

---

# XAI Saliency Maps Can Be Used for OOD Detection

---

**Jonatan Hoffmann Hanssen**

Department of Informatics  
University of Oslo  
jonatahh@uio.no

**Hugo Lewi Hammer**

Department of Computer Science  
OsloMet  
Department of Holistic systems  
SimulaMet  
hugo.hammer@oslomet.no

## Abstract

When neural networks are used in high-impact settings such as cancer detection or autonomous driving, we must not only require that they make predictions with high accuracy, but also that they are aware of their shortcomings, and alert us when faced with unusual data. The need for models which "know what they do not know" has led to the field of Out-of-Distribution (OOD) detection, which attempts to detect when models are exposed to data points that are far outside of their training data and thus unlikely to be classified correctly. OOD detection is a young and developing field, and there are, to date, no methods which achieve superior performance on all benchmarks. Thus, there is a need for novel techniques which push the field forward. In this paper, we show that simply aggregating the saliency values produced by Explainable Artificial Intelligence (XAI) methods such as Integrated Gradients, LIME or GBP is an effective way to detect OOD data points. Furthermore, we show that these aggregates are only weakly correlated with the model's confidence in the predicted class, which allows us to combine the Maximum Logit Score with saliency aggregates, achieving AUROC scores which are close to State-of-the-Art (SoTA) OOD detection methods. Code available on GitHub. TODO: make link.

## 1 Introduction

Machine Learning (ML) generally, and Deep Learning (DL) specifically, has seen a tremendous increase in performance in recent years, performing comparable to humans in tasks such as image classification, speech and handwriting recognition, as well as many others [3]. Consequently, DL methods have been deployed in a multitude of fields, and have become a part of our daily lives through their role in web search, text translation, computer vision, as well as in many other technologies which are taken for granted. In the medical field, DL has the potential to provide faster and more accurate detection of diseases by being trained on cases from thousands of previous patients [10]. Despite this, the adoption DL in high impact fields, such as medicine, has been slow, with Rajkomar et al. [11] stating that: "surprisingly little in health care is driven by machine learning".

To explain this discrepancy, we should consider that despite their impressive performance, the application of DL methods is not without drawbacks. Firstly, deep neural networks are inherently unexplainable due to the large number of parameters that any non-trivial network has. State-of-the-Art (SoTA) models will perform millions of operations to evaluate a single data point, and it is therefore impossible for humans to comprehend and explain the entire process which led the model to make a particular decision. In medicine, this is a major limitation of DL methods, as both doctors and patients expect to be able to understand why a decision was made [21]. In other high-impact fields, such as autonomous driving, this lack of transparency also has serious practical and legal ramifications.

Secondly, although neural networks may attain high accuracy on test data and appear to have learned great insights about the tasks they are employed in, they often lack robustness and can suffer large drops in performance on data points which are slightly different from the training data. As Szegedy et al. [18] has shown, it is possible to create data points which are imperceptibly different from normal data points, yet still fool otherwise high performing models. More problematically, unlike humans, who recognize when they are faced with a novel situation where their expertise might be lacking, DL methods will predict equally confidently on data points which are far outside the data they have been trained on [21].

These two problems lead to the fields of Explainable Artificial Intelligence (XAI) and Out-of-Distribution (OOD) detection. XAI attempts to explain the reasons why a model came to a decision, which helps to remedy the black-box nature of complicated DL models. In a healthcare setting, such explanations can be inspected by medical practitioners to confirm the diagnosis, and can be used to give patients information about why decisions regarding their health were made. In autonomous driving or other automated high impact fields, XAI can be used to detect failure modes or to understand and improve trained models. OOD detection attempts to uncover when a data point is too different from the training data to be classified reliably. These methods could alert medical practitioners when such data points occur, thus avoiding potentially fatal misclassifications. In autonomous driving, the system could detect novel situations and cede control back to the user, avoiding accidents.

Both of these fields have seen increased interest in recent years, and are vital parts of any integration of DL in high impact settings. As two vibrant fields of study, there is great potential to combining insights from one field to improve performance in the other; an area which is underexplored. This paper investigates the possibility of using XAI methods to aid OOD detection. The overarching intuition is that when saliency maps are generated on semantically shifted data, they are creating an explanation for a wrong prediction, which leads to systematic differences when compared to explanations generated on In-Distribution (ID) data. This methodology is essentially non-existent in the literature: OpenOOD [22, 25], the standard OOD detection benchmarking framework which includes over 40 different methods, contains no method which uses XAI as part of its functioning.

In this paper, we show that the magnitude of values generated by XAI saliency mapping methods are significantly different between ID and OOD data. By aggregating all saliency values for a given input image, we get a single number which can be used for OOD detection, achieving Area Under Receiver Operating Characteristic (AUROC) scores which are comparable to baseline methods such as the Maximum Softmax Probability (MSP) or the Maximum Logit Score (MLS). Furthermore, we show that saliency aggregates are only somewhat correlated with these baselines. Inspired by the work of Shekhar et al. [14], we combine saliency aggregations with the MLS, to exploit the information gained both by the explanation and the model prediction for OOD detection.

The key contributions of this work are as follows:

- We show that XAI saliency mapping methods such as Local Interpretable Model-Agnostic Explanations (LIME), Gradient Class Activation Mapping (GradCAM) and Guided Back-propagation (GBP) capture valuable magnitude information which can be used for OOD detection. By forgoing normalization, which is common when displaying saliency maps, we achieve higher AUROC scores than previous works which have attempted to use XAI for OOD detection, such as Martinez et al. [9].
- Further illustrating the importance of the magnitude information in saliency maps, we show how aggregations such as the Coefficient of Variation (CV), Relative Mean Absolute Difference (RMD) and Quartile Coefficient of Determination (QCD), which are magnitude invariant, separate ID and OOD data points far worse than aggregations which capture the magnitude.
- We show that for many XAI saliency mapping methods, saliency aggregates are only weakly correlated with the confidence of the prediction (the MLS). This allows us to combine the discriminative power of XAI saliency maps with traditional OOD detection. We introduce the Saliency plus Logit framework, achieving results which are close to SoTA OOD detection methods.

## 2 Related Work

While the combination of XAI and OOD detection has been explored in many previous works, the majority of these focus on explaining why a data point was marked as OOD, as opposed to using XAI to aid the detection itself. Delaney et al. [4], Sipple et al. [15] and Tallón-Ballesteros et al. [19] used XAI methods to explain OOD detection decisions. Within network security, XAI has been as part of anomaly detection systems to detect malicious or faulty network traffic. Here, it has been used to explain detections [1, 8], but also to aid in detection itself by inspecting the explanations of the detection system [20, 23]. These methods thus use XAI to aid OOD detection in a similar manner to our work, however, they are strictly focused on sequential network traffic data as opposed to images, and are mostly concerned with detection "unnatural" data samples such as intentionally malicious traffic or that generated by faulty equipment, as opposed to natural OOD data caused by semantic or covariate shift occurring when a model is deployed.

Martinez et al. [9] is the most relevant previous work. Here, the authors explicitly aim to use XAI to improve OOD detection on images. They do this by looking at saliency maps produced by GradCAMPlusPlus [2] during inference, i.e the heatmaps that explain which parts of the image was most influential to classify the image as a specific class. Using these heatmaps, they perform distance-based OOD detection: By collecting all explanations for each image in the ID dataset, they are able to construct archetypical explanations, and can make clusters of explanations. To perform OOD detection, they simply compare the explanation of a new data point to the clusters of archetypical explanations, and mark it as OOD if it has a distance which is over a certain threshold. In contrast to our work, Martinez et al. consider normalized heatmaps, not raw saliency values.

This method performs decently on toy benchmarks, achieving scores similar to SoTA methods when using *Fashion MNIST* as ID and *MNIST* as OOD. However, this method fails in more complicated scenarios, achieving an AUROC score of only 52% on *CIFAR10* vs *SVHN*, which is far below most other OOD detection methods. The paper thus ends with the authors concluding that "OoD detection approaches that are specifically designed for the purpose achieve in general better detection scores at the cost of an additional computational burden in the model's construction" [9].

For more potential related work, we can look to OpenOOD [22, 25], which aims to provide a comprehensive benchmark of all relevant methods in the field of OOD detection. Out of all 41 OOD detection methods included in this benchmark, there are no methods which use XAI. However, as many XAI methods utilize the gradients of the network to generate saliency values, we could also consider OOD detection models which utilize gradients in some form as tangentially related to our work. In this regard, GradNorm [6] is somewhat related, as it utilizes the norm of the gradients of the network with respect to the Kullback-Leibler distance between the outputs and a uniform distribution to perform OOD detection.

## 3 Preliminaries

Let  $\mathcal{X}$  be an input space, with  $\mathcal{Y} = \{1, 2, 3, \dots, C\}$  as the corresponding output space. A classifier  $f$  is trained on a training data set  $\mathcal{D}_{in} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , which contains samples drawn from the joint distribution  $P_{XY}$ . We denote the marginal distribution over  $\mathcal{X}$  as  $\mathcal{P}_{in}$ , which becomes the in-distribution of our classifier  $f$ . During inference, the classifier may be exposed to samples which are not from  $\mathcal{P}_{in}$ , which are unlikely to be classified correctly by the classifier. The aim of OOD detection is to detect such samples. Thus, OOD detection is a binary classification problem, where the goal is to design a scoring function  $g$  and a corresponding threshold  $\delta$  such that

$$\text{OOD}(\mathbf{x}) = \begin{cases} \text{in} & g(\mathbf{x}) \geq \delta \\ \text{out} & g(\mathbf{x}) < \delta \end{cases} \quad (1)$$

## 4 SalAgg+MLS

In this section, we introduce our framework for using XAI saliency methods for OOD detection. This framework aggregates the saliency values outputted by methods such as GradCAM, GBP or LIME, and uses this value to determine OOD-ness. Given the large number of existing XAI saliency

methods and possible aggregate functions (such as the mean, maximum, third quartile or median), the framework opens up for a range of new OOD detection methods. In addition to the information extracted by saliency aggregation, we also utilize the Maximum Logit score. As Shekhar et al. [14] have shown, SoTA OOD detection performance can be achieved by combining different methods which extract information from the network in different ways. Inspired by this work, we present the SalAgg+MLS (Saliency Aggregation plus Maximum Logit Score) framework:

Under this framework, the OOD score is essentially a sum of a saliency aggregate and the maximum logit score, for a given prediction. However, due to the fact that both the logits and saliency aggregates can be of arbitrary magnitude, we must normalize them before summing if we to control how much each part contributes to the final score. Thus, we can sum the Z-scores of each metric instead. This ensures that the values of the maximum logit and the saliency aggregate are distributed similarly. To calculate the Z-scores, we subtract the mean and divide by the standard deviation over an entire ID validation dataset, for each metric. Thus, we calculate the mean and standard deviations of the maximum logit over an ID validation set  $\mu_{\text{MLS}}^{\text{id}}$  and  $\sigma_{\text{MLS}}^{\text{id}}$ , as well as the mean and standard deviation of the aggregate of saliencies  $\mu_{\text{Agg}}^{\text{id}}$  and  $\sigma_{\text{Agg}}^{\text{id}}$ .

To align ourselves with convention in the field of OOD detection, we must define the OOD detection score as one which is higher for ID data than for OOD data. However, given that our framework does not place any limits on the choice of aggregation, we cannot know in advance whether a specific aggregation will have higher or lower values for ID data. To keep the framework as general as possible, we therefore include a *sign* factor, which multiplies the saliency aggregate by 1 or  $-1$ , depending on whether the ID aggregates are higher or lower, respectively. To do this, we calculate the mean value of the aggregation metric over a validation OOD dataset. We denote this value as  $\mu_{\text{Agg}}^{\text{ood}}$ . The sign factor can then be calculated by  $\text{sign}(\mu_{\text{Agg}}^{\text{id}} - \mu_{\text{Agg}}^{\text{ood}})$ , which we denote as  $S$ .

We have now defined the necessary variables required to describe this framework mathematically. We assume we have a model  $f : \mathbf{x} \rightarrow \mathbb{R}^C$ , an XAI saliency mapping method  $s : (f, \mathbf{x}) \rightarrow \mathbb{R}^{K \times N \times M}$ , and an aggregation function  $A : \mathbb{R}^{K \times N \times M} \rightarrow \mathbb{R}$ . The saliency mapping method is defined as one which does not normalize or rectify its outputs, meaning that if we use methods such as GradCAM, we must modify them to remove the normalization step. An OOD detector under this framework then has the following form, given a threshold  $\delta$ :

$$g(\mathbf{x}; s, A, \delta) = \begin{cases} \text{in} & S \cdot \frac{A(s(\mathbf{x}, f)) - \mu_{\text{Agg}}^{\text{id}}}{\sigma_{\text{Agg}}^{\text{id}}} + \frac{\max_i f(\mathbf{x}) - \mu_{\text{MLS}}^{\text{id}}}{\sigma_{\text{MLS}}^{\text{id}}} \geq \delta \\ \text{out} & S \cdot \frac{A(s(\mathbf{x}, f)) - \mu_{\text{Agg}}^{\text{id}}}{\sigma_{\text{Agg}}^{\text{id}}} + \frac{\max_i f(\mathbf{x}) - \mu_{\text{MLS}}^{\text{id}}}{\sigma_{\text{MLS}}^{\text{id}}} < \delta \end{cases} \quad (2)$$

In fact, this detector can be simplified somewhat. Consider the following:

$$S \cdot \frac{A(s(\mathbf{x}, f)) - \mu_{\text{Agg}}^{\text{id}}}{\sigma_{\text{Agg}}^{\text{id}}} + \frac{\max_i f(\mathbf{x}) - \mu_{\text{MLS}}^{\text{id}}}{\sigma_{\text{MLS}}^{\text{id}}} = \quad (3)$$

$$S \left( \frac{A(s(\mathbf{x}, f))}{\sigma_{\text{Agg}}^{\text{id}}} - \frac{\mu_{\text{Agg}}^{\text{id}}}{\sigma_{\text{Agg}}^{\text{id}}} \right) + \frac{\max_i f(\mathbf{x})}{\sigma_{\text{MLS}}^{\text{id}}} - \frac{\mu_{\text{MLS}}^{\text{id}}}{\sigma_{\text{MLS}}^{\text{id}}} = \quad (4)$$

$$S \cdot \frac{A(s(\mathbf{x}, f))}{\sigma_{\text{Agg}}^{\text{id}}} + \frac{\max_i f(\mathbf{x})}{\sigma_{\text{MLS}}^{\text{id}}} - \left( S \cdot \frac{\mu_{\text{Agg}}^{\text{id}}}{\sigma_{\text{Agg}}^{\text{id}}} + \frac{\mu_{\text{MLS}}^{\text{id}}}{\sigma_{\text{MLS}}^{\text{id}}} \right) \quad (5)$$

All the values in the third term of the above summation are constants; they do not depend on  $\mathbf{x}$ . Thus, we can disregard these terms, as all they do is shift all outputs by a constant value. The final OOD detector thus has the following form:

$$g(\mathbf{x}; s, A, \delta) = \begin{cases} \text{in} & S \cdot \frac{A(s(\mathbf{x}, f))}{\sigma_{\text{Agg}}^{\text{id}}} + \frac{\max_i f(\mathbf{x})}{\sigma_{\text{MLS}}^{\text{id}}} \geq \delta \\ \text{out} & S \cdot \frac{A(s(\mathbf{x}, f))}{\sigma_{\text{Agg}}^{\text{id}}} + \frac{\max_i f(\mathbf{x})}{\sigma_{\text{MLS}}^{\text{id}}} < \delta \end{cases} \quad (6)$$

To use a method under this framework, one simply chooses an XAI saliency mapping method and an aggregate. In our testing, we have chosen LIME [12], Occlusion [24], GradCAM [13], Integrated Gradients [17] and GBP [16] as the XAI saliency mapping methods. For aggregate functions, six aggregates which capture the magnitude of the saliency values in different ways have been chosen. These are the mean, median, vector norm, range, maximum value and third quartile. In addition, three aggregate functions which are (more or less) scale invariant have also been chosen, which measure the statistical spread of the saliencies without considering the magnitude. These aggregates are the Coefficient of Variation (CV), the Relative Mean Absolute Difference (RMD) and the Quartile Coefficient of Determination (QCD). By including these aggregate functions, we can compare the performance of magnitude variant and invariant functions and investigate the importance of saliency magnitudes in an OOD detection context.

## 5 Results

This section is divided into three parts: First, we show that the magnitudes of saliency values, by themselves, are discriminative in an OOD detection context, even without combining them with other OOD detection methods. We do this by using different saliency aggregates to perform OOD detection directly on the ImageNet200 [5] and CIFAR10 [7] OpenOOD benchmarks [22]. In addition, we show that aggregates which are magnitude invariant exhibit very low discriminative power, showcasing the importance of using raw saliency values as opposed to normalized heatmaps. Next, we show that saliency aggregates are not strongly correlated with baseline OOD detection metrics such as the MLS, justifying their combination in the SalAgg+MLS framework. Finally, we test the performance of a selection of methods under this framework, and show that XAI based OOD detection can be competitive with SoTA methods.

In all cases, the experiments are conducted on pretrained ResNet architectures with weights provided as part of OpenOOD. This ensures that the results are comparable to other methods tested under this framework.<sup>1</sup>

### 5.1 The discriminative performance of saliency aggregates

Figure 1 shows the combined Near- and Far-OOD detection performance on both ImageNet200 and CIFAR10 for all combinations of XAI methods and aggregate functions. As we can see, simply aggregating the saliency maps generated on ID and OOD data points allows for efficient detection of OOD data, with the best performing combination (GradCAM saliency mapping with vector norm aggregation) achieving an AUROC score of 88.68%. Among the other XAI saliency mapping methods, we also see high AUROC scores, given a suitable choice of aggregation. In comparison, the baseline methods MLS and MSP achieve average AUROC scores of 88.23 and 88.05, respectively. As we see, simply aggregating saliency maps can lead to results which are above baseline models which use the model confidence.

Another interesting observation is that while the first six aggregations (those which capture magnitude) perform well, the last three (which are magnitude invariant) perform very poorly. This suggests that it is primarily the magnitude information of saliency values which is important for OOD detection, not their spread or relation to each other.

In most cases, ID saliencies have higher magnitudes than OOD saliencies, similarly to how ID data points more often have higher maximum logit scores than OOD data points. However, depending on the choice XAI saliency mapping method and aggregate function, this is not always the case, as can be seen in the appendix.

### 5.2 Correlation between saliency aggregates and the maximum logit

Next, we show that the discriminatory power of saliency aggregations is not simply due to a high correlation with the model confidence in a potential OOD sample. Like in the previous section, we use the ImageNet200 and CIFAR10 OpenOOD benchmarks for our analysis, and average the results over both benchmarks. The full results can be found in the appendix.

<sup>1</sup><https://zjysteven.github.io/OpenOOD/>

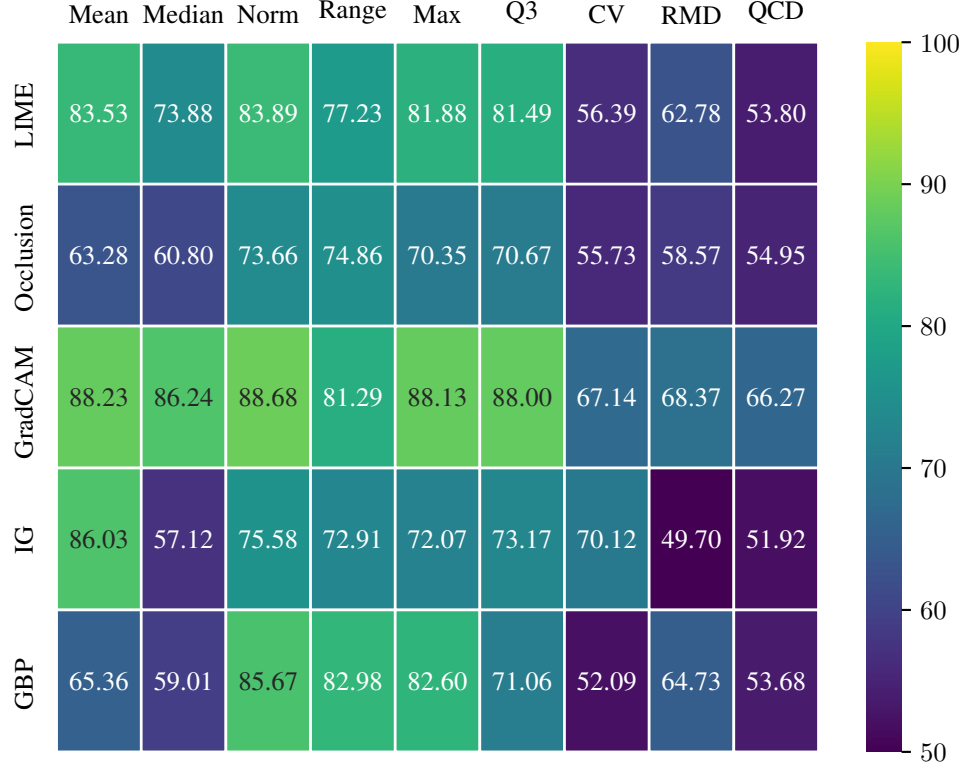


Figure 1: Heatmap of overall AUROC performance for all combinations of XAI methods and aggregations on ImageNet200 and CIFAR10

Figure 2 shows the correlation between saliency aggregates and the MLS, averaged over both the ImageNet200 and CIFAR10 benchmarks. As we can see, GradCAM is very strongly correlated with the MLS. In fact, it can be shown that performing mean aggregation on GradCAM saliencies is equivalent to MLS OOD detection, given certain conditions (which happen to be satisfied in our case).<sup>2</sup> However, the other saliency methods are less correlated, showing that the discriminatory performance of saliency aggregation is not just due to their connection with the logit of the predicted class.

### 5.3 Performance of SalAgg+MLS

Finally, we present the results of using the proposed framework, which combines saliency aggregations and the maximum logit score to predict whether a data point is OOD or not. To do this, we select five combinations of saliency mapping methods and aggregations to test and compare against the baseline MLS. These combinations are selected by choosing the best performing aggregation for each XAI method, as reported in Figure 1. To ensure unbiased results, the OpenOOD benchmarks have been split into validation and testing benchmarks, with the choice of aggregation being done on the validation benchmarks, while the final results are reported on the testing benchmarks. In addition, the testing benchmarks have been bootstrapped ten times, allowing for statistical comparisons between the baseline MLS and the methods which complement MLS by adding a saliency aggregate.

<sup>2</sup>See the appendix.

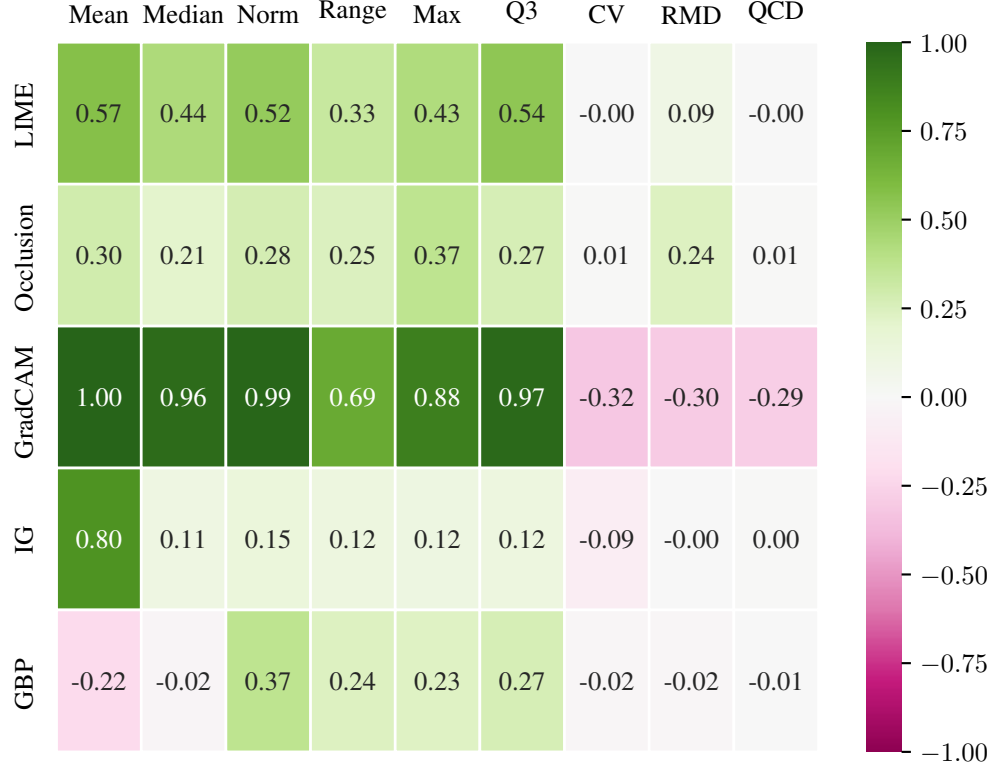


Figure 2: Heatmap showing the correlation between the saliency aggregate of a given XAI method and the MLS, averaged over all four OpenOOD benchmark datasets

As we can see from Table 1, in all benchmarks except for CIFAR100 Near-OOD, there are statistically significant improvements over just using the MLS for multiple methods, showing that XAI saliency aggregations add relevant information to traditional OOD detection methods. In addition, we see that combining the vector norm of GBP saliencies with the MLS is particularly effective, being the best performing method in every benchmark except for CIFAR100 Near-OOD.

Finally, we compare the best performing method under the SalAgg+MLS framework, GBPNorm, to SoTA OOD detection methods, in Table 2. In this case, the testing is done on the entire benchmark as developed by OpenOOD, as opposed to a testing split. This is done to ensure accurate comparison with the results reported by OpenOOD [25], which use the entire dataset. From this table we see that although GBPNorm does not achieve the highest AUROC in any benchmark, its performance is quite close to the SoTA in many cases. This demonstrates the clear potential of using XAI saliency maps for OOD detection.

## 6 Conclusion

This paper presents the first comprehensive investigation of the discriminative power of XAI saliency aggregation in an OOD detection context. We show that many commonly used XAI saliency mapping methods generate outputs which have differing magnitudes between ID and OOD samples. In addition, we show that these saliency magnitudes are not strongly correlated with the MLS, meaning that XAI saliency maps extract different aspects of the model’s behaviour. Finally, we introduce the SalAgg+MLS framework, which describes a method of combining saliency aggregates with the MLS.

Dataset	MLS	LIMENorm	OccRange	GradCAMNorm	IGMean	GBPNorm
ImageNet200 Near-OOD	82.75	83.46**	80.21	83.15**	83.15**	<b>83.70**</b>
ImageNet200 Far-OOD	91.28	93.12**	92.20**	92.09**	91.04	<b>94.03**</b>
ImageNet1K Near-OOD	76.33	74.23	70.95	76.68**	75.37	<b>78.98**</b>
ImageNet1K Far-OOD	89.51	92.17**	88.98	91.25**	88.85	<b>94.09**</b>
CIFAR10 Near-OOD	86.97	87.80**	88.00**	86.92	87.21**	<b>89.85**</b>
CIFAR10 Far-OOD	91.82	93.35**	90.54	91.81	92.47**	<b>94.79**</b>
CIFAR100 Near-OOD	<b>81.24</b>	78.92	78.56	81.24	80.58	78.73
CIFAR100 Far-OOD	79.86	82.13**	81.39**	79.93**	81.95**	<b>84.18**</b>

Table 1: Average AUROC scores over ten bootstraps for the five example methods under the SalAgg+MLS framework, as well as the MLS. The asterisks denote the Bonferroni corrected statistical significance of a Wilcoxon signed-rank test done against the null hypothesis that the developed methods are no better than the MLS. For each benchmark, the best performing method is highlighted in bold.

Dataset	RotPred	CombOOD	OE+MSP	AdaScale-L	GBPNorm
ImageNet200 Near-OOD	81.59	<b>95.74</b>	85.73	84.84	83.37
ImageNet200 Far-OOD	92.56	92.57	89.02	<b>94.86</b>	93.87
ImageNet1K Near-OOD	76.52	<b>95.22</b>	N/A	84.27	78.90
ImageNet1K Far-OOD	90.00	90.24	N/A	<b>97.28</b>	94.03
CIFAR10 Near-OOD	92.68	91.13	<b>94.82</b>	74.99	89.90
CIFAR10 Far-OOD	96.62	94.65	<b>96.00</b>	79.02	94.92
CIFAR100 Near-OOD	76.43	78.77	<b>88.30</b>	80.54	78.73
CIFAR100 Far-OOD	88.40	<b>85.87</b>	81.41	83.38	84.18

Table 2: The GBPNorm method under the SalAgg+MLS compared to the performance of SoTA OOD detection methods as reported by Zhang et al. [25]. The best performing method in each case is highlighted in bold.

When using GBP and vector norm aggregation within this framework, we achieve results which are close to the SoTA on OpenOOD benchmarks, especially on Far-OOD.

## References

- [1] Osvaldo Arreche, Tanish R. Guntur, Jack W. Roberts, and Mustafa Abdallah. E-xai: Evaluating black-box explainable ai frameworks for network intrusion detection. *IEEE Access*, 12:23954–23988, 2024.
- [2] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, March 2018.
- [3] Shaveta Dargan, Munish Kumar, Maruthi Rohit Ayyagari, and Gulshan Kumar. A survey of deep learning and its applications: A new paradigm to machine learning. *Archives of Computational Methods in Engineering*, 27(4):1071–1092, 09 2020.
- [4] Eoin Delaney, Derek Greene, and Mark T. Keane. Uncertainty estimation and out-of-distribution detection for counterfactual explanations: Pitfalls and solutions, 2021.
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [6] Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild, 2021.



- [7] Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009.
- [8] Basim Mahbooba, Mohan Timilsina, Radhya Sahal, and Martin Serrano. Explainable artificial intelligence (xai) to enhance trust management in intrusion detection systems using decision tree model. *Complexity*, 2021(1):6634811, 2021.
- [9] Aitor Martinez-Seras, Javier Del Ser, and Pablo Garcia-Bringas. Can post-hoc explanations effectively detect out-of-distribution samples? In *2022 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE)*, pages 1–9, 2022.
- [10] Sajid Nazir, Diane M. Dickson, and Muhammad Usman Akram. Survey of explainable artificial intelligence techniques for biomedical imaging with deep neural networks. *Computers in biology and medicine*, 156:106668, 2023.
- [11] Alvin Rajkomar, Jeffrey Dean, and Isaac Kohane. Machine learning in medicine. *New England Journal of Medicine*, 380(14):1347–1358, 2019.
- [12] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier, 2016.
- [13] 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. *International Journal of Computer Vision*, 128(2):336–359, October 2019.
- [14] Shashi Shekhar, Vagelis Papalexakis, Jing Gao, Zhe Jiang, and Matteo Riondato. *Proceedings of the 2024 SIAM International Conference on Data Mining (SDM)*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2024.
- [15] John Sipple and Abdou Youssef. A general-purpose method for applying explainable ai for anomaly detection. In Michelangelo Ceci, Sergio Flesca, Elio Masciari, Giuseppe Manco, and Zbigniew W. Raś, editors, *Foundations of Intelligent Systems*, pages 162–174, Cham, 2022. Springer International Publishing.
- [16] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net, 2015.
- [17] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks, 2017.
- [18] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks, 2014.
- [19] AJ Tallón-Ballesteros and C Chen. Explainable ai: Using shapley value to explain complex anomaly detection ml-based systems. *Machine Learning and Artificial Intelligence: Proceedings of MLIS 2020*, 332:152, 2020.
- [20] Erzhena Tcydenova, Tae Woo Kim, Changhoon Lee, and Jong Hyuk Park. Detection of adversarial attacks in ai-based intrusion detection systems using explainable ai. *Human-Centric Comput Inform Sci*, 11, 2021.
- [21] Jordan Zheng Ting Sim, Qi Wei Fong, Weimin Huang, and Cher Heng Tan. Machine learning in medicine: what clinicians should know. *Singapore Med J*, 64(2):91–97, May 2021.
- [22] Jingkan Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, WenXuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyun Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Dan Hendrycks, Yixuan Li, and Ziwei Liu. OpenOOD: Benchmarking generalized out-of-distribution detection. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- [23] Tahmina Zebin, Shahadate Rezvy, and Yuan Luo. An explainable ai-based intrusion detection system for dns over https (doh) attacks. *IEEE Transactions on Information Forensics and Security*, 17:2339–2349, 2022.

- [24] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 818–833, Cham, 2014. Springer International Publishing.
- [25] Jingyang Zhang, Jingkang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyu Sun, Xuefeng Du, Yixuan Li, Ziwei Liu, Yiran Chen, and Hai Li. Openood v1.5: Enhanced benchmark for out-of-distribution detection. *arXiv preprint arXiv:2306.09301*, 2023.

## Appendix

### Proof of Equivalence between GradCAM Saliency Aggregation and MLS

The following section proves that taking the mean value of the saliency values outputted by GradCAM is equivalent to MLS in an OOD detection context, under certain conditions. This shows that XAI saliency mapping methods extract information which can be used for OOD detection.

In this case, we assume that the classification stage of the Convolutional Neural Network (CNN) is a simple Global Average Pooling (GAP) over the feature map followed by a single linear layer. Such a structure is the classification head of all ResNet models. Furthermore, we choose to perform GradCAM on the final layer of the network, which is recommended by Selvaraju et al. [13]. These two conditions give rise to the following theorem:

**Theorem 1.** Assume we have a network where  $y^c = \sum_k \text{mean}(F_k) \cdot W_{ck}$ , where  $F$  is a convolutional feature map and  $W$  is a linear layer of size  $\text{channels} \times \text{classes}$ . Taking the mean value of the unrectified output of GradCAM for this network and a given input is equal to the MLS up to a constant  $a$ , and thus equivalent to MLS in an OOD detection context:

$$\text{mean}(\text{GradCAM}_F(\mathbf{x})) = a \cdot \text{MLS}(\mathbf{x})$$

*Proof.* Following Selvaraju et al. [13], the saliency map generated by GradCAM has the following definition:

$$\text{GradCAM}_F(\mathbf{x}) = \text{ReLU} \left( \sum_k \left( \frac{1}{N \cdot M} \sum_i \sum_j \frac{\delta y^c}{\delta F_{ij}^k} \right) F^k \right). \quad (7)$$

Here,  $F^k$  is the  $k$ 'th channel of the final convolutional feature map, while  $N$  and  $M$  are its dimensions. The above equation simply describes averaging the gradients of the logit of class  $c$  for each channel, and using these values to perform a weighted sum of the channels in the feature map. While  $c$  can be any class index, in this case, we define  $c = \max_i f_i(\mathbf{x})$ , i.e we calculate the saliency map for the predicted class. We remove the Rectified Linear Unit (ReLU) activation function from the GradCAM definition and use the raw saliency values, given that this is specified in the theorem. Taking the mean of this function, we get the following equation:

$$\text{mean} \left( \sum_k \left( \frac{1}{N \cdot M} \sum_i \sum_j \frac{\delta y^c}{\delta F_{ij}^k} \right) F^k \right). \quad (8)$$

Given the assumptions of the theorem, we have that the logit  $y^c$  for class  $c$  is calculated in the following manner:

$$y^c = \sum_k \text{mean}(F_k) \cdot W_{ck} \quad (9)$$

$$= \sum_k \left( \frac{\sum_i \sum_j F_{ij}^k}{N \cdot M} \cdot W_{ck} \right). \quad (10)$$

This equation simply describes GAP (all channels are averaged to a single number) followed by a single linear layer (each logit is a weighted sum of the average pooled channels, with the weights defined by a specific row/column in the weight matrix  $W$ ). Given our definition of  $c = \max_i f_i(\mathbf{x})$ ,  $y^c = \text{MLS}(\mathbf{x})$ . We return to Equation 8

$$\text{mean} \left( \sum_k \left( \frac{1}{N \cdot M} \sum_i \sum_j \frac{\delta y^c}{\delta F_{ij}^k} \right) F^k \right). \quad (11)$$

Given equation 10,

$$\frac{\delta y^c}{\delta F_{ij}^k} = \frac{W_{ck}}{N \cdot M}. \quad (12)$$

As we can see, the indices  $i$  and  $j$  have disappeared. This is to be expected, as global average pooling means that all values in each channel are multiplied by the same value when calculating the logit of a specific class. We may now substitute this derivative in equation 11:

$$\text{mean} \left( \sum_k \left( \frac{1}{N \cdot M} \sum_i \sum_j \frac{W_{ck}}{N \cdot M} \right) F^k \right). \quad (13)$$

We now perform some simple algebra, exploiting the fact that  $\text{mean}(a \cdot \mathbf{x}) = a \cdot \text{mean}(\mathbf{x})$  and that  $\sum_i c \cdot x_i = c \sum_i x_i$ :

$$\text{mean} \left( \sum_k \left( \frac{1}{N \cdot M} \sum_i \sum_j \frac{W_{ck}}{N \cdot M} \right) F^k \right) \quad (14)$$

$$\text{mean} \left( \sum_k \left( \sum_i \sum_j \frac{W_{ck}}{N \cdot M} \right) F^k \right) \quad (15)$$

$$\frac{1}{N \cdot M} \text{mean} \left( \sum_k \left( (N \cdot M) \frac{W_{ck}}{N \cdot M} \right) F^k \right) \quad (16)$$

$$\frac{1}{N \cdot M} \text{mean} \left( \sum_k W_{ck} \cdot F^k \right) \quad (17)$$

$$\frac{1}{N \cdot M} \cdot \left( \sum_k W_{ck} \cdot \text{mean}(F^k) \right). \quad (18)$$

Let us denote  $1/(N \cdot M)$  as  $a$ . We recognize the second factor as  $y^c = \text{MLS}(\mathbf{x})$  as described in equation 9. We then have

$$\text{mean}(\text{GradCAM}_F(\mathbf{x})) = a \cdot \text{MLS}(\mathbf{x}). \quad (19)$$

□

This theorem shows that XAI saliency mapping methods, although they have been developed for an entirely different purpose than OOD detection, may also collect information which can be used for OOD detection.

### Further results

Figure 1 aggregated AUROC scores over both Near- and Far-OD and over the two OpenOOD benchmarks CIFAR10 and ImageNet200. Below follows the four individual results; Near- and Far-OD for CIFAR10 and ImageNet200.

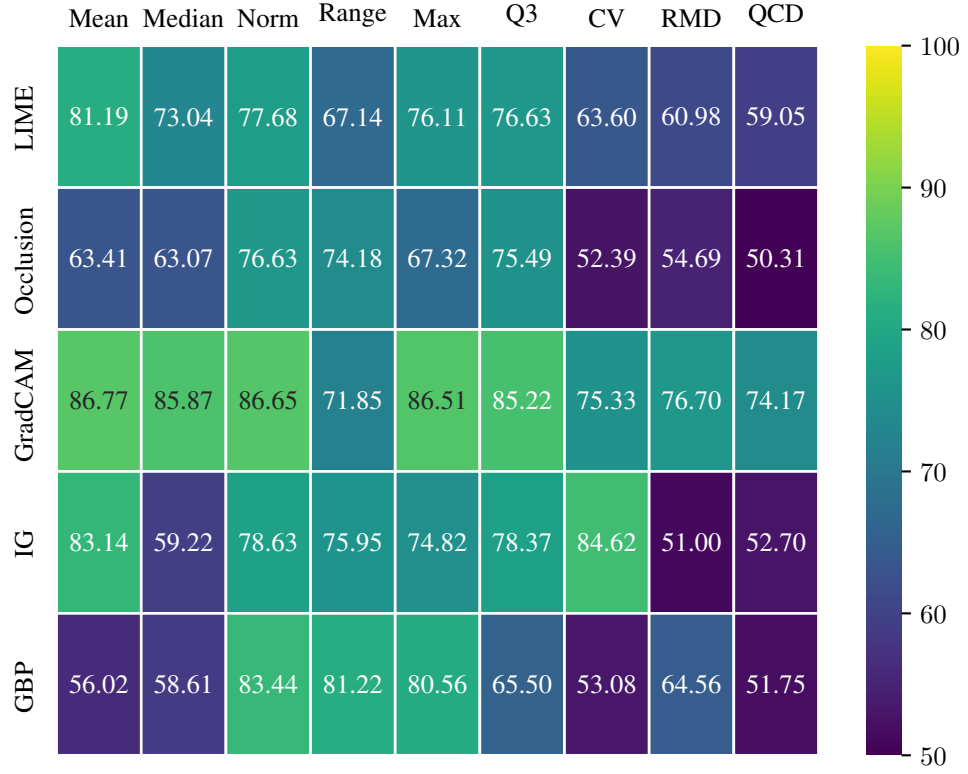


Figure 3: Near-OOD AUROC performance on CIFAR10 for all combinations of XAI methods and aggregations. Baseline performance was 86.75% for MLS and 87.69% for MLS. Due to numerical imprecision, the AUROC score of the mean of GradCAM is not exactly equal to the AUROC of MLS.

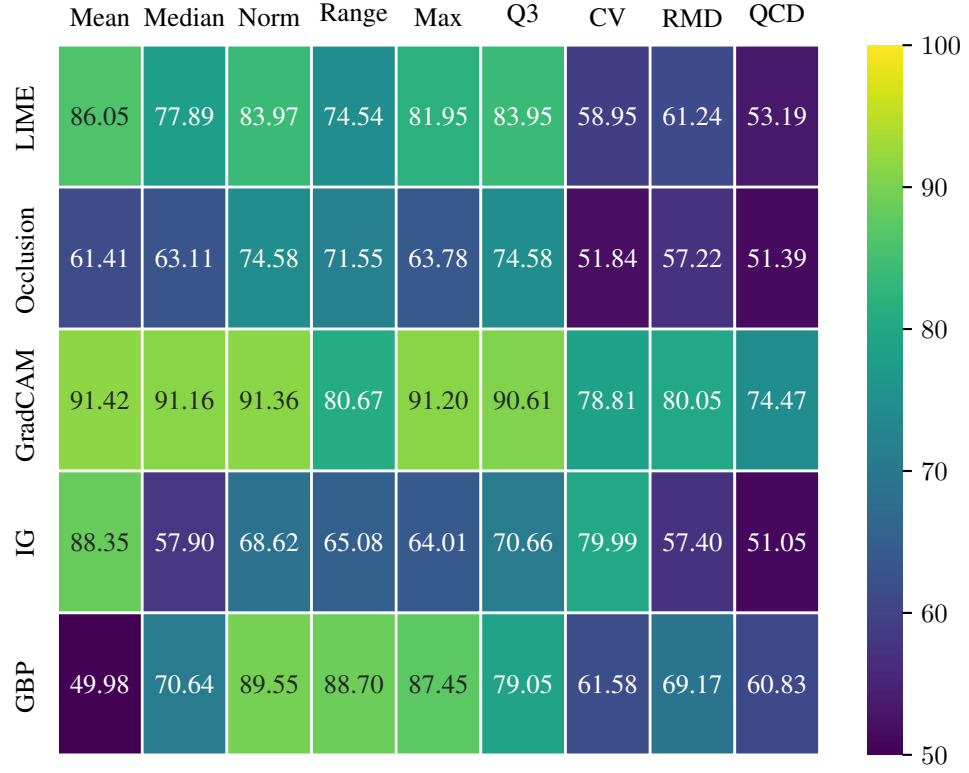


Figure 4: Far-OOD AUROC performance on CIFAR10 for all combinations of XAI methods and aggregations. Baseline performance was 91.39% for MLS and 90.78% for MLS.

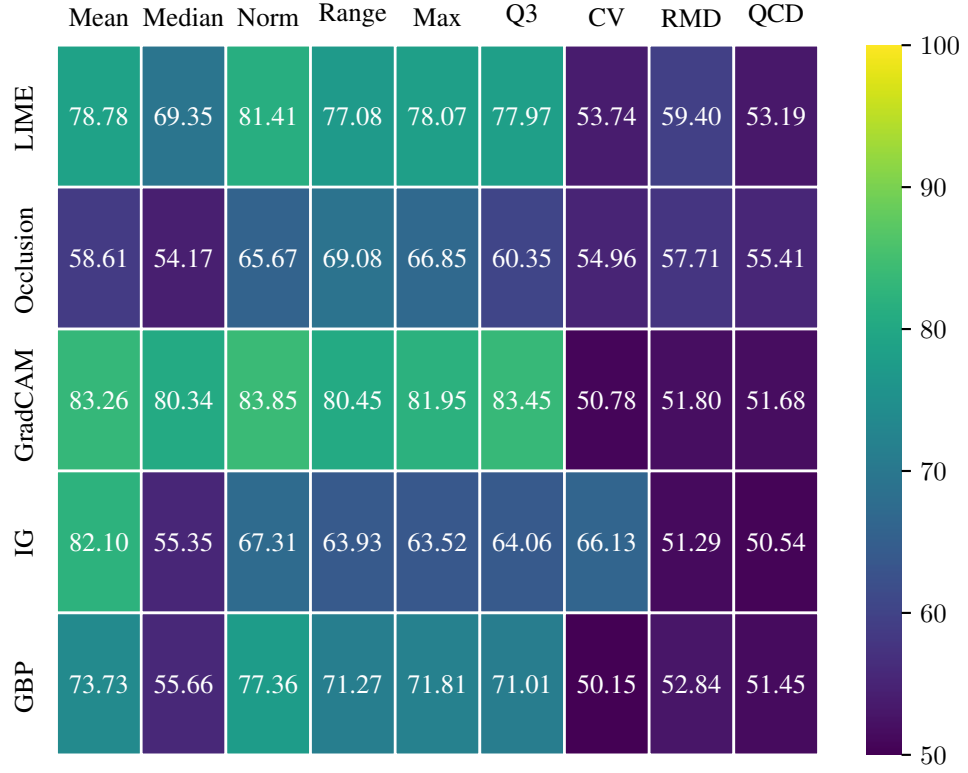


Figure 5: Near-OOD AUROC performance on ImageNet200 for all combinations of XAI methods and aggregations. The baseline performance was 83.28% for MLS and 83.43% for MLS.

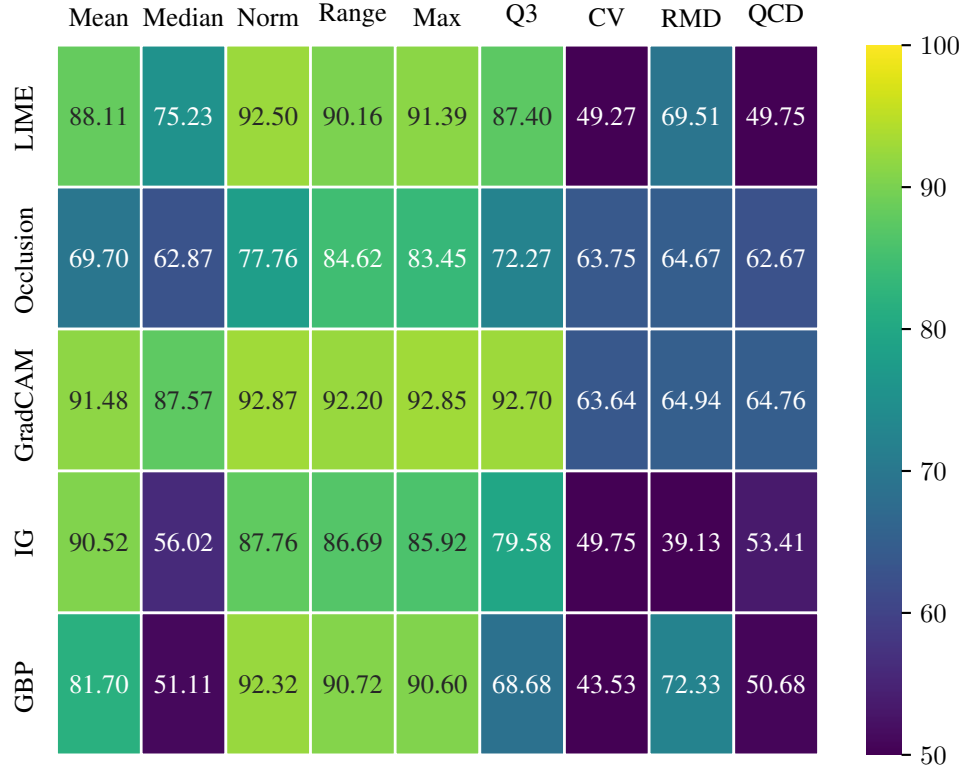


Figure 6: Far-OOD AUROC performance on ImageNet200 for all combinations of XAI methods and aggregations. The baseline performance was 91.48% for MLS and 90.30% for MLS.



### **Hyperparameters used**

As explained in section 5, the testing is done on pretrained ResNet classifiers supplied by [22]. Thus, we do not list the hyperparameters used during training, as this is outside of the scope of this paper. However, the XAI saliency mapping methods also accept hyperparameters, and these are listed below:

#### **LIME**

We use simple rectangular segmentation, for LIME, masking out regions by setting their pixel values to zero. Each image is split into 16 equally sized rectangular regions, by dividing the height and width by 4. Thus, for the CIFAR benchmarks, the rectangular regions are  $32/4 \times 32/4 = 8 \times 8$ , while the regions are  $224/4 \times 224/4 = 56 \times 56$  on the ImageNet benchmarks. The masking probability was 50% and the kernel width was 0.25.

#### **Occlusion**

Like with LIME, we have used a rectangular segmentation approach with the same number of rectangular regions, and zero based masking.

#### **GradCAM**

Following the recommendations of Selvaraju et al. [13], we have calculated gradients from the final layer when calculating saliencies using GradCAM.

#### **Guided Backpropagation**

GBP is non-parametric.

#### **Integrated Gradients**

As recommended by Sundararajan et al. [17], we have chosen the zero vector as the baseline for the calculation of the integrated gradients.