

Predicting Next Local Appearance for Video Anomaly Detection

Pankaj Raj Roy¹ Guillaume-Alexandre Bilodeau¹ Lama Seoud²
¹LITIV ²Institute of Biomedical Engineering
Polytechnique Montréal, Montréal, Canada
{pankaj-raj.roy,gabilodeau,lama.seoud}@polymtl.ca

Abstract

We present a local anomaly detection method in videos. As opposed to most existing methods that are computationally expensive and are not very generalizable across different video scenes, we propose an adversarial framework that learns the temporal local appearance variations by predicting the appearance of a normally behaving object in the next frame of a scene by only relying on its current and past appearances. In the presence of an abnormally behaving object, the reconstruction error between the real and the predicted next appearance of that object indicates the likelihood of an anomaly. Our method is competitive with the existing state-of-the-art while being significantly faster for both training and inference and being better at generalizing to unseen video scenes.

1 Introduction

Video anomaly detection (VAD) is one of the behavioral recognition tasks that can help ensuring a safer society by detecting quickly various potential or ongoing incidents. Generally, VAD consists in identifying video frames containing spatio-temporal regions deviating from the normal behavior expected for a given scene. Since anomalies are rare and unpredictable, many researchers (*e.g.* [1–8]) aim at identifying them as novelties in videos by relying only on normal known observations during training and by detecting anomalies as observations that fall outside of the known normal boundary.

Most of the recent VAD methods [1–5, 7] are holistic. The normal pixel-level behaviors in videos are learned by training unsupervised generative neural networks minimizing the reconstruction and/or the prediction of the appearance/motion cues of the whole video frames. Color intensities and/or dense optical flows of video frames are widely used as appearance and/or motion cues respectively. In reconstruction-based methods [1, 4], these normal appearance and/or motion cues are learned through a convolutional auto-encoder (CAE) trained by minimizing the reconstruction error between the input and the predicted output. In the case of prediction-based approaches, which generally perform better compared to the reconstruction-based ones, CAEs [3] or convolutional LSTMs [2] with two decoders can reconstruct the input and predict

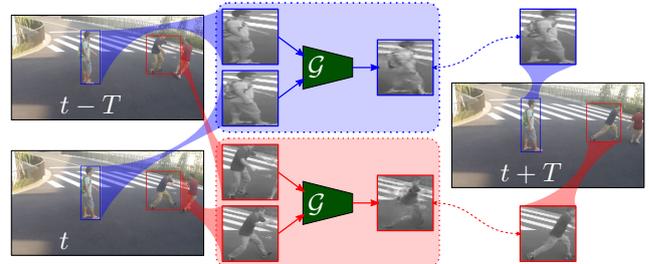


Figure 1. Overview of our method. For each object, the generator \mathcal{G} predicts the appearance of that object in the next frame $t+T$ by using its appearances in the past $t-T$ and current t frames. For an anomaly (person fighting), \mathcal{G} is expected to produce bad reconstruction (as shown in red).

the appearance of a different video frame by solely using some consecutive appearance cues, whereas U-Net frameworks [5, 7] use both the appearance and motion cues. Even though they can perform VAD in real-time, most of these holistic methods are computationally heavy and tend to overfit on the background appearance in the training set. To alleviate the latter issue, model agnostic meta-learning (MAML [23]) was proposed in [8] to extend a holistic VAD model beyond known video scenes. However, meta-learning scheme can be computationally very expensive when working with a high number of video scenes.

To overcome the issues of holistic methods, some researchers [6, 9] instead consider object-centric pixel-level features because they allow VAD systems to generalize better across different scenes. Indeed, each of these object-centric features corresponds to a spatio-temporal region occupied by one of the objects of interest (*e.g.* pedestrians, cyclists, cars, *etc.*), ignoring the non-object background information in videos that prevents VAD to generalize normal local behaviors to unseen videos. These approaches focus on the detection of object-centric local anomalies in videos, independently of other objects [10], but actually they can also capture close-range interactions between objects allowing them to capture a broad variety of anomalies.

Because of their advantages over holistic models, we propose to use object-centric appearances extracted with a pretrained object detector for training a VAD framework that can detect local anomaly. We also want

a light-weight VAD solution that can quickly and effectively detect local anomalies with just a few video frames. As opposed to some existing methods [5–7], we do not use any motion features (*e.g.* optical flows) mainly due to the fact that they are noisy for low resolution/poor quality videos. Moreover, the explicit use of motion adds another layer of complexity that can make the VAD framework less generalizable to other videos with different viewpoints.

Unlike the framework proposed in [6] where each CAE model learns normal features independently from each other, we propose a next local appearance prediction network (NLAPnet) comprising of a single generator that follows a U-Net architecture with skip connections to learn to predict the next local appearance of a normally behaving object by using only the past and the current local appearances of the corresponding object in the video frames. Thus, contrarily to the CAE models in [6], our generator explicitly learns the short temporal appearance variations of a given object which enables better characterization of normal behavior. Inspired by [5, 7] in which the generator is trained as a generative adversarial network (GAN [11]), an adversarial loss is also added during training in which a discriminator learns jointly with the generator to separate the real against the generated images. During inference for a given object, we rely on the structural similarity index measure (SSIM [17]) between the real and the predicted image for producing an anomaly score.

Our contributions are: 1) Unlike previous methods, ours explicitly learns the temporal appearance variations by predicting the next local appearance using just a few frames and we show that this approach gives competitive results, 2) our method can perform real-time VAD when using a pretrained DLA backbone for object detection, and 3) because it does not rely on optical flow and it is object-centric, our method shows very good generalization capabilities on new scene as demonstrated by a few-shot scene adaptation study.

2 Proposed Method

Assuming that local video anomalies are caused by objects in the scene, we propose to adversarially train, with a discriminator, a generative model that learns to predict the appearance of objects behaving normally in the next frame given its appearances in the past and current frames. To reduce the computation cost, we define past (a_i^{t-T}), current (a_i^t) and next (a_i^{t+T}) object-centric appearances as the gray-scaled pixel-level intensity images of an object (i) in the past ($t-T$), current (t) and next ($t+T$) frames respectively. Hence, each object in a given scene at a given time has an associated appearance triplet ($\langle a_i^{t-T}, a_i^t, a_i^{t+T} \rangle$) which slides temporally through the video frames. The frame gap $T = 3$ is empirically selected to ensure significant local changes in appearance while having a small spatial displacement to avoid the need to track the ob-

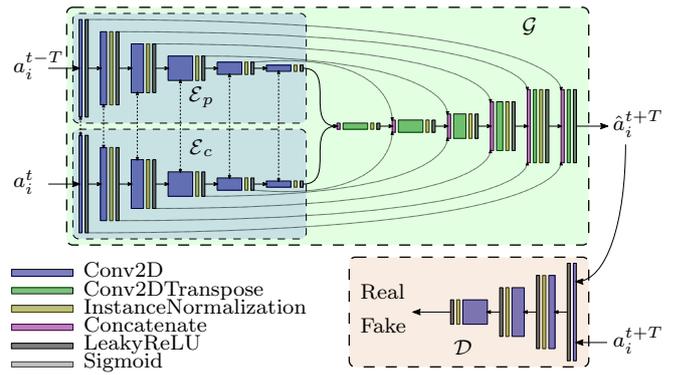


Figure 2. Overview of our proposed next local appearance prediction network (NLAPnet). The encoders \mathcal{E}_p and \mathcal{E}_c of the generator \mathcal{G} take the past and current appearances a_i^{t-T} and a_i^t respectively. The decoder learns to predict the next appearance a_i^{t+T} . The discriminator \mathcal{D} learns to discriminate the real and the predicted images.

jects. After learning with normal behaviors, the neural network will produce bad reconstructions for the next appearance of objects with abnormal behavior, thus indicating the likelihood of anomalies.

2.1 Object-Centric Appearance Extraction

Since our method requires object-centric appearances in videos, we need a pretrained multi-class object detector (MOD) to estimate the bounding box locations of all the objects of interest for each of the video frames. For this, we used the CenterNet [12] detector, for his state-of-the-art (SOTA) performance, with pretrained weights on MS-COCO [13] for detecting objects within 81 different classes. To extract the appearance images of a given object, we simply crop the past ($t-T$), the current (t) and the next ($t+T$) video frames using the bounding box coordinates of that object in the current (t) frame.

2.2 Next Local Appearance Prediction Network

NLAPnet predicts the next local appearance of an object using a generator \mathcal{G} that learns to predict a_i^{t+T} from a pair of object crops $\langle a_i^{t-T}, a_i^t \rangle$. In order to solve the issue of vanishing gradients while capturing better fine-grained details, our generator \mathcal{G} is formed by two encoders that are linked to one decoder with skip-connections [14]. The encoders share the same weights and are composed of multiple 2-strided convolutional layers, while the decoder is composed of multiple 2-strided transpose convolutional layers. A discriminator \mathcal{D} is also added to impose a generative adversarial constraint by learning to differentiate the real a_i^{t+T} against the generated one \hat{a}_i^{t+T} by \mathcal{G} . To increase

generative performance, our model follows the adversarial framework in [5] by using PatchGAN [15] with the Least Square GAN [16] which makes \mathcal{D} a discriminative feature generator that is composed of multiple 2-strided convolutional layers and learns to output a feature matrix of confidence values ranging from one (for real images) to zero (for generated images).

For a given pair of $\langle a_i^{t-T}, a_i^t \rangle$ of an object i in a video frame t , the quality of the prediction by \mathcal{G} is evaluated through the generative reconstruction loss $\mathcal{L}_{\mathcal{G}}$, the generative adversarial loss $\mathcal{L}_{adv}^{\mathcal{G}}$ and the discriminative adversarial loss $\mathcal{L}_{adv}^{\mathcal{D}}$, which are given by:

$$\mathcal{L}_{\mathcal{G}} = \frac{1}{2} (1 - \text{SSIM}(a_i^{t+T}, \mathcal{G}(a_i^{t-T}, a_i^t))) \quad (1)$$

$$\mathcal{L}_{adv}^{\mathcal{G}} = \sum_{x,y} \left(1 - \mathcal{D}(\mathcal{G}(a_i^{t-T}, a_i^t))_{x,y}\right)^2 \quad (2)$$

$$\mathcal{L}_{adv}^{\mathcal{D}} = \frac{1}{2} \left[\sum_{x,y} \left(1 - \mathcal{D}(a_i^{t+T})_{x,y}\right)^2 + \sum_{x,y} \left(\mathcal{D}(\mathcal{G}(a_i^{t-T}, a_i^t))_{x,y}\right)^2 \right], \quad (3)$$

where x and y denotes the patch coordinates of the discriminative feature for computing the mean squared errors (MSE). During training, \mathcal{D} and \mathcal{G} update their weights alternately: \mathcal{D} first minimizes $\mathcal{L}_{adv}^{\mathcal{D}}$ and then \mathcal{G} minimizes $\mathcal{L}_{\mathcal{G}} + \mathcal{L}_{adv}^{\mathcal{G}}$. For $\mathcal{L}_{\mathcal{G}}$, we found that the structural similarity index measure (SSIM [17]) between gray-scaled intensity images yields better overall VAD performance.

2.3 Video Anomaly Detection

During inference, we simply use $\mathcal{L}_{\mathcal{G}}$ which employs SSIM as described in eq. 2 for producing a region-level anomaly score s_i^t for a triplet $\langle a_i^{t-T}, a_i^t, a_i^{t+T} \rangle$ of an object i at a time t . Anomaly detection decisions are taken at frame-level with provisions to reduce noise. To get the frame-level anomaly score s^t , we take the highest corresponding region-level anomaly score. To mitigate the issue of missing detection due to occlusion, we temporally smooth the frame-level anomaly scores with a Gaussian filter.

3 Experiments

3.1 Experimental setup

To validate our proposed method, we used three publicly available VAD benchmarks with training videos assumed as being normal: UCSD Pedestrian [18], CUHK Avenue [19] and ShanghaiTech [20]. The UCSD Pedestrian benchmark is divided into two datasets (Ped1 and Ped2), where both share the same

Table 1. Frame-level AUC results (in %) of our method on four datasets. * Means that the results were obtained by our implementation of the method. The best results are shown in boldface.

Method	Ped1	Ped2	Ave	ST
CAE [1]	81.0	90.0	70.2	60.9
ConvLSTM-AE [2]	75.5	88.1	77.0	-
STAE-OF [3]	87.1	88.6	80.9	-
Deep CAEs [4]	56.9	84.7	77.2	-
FFP [5]	83.1	95.4	84.9	72.8
MAC [7]	-	96.2	86.9	-
Few-shot [8]	86.3	96.2	85.8	77.9
OC-CAEs-dla* [6]	77.4	95.5	80.3	79.3
OC-CAEs-hg* [6]	79.6	96.4	82.2	80.5
NLAPnet-dla	81.1	96.3	81.3	82.0
NLAPnet-hg	82.3	97.2	85.4	82.5

abnormalities in videos consisting of vehicles, cyclist, skateboarders and wheelchairs going on pedestrian pathways, the CUHK Avenue dataset contains anomalies that are based on some irregular human activities like running, loitering and throwing/leaving carried objects, and ShanghaiTech (ST) is composed of multiple videos having different scenes with varying types of abnormal activities like cycling, fighting, robbing, *etc.*

To measure performance, we employ the frame-level area under the ROC curve, following the evaluation protocol in [5]. To study the effect of MOD on local VAD, we consider two different CenterNet [12] backbones, namely 1) deep layer aggregation with 34 layers (DLA [21]) and 2) hourglass with 104 layers (HG [22]). The former backbone is less accurate than the latter one, but enables fast real-time detection. When conducting experiments related to scene adaptation and for the ablation study, we used the HG backbone.

3.2 Results and Discussion

Comparison with SOTA methods: We first compared our proposed method using two different MOD backbones (NLAPnet-dla and NLAPnet-hg) with several existing SOTA methods [1–8], as shown in table 1. Note that we present only the results from our version of the implementation of [6] using the same two MOD backbones (OC-CAEs-dla and OC-CAEs-hg) so that all results are obtained following the same evaluation protocol for a rigorous comparison. Overall, results show that our proposed method, using either one of the two MOD backbones, significantly outperforms other existing approaches on the most challenging dataset, ST, as well as on Ped2, while being competitive on other datasets. However, our method does not perform as well on Ped1 which has a significantly lower resolution, which may affect the learning of normal behavior for predicting frames. Nevertheless, we can see that a better performing MOD backbone helps increasing the performance on these datasets, es-

Table 2. Comparison between frame-level AUC results (in %) on three target datasets using pre-trained weights on ST. Met.: Method, Tr.: Training scheme. PO: use pretrained weights only, FT: fine-tune the pretrained weights on the target dataset and MT: meta-training on ST. *: results were obtained by our implementation of the method. †: shots are not used with PO. Best results are shown in boldface.

Target	Met.	Tr.	1-shot	5-shot	10-shot
Ped1	[8]	PO	73.1†		
	[8]	FT	77.0	77.9	78.2
	[8]	MT	80.6	81.4	82.4
	* [6]	PO	71.6†		
	* [6]	FT	72.1	74.3	77.2
	(ours)	PO	74.4†		
Ped2	[8]	PO	82.0†		
	[8]	FT	85.6	89.7	91.1
	[8]	MT	91.2	91.8	92.8
	* [6]	PO	90.3†		
	* [6]	FT	92.2	93.1	94.2
	(ours)	PO	95.9†		
Ave	[8]	PO	71.4†		
	[8]	FT	75.4	76.5	77.8
	[8]	MT	76.6	77.1	78.8
	* [6]	PO	74.7†		
	* [6]	FT	75.7	76.5	77.1
	(ours)	PO	78.8†		
	(ours)	FT	79.5	81.1	82.3

pecially on Avenue, with heavily occluded areas. Still, our results with NLAPnet-dla are competitive and are obtained at real-time speed. With NLAPnet-hg, our performances are similar or better than holistic methods [1–5, 7, 8], except for Ped1. This shows the benefit of an object-centric approach in VAD.

Performance in Scene Adaptation: Since the non-object background information is not used in our proposed framework, we are more capable of adapting NLAPnet to a different unseen video scene with similar anomalies. To demonstrate that, we employed a similar protocol as [8] by fine-tuning our pretrained model with only K randomly selected frames per video (K -shot) in the training set of the target datasets: Ped1, Ped2 and Avenue (Ave). As in [8], we use ST as the source dataset for pretraining our model since it has various scenes with similar anomalies as in other datasets. Moreover, to further compare the effectiveness of our proposed framework for scene adaptation, we include results with the existing object-centric method [6] by fine-tuning the CAE models on the target dataset. Results are shown in table 2.

Results show that our method outperforms the model agnostic meta-learning (MAML [23]) framework that was used as an holistic method for the Ped2 and

Table 3. Ablation study AUC results (in %) on ST. \mathcal{E}_p and \mathcal{E}_c : encoders of the generator \mathcal{G} for past and current local images respectively. SC: with skip-connections. *adv*: \mathcal{G} is trained adversarially with a discriminator \mathcal{D} . Best results are shown in boldface.

\mathcal{E}_p	✓	—	✓	✓	✓	✓
\mathcal{E}_c	—	✓	✓	✓	✓	✓
SC	—	—	—	✓	—	✓
<i>adv</i>	—	—	—	—	✓	✓
AUC	75.3	75.5	76.7	80.8	76.2	82.2

Table 4. Training (Tr) time in hours and inference (Inf) speed in frame per seconds of our method and others on ST.

	[5]	dla [6]	hg [6]	dla (ours)	hg (ours)
Tr	48	10	10	3	3
Inf	25	25	11	34	13

Avenue datasets. We also note a significant increase in performance when using our framework compared to the object-centric method of [6]. Interestingly, we can see that using pretrained weights on ST can even improve the performance on Ped2 compared to the one that is trained only on Ped2 (see table 1). However, for Ped1, since the appearances of crowded objects are of low pixel resolution, similarly to [6], our method does not perform as well compared to the holistic model [8].

Ablation Study: In this experiment, we validate the crucial parts of NLAPnet using the ST dataset. We test the effects of the two encoders (for past and current images), the skip-connections and the adversarial loss $\mathcal{L}_{adv}^{\mathcal{D}}$. As we can see in table 3, results illustrate the importance of using both encoders with skip-connections when training the model in an adversarial manner.

Running Time: We used Python 3 and Keras 2 with TensorFlow binding to implement our proposed method¹ on a Intel i7-8700 machine with 16 GB RAM using Nvidia RTX 2080 GPU. Table 4 gives the approximated training time and inference speed of NLAPnet, of [5] and of our implementation of [6] using DLA or HG MOD backbone. We can see that our method is significantly faster than [6] while having competitive VAD performances.

4 Conclusion

To conclude, this paper proposes an adversarial framework that learns to predict the local appearance of a normally behaving object in the next video frame by using only the appearances of that object in the past and current frames. Results on four public benchmarks demonstrate the effectiveness of our method with competitive performance in VAD while at the same being faster at inference, light-weight and capable to better adapt to unseen video scenes than other methods.

¹Code: https://github.com/proy3/NLAP-net_VAD.

Acknowledgments

This research was supported by grants from IVADO and NSERC funding programs.

References

- [1] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis, "Learning Temporal Regularity in Video Sequences," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2016.
- [2] W. Luo, W. Liu, and S. Gao, "Remembering history with convolutional LSTM for anomaly detection," in *2017 IEEE International Conference on Multimedia and Expo (ICME)*, IEEE, 2017.
- [3] Y. Zhao, B. Deng, C. Shen, Y. Liu, H. Lu, and X.-S. Hua, "Spatio-Temporal AutoEncoder for Video Anomaly Detection," in *Proceedings of the 2017 ACM on Multimedia Conference - MM '17*, New York, New York, USA: ACM Press, 2017.
- [4] M. Ribeiro, A. E. Lazzaretti, and H. S. Lopes, "A study of deep convolutional auto-encoders for anomaly detection in videos," *Pattern Recognition Letters*, 2018.
- [5] W. Liu, W. Luo, D. Lian, and S. Gao, "Future Frame Prediction for Anomaly Detection - A New Baseline," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, IEEE, 2018.
- [6] R. T. Ionescu, F. S. Khan, M.-I. Georgescu, and L. Shao, "Object-Centric Auto-Encoders and Dummy Anomalies for Abnormal Event Detection in Video," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2019.
- [7] T. N. Nguyen and J. Meunier, "Anomaly Detection in Video Sequence With Appearance-Motion Correspondence," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, IEEE, 2019.
- [8] Y. Lu, F. Yu, M. K. K. Reddy, and Y. Wang, "Few-shot Scene-adaptive Anomaly Detection," in *Proceedings of the European Conference on Computer Vision*, 2020.
- [9] R. Hinami, T. Mei, and S. Satoh, "Joint Detection and Re-counting of Abnormal Events by Learning Deep Generic Knowledge," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [10] C.-Y. Chen and Y. Shao, "Crowd Escape Behavior Detection and Localization Based on Divergent Centers," *IEEE Sensors Journal*, 2015.
- [11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, 2020.
- [12] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as Points," arXiv.org, Tech. Rep., 2019.
- [13] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2014.
- [14] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," *MICCAI2015*, 2015.
- [15] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017.
- [16] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. P. Smolley, "Least Squares Generative Adversarial Networks," in *2017 IEEE International Conference on Computer Vision (ICCV)*, IEEE, 2017.
- [17] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image Quality Assessment: From Error Visibility to Structural Similarity," *IEEE Transactions on Image Processing*, 2004.
- [18] V. Mahadevan, W. Li, V. Bhalodia, and N. Vasconcelos, "Anomaly detection in crowded scenes," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, IEEE, 2010.
- [19] C. Lu, J. Shi, and J. Jia, "Abnormal Event Detection at 150 FPS in MATLAB," in *2013 IEEE International Conference on Computer Vision*, IEEE, 2013.
- [20] W. Luo, W. Liu, and S. Gao, "A Revisit of Sparse Coding Based Anomaly Detection in Stacked RNN Framework," in *2017 IEEE International Conference on Computer Vision (ICCV)*, IEEE, 2017.
- [21] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, "Deep Layer Aggregation," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, IEEE, 2018.
- [22] H. Law and J. Deng, "CornerNet: Detecting Objects as Paired Keypoints," *International Journal of Computer Vision*, 2020.
- [23] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *34th International Conference on Machine Learning, ICML 2017*, JMLR.org, 2017.