.728** |
| Ours | **0.954** | 0.725 |



Figure 5: On the left, AUROC performances of our model w.r.t. state-of-the-art competitors. On the right, novelty scores and localizations maps for samples drawn from UCSD Ped2 and ShanghaiTech. For each example, we report the trend of the assessed score, highlighting with a different color the time range in which an anomalous subject comes into the scene.

the scenario by estimating the foreground for each frame of a clip with a standard MOG-based approach and removing the background. We fed the model with 16-frames clips, but ground-truth anomalies are labeled at frame level. In order to recover the novelty score of each frame, we compute the mean score of all clips in which it appears. We then merge the two terms of the loss function following the same strategy illustrated in Eq. 9, computing however normalization coefficients in a per-sequence basis, following the standard approach in the anomaly detection literature. The scores for each sequence are then concatenated to compute the overall AUROC of the model. Additionally, we envision localization strategies for both datasets. To this aim, for UCSD, we denote a patch exhibiting the highest novelty score in a frame as anomalous. Differently, in ShanghaiTech, we adopt a sliding-window approach [47]: as expected, when occluding the source of the anomaly with a rectangular patch, the novelty score drops significantly. Fig. 5 reports results in comparison with prior works, along with qualitative assessments regarding the novelty score and localization capabilities. Despite a more general formulation, our proposal scores on-par with the current state-of-the-art solutions specifically designed for video applications and taking advantage of optical flow estimation and motion constraints. Indeed, in the absence of such hypotheses (FFP entry in Fig. 5), our method outperforms future frame prediction on UCSD Ped2.

## 4.3. Model Analysis

**CIFAR-10 with semantic features.** We investigate the behavior of our model in the presence of different assumptions regarding the expected nature of novel samples. We expect that, as the correctness of such assumptions increases, novelty detection performances will scale accordingly. Such a trait is particularly desirable for applications in which prior beliefs about novel examples

can be envisioned. To this end, we leverage the CIFAR-10 benchmark described in Sec. 4.1 and change the type of information provided as input. Specifically, instead of raw images, we feed our model with semantic representations extracted by ResNet-50 [12], either pre-trained on Imagenet (i.e., assume semantic novelty) or CIFAR-10 itself (i.e., assume data-specific novelty). The two models achieved respectively 79.26 and 95.4 top-1 classification accuracies on the respective test sets. Even though this procedure is to be considered unfair in novelty detection, it serves as a sanity check delivering the upper-bound performances our model can achieve when applied to even better features. To deal with dense inputs, we employ a fully connected autoencoder and MFC layers within the estimation network.

Fig. 6-(a) illustrates the resulting ROC curves, where semantic descriptors improve AUROC w.r.t. raw image inputs (entry "Unsupervised"). Such results suggest that our model profitably takes advantage of the separation between normal and abnormal input representations and scales accordingly, even up to optimal performances for the task under consideration. Nevertheless, it is interesting to note how different degrees of supervision deliver significantly different performances. As expected, dataset-specific supervision increases the AUROC from 0.64 up to 0.99 (a perfect score). Surprisingly, semantic feature vectors trained on Imagenet (which contains all CIFAR classes) provide a much lower boost, yielding an AUROC of 0.72. Such result suggests that, even in the rare cases where the semantic of novelty can be known in advance, its contribution has a limited impact in modeling the normality, mostly because novelty can depend on other cues (e.g., low-level statistics).

**Autoregression via recurrent layers.** To measure the contribution of the proposed MFC and MSC layers described in Sec. 3, we test on CIFAR-10 and UCSD
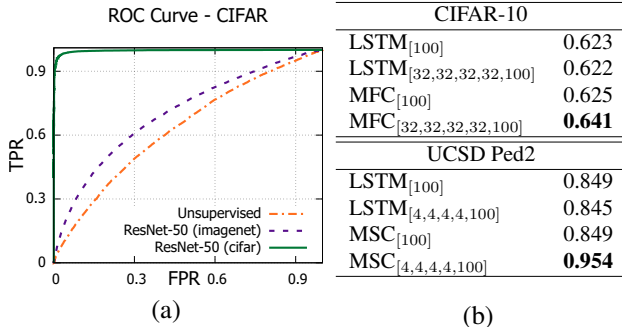
Figure 6: (a) CIFAR-10 ROC curves with semantic input vectors. Each curve is an interpolation among the ten classes. (b) Comparison of different architectures for the autoregressive density estimation in feature space. We indicate with $LSTM_{[F_1, F_2, ..., F_N]}$ - same goes for MFC and MSC - the output shape for each of the $N$ layers composing the estimator. Results are reported in terms of test AUROC.

Ped2, alternative solutions for the autoregressive density estimator. Specifically, we investigate recurrent networks, as they represent the most natural alternative featuring autoregressive properties. We benchmark the proposed building blocks against an estimator composed of LSTM layers, which is designed to sequentially observe latent symbols $\mathbf{z}_{<i}$ and output the CPD of $z_i$ as the hidden state of the last layer. We test MFC, MSC and LSTM in single-layer and multi-layer settings, and report all outcomes in Fig. 6-(b).

It emerges that, even though our solutions perform similarly to the recurrent baseline when employed in a shallow setting, they significantly take advantage of their depth when stacked in consecutive layers. MFC and MSC, indeed, employ disentangled parametrizations for each output CPD. This property is equivalent to the adoption of a specialized estimator network for each $z_i$, thus increasing the proficiency in modeling the density of its designated CPD. On the contrary, LSTM networks embed all the history (i.e., the observed symbols) in their memory cells, but manipulate each input of the sequence through the same weight matrices. In such a regime, the recurrent module needs to learn parameters shared among symbols, losing specialization and eroding its modeling capabilities.

### 4.4. Novelty in cognitive temporal processes

As a potential application of our proposal, we investigate its capability in modeling human attentional behavior. To this end, we employ the DR(eye)VE dataset [33], introduced for the prediction of focus of attention in driving contexts. It features 74 driving videos where frame-wise fixation maps are provided, highlighting the region of the scene attended by the driver. In order to capture the dynamics of attentional patterns, we purposely discard the visual content of
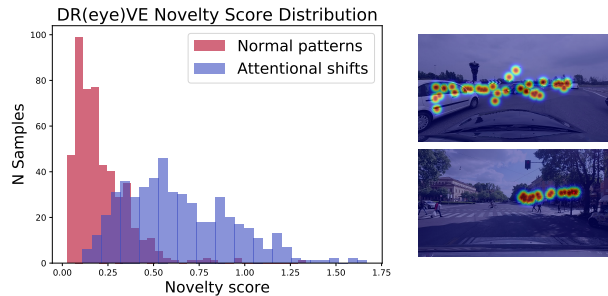


Figure 7: Left, the distribution of novelty scores assigned to normal patterns against attentional shifts labeled within the DR(eye)VE dataset. Right, DR(eye)VE clips yielding the highest novelty score (i.e., clips in which the attentional pattern shifts from the expected behavior). Interestingly, they depict some peculiar situations such as waiting for the traffic light or approaching a roundabout.

the scene and optimize our model on clips of fixation maps, randomly extracted from the training set. After training, we rely on the novelty score of each clip as a proxy for the uncommonness of an attentional pattern. Moreover, since the dataset features annotations of peculiar and unfrequent patterns (such as distractions, recording errors), we can measure the correlation of the captured novelty w.r.t. those. In terms of AUROC, our model scores 0.926, highlighting that novelty can arise from unexpected behaviors of the driver, such as distractions or other shifts in attention. Fig. 7 reports the different distribution of novelty scores for ordinary and peculiar events.

## 5. Conclusions

We propose a comprehensive framework for novelty detection. We formalize our model to capture the twofold nature of novelties, which concerns the incapability to remember unseen data and the surprisal aroused by the observation of their latent representations. From a technical perspective, both terms are modeled by a deep generative autoencoder, paired with an additional autoregressive density estimator learning the distribution of latent vectors by maximum likelihood principles. To this aim, we introduce two different masked layers suitable for image and video data. We show that the introduction of such an auxiliary module, operating in latent space, leads to the minimization of the encoder's differential entropy, which proves to be a suitable regularizer for the task at hand. Experimental results show state-of-the-art performances in one-class and anomaly detection settings, fostering the flexibility of our framework for different tasks without making any data-related assumption.

# References

[1] 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*, 30(3):555–560, 2008. 2

[2] S. Bai, J. Z. Kolter, and V. Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv:1803.01271*, 2018. 4

[3] J. Ballé, V. Laparra, and E. P. Simoncelli. End-to-end optimized image compression. *International Conference on Learning Representations*, 2017. 1

[4] A. Barto, M. Mirolli, and G. Baldassarre. Novelty or surprise? *Frontiers in psychology*, 4:907, 2013. 1

[5] A. Basharat, A. Gritai, and M. Shah. Learning object motion patterns for anomaly detection and improved object detection. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. 2

[6] M. Bauer and A. Mnih. Resampled priors for variational autoencoders. *International Conference on Artificial Intelligence and Statistics*, 2019. 2

[7] S. Calderara, U. Heinemann, A. Prati, R. Cucchiara, and N. Tishby. Detecting anomalies in peoples trajectories using spectral graph analysis. *Computer Vision and Image Understanding*, 115(8):1099–1111, 2011. 1

[8] A. Chan and N. Vasconcelos. Ucsd pedestrian database. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2008. 6

[9] Y. Cong, J. Yuan, and J. Liu. Sparse reconstruction cost for abnormal event detection. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 3449–3456. IEEE, 2011. 1, 2

[10] M. Germain, K. Gregor, I. Murray, and H. Larochelle. Made: Masked autoencoder for distribution estimation. In *International Conference on Machine Learning*, pages 881–889, 2015. 3

[11] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis. Learning temporal regularity in video sequences. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 733–742. IEEE, 2016. 1, 2, 7

[12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 7

[13] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Neural Information Processing Systems*, pages 6626–6637, 2017. 12

[14] R. Hinami, T. Mei, and S. Satoh. Joint detection and recounting of abnormal events by learning deep generic knowledge. In *IEEE International Conference on Computer Vision*, pages 3639–3647, 2017. 1, 2, 7

[15] R. T. Ionescu, S. Smeureanu, B. Alexe, and M. Popescu. Unmasking the abnormal events in video. *IEEE International Conference on Computer Vision*, 2017. 7

[16] L. Itti and P. Baldi. Bayesian surprise attracts human attention. *Vision research*, 49(10):1295–1306, 2009. 1

[17] J. Kim and K. Grauman. Observe locally, infer globally: a space-time mrf for detecting abnormal activities with incremental updates. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 2921–2928. IEEE, 2009. 2, 7

[18] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015. 11

[19] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *International Conference on Learning Representations*, 2014. 2, 5, 11

[20] T. Kohonen. *Self-organization and associative memory*, volume 8. Springer Science & Business Media, 2012. 1

[21] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning*. The MIT Press, 2009. 13

[22] A. Kumar. Computer-vision-based fabric defect detection: A survey. *IEEE Transactions on Industrial Electronics*, 55(1):348–363, 2008. 1

[23] J. Kwon and K. M. Lee. A unified framework for event summarization and rare event detection from multiple views. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1737–1750, 2015. 2

[24] H. Larochelle and I. Murray. The neural autoregressive distribution estimator. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 29–37, 2011. 3

[25] W. Li, V. Mahadevan, and N. Vasconcelos. Anomaly detection and localization in crowded scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(1):18–32, 2014. 2

[26] W. Liu, W. Luo, D. Lian, and S. Gao. Future frame prediction for anomaly detection – a new baseline. *IEEE International Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2, 7

[27] C. Lu, J. Shi, and J. Jia. Abnormal event detection at 150 fps in matlab. In *IEEE International Conference on Computer Vision*, pages 2720–2727. IEEE, 2013. 2

[28] P. Luc, N. Neverova, C. Couprie, J. Verbeek, and Y. LeCun. Predicting deeper into the future of semantic segmentation. In *IEEE International Conference on Computer Vision*, pages 648–657, 2017. 3

[29] W. Luo, W. Liu, and S. Gao. Remembering history with convolutional lstm for anomaly detection. In *IEEE International Conference on Multimedia and Expo*, pages 439–444. IEEE, 2017. 7

[30] W. Luo, W. Liu, and S. Gao. A revisit of sparse coding based anomaly detection in stacked rnn framework. *IEEE International Conference on Computer Vision*, 2017. 1, 2, 6, 7

[31] V. Mahadevan, W. Li, V. Bhalodia, and N. Vasconcelos. Anomaly detection in crowded scenes. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 1975–1981. IEEE, 2010. 2, 7

[32] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016. 3

[33] A. Palazzi, D. Abati, S. Calderara, F. Solera, and R. Cucchiara. Predicting the driver's focus of attention: the dr(eye)ve project. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018. 8

[34] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-supervised prediction. In *International Conference on Machine Learning*, volume 2017, 2017. 1

[35] M. Ravanbakhsh, E. Sangineto, M. Nabi, and N. Sebe. Training adversarial discriminators for cross-channel abnormal event detection in crowds. *arXiv preprint arXiv:1706.07680*, 2017. 2

[36] M. Sabokrou, M. Fayyaz, M. Fathy, Z. Moayed, and R. Klette. Deep-anomaly: Fully convolutional neural network for fast anomaly detection in crowded scenes. *Computer Vision and Image Understanding*, 2018. 2

[37] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli. Adversarially learned one-class classifier for novelty detection. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 3379–3388, 2018. 1, 2

[38] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International Conference on Information Processing in Medical Imaging*, pages 146–157. Springer, 2017. 1, 2, 5

[39] B. Schölkopf, R. C. Williamson, A. J. Smola, J. Shawe-Taylor, and J. C. Platt. Support vector method for novelty detection. In *Neural Information Processing Systems*, 2000. 5

[40] L. Theis, A. v. d. Oord, and M. Bethge. A note on the evaluation of generative models. *International Conference on Learning Representations*, 2016. 5

[41] J. M. Tomczak and M. Welling. Vae with a vamp prior. *International Conference on Artificial Intelligence and Statistics*, 2018. 2

[42] M. Tribus. *Thermostatics and thermodynamics: an introduction to energy, information and states of matter, with engineering applications*. van Nostrand, CS7, 1961. 1, 2

[43] B. Uria, I. Murray, and H. Larochelle. Rnade: The real-valued neural autoregressive density-estimator. In *Advances in Neural Information Processing Systems*, pages 2175–2183, 2013. 3

[44] B. Uria, I. Murray, and H. Larochelle. A deep and tractable density estimator. In *International Conference on Machine Learning*, pages 467–475, 2014. 3

[45] A. van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, and A. Graves. Conditional image generation with pixelcnn decoders. In *Neural Information Processing Systems*, 2016. 3, 5

[46] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu. Pixel recurrent neural networks. *International Conference on Machine Learning*, 2016. 3, 11

[47] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833. Springer, 2014. 7

[48] B. Zhao, L. Fei-Fei, and E. P. Xing. Online detection of unusual events in videos via dynamic sparse coding. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 3313–3320. IEEE, 2011. 2

[49] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International Conference on Learning Representations*, 2018. 1, 2

# Supplementary material

## 6. On the implementation details

Architectures and hyperparameters employed for each experiment are reported in Tab. 2, in terms of the type of blocks, autoregressive layers, mini-batch size, learning rate and weight of the log-likelihood objective. All intermediate layers are Leaky ReLU activated. The objective function is optimized using Adam [18]. All hyperparameters are tuned on a held-out validation set, by minimizing the raw objective (Eq. 4 with $\lambda = 1$).

## 7. On the log-likelihood objective

In this section, we detail how the log-likelihood term (Eq. 4 in the main paper) has been computed and optimized. Importantly, as mentioned in the main paper, we model each CPD through a multinomial. To this aim, we firstly need

|  | MNIST | CIFAR-10 | UCSD Ped2 | ShanghaiTech | DR(eye)VE |
|---|---|---|---|---|---|
| Input Shape | 1,28,28 | 3,32,32 | 1,8,32,32* | 3,16,256,512 | 1,16,160,256 |
| Encoder Network | $D_{2,2}^{32}$ $D_{2,2}^{64}$ $FC^{64}$ $FC^{64}$ | 2D Conv$_{3X3}^{32}$ $R^{32}$ $D_{2,2}^{64}$ $D_{2,2}^{128}$ $D_{2,2}^{256}$ $FC^{256}$ $FC^{64}$ | $D_{1,2,2}^{8}$ $D_{2,1,1}^{12}$ $D_{1,2,2}^{18}$ $D_{2,1,1}^{27}$ $D_{1,2,2}^{40}$ $TFC^{64}$ | $D_{1,2,2}^{8}$ $D_{1,2,2}^{16}$ $D_{2,2,2}^{32}$ $D_{1,2,2}^{64}$ $D_{2,2,2}^{64}$ $TFC^{512}$ $TFC^{64}$ | $D_{1,2,2}^{8}$ $D_{1,2,2}^{16}$ $D_{2,2,2}^{32}$ $D_{1,2,2}^{64}$ $D_{2,2,2}^{64}$ $TFC^{512}$ $TFC^{64}$ |
| Decoder Network | $FC^{64}$ $FC^{64}$ $U_{2,2}^{32}$ $U_{2,2}^{16}$ 2D Conv$_{1x1}^{3}$ | $FC^{256}$ $FC^{256}$ $U_{2,2}^{128}$ $U_{2,2}^{64}$ $U_{2,2}^{32}$ $R^{32}$ 2D Conv$_{1x1}^{3}$ | $TFC^{64}$ $U_{1,2,2}^{40}$ $U_{1,2,2}^{27}$ $U_{2,1,1}^{18}$ $U_{1,2,2}^{12}$ $U_{2,1,1}^{8}$ 3D Conv$_{1x1}^{1}$ | $TFC^{64}$ $TFC^{512}$ $U_{2,2,2}^{64}$ $U_{1,2,2}^{32}$ $U_{2,2,2}^{16}$ $U_{1,2,2}^{8}$ $U_{1,2,2}^{8}$ 3D Conv$_{1x1}^{3}$ | $TFC^{64}$ $TFC^{512}$ $U_{2,2,2}^{64}$ $U_{1,2,2}^{32}$ $U_{2,2,2}^{16}$ $U_{1,2,2}^{8}$ $U_{1,2,2}^{8}$ 3D Conv$_{1x1}^{3}$ |
| Estimator Network | $MFC^{32}$ $MFC^{32}$ $MFC^{32}$ $MFC^{32}$ $MFC^{100}$ | $MFC^{32}$ $MFC^{32}$ $MFC^{32}$ $MFC^{32}$ $MFC^{100}$ | $MSC^{4}$ $MSC^{4}$ $MSC^{4}$ $MSC^{4}$ $MSC^{100}$ | $MSC^{4}$ $MSC^{4}$ $MSC^{100}$ | $MSC^{4}$ $MSC^{4}$ $MSC^{4}$ $MSC^{4}$ $MSC^{100}$ |
| Mini Batch | 256 | 256 | 2760 | 8 | 16 |
| Learning Rate | $10^{-4}$ | $10^{-3}$ | $10^{-3}$ | $10^{-3}$ | $10^{-3}$ |
| $\lambda$ | 1 | 0.1 | 0.1 | 1 | 1 |

*Patches extracted from input clips having shape 1,16,256,384.

Table 2: Architectural and optimization hyperparameters of each setting. We denote with $D_S^C$ (downsampling), $U_S^C$ (upsampling) and $R^C$ (residual) the parametrizations for the employed building blocks (see Fig. 1ii in the main paper). On the one hand, $C$ is the number of output channels, whereas $S$ is the stride of the first convolution in the block. Additionally, $FC^C$ and $TFC^C$ denote dense layers and temporally-shared full connections respectively (in this case, $C$ is the number of output features). Finally, we refer to $MFC^C$ and $MSC^C$ for the proposed autoregressive layers, illustrated in Fig. 3 in the manuscript. For a comprehensive description of each type of layer, please refer to Sec. 3.1 of the main paper.

that the encoder acts as a bounded function. To achieve such desideratum, we simply employ a sigmoidal activation, ensuring that latent representations $\mathbf{z} = f(\mathbf{x}; \theta_f)$ reside in $[0,1]^d$. Therefore, for each $z_j$ with $j = 1, 2, \ldots, d$, we perform a linear quantization of the space $[0,1]$ in $B$ bins (where $B$ is a hyperparameter). This latter step provides for $z_j$ a $B$-dimensional categorical distribution $\phi(z_j)$, highlighting the correct bin to which $z_j$ belongs. For each CPD, such distribution will serve as ground truth for the estimator $h(\mathbf{z}; \theta_h)$, the latter coherently predicting $d$ distributions $p(z_j|\mathbf{z}_{<j})$ across the $B$ bins, employing a softmax activation. This way, as shown in Eq. 11, the $\mathcal{L}_{\text{LLK}}$ loss turns out to be a valid likelihood term, defined as the cross-entropy loss between each one of the estimated CPD and their categorical counterparts:

$$\mathcal{L}_{\text{LLK}}(\theta_f, \theta_h) = \mathbb{E}_{\mathbf{x} \sim P}\left[ -\sum_{j=1}^{d}\sum_{k=1}^{B} \phi(z_j)_k \log(p(z_j|\mathbf{z}_{<j})_k) \right].$$
(11)

It is worth noting that multinomials are just one of the plausible models for the CPDs. Indeed, if we replace them with Gaussians, the overall framework would leave standing. However, as we observed in different trials, this choice does not yield considerable improvements but rather numerical instabilities, as described in prior works [46].

## 8. On the relations to Variational Autoencoders

Our framework yields some similarities with the Variational Autoencoder (VAE) [19]. Indeed, they both approximate the integral of Eq. 1 in the main paper through the minimization of the reconstruction error under a regularization constraint involving a prior distribution on latent vectors. However, it is worth noting several fundamental distinctions. Firstly, our model does not provide an explicit strategy to sample from the posterior distribution, thus resulting in a deterministic mapping from the input to the hidden representation. Secondly, while VAE specifies an explicit and adamant form for modeling the prior $p(\mathbf{z})$, in our formulation its landscape is free from any assumption and directly learnable as a result of the estimator's autoregressive nature. On this point, our proposal leads to two beneficial aspects. First, as the VAE forces the codes' distribution to match the prior, their differential entropy converges to be the same as the prior. This behavior results in approximately stationary entropies across different settings (appreciable in Fig. 2 in the main paper, where we discuss the intuition behind the entropy minimization within a novelty detection task). Secondly, the employment of a too simplistic prior may lead to over-regularized representations, whereas our proposal is less prone to such risk. Empirical evidence of
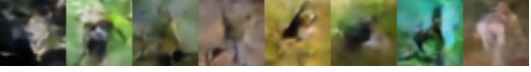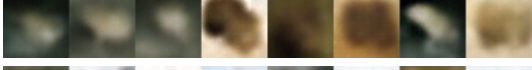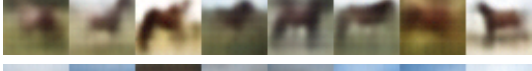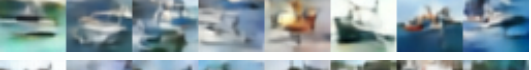
| FID | VAE Samples | Our Samples | FID |
|---|---|---|---|
| 149.72 | | | 72.96 |
| 172.02 | | | 72.53 |
| 181.56 | | | 76.27 |
| 188.37 | | | 67.33 |
| 202.06 | | | 68.33 |
| 207.47 | | | 73.92 |
| 186.48 | | | 62.26 |
| 220.79 | | | 64.38 |
| 164.36 | | | 52.53 |
| 204.84 | | | 67.17 |



Figure 8: For all CIFAR-10 classes (organized in different rows), we report images sampled from VAEs (left) and the proposed autoencoders with autoregressive priors. As can be seen, our samples visually exhibit fine-grained details and sharpness, differently from the heavily blurred ones coming from VAEs. Finally, the over-regularization arising from VAE is confirmed when looking at FID scores (at the extremes of the figure, the lower, the better).

such behavior can also be appreciated in Fig. 8, where we draw new samples from VAE and our model, both of which has been trained on CIFAR-10. All settings being equal, our hallucinations are visually much more realistic than the ones coming from VAEs, the latter leading to over-smooth shapes and lacking any details, as further confirmed by the substantial differences in Fréchet Inception Distance (FID) scores [13].

## 9. On the dual nature of novelty

In this section, we stress how significant is the presence of both terms for obtaining a highly discriminative novelty score (NS, Eq. 9 in the main paper): namely the reconstruction error (REC), modeling the memory capabilities, and the log-likelihood term (LLK), capturing the surprisal inducted from latent representations. Aiming to reinforce this latter point, just briefly illustrated in Fig. 4 of the manuscript, we report in Tab. 3 performances - expressed in AUROC - delivered by different scoring strategies on each setting mentioned in the main paper. Except for ShanghaiTech, we systematically observe a reward in accounting for both as-

pects. Furthermore, for MNIST and CIFAR-10, we find particularly interesting the gap in performance arising from our reconstruction error w.r.t. the one arising from the denoising autoencoder (DAE) variants (0.942 and 0.590 for the two datasets respectively, as reported in Tab. 1 of the main paper). In this respect, we gather new evidence supporting that surprisal minimization acts as a novelty-oriented

| | LLK | REC | NS |
|---|---|---|---|
| MNIST | 0.926 | 0.949 | **0.975** |
| CIFAR-10 | 0.627 | 0.603 | **0.641** |
| UCSD Ped2 | 0.933 | 0.909 | **0.954** |
| ShanghaiTech | 0.695 | **0.726** | 0.725 |
| DR(eye)VE | 0.917 | 0.863 | **0.926** |

Table 3: For each setting, AUROC performances under three different novelty scores: i) the log-likelihood term (LLK), ii) the reconstruction term (REC), and iii) the proposed scheme accounting for both (NS).
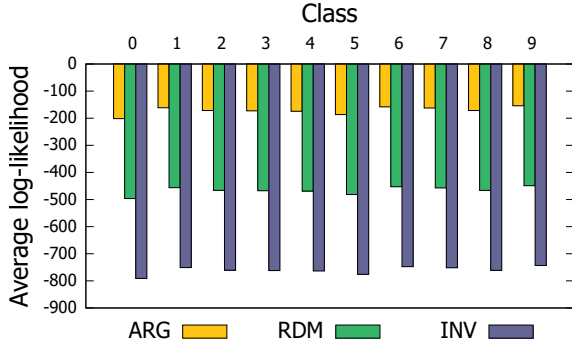
Figure 9: Sample training log-likelihood of a Bayesian Network modeling the distribution of latent codes produced by the encoder of our model trained on MNIST digits. When the BN structure resembles the autoregressive order imposed during training, a much higher likelihood is achieved. This behavior is consistent in all classes and supports the capability of the encoder to produce codes that respect a pre-imposed autoregressive structure.

regularizer for the overall architecture, as it improves the discriminative capability of the reconstruction (as already conjectured in Sec. 4.1 of the main paper).

## 10. On the causal structure of representations

We now investigate the capability of our encoder to produce representations that respect the autoregressive causal structure imposed by the LLK loss (mentioned in Sec. 3 of the main paper). To this aim, we extract representations out of the ten models trained on MNIST digits and fit their distribution using a structured density estimator. Specifically, we employ Bayesian Networks (BNs) with different autoregressive structures. In this respect, each BN is modeled with Linear Gaussians [21], s.t. each CPD $p(z_i|Pa(z_i))$ with $i = 1, 2, \ldots, d$ is given by:

$$p(z_i|Pa(z_i)) = \mathcal{N}(z_i \mid w_0^{(i)} + \sum_{z_j \in Pa(z_i)} w_j^{(i)} z_j, \sigma_i^2), \quad (12)$$

where each $w_j^{(i)}$, $\sigma_i^2$ are learnable parameters. We indicate with $Pa(z_i)$ the parent variables of $z_i$ in the BN. The previous equation holds for all nodes, except for the root one, which is modeled through a Gaussian distribution. Concerning the BN structure, we test:

- Autoregressive order: the BN structure follows the autoregressive order imposed during training, namely $Pa(z_i) = \{z_j \mid j = 1, 2, \ldots, i-1\}$

- Random order: the BN structure follows a random autoregressive order.

- Inverse order: the BN structure follows an autoregressive order which is the inverse with respect to the one

imposed during training, namely $Pa(z_i) = \{z_j \mid j = i+1, i+2, \ldots, d\}$

It is worth noting that, as the three structures exhibit the same number of edges and independent parameters, the difference in their fitting capabilities is only due to the causal order imposed over variables.

Fig. 9 reports the sample training log-likelihood of all BN models. Remarkably, the autoregressive order delivers a better fit, supporting the capability of the encoder network to extract features with learned autoregressive properties. Moreover, to show that this result is not due to overfitting or other lurking behaviors, we report in Tab. 4 log-likelihoods for training, validation and test set.

## 11. On the entropy minimization

To provide an additional grasp about the role of the representation's entropy minimization, we focus on a single MNIST digit (class 7) and report in Fig. 10 some randomly sampled reconstructions from the training set. Such reconstructions are learned under three different regularization regimes, represented by different weights on the log-likelihood objective ($\lambda$, Eq. 4 in the main paper). As shown in Fig. 10, higher degrees of regularization (i.e., stricter constraints on entropy) deliver near mode-collapsed reconstructions, losing sharp variations in favor of capturing fewer prototypes for the input distribution.

## 12. On the complexity of autoregressive layers

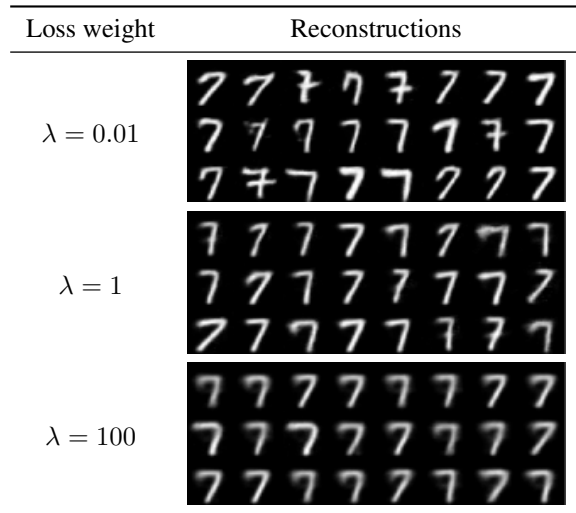In this section, we briefly discuss the complexity of Masked Fully Connected (MFC) and Masked Stacked Convolution

| Loss weight | Reconstructions |
|---|---|
| $\lambda = 0.01$ |  |
| $\lambda = 1$ |  |
| $\lambda = 100$ |  |

Figure 10: MNIST reconstructions delivered by different values of $\lambda$, the latter controlling the impact of the differential entropy minimization.

| | | Classes | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| ARG | Train | -201.60 | -161.60 | -171.43 | -172.73 | -174.17 | -186.48 | -158.22 | -162.37 | -171.65 | -154.11 |
| | Val | -200.96 | -160.38 | -170.10 | -172.29 | -173.85 | -185.25 | -157.22 | -162.20 | -171.42 | -154.02 |
| | Test | -200.89 | -159.73 | -169.64 | -170.75 | -172.40 | -184.27 | -157.74 | -161.65 | -170.10 | -152.70 |
| RDM | Train | -496.33 | -456.34 | -466.16 | -467.47 | -468.90 | -481.21 | -452.95 | -457.10 | -466.39 | -448.84 |
| | Val | -495.69 | -455.11 | -464.83 | -467.02 | -468.58 | -479.98 | -451.95 | -456.93 | -466.15 | -448.75 |
| | Test | -495.62 | -454.47 | -464.37 | -465.48 | -467.13 | -479.00 | -452.48 | -456.38 | -464.83 | -447.43 |
| INV | Train | -791.06 | -751.07 | -760.89 | -762.20 | -763.63 | -775.94 | -747.68 | -751.83 | -761.12 | -743.57 |
| | Val | -790.42 | -749.84 | -759.56 | -761.75 | -763.31 | -774.71 | -746.68 | -751.66 | -760.88 | -743.48 |
| | Test | -790.35 | -749.20 | -759.11 | -760.22 | -761.86 | -773.73 | -747.21 | -751.12 | -759.56 | -742.16 |

Table 4: Sample log-likelihood obtained by different BN structures when fitting MNIST representations. Each BN is trained on latent codes computed from the training set of a single class, following either the autoregression order (ARG), a random order (RDM) or the order inverse to autoregression (INV). We report the log-likelihood also on the validation and test set. For train-val-test split, see Sec 4.1 of the paper. Only "normal" test samples are used in this evaluation.

(MSC) layers (Fig. 3 of the main paper)[2]: adhering to the notation introduced in Sec. 3 from the main paper, MFC exhibits $\frac{d^2+d}{2} \cdot ci \cdot co + d \cdot co$ trainable parameters and a computational complexity $\mathcal{O}(d^2 \cdot ci \cdot co)$. MSC, instead, features $\frac{3d^2+d}{2}ci \cdot co + d \cdot c_o$ free parameters and a time complexity $\mathcal{O}(d^2 \cdot ci \cdot co \cdot t)$.

## 13. On the localizations and novelty scores in video anomaly detection

We show in Fig. 11 other qualitative evidence of the behavior of our model in video anomaly detection settings, namely UCSD Ped2 and ShanghaiTech.

---

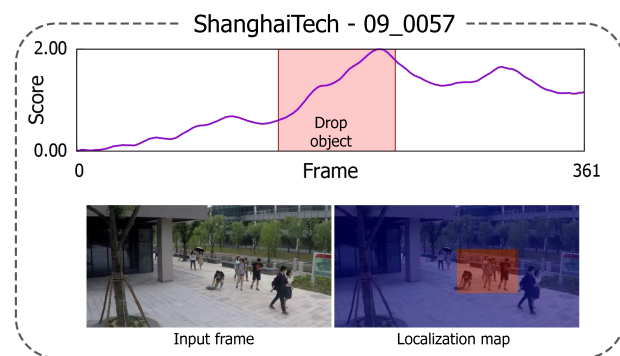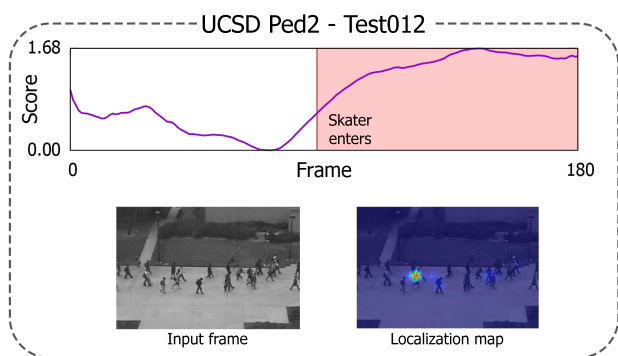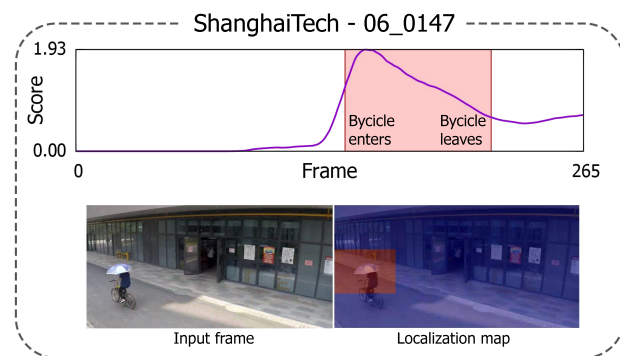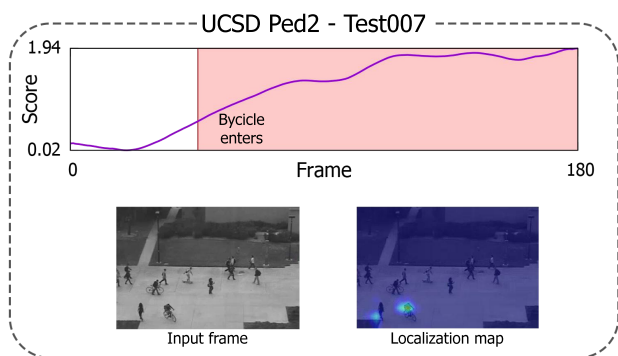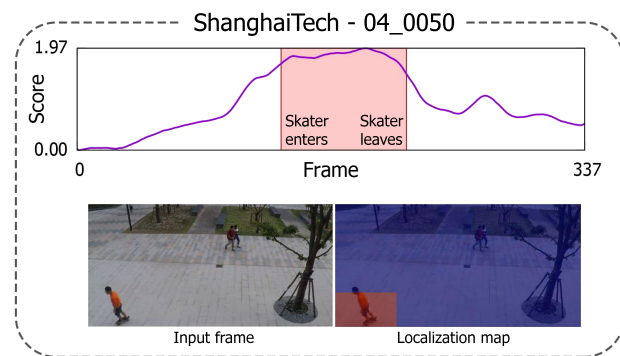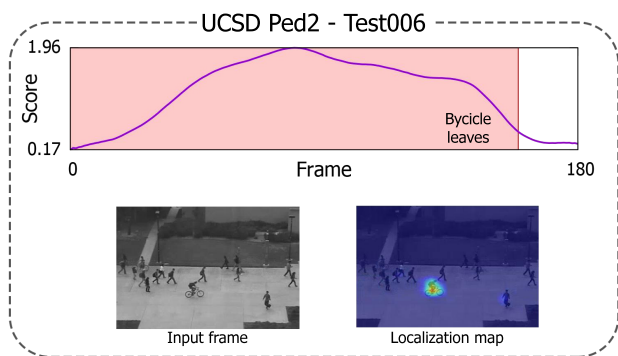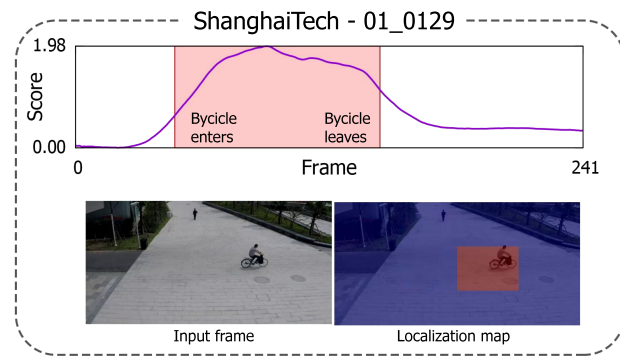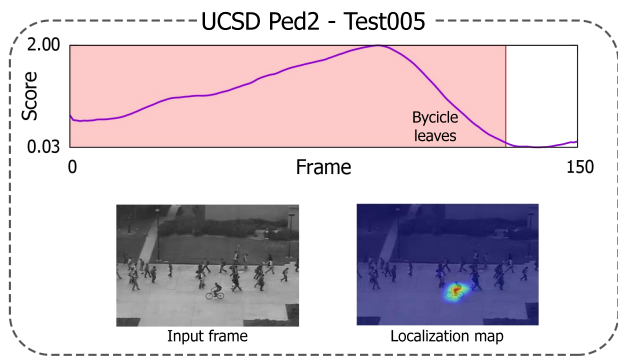[2]We refer to the type 'B' of both layers, since it is an upper bound to the type 'A'

Figure 11: Novelty scores and localizations maps for several test clips from UCSD Ped2 (left) and ShanghaiTech (right).