



# From CNN to CNN + RNN: Adapting Visualization Techniques for Time-Series Anomaly Detection

Fabien Poirier

## ► To cite this version:

Fabien Poirier. From CNN to CNN + RNN: Adapting Visualization Techniques for Time-Series Anomaly Detection. 2024. <hal-04855169>

**HAL Id: hal-04855169**

**<https://hal.science/hal-04855169v1>**

Preprint submitted on 24 Dec 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY-NC-SA 4.0 - Attribution - Non-commercial use - ShareAlike - International License

# From CNN to CNN + RNN: Adapting Visualization Techniques for Time-Series Anomaly Detection

Fabien Poirier

Université Paris 8 - IUT de Montreuil, LIASD, France

December 24, 2024

## Abstract

Deep neural networks, while effective in solving complex problems, are often perceived as "black boxes", limiting their adoption in contexts where transparency and explainability are crucial. This lack of visibility raises ethical and legal concerns, particularly in critical areas such as security, where automated decisions can have significant consequences. The implementation of the General Data Protection Regulation (GDPR) emphasizes the importance of justifying decisions made by these systems. In this work, we explore the use of visualization techniques to improve the understanding of anomaly detection models based on convolutional recurrent neural networks (CNN + RNN), integrating a TimeDistributed layer. Our model combines VGG19 for convolutional feature extraction and a GRU layer for sequential analysis to process real-time video data. Although this approach is suitable for models dealing with temporal data, it complicates gradient propagation since each element of the sequence is processed independently. The TimeDistributed layer applies a model or layer to each element of the sequence, but unfortunately, this structure dissociates temporal information, which can make it difficult to associate gradients with specific elements of the sequence. Therefore, we attempt to adapt visualization techniques, such as saliency maps and Grad-CAM, to make them applicable to models incorporating a temporal dimension. This article highlights current challenges in the visual interpretation of models that handle video data. While specific methods for video are still limited, it demonstrates the possibility of adapting visualization techniques designed for static images to neural network architectures processing video sequences. This approach extends the use of common interpretation methods from classical convolutional networks to recurrent convolutional networks and video data, thus offering a transitional solution in the absence of dedicated techniques.

**Keywords:** Deep learning, Explicability, Visualization, Time Distributed Convolution, Saliency, Grad Cam

# 1 Introduction

Deep neural networks play a central role in many fields, particularly in solving complex problems such as real-time anomaly detection in videos. However, their effectiveness comes with significant opacity: they are often seen as "black boxes", unable to explain or justify their decisions. This opacity presents critical challenges, especially in sensitive areas like security, where anomalies such as fights, gunshots, or accidents need to be detected quickly and reliably. The lack of transparency raises ethical concerns about the responsibility for automated decisions, as well as legal issues, particularly with the implementation of the General Data Protection Regulation (GDPR) in May 2018. This regulatory framework, particularly Article 22 [1], requires that significant decisions not be based solely on automated processing, making the explainability of AI systems essential. In this context, the explainability of deep learning models, particularly in anomaly detection systems, poses a major challenge. While visualization techniques are well established for convolutional networks used on images, their adaptation to models handling video data remains limited. Videos introduce a temporal dimension that complicates traditional visualization techniques.

Thus, a key question arises: how can we make the decisions of a deep learning model designed to detect anomalies in video sequences visible and interpretable? To address this issue, we explore visualization techniques, such as saliency maps and Grad-CAM, applied to a Time Distributed CNN + RNN model designed for anomaly detection in videos. Our model, developed using the Keras library<sup>1</sup>, combines convolutional layers (VGG19) for visual feature extraction and GRU layers for sequential analysis, allowing for real-time video data processing. The "Time Distributed" architecture enables each image in a video sequence to be analyzed individually, while preserving the temporal continuity of the data, making it particularly suited for complex video processing applications [20]. However, while the Time Distributed structure is useful for processing sequential data, it introduces specific challenges. In fact, it complicates the propagation of the gradient, as each element in the sequence is processed independently. The TimeDistributed layer applies a model or layer to each element of the sequence, but unfortunately, this structure dissociates temporal information, making it difficult to associate gradients with specific elements in the sequence. We thus attempt to adapt visualization techniques, such as saliency maps and Grad-CAM, to make them applicable to models incorporating a temporal dimension.

This study presents two major contributions: on the one hand, it analyzes the limitations of current visualization techniques, highlighting the specific constraints imposed by the "Time Distributed" architecture, particularly concerning gradient propagation and the generation of interpretable maps; on the other hand, it adapts visualization techniques from convolutional networks to make them applicable to models processing video data, while integrating the temporal dimension. Furthermore, this study shows that explainability can also guide our choices in AI model design, particularly with respect to hyperparameters.

---

<sup>1</sup>see [keras.io](https://keras.io)

## 2 Related Works

Surveillance videos, often devoid of audio, led us to focus our exploration on technologies applicable to images, which align with the type of data processed by the CNN component of our model. Moreover, the decisions made by the CNN can be visually represented, making them more intuitive and accessible, even to non-expert audiences. In contrast, the RNN relies on abstract temporal relationships, making the explanation of its decisions significantly more complex. In this context, the book by Christoph Molnar, "A Guide for Making Black Box Models Explainable", published in 2019, serves as a key reference. It synthesizes and analyzes various explainability and visualization approaches used to interpret complex machine learning models [14]. This book marks a significant milestone in popularizing and structuring explainability techniques, highlighting their strengths, limitations, and areas of application, providing a valuable framework for guiding research in this domain.

On one hand, some technologies, such as LIME (Local Interpretable Model-Agnostic Explanations), proposed on August 9, 2016, by Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin [22], or SHAP (SHapley Additive exPlanations), introduced in November 2017 by Scott M. Lundberg and Su-In Lee [13], stand out due to their independence from the internal structure of the models. These methods work by analyzing the inputs/outputs of a model to assign local or global importance to input features. For instance, SHAP, based on Shapley value theory, provides a coherent and unified explanation of each feature impact on the prediction, while LIME generates locally simplified models to interpret specific predictions. These technologies are particularly useful in environments where models are complex or proprietary, and their flexibility allows them to be applied to a variety of models, whether linear, ensemble-based, or neural networks. However, their implementation often requires the installation of specific libraries, along with pre- or post-processing steps to make the models compatible with these tools.

On the other hand, some technologies are specific to certain types of models, particularly convolutional neural networks (CNNs), which primarily process structured data such as images or videos. For these models, visualization techniques are used that directly exploit their internal architecture. For example:

### 2.1 Activation Maps and Their Variants

- Grad-CAM: Generates heatmaps by weighting the activations of the last layers with gradients to highlight the regions important for classification [23].
- Grad-CAM++: An improvement on Grad-CAM that considers higher-order interactions for more precise maps, particularly in cases where multiple objects or classes are present in the image [2].
- LayerCAM: Another method that generates activation maps using the activations from the network's shallow layers to provide more accurate

localization information [10].

## 2.2 Saliency Maps and Their Improvements

- Saliency maps identify the regions of an image that have the greatest impact on the model’s decision [25]. Several techniques have been developed to improve their quality:
- SmoothGrad: Reduces noise by averaging gradients calculated over slightly perturbed inputs [26].
- Guided Backpropagation: Combines gradients of activations with forward propagation information for sharper maps [27].
- Layer-wise Relevance Propagation (LRP): Assigns relevance to each neuron by following the decision flow through the layers of the network to form a saliency map [3].
- Integrated Gradients: Creates saliency maps by calculating the importance of pixels by integrating gradients between a reference input and the actual input [28].
- DeepLIFT: Compares a neuron’s activations to a reference activation to generate more understandable saliency maps [24].
- Poly-CAM: This method generates saliency maps by recursively merging activation maps from early layers with those from deeper layers for more detailed visualizations [7].
- Eigen-CAM: Uses the decomposition of activations into principal components to generate saliency maps while integrating gradients to refine the discrimination between classes [15].
- Fast-CAM: A fast method for generating saliency maps, using Grad-CAM or Grad-CAM++ for better computational efficiency [16].
- RISE (Randomized Input Sampling for Explanation): Generates saliency maps by perturbing inputs and calculating average predictions for each perturbed sample, assigning importance to pixels [19].

## 2.3 Attention Maps and Their Improvements

- Attention Maps: These maps visualize the areas of the input that a model pays the most attention to, typically used in models with an attention mechanism such as Transformers. They help identify influential regions in the model’s decisions by exploiting self-attention weights calculated in deep layers [5].

- **Class-Discriminative Attention Maps (CDAM):** An extension of classic attention maps, CDAM adjusts the attention scores to make them specific to a target class. This is achieved by weighting attention scores by the relevance of the "tokens" associated with a class, using gradients. CDAM is particularly effective for providing compact and semantic explanations by targeting specific concepts [5].
- **Smooth CDAM:** A variant of CDAM that adds Gaussian noise to the input "tokens" to generate averaged explanations, thus improving the robustness of attention maps [5].
- **Integrated CDAM:** Another improvement of CDAM that interpolates between a reference input and the actual input to produce integrated explanations, inspired by integrated gradients. This enhances the stability of explanations, particularly for latent or complex concepts [5].

## 2.4 Other Techniques

- **3D Attention, Activation, and Saliency Maps:** These maps are designed for volumetric or spatio-temporal data, such as videos or 3D medical images. They generalize classical attention maps by integrating spatial and temporal dimensions, allowing for the localization of important regions in complex data such as video sequences or MRI scans [12].
- **Occlusion Sensitivity:** A technique where portions of the input are masked to observe how this affects the model prediction, providing insight into important regions [29].
- **DeconvNets:** Use a deconvolution process to reverse the activation process of a network, allowing more intuitive visualizations of intermediate layer activations [30].

Many libraries exist to help users visualize the features learned by neural network models, particularly in the fields of computer vision and medical learning. In 2017, Kotikalapud Raghavendra proposed `keras-vis` [11], a public library that allows users to visualize the convolution filters of each layer, their evolution during training, and the activation maps. Later, in 2020, Philippe Remy developed `keract` [21], another library to perform similar operations. In 2020, Gotkowski, Karol et al. proposed a library allowing users to visualize both 2D and 3D attention maps in deep learning models for volumetric or spatio-temporal data, particularly in medical applications [8]. Additionally, the Keras library, developed by François Chollet in 2015 [6], now includes some of these visualization techniques. Another notable library is M3d-CAM, which generates 3D attention maps for medical deep learning applications, facilitating the visualization of relevant areas in 3D images such as MRIs or CT scans [9].

### 3 Approach

The model used in this study has already been the subject of a previous publication [20]. It is based on a CNN + RNN architecture developed using the Keras library, specifically designed for anomaly detection in video sequences. The architecture consists of two main parts: a convolutional component for feature extraction and a recurrent component for temporal relationship analysis. For the extraction of visual features from images, we chose the VGG19 model (Visual Geometry Group). To integrate the temporal dimension in video processing, we encapsulated VGG19 in a Time Distributed layer. This layer allows us to apply VGG19 independently to each image in a video sequence while preserving the temporal structure of the data. In other words, it applies the same processing to each image of the sequence in an identical manner while maintaining the temporal order. The recurrent part of the model uses a GRU (Gated Recurrent Unit). This GRU consists of 1024 neurons and includes a dropout rate of 50% to prevent overfitting. After this recurrent part, an MLP (Multi-Layer Perceptron) consisting of three successive blocks is integrated. Each block contains a dense layer and a dropout layer with a rate of 50%, helping to stabilize predictions while reducing the risk of overfitting. The details of the model’s hyperparameter choices, illustrated in figure 1, will be detailed in Section 4.

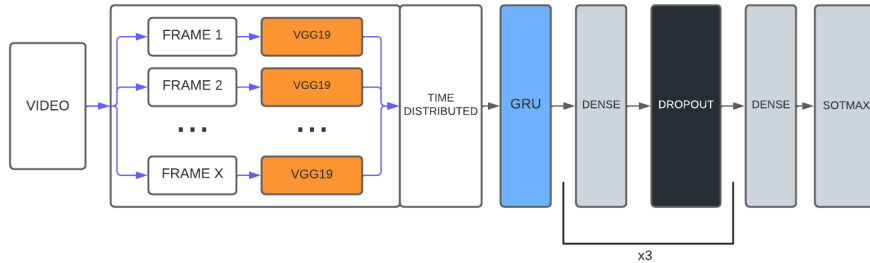


Figure 1: Model Architecture

2D convolutions are commonly used in image processing. They apply a filter to an image to extract local features such as edges, textures, or patterns. Each filter is applied to a spatial dimension of the image (width and height), without considering the temporal dimension. Thus, a 2D convolution treats an image as a single input, without temporal relation between multiple images in a sequence.

3D convolution, on the other hand, is specifically designed for processing volumetric data or videos, where the input includes both spatial dimensions (width and height) and a temporal dimension (depth or the number of images in a sequence). Unlike 2D convolutions, which focus on processing each image inde-

pendently, 3D convolutions apply filters simultaneously on both the spatial and temporal dimensions. This allows the model to capture features not only within each image but also across relationships between consecutive images. Thanks to this capability, 3D convolutions are particularly suited for identifying spatio-temporal patterns, tracking the evolution of visual information over time, which is crucial for understanding movements, transitions, or dynamic changes in a sequence.

The Time Distributed layer allows the same layer, such as VGG19, to be applied independently to each image of a video sequence while maintaining the temporal integrity of the data. In other words, it "distributes" the processing of each image independently while respecting the temporal order of the data. The difference here with 2D or 3D convolutions is that, in the case of Time Distributed, each image in the sequence is treated identically, but the model retains the temporal information through the sequence structure, without applying convolution on the temporal dimension itself. This allows for independent extraction of visual features while preserving the order of the images for further processing by recurrent layers like GRU.

To better understand the internal workings of our model, we opted for visualization techniques suited to CNNs, such as saliency maps (cf. Figure 2), activation maps (cf. Figure 3), and filter visualizations. These methods, based on gradient propagation, allow for direct exploration of the network's activations. Saliency maps identify the most influential pixels for a prediction, while activation maps reveal the specific regions of an image that contribute to the final decisions. By avoiding external dependencies, we ensure consistency and reproducibility in our analyses.

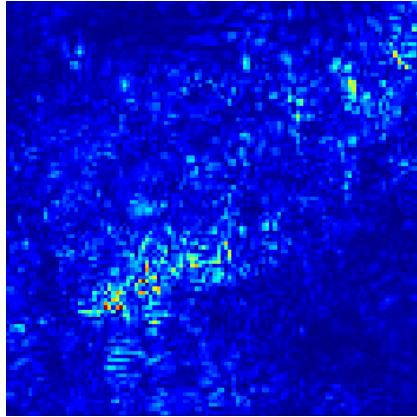


Figure 2: Saliency maps



Figure 3: Activation maps  
Source "tap rack gajoubert"



Unfortunately, due to our specific architecture, we found that visualization libraries designed for classic CNNs were not suitable for our model. This limitation is directly related to the structure and peculiarities of our approach.

The first challenge lies in the organization of our model. Unlike a standard network using 2D or 3D convolutions, where the layers are directly connected and the outputs of one layer serve as inputs to the next, our architecture encapsulates convolutions within a Time Distributed layer (see Figure 4). This layer processes each image in the video sequence independently, applying a sub-model to each image. While this preserves the temporal structure of the data, it complicates associating gradients with specific elements of the sequence. As a result, some important temporal information may become fragmented or difficult to link.

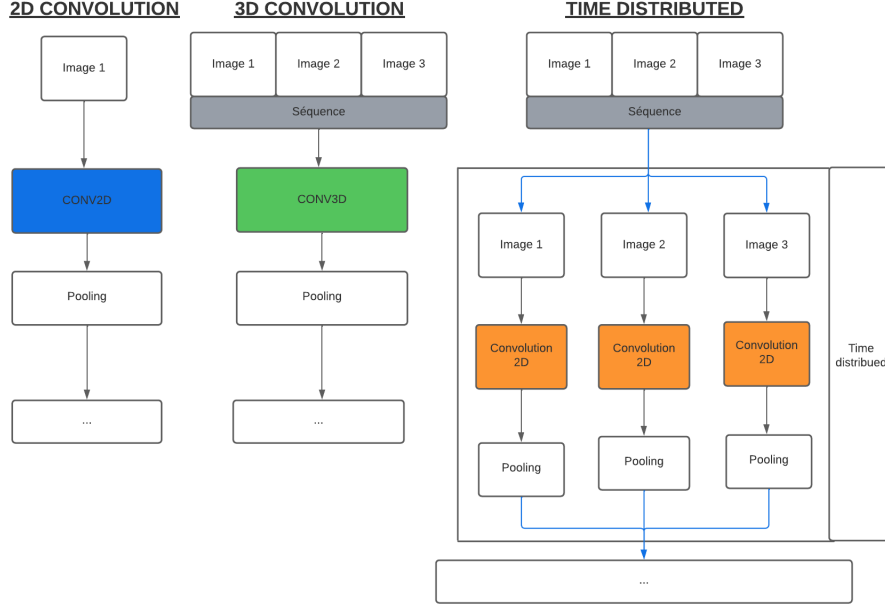


Figure 4: Structure of 2D convolution / 3D convolution / Time distributed convolution

The second challenge comes from the addition of the temporal dimension. Unlike models based on 2D or 3D data, which generate a single output visualization for a single input (such as an image or volumetric object), our model processes video sequences, which involve multiple input images for a single output sequence. This paradigm shift requires rethinking how visualizations are generated and

interpreted, as each input image contributes to the overall prediction of the sequence.

These limitations directly affect the visualization methods. To generate saliency maps, we had to compute gradients for each image in the sequence, resulting in a series of gradients corresponding to the length of the sequence. These gradients are then displayed individually, generating a distinct saliency map for each image. For activation maps, the solution we implemented relies on using the outputs from the Time Distributed layer. This layer applies identical processing to each image in the sequence, generating a list of outputs, each corresponding to an individual image. The calculated gradients have the same dimensions as these outputs, which allows them to be associated image by image and projected onto the corresponding image. However, with this architecture, activation maps can only be generated for the output layer of our submodel, limiting the depth of analysis of the internal layers.

## 4 Experimentation

In this section, we will present our results along with the advantages and disadvantages of each approach explained previously. All of our visualizations were obtained from a proprietary dataset<sup>2</sup>. By examining the activation maps shown in Figure 5, which represent a sequence of images from a video illustrating a gunshot, we can observe that our model does not consistently focus on the same areas of an image, even when the images are successive. This is due to the operation of the convolutional layers, which process each image independently, without explicitly linking important areas of one image to those of others in the sequence.

---

<sup>2</sup>The data sets used are confidential and cannot be shared publicly.

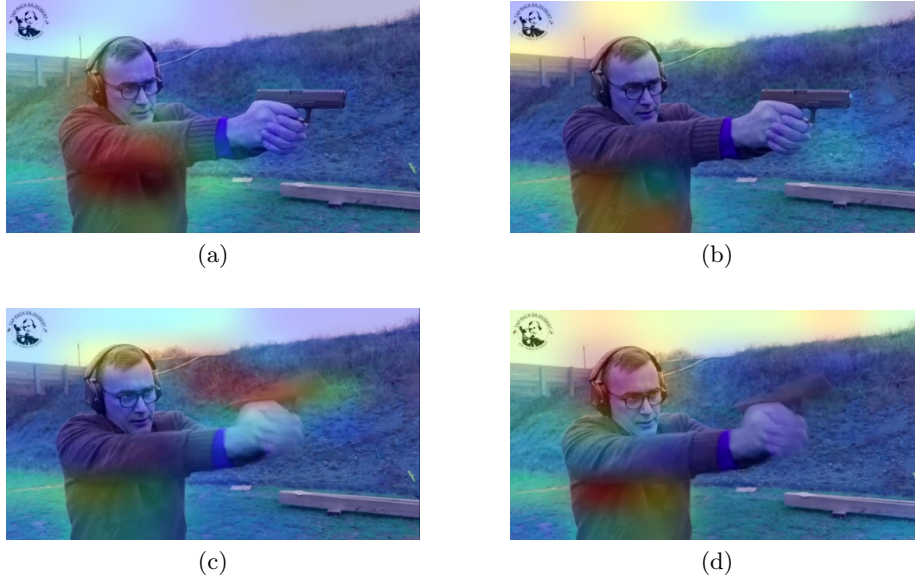


Figure 5: Class activation maps (Grad-CAM) for a shooting video  
Source: “tap rack gajoubert”

To facilitate the interpretation of sequences, we used OpenCV to extract the contours from the activation maps. This new visualization is presented in Figure 6. Images (a) and (c) represent the results of the activation maps: image (a) is from a video illustrating the “gunshot” anomaly, while image (c) corresponds to a video representing the “fight” anomaly. Images (b) and (d), on the other hand, show our contour visualization applied to the activation areas of images (a) and (c), providing a clearer representation of the regions of interest identified by the model.



Figure 6: Example of visualization for the classes: fight and shooting.  
Source: “Le Parisien newspaper”

This new visualization allowed us to identify low-activation areas that would be difficult to detect using traditional activation maps. However, this method has some limitations. First, the contour detection lacks precision: contours can be nested, especially when a major activation area is surrounded by minor activation areas. Second, this approach does not yet allow us to assess the intensity of activations, which limits its use for more in-depth analysis. Figures 7 illustrate both the advantages and disadvantages of this technique.

For example, in image (a), the gun is identified as a low-activation area, which is difficult to distinguish on the activation map but clearly visible through contour detection in image (b). We also observe a major activation area surrounded by a minor area on the left side of the image. Similarly, image (c) shows a person hitting another. Although the activation map seems to indicate that the model correctly perceives the action, it incorrectly labels it as “Normal.” However, by analyzing the contours (image d), we realize that the model completely misses the actual action, highlighting the limitations of the activation maps. Activation maps played a crucial role in the analysis and optimization of our architecture. By observing the effect of the different layers (GRU, Dense, Dropout) on the features learned by the convolutional part, we were able to adjust our parameter choices. These adjustments, validated by metrics such as accuracy and recall, confirm the consistency between visual observations and quantitative results. First, we analyzed the impact of the number of neurons in the GRU layer.

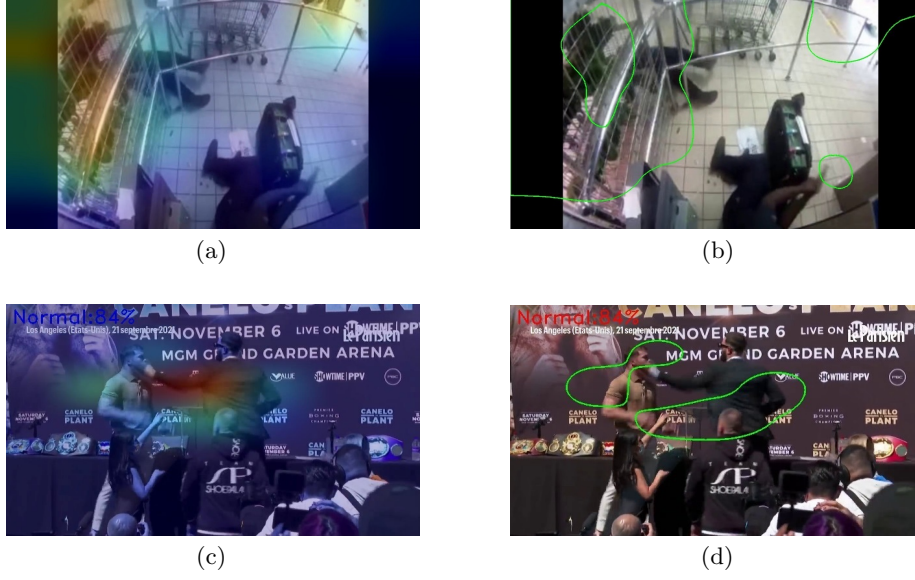


Figure 7: Activation map and contour visualization applied to images extracted from shooting and fight videos.

Source: UCF Crime + “Le Parisien newspaper“

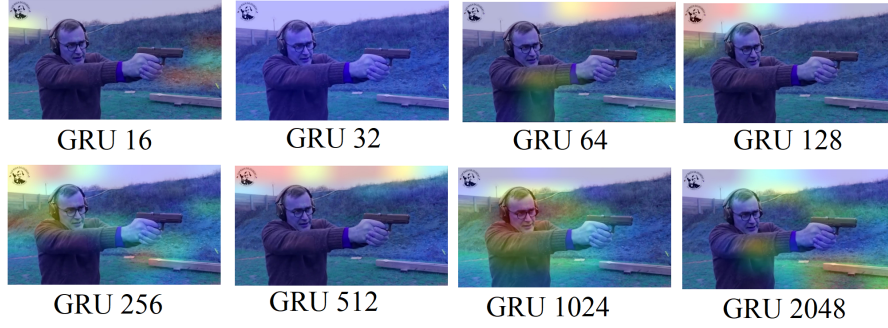


Figure 8: Activation map with variation of the number of neurons in the GRU Layer (0 dropout).

Source “tap rack gajoubert“

The activation maps shown in Figure 8 show a significant improvement in the concentration of activation areas for the 1024-neuron configuration. This is the only configuration where the activation area is clearly focused on the person in the video, which is crucial for anomaly detection. This visual observation is supported by the data in Table 1, where this configuration achieves an optimal

balance between accuracy (88.4%) and recall (86.3%).

Neurons (nb)	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)
8	83,3	19,6	84,2	31,7
16	83,1	89,1	83,1	85.9
32	82,9	91,7	82,9	87
64	85,4	89,6	85,4	87.4
128	87,4	80,7	87,4	83.9
256	85,3	81,1	85,3	83.1
512	85,1	81,4	84,9	83.2
1024	86,3	88,4	86,3	87.3
2048	83,4	84,0	83,4	83.6

Table 1: Impact of the number of neurons on performance with dropout rate set to 50%.

In contrast, lower or higher numbers of neurons lead to either a loss of information (dispersed activation) or overfitting (activation too focused on irrelevant areas).

Next, we studied the effect of the dropout rate. Although the images shown in Figure 9 are from a configuration with 8 neurons in the GRU layer and not 1024, they help visualize how the dropout rate impacts the model’s activation areas<sup>3</sup>.



Figure 9: Activation map showing the variation of the dropout with an 8 neurons GRU.

Source “tap rack gajoubert”

These observations are consistent with the results in Table 2. The absence of dropout (0%) produces dense activations but exposes the model to overfitting, while a dropout rate of 50% offers an interesting compromise. This rate maintains high accuracy (89.6%) and recall (87.4%) while improving the robustness and generalization of predictions.

<sup>3</sup>Due to confidentiality clauses, the results related to the dropout rate for the 1024-neuron configuration are no longer accessible.

<b>Dropout (%)</b>	<b>Accuracy (%)</b>	<b>Precision (%)</b>	<b>Recall (%)</b>	<b>F1-score (%)</b>
0	92,8	83,0	92,8	87.6
25	87,3	84,7	87,3	85.9
50	87,4	89,6	87,4	88.4
75	84,6	78,8	84,6	81.5

Table 2: Performance en fonction du taux de dropout pour 1024 neurones.

Finally, we analyzed the effect of the number of output layers (Figure 10). The activation maps reveal significant differences between configurations. In the 2-layer dense configuration, the model focuses mainly on the face and posture of the person in the frame, indicating a limited focus on specific aspects. In contrast, with 3 dense layers, the activations are more evenly distributed across the person’s overall posture, reflecting a more comprehensive understanding of the scene.



Figure 10: Activation map showing the variation of dense layer in the MLP using a GRU with 50% dropout and 1024 neurones.

Source “tap rack gajoubert”

This better concentration leads to improved performance, as shown in Table 3: the 3-layer configuration achieves an accuracy of 84.8% and a recall of 92.2%, outperforming the other configurations. The 4-layer configuration seems to introduce excessive complexity, reducing the coherence of activation and resulting in a slight performance degradation.

<b>Accuracy (%)</b>	<b>Precision (%)</b>	<b>Recall (%)</b>	<b>F1-score (%)</b>	<b>Blocks (Nb.)</b>
89,8	84,2	89,8	86,9	2
92,2	84,8	92,2	88,3	3
83,2	89,6	83,2	86,2	4

Table 3: Performance en fonction du nombre de blocks entre le GRU et la sortie. Chaque block est composé d’une couche dense et d’une couche de dropout à 50%.

This iterative approach, combining visualization and quantitative validation, allowed us to define an optimal configuration: 1024 neurons in the GRU, a



dropout rate of 50%, and 3 dense layers in the output. These adjustments ensure a high-performance, robust model that is consistent with the visual observations. These visualizations also allow us to examine the features associated with the "normal" class. In a supervised learning context, this class represents the absence of anomalies, encompassing a wide range of actions such as working, walking, exercising, etc. It is interesting to note that movements in this class are often slow, unlike anomalies, which are characterized by sudden and rapid movements. While a human would by default assign the normal class in the absence of anomaly indicators, our model does not follow this logic. To predict the normal class, it must detect specific features that represent it. By visualizing examples of videos belonging to this class through activation and saliency maps, we found that our model primarily relies on the posture of the people on screen to predict this class. Observing the activation maps from successive frames (see Figure 11), we can see that the activations target the hands of the person approaching the receptionist.



Figure 11: Activation Map for a normal video (images are sequential within the video)

On the saliency map (Figure 12), the prominent pixels are mostly the outstretched arms in Figure (a), visible in the activation map, but we can also note the beginning of the shoulders and the lower body in Figure (b).



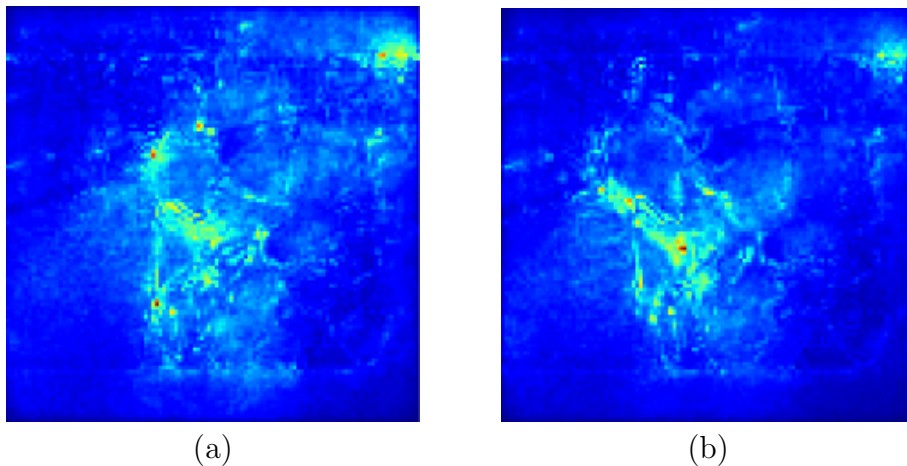


Figure 12: Saliency map for a normal video (images are sequential within the video).

## 5 Conclusions and Future Works

Through this article, we have demonstrated how to adapt visualization techniques, such as activation maps and saliency maps, which are commonly used in convolutional neural networks (CNNs), to be applicable with convolutional recurrent networks (CNN + RNN) that integrate convolutions into a "Time Distributed" layer. These techniques can be extended to other advanced variants, such as SmoothGrad, Grad-CAM++, and other similar methods discussed in the state of the art. This approach allowed us to meet the requirements of the GDPR, which stipulates that no decision should solely result from automated processing. By providing the ability for non-experts to verify the predicted class based on the action area targeted by the model and generating visualizations alongside the prediction, we have made the model's decisions more transparent and understandable. However, several research ways remain to be explored:

- The visualization of contours could be improved by incorporating a representation of the intensity of the activation areas.
- Although the visualization is performed in parallel with the model's prediction and can be disabled, it may not be suitable for real-time tasks where immediate feedback is required.
- Another promising approach would be to use object detection models to precisely locate anomalies, which could facilitate real-time visualizations and be a major asset for time-sensitive applications. However, although vision transformers allow for visualizations using attention maps, this process remains slow and computationally expensive, making it difficult to use in applications requiring immediate feedback.

- The development of visualization techniques specifically tailored to video data could significantly enhance our analytical capabilities.

With the rapid advancement of Artificial Intelligence (AI), explainability has become a major issue that has sparked significant interest within the scientific community. This interest has led to numerous works each year, particularly in the field of convolutional neural networks and image processing [4, 17, 18, 31]. These ongoing research efforts contribute to clarifying and improving our understanding of this field, paving the way for increasingly transparent and interpretable AI systems.

## References

- [1] Gdpr - article 22 :automated individual decision-making, including profiling. <https://gdpr-info.eu/art-22-gdpr/>, 2018. Adopted on April 27, 2016, enforced on May 25, 2018.
- [2] C Aditya, S Anirban, D Abhishek, and H Prantik. Grad-cam++: Improved visual explanations for deep convolutional networks. arxiv 2018. *arXiv preprint arXiv:1710.11063*.
- [3] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.
- [4] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Agata Lapedriza, Bolei Zhou, and Antonio Torralba. Understanding the role of individual units in a deep neural network. *Proceedings of the National Academy of Sciences*, 2020.
- [5] Lennart Brocki, Jakub Binda, and Neo Christopher Chung. Class-discriminative attention maps for vision transformers. *arXiv preprint arXiv:2312.02364*, 2023.
- [6] Francois Chollet et al. Keras, 2015.
- [7] Alexandre Englebert, Olivier Cornu, and Christophe De Vleeschouwer. Poly-cam: high resolution class activation map for convolutional neural networks. *Machine Vision and Applications*, 35(4):89, 2024.
- [8] Karol Gotkowski, Camila Gonzalez, Andreas Bucher, and Anirban Mukhopadhyay. M3d-cam: A pytorch library to generate 3d data attention maps for medical deep learning. *arXiv preprint arXiv:2007.00453*, 2020.
- [9] Karol Gotkowski, Camila Gonzalez, Andreas Bucher, and Anirban Mukhopadhyay. M3d-cam: A pytorch library to generate 3d attention maps for medical deep learning. In *Bildverarbeitung für die Medizin 2021: Proceedings, German Workshop on Medical Image Computing, Regensburg, March 7-9, 2021*, pages 217–222. Springer, 2021.

- [10] Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. Layercam: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing*, 30:5875–5888, 2021.
- [11] Raghavendra Kotikalapudi and contributors. keras-vis. <https://github.com/raghakot/keras-vis>, 2017.
- [12] Gabriele Lozupone, Alessandro Bria, Francesco Fontanella, Frederick JA Meijer, and Claudio De Stefano. Axial: Attention-based explainability for interpretable alzheimer’s localized diagnosis using 2d cnns on 3d mri brain scans. *arXiv preprint arXiv:2407.02418*, 2024.
- [13] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
- [14] Christoph Molnar. *Interpretable Machine Learning*. 2019.
- [15] Mohammed Bany Muhammad and Mohammed Yeasin. Eigen-cam: Class activation map using principal components. In *2020 international joint conference on neural networks (IJCNN)*, pages 1–7. IEEE, 2020.
- [16] T Nathan Mundhenk, Barry Y Chen, and Gerald Friedland. Efficient saliency maps for explainable ai. *arXiv preprint arXiv:1911.11293*, 2019.
- [17] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017. <https://distill.pub/2017/feature-visualization>.
- [18] Chris Olah, Arvind Satyanarayan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine Ye, and Alexander Mordvintsev. The building blocks of interpretability. *Distill*, 2018. <https://distill.pub/2018/building-blocks>.
- [19] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models, 2018.
- [20] Fabien Poirier, Rakia Jaziri, Camille Srouf, and Gilles Bernard. Enhancing anomaly detection in videos using a combined yolo and a vgg gru approach. In *2023 20th ACS/IEEE International Conference on Computer Systems and Applications (AICCSA)*, pages 1–6, 2023.
- [21] Philippe Remy. Keract: A library for visualizing activations and gradients. <https://github.com/philipperemy/keract>, 2020.
- [22] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ” why should i trust you?” explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- [23] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.

- [24] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *International conference on machine learning*, pages 3145–3153. PMIR, 2017.
- [25] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *In Workshop at International Conference on Learning Representations*. Citeseer, 2014.
- [26] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- [27] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- [28] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
- [29] Tomoki Uchiyama, Naoya Sogi, Satoshi Iizuka, Koichiro Niinuma, and Kazuhiro Fukui. Adaptive occlusion sensitivity analysis for visually explaining video recognition networks, 2023.
- [30] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks, 2013.
- [31] Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu. Interpretable convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8827–8836, 2018.