# TorchEsegeta: Framework for Interpretability and Explainability of Image-based Deep Learning Models

Soumick Chatterjee[a,b,c,**], Arnab Das[a], Chirag Mandal[a], Budhaditya Mukhopadhyay[a], Manish Vipinraj[a], Aniruddh Shukla[a], Rajatha Nagaraja Rao[a], Chompunuch Sarasaen[c,d], Oliver Speck[c,e,f,g], Andreas Nürnberger[a,b,f]

[a]*Faculty of Computer Science, Otto von Guericke University Magdeburg, Germany*
[b]*Data and Knowledge Engineering Group, Otto von Guericke University Magdeburg, Germany*
[c]*Biomedical Magnetic Resonance, Otto von Guericke University Magdeburg, Germany*
[d]*Institute for Medical Engineering, Otto von Guericke University Magdeburg, Germany*
[e]*German Center for Neurodegenerative Disease, Magdeburg, Germany*
[f]*Center for Behavioral Brain Sciences, Magdeburg, Germany*
[g]*Leibniz Institute for Neurobiology, Magdeburg, Germany*

## Abstract

Clinicians are often very sceptical about applying automatic image processing approaches, especially deep learning based methods, in practice. One main reason for this is the black-box nature of these approaches and the inherent problem of missing insights of the automatically derived decisions. In order to increase trust in these methods, this paper presents approaches that help to interpret and explain the results of deep learning algorithms by depicting the anatomical areas which influence the decision of the algorithm most. Moreover, this research presents a unified framework, TorchEsegeta, for applying various interpretability and explainability techniques for deep learning models and generate visual interpretations and explanations for clinicians to corroborate their clinical findings. In addition, this will aid in gaining confidence in such methods. The framework builds on existing interpretability and explainability techniques that are currently focusing on classification models, extending them to segmentation tasks. In addition, these methods have been adapted to 3D models for volumetric analysis. The proposed framework provides methods to quantitatively compare visual explanations using infidelity and sensitivity metrics. This framework can be used by data scientists to perform post-hoc interpretations and explanations of their models, develop more explainable tools and present the findings to clinicians to increase their faith in such models. The proposed framework was evaluated based on a use case scenario of vessel segmentation models trained on Time-of-fight (TOF) Magnetic Resonance Angiogram (MRA) images of the human brain. Quantitative and qualitative results of a comparative study of different models and interpretability methods are presented. Furthermore, this paper provides an extensive overview of several existing interpretability and explainability methods.

*Keywords:* Deep Learning, Black-box, Interpretability, Explainability, Model introspection, MRA segmentation

## 1. Introduction

The use of artificial intelligence is widely prevalent today in medical image analysis. It is, however, imperative to do away with the black-box nature of Deep Learning techniques to gain the trust of radiologists and clinicians, as a model's erroneous output might have a high impact in the medical domain.

*Interpretability vs Explainability:* Interpretability and Explainability methods aim to unravel the black-box nature of decision-making in machine learning and deep learning models. Explainability focuses on explaining the internal working mechanisms of the model. In contrast, interpretability focuses on observing the effects of changes in model parameters and inputs on the model prediction, hence attributing properties to the input, and it also requires a context, considering the audience who is to interpret the properties and characteristics for the model outcomes (Marcinkevičs and Vogt, 2020;

Chakraborty et al., 2017). It is to be noted that earlier work failed to draw a clear line of distinction between interpretability and explainability as these terms are subjective and pertain to the stakeholders, considering the audience who must understand the model and its outcomes. Interpretability and explainability nurture a sense of human-machine trust as they help the users of the machine/deep-learning models understand how certain decisions are made by the model and are not limited to statistical metric-based approaches such as accuracy or precision. This in turn, allows employing such models in mission-critical problems such as medicine, autonomous driving, legal systems, banking etc. The basis of measuring model transparency is said to comprise simulatability, decomposability, and algorithmic transparency (Chakraborty et al., 2017); and model functionality which consists of textual descriptions of the model's output and visualisations of the model's parameters. The goals to be attained through interpretability of models (Marcinkevičs and Vogt, 2020) are trust, reliability, robustness, fairness, privacy, and causality. Explainability can be formulated as the explanation of the decisions made internally by the model that in turn

generates the observable, external conclusions arrived at by the model. This promotes human understanding of the model's internal mechanisms and rationale, which in turn helps build user trust in the model (Emmert-Streib et al., 2020). The evaluation criteria for model explainability include comprehensibility by humans, fidelity and scalability of the extracted representations, and scalability and generality of the method (Belle and Papantonis, 2020).

Some models are transparent and hence inherently understandable, such as Decision Trees, K-Nearest Neighbours, Rule-based and Bayesian models, but some are opaque and require post-hoc explanations. A post-hoc model explanation can be through model-agnostic and model-specific techniques that aim towards explaining a black-box model in human-understandable terms. The European Union has also incorporated transparency and accountability of models in 2016 as a criterion in the ethical guidelines of trustworthy AI (Arrieta et al., 2020; Choo and Liu, 2018). Through interpretability and explainability, to some degree, the following goals can be achieved: trustworthiness, causality, confidence, fairness, informativeness, transferability and interactivity, which in turn can help improve the model (Arrieta et al., 2020).

Data scientists can generate interpretability-explainability results for their opaque models and present the model's reasoning or the model's inner mechanism to the domain experts (e.g. clinicians). If the reasoning is the same as a human domain expert would have done, then the experts can have more faith in that model, which in turn implies that the model can be incorporated into real-life workflows. Apart from building trust in the models, interpretability and explainability techniques can be used by the data scientists to improve their models as well, by discussing the outcomes with the domain experts and improving the model based on expert feedback. Moreover, accurate (verified by experts) interpretability-explainability results can be used as part of automatic or semi-automatic teaching programs for trainees.

For classification models, there are multiple interpretability ad explainability techniques supported by various libraries such as Captum[1], Torchray (Fong et al., 2019) and CNN Visualization library [2] and many of the techniques are introduced in the following section of this paper. However, this is more challenging for segmentation as the output is more complex than just a simple class prediction. In this contribution, various interpretability and explainability techniques used for classification models have been adapted to work with segmentation models. On applying the interpretability techniques, essential features or areas of the input image can be visualised, on which the model's output is critically based. By applying explainability techniques, a better understanding of the knowledge represented in the model's parameters can be achieved.

## 1.1. Contributions

This paper proposes a unified, flexible and scaleable interpretability-explainability pipeline for PyTorch, Torch-

---

[1]`https:captum.ai`
[2]`https://github.com/utkuozbulak/pytorch-cnn-visualizations`

Esegeta, which leverages post-hoc interpretability and explainability methods and can be applied on 2D or 3D deep learning models working with images. Apart from implementing the existing methods for classification models, this research extends them for segmentation models. This pipeline can be applied to trained models with little to no modification, using either a graphical user interface for easy access or directly using a python script. It also provides an easy platform for incorporating new interpretability or explainability techniques. Moreover, this pipeline provides features to evaluate the interpretability and explainability methods using two different methods. First, using cascading randomisation of the model's weights, which can help evaluate how much the results are dependent on the actual weights, and second, using quantitative metrics: infidelity and sensitivity. Finally, the pipeline was applied on models trained to segment vessels from magnetic resonance angiograms (MRA), and the interpretability results are shown here. Furthermore, apart from the technical contributions with the TorchEsegeta pipeline, this paper also provides a comprehensive overview of several post-hoc interpretability and explainability techniques.

## 2. State-of-the-art methods

Various interpretability and explainability techniques have been proposed over the years to address the black-box problem of deep learning. Some of those techniques are explored here in this section.

### 2.1. Interpretability

Interpretability techniques help to understand the focus area of a model - can help understand the reasoning done by the model. These techniques can be categorised into two groups: model attribution and layer attribution.

#### 2.1.1. Model Attribution Techniques

Model attribution techniques are the techniques to assess the contribution of each attribute to the prediction of the model. There is a long list of methods available in the literature under the rubric model attribution techniques, hence for better understanding, they are further divided here into two groups: feature-based and gradient-based.

1. *Feature-based Techniques:*
   Methods under this group bring into play input and/or output feature space to compute local or global model attributions.
   (a) *Feature Permutation* - This is a perturbation based technique (Breiman, 2001; Fisher et al., 2019) in which the value of an input or group of inputs is changed utilising random permutation of values in a batch of inputs and calculating its corresponding change in output. Hence, meaningful feature attributions are calculated only when a batch of inputs is provided.

(b) *Shapley Value Sampling* - The proposed method defines an input baseline to which all possible permutations of the input values are added one at a time, and the corresponding output values and hence each feature attribution is calculated. For permutation O, given player set N, the set of all possible permutations $\pi(N)$ and the predecessors of player i $Pre^i(O)$, the Shapley value is given by

$$Sh_i(v) = \sigma_{O \in \pi(N)} \frac{1}{n!} (v(Pre^i(O) \bigcup i) - v(Pre^i(O)))$$
(1)

As all possible permutations are considered, this technique is computationally expensive, and this can be overcome by sampling the permutations and averaging their marginal contributions instead of considering all possible permutations such as the ApproShapley sampling technique (Castro et al., 2009)

(c) *Feature Ablation* - This method works by replacing an input, or a group of inputs with another value defined by a range or reference value and the feature attribution of the input or group of inputs such as a segment are computed. This method works based on perturbation.

(d) *Occlusion* - Similar to the Feature Ablation method, Occlusion (Zeiler and Fergus, 2014) works as a perturbation based approach wherein the continuous inputs in a rectangular area are replaced by a value defined by a range or a reference value. Using this approach, the change in the corresponding outputs is calculated in order to find the attribution of the feature.

(e) *RISE Randomised Input Sampling for Explanation of Black-box Models* - It generates an importance map indicating how salient each pixel is for the model's prediction (Petsiuk et al., 2018). RISE works on black-box models since it does not consider gradients while making the computations. In this approach, the model's outputs are tested by masking the inputs randomly and calculating the importance of the features.

(f) *Extremal Perturbations* - They are regions of an image that maximally affect the activation of a certain neuron in a neural network for a given area in an image (Fong et al., 2019). Extremal Perturbations lead to the largest change in the prediction of the deep neural network, when compared to other perturbations defined by

$$m_a = \underset{m:\|m\|_1 a|\Omega|, m \in M}{\arg\max} \quad \Phi(m \otimes x)$$
(2)

for the chosen area a. The paper also introduces area loss to enforce the restrictions while choosing the perturbations.

(g) *Score-weighted Class Activation (Score CAM)* - It is a gradient-independent interpretability technique based on class activation mapping (Wang et al., 2020). Activation maps are first extracted, and each activation then works as a mask on the original image, and its forward-passing score on the target class is obtained. Finally, the result can be generated by the linear combination of score-based weights and activation maps[3]. Given a convolutional layer l, class c, number of channels k and activations A:

$$L^c_{Score-CAM} = ReLU(\sum_k \alpha^c_k A^k_l)$$
(3)

*Gradient-based Techniques:*. This group of methods mainly use the model's parameter space to generate the attribute maps - calculated with the help of the gradients.

(a) *Saliency* - It was initially designed for visualising the image classification done by convolutional networks and the saliency map for a specific image (Simonyan et al., 2013). It is used to check the influence of each pixel of the input image in assigning a final class score to the image $I$ using the linear score model $S_c(I) = w_c^T I + b_c$ for weight w and bias b.

(b) *Guided Backpropagation* - In this approach, during the backpropagation, the gradients of the outputs are computed with respect to the inputs (Springenberg et al., 2014). RELU activation function is applied to the input gradients, and direct backpropagation is performed, ensuring that the backpropagation of non-negative gradients does not occur.

(c) *Deconvolution* - This approach is similar to the guided backpropagation approach wherein it applies RELU to the output gradients instead of the input gradients and performs direct backpropagation (Zeiler and Fergus, 2014). Similarly, RELU backpropagation is overridden to allow only non-negative gradients to be backpropagated.

(d) *Input X Gradient* - It extends the saliency approach such that the contribution of each input to the final prediction is checked by multiplying the gradients of outputs and their corresponding inputs in a setting of a system of linear equations AX = B where A is the gradients and B is the calculated final contribution of input X (Shrikumar et al., 2016).

(e) *Integrated gradients* - This technique (Sundararajan et al., 2017) calculates the path integral of the gradients along the straight line path from the baseline $x'$ to the input x. It satisfies the axioms of Completeness i.e. the attributions must account for the difference in output for the baseline $x'$ and input x, Sensitivity i.e. a non-zero attribution must be provided even to inputs that flatten the prediction function where the input differs from the baseline, and Implementation Invariance of gradients i.e. two functionally equivalent networks must have identical feature attributions. It requires no instrumentation of the deep neural network for its application. All the gradients along the

---

[3]`https://github.com/haofanwang/Score-CAM`

straight line path from $x'$ to x are integrated along the $i^{th}$ dimension as:

$$IG_i(x) ::= (x_i - x'_i) \times \int_{\alpha=0}^{1} \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha \tag{4}$$

(f) *Grad Times Image* - In this technique (Shrikumar et al., 2017b) the gradients are multiplied with the image itself. It is a predecessor of the DeepLift method as the activation of each neuron for a given input is compared to its reference activation, and contribution scores are assigned according to the difference for each neuron.

(g) *DeepLift* - It is a method (Shrikumar et al., 2017a) that considers not only the positive but also the negative contribution scores of each neuron based on its activation concerning its reference activation. The difference in the output of the activation function of the neuron t under observation is calculated based on the difference in input with respect to a reference input. Contribution scores for each of the preceding neurons $x_1, ..x_n$ that influence t are assigned $C_{\triangle x_{i_{\triangle t}}}$ that sums up to the difference in t's activation:

$$\sum_{i=1}^{n} C_{\triangle x_{i_{\triangle t}}} = \triangle t \tag{5}$$

(h) *DeepLiftShap* - This method extends the DeepLift method and computes the SHAP values on an equivalent, linear approximation of the model (Lundberg and Lee, 2017). It assumes the independence of input features. For model f, the effective linearization from SHAP for each component is computed as:

$$\phi_i(f_3, y) \approx m_{y_i f_3}(y_i - E[y_i]) \tag{6}$$

(i) *GradientShap* - This technique also assumes the independence of the inputs and computes the game-theoretic SHAP values of the gradients on the linear approximation of the model (Lundberg and Lee, 2017). Gaussian noise is added to randomly selected points, and the gradients of their corresponding outputs are computed.

(j) *Guided GradCAM* - The Gradient-weighted Class Activation Mapping approach provides a means to visualise the regions of an image input that are predominant in influencing the predictions made by the model. This is done by visualising the gradient information pertaining to a specific class in any layer of the user's choice. It can be applied to any CNN-based model architecture. The guided backpropagation and GradCAM approaches are combined by computing the element-wise product of guided backpropagation attributions with upsampled GradCAM attributions (Selvaraju et al., 2017).

(k) *Grad-Cam++* - Generalised Gradient-based Visual Explanations for Deep Convolutional Networks is a method (Chattopadhay et al., 2018) that claims to provide better predictions than the Grad-CAM and other state-of-the-art approaches for object localisation and explaining occurrences of multiple object instances in a single image. This technique uses a weighted combination of the positive partial derivatives of the last convolutional layer's feature maps concerning a specific class score as weights to generate a visual explanation for the corresponding class label. The importance $w_k^c$ of an activation map $A^k$ over class score $Y^c$ is given by

$$w_k^c = \sum_i \sum_j \alpha_{ij}^{kc} . relu(\frac{\partial Y^c}{\partial A_{ij}^k}) \tag{7}$$

Grad-CAM++ provides human-interpretable visual explanations for a given CNN architecture across multiple tasks, including classification, image caption generation and 3D action recognition.

(l) *Vanilla Backpropagation and layer visualisation* - It is the standard gradient backpropagation technique through the deep neural network wherein the gradients are visualised at different layers and what is being learnt by the model is observed, given a random input image.

(m) *Smooth Grad* - In this method (Wang et al., 2020; Smilkov et al., 2017) random Gaussian noise $N(0, \sigma^2)$ is added to the given input image and the corresponding gradients are computed to find the average values over n image samples

$$\hat{M}_c(x) = \frac{1}{n} \sum_1^n M_c(x + N(0, \sigma^2)) \tag{8}$$

Vanilla and guided backpropagation techniques can be used to calculate the gradients.

### 2.1.2. Layer Attribution Techniques

Layer Attribution methods evaluate the effect of each individual neuron in a particular layer on the model output. There are various types of layer attribution methods that have been explored as part of this research work. Perturbation-Based approaches (Wang et al., 2020) perturb the original input and observe the change in the prediction of the model, and are highly time-consuming.

1. *Inverted representation* - The aim of this technique is to generate the given input image after a number of n target layers. An inverse of the given input image representation is computed (Mahendran and Vedaldi, 2014) in order to find an image $\Phi(x_0)$ whose representation best matches the input image $\Phi_0$ while minimising the given loss l such that

$$x^* = \underset{x \in \mathbb{R}^{HxWxC}}{\arg \min} l(\Phi(x), \Phi_0) + \lambda R(x) \tag{9}$$

2. *Layer Activation with Guided Backpropagation* - This method (Springenberg et al., 2015) is quite similar to guided backpropagation, but instead of guiding the signal

from the last layer and a specific target, it guides the signal from a specific layer and filter. The guided backpropagation method adds an additional guidance signal from the higher layers to the usual backpropagation.

3. *Layer DeepLift* - This is the DeepLift method as mentioned in the model attribution techniques (Lundberg and Lee, 2017), but applied with respect to the particular hidden layer in question.

4. *Layer DeepLiftShap* - This method is similar to the DeepLIFT SHAP technique mentioned in the model attribution techniques (Lundberg and Lee, 2017), applied for a particular layer. The original distribution of baselines is taken, and the attribution for each input-baseline pair is calculated using the Layer DeepLIFT method, and the resulting attributions are averages per input example. Assuming model linearity, $\phi(f_3, y) \approx m_y i$

5. *Layer GradCam* – The GradCam attribution for a given layer is provided by this method (Selvaraju et al., 2017). The target output's gradients are computed concerning the particularly given layer. The resultant gradients are averaged for each output channel (dimension 2 of output). The average gradient for each channel is then multiplied by the layer activations. The results are then added over all the channels. For the class feature weights w, global average pooling is performed over the feature maps A such that

$$S^c = \sum_k w_k^c \frac{1}{Z} \sum_i \sum_j A_{ij}^k \qquad (10)$$

6. *Layer GradientShap* - This is analogous with GradientSHAP (Lundberg and Lee, 2017) method as mentioned in the model attribution techniques but applied for a particular layer. Layer GradientSHAP adds Gaussian noise to each input sample multiple times, wherein a random point along the path between baseline and input is selected, and the gradient of the output with respect to the identified layer is computed. The final SHAP values approximate the expected value of gradients * (layer activation of inputs - layer activation of baselines).

7. *Layer Conductance* – This method (Dhamdhere et al., 2018) provides the conductance of all neurons of a particular hidden layer. Conductance of a particular hidden unit refers to the flow of Integrated Gradients attribution through this hidden unit. The main idea behind using this method is to decompose the computation of the Integrated Gradients via the chain rule. One property that this method of conductance satisfies is that of completeness. Completeness means that the conductances of a particular layer add up to the prediction difference of F(x) − F(x′) for input x and baseline input x′. Other properties that are satisfied by this method are that of Linearity and insensitivity.

8. *Internal Influence* - This method calculates the contribution of a layer to the final prediction of the model by integrating the gradients with respect to the particular layer under observation. The internal representation is influenced by an element j as defined by (Leino et al., 2018)

such that

$$\chi_i^s(f, P) = \int_\chi \frac{\partial g}{\partial z_j}\Big|_{h(x)} P(x)dx \qquad (11)$$

wherein $s = \langle g, h \rangle$ for the slice of the network represented by s as a function of tuple g,h. It is similar to the integrated gradients approach where instead of the input, the gradients are integrated with respect to the layer.

9. *Contrastive Excitation Backpropagation/Excitation Backpropagation* - This approach is used to generate and visualise task-specific attention maps. Excitation Backprop method as proposed by (Zhang et al., 2018) is to pass along top-down signals downwards in the network hierarchy via a probabilistic Winner-Take-All process wherein the most relevant neurons in the network are identified for a given top-down signal. Both top-down and bottom-up information is integrated to compute the winning probability of each neuron as defined by

$$y_i = \sum_{j=1}^N w_{ij} x_j \qquad (12)$$

for input x and weight matrix w. The contrastive excitation backpropagation is used to make the top-down attention maps more discriminative.

10. *Layer Activation* – It computes the activation of a particular layer for a particular input image (Liu et al., 2020). It helps to understand how a given layer reacts to an input image. One can get an excellent idea of what part or features of the image at which a particular layer looks.

11. *Linear Approximator* - This is a technique to overcome inconsistencies of post-hoc model interpretation; linear approximator combines a piecewise linear component and a nonlinear component (Guo et al., 2020; Fong et al., 2019). The piecewise linear component describes the explicit feature contributions by piecewise linear approximation that increases the expressiveness of the deep neural network. The nonlinear component uses a multi-layer perceptron to capture feature interactions and implicit nonlinearity, which in turn increases the prediction performance. Here, the interpretability is obtained once the model is learned in the form of feature shapes and has high accuracy.

12. *Layer Gradient X Activation* – This method (Ancona et al., 2017) computes the element-wise product of a given layer's gradient and activation. It is a combination of the gradient and activation methods of layer attribution. The output attributions are returned as a tuple if the layer input/output contains multiple tensors or a single tensor is returned.

*2.2. Explainability*

Explainability techniques aim to unravel the internal working mechanism of a model - tries to explain the knowledge represented inside the model's parameters.

### 2.2.1. DeepDream

As the network is trained using many examples, it is essential to check what has been learnt from the input image. Hence, visualising what has been learnt at different layers of the CNN based network gives rise to repetitive patterns of different levels of abstraction, enabling to interpret of what has been learnt by the layers. This visualisation can be realised using Deep-Dream. Given an arbitrary input image, any layer from the network can be picked, and the detected features at that layer can be enhanced. It is observed that the initial layers are sensitive to basic features in the input images, such as edges, and the deeper layers identify complex features from the input image. As an example, InceptionNet [4] was trained on animal images, and it was observed that the different layers of the network could interpret an interesting remix of the learnt animal features in any given image.

### 2.2.2. LIME

Local Interpretable Model-agnostic Explanations (LIME) is a technique for providing post-hoc model explanations. LIME constructs an approximated surrogate, locally linear model or decision tree of a given complex model that helps to explain the decisions of the original model, and due to its model-agnostic nature, it can be employed to explain any model even when its architecture is unknown (Choo and Liu, 2018; Samek et al., 2020; Arrieta et al., 2020). LIME also performs the transformation of input features to obtain a representation that is interpretable to humans (Belle and Papantonis, 2020). In the original work (Ribeiro et al., 2016), the authors propose LIME, an algorithm that provides explanations for a model f, that are locally faithful in the locality $\Pi_x$, for an individual prediction of the model, such that users can ensure that they can trust the prediction before acting on it. LIME provides an explanation for f in the form of a model g and $g \in G$ where G is a set of all possible interpretable models and $\Omega(g)$ which is the complexity of the interpretable model, is minimised. $L(f, g, \Pi_x)$ is the fidelity function that measures the unfaithfulness of g in estimating f in the locality $\Pi_x$. Hence, the explanation provided by LIME is (Ribeiro et al., 2016):

$$\xi(x) = argmin_{g \in G} L(f, g, \Pi_x) + \Omega(g) \tag{13}$$

### 2.2.3. SHAP

SHapley Additive exPlanation (Lundberg and Lee, 2017) is similar to LIME such that it approximates an interpretable, explanation model g of the original, complex model f, in order to explain a prediction made by the model f(x). SHAP provides post-hoc model explanations for an individual output and is model-agnostic. It calculates the contribution of each feature in producing the final prediction and is based on the principles of Game Theory such as Shapley Values (Belle and Papantonis, 2020). SHAP works on simplified, interpretable input data $x^{\cdot}$, which is analogous to the original input data such that

$x = h_x(x)$. SHAP must ensure consistency theorems (Lundberg and Lee, 2017),local accuracy such that $g(x^{\cdot})$ matches $f(x)$ when $x = h_x(x)$ and missingness such that if an attribute's value is 0, it does not not imply that its corresponding contribution to the prediction is 0. Hence, the contribution of each attribute x can be calculated as follows (Lundberg and Lee, 2017):

$$\phi_i(f, x) = \Sigma_{z^{\cdot} \subseteq x^{\cdot}} \frac{|z^{\cdot}|!(M - |z^{\cdot}| - 1)!}{M!} [f_x(z^{\cdot}) f_x(z^{\cdot} \setminus i)] \tag{14}$$

where $z^{\cdot} \in \{0, 1\}^M$ and $|z^{\cdot}|$ is the number of non-zero entries in $z^{\cdot}$ and $z^{\cdot} \subseteq x^{\cdot}$ represents all $z^{\cdot}$ vectors where the non-zero entries are a subset of the non-zero entries in $x^{\cdot}$.

### 2.2.4. Lucent

Based on Lucid [5], Lucent is a library implemented using PyTorch for the explainability of deep learning models. The implementation provides optivis, which is the main framework for providing the visualisations of parameters learnt by the different layers of the deep neural network. It can be used to visualise torchvision models with no overhead for setup. Activation atlas methods, feature visualisation methods, building blocks and differentiable image parameterisations can be used to visualise the different features learnt by the network. Activation atlas methods comprise different methods which show the network activations for a particular class or average activations in a grid cell of the image. Feature visualisation methods help understand the crucial features for a neuron, entire channel, or layer. Building blocks help to visualise the activation vector and its components for the given image. Differentiable image parameterisations find the types of image generation processes which can be backpropagated through, and this helps to perform appropriate preconditioning of the input image to improve the optimisation of the neural network.

## 3. Architecture of TorchEsegeta

One of the main contributions of this research work is the TorchEsegeta framework which integrates various interpretability and explainability techniques available in different libraries and extends these techniques for segmentation models. It is noteworthy that the development of this pipeline started with the exploration of various interpretability techniques for classifying COVID-19 and other types of pneumoniæ (Chatterjee et al., 2020b). An initial pipeline was developed under that project for classification models, but only for 2D images. This research further extends that for 3D volumetric images, as well as for segmentation models. Apart from incorporating these features, the original pipeline was further streamlined and improved during this research - to create the first version of Torch-Esegeta.

---

[4]https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html

[5]`https://github.com/tensorflow/lucid`

### 3.1. Incorporated libraries

The following interpretability and explainability based libraries have been explored in this research work - Captum, CNN visualisation, TorchRay, DeepDream, Lucent, LIME and SHAP

Captum is a library built on PyTorch and is used to provide the interpretability of machine/deep-learning models. It provides many algorithms that evaluate the contribution of different features in providing a model's prediction and thus helps in improving the model.

CNN visualisations (Ozbulak, 2019) provides different implementations for the different interpretability techniques and visualisations for CNN based models architectures.

TorchRay [6] is a package used for several visualisation methods for convolutional neural network architectures using PyTorch. It focuses on interpretability wherein it attempts to determine which regions of the input image influence the final prediction made by the model (Fong et al., 2019).

LIME (Local Interpretable Model-agnostic Explanations) is a model-agnostic technique that provides post-hoc model explanations to explain the decisions of the deep learning model, and due to its model-agnostic nature, LIME flexibly explains any unknown model (Choo and Liu, 2018; Samek et al., 2020; Arrieta et al., 2020)

SHAP [7] (SHapley Additive exPlanation) is a post-hoc, game-theoretical approach that computes shapley values in order to explain the prediction made by the deep learning model (Lundberg and Lee, 2017).

Lucent [8] is the PyTorch implementation of lucid for the explainability of deep learning models. It aims to explain the decision made by the deep neural network by explaining what is being learnt by the various layers of the network.

DeepDream [9] is also an explainability technique that is used to visualise the parameters learnt by the convolutional neural network.

### 3.2. The Pipeline: TorchEsegeta

The pipeline architecture is shown in Figure 1. The pipeline is implemented keeping in line with the object-oriented methodology. One of the key features of the pipeline is that it is easily scalable, according to the customised needs of the end-user because of the JSON based configuration. The pipeline is a plug and play software and can be easily used by the end-users. This pipeline is publicly available on GitHub[10].

### 3.3. Features of TorchEsegeta

3D input data is especially important for medical images like MRI and CT scans. Hence, the framework has been built to support both 2D and 3D input data. One of the most important features of the pipeline is scalability. New methods can be added to the pipeline seamlessly by the users and used as per requirement.

The JSON based configuration will allow the user to execute the pipeline according to their specific needs and is easy to use the feature. The logging facility allows the users to track the pipeline execution and make changes in case of any error because of wrong parameter usage.

The timeout facility makes sure that no method execution is above a certain threshold of time. That way, it saves both computation resources as well as valuable time for the end-user.

#### 3.3.1. Parallel Execution

Another way of saving much valuable time for the user is by using Multi-GPU, which helps in the parallel execution of multiple methods simultaneously. The multi-threading feature enables the execution of multiple methods on the same GPU allowing the users to make optimum use of resources.

#### 3.3.2. Patch-based Execution

Patch-based models are supported in the pipeline and thus helps running computationally intensive methods. Existing interpretability and explainability methods do not give an option to be directly used on patch-based models.

#### 3.3.3. Automatic Mixed Precision

Automatic Mixed Precision (Micikevicius et al., 2017) facility is also incorporated in the framework aiding in the faster execution of the code. Due to the use of this technique, the overall memory consumption while executing the methods is greatly reduced. The execution time is also reduced as a result.

#### 3.3.4. Wrapper for Segmentation Models

The interpretability techniques that are available publicly are basically for classification problems. The framework is an extension to the segmentation problem with the help of a task-specific wrapper functionality -

1. Pixel-wise multi-class classification
2. Threshold-based pixel classification

*Pixel-wise multi-class classification.* In this method, the pixel scores are summed up for all pixels predicted as each class. Two main steps are performed:

a. Pixel-wise Class Assignment - In this step, for every pixel, the argmax class is computed. Let y be the image with dimensions

$$y_{ij} = \arg\max_k(y_{ijk}) \qquad (15)$$

b. Final Output Calculation - In this step, the sum of the pixel scores for all pixels is predicted as each class is computed.

$$Out_{in} = count\{y_{ij} \mid y_{ij} \in class\ m\} \qquad (16)$$
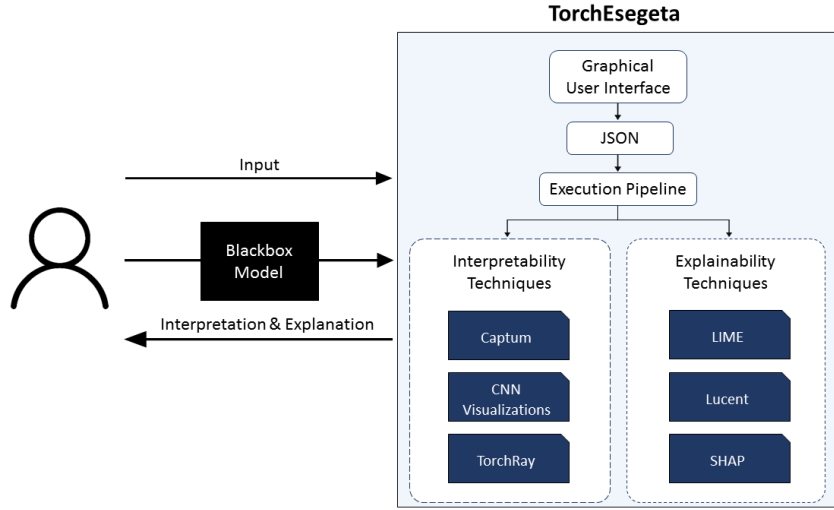
where $0 \leq m \leq N_c$

---

Figure 1: TorchEsegeta pipeline architecture

*Threshold based pixel classification.* This method performs class identification by Otsu threshold and then sums up the pixels for each class. This task is also performed in two steps:

a. Normalisation - In this step, the input image is normalised by the following function:

$$y_{ij_{norm}} = \frac{y_{ij} - min(y)}{max(y) - min(y)} \qquad (17)$$

b. Pixel-wise binarisation - The pixel-wise binarisation is performed with the help of Otsu thresholding.

$$y_{ij} = \begin{cases} 1 & where \quad y_{ij_{norm}} > th \\ 0 & elsewhere \end{cases} \qquad (18)$$

where th = otsu($y_{ij_{norm}}$)

The output for both the processes is a tensor with a single value for each class.

### 3.3.5. Graphical User Interface

A Graphical User Interface (GUI) has been created to provide the end-user with an intuitive graphical layout to select the parameters and interpretability methods according to their requirement. The first screen shown in Figure 2 provides the user with a dialogue box for parameter selection. The user can choose the model nature, model name, dataset and many more run-time parameters. Once these selections are made, the 'Select Methods' button will display the following dialogue box as shown in Figure 3. The users have a wide range of interpretability methods to choose from in this dialogue box. Once the interpretability methods are chosen, the method-specific parameter dialogue box is displayed to the user as in Figure 4. The users can then choose the visualisation method, device id and many other parameters. On clicking the 'Next' button, the code will be executed in the back-end, and the output will be generated in the output path specified in Figure 2.
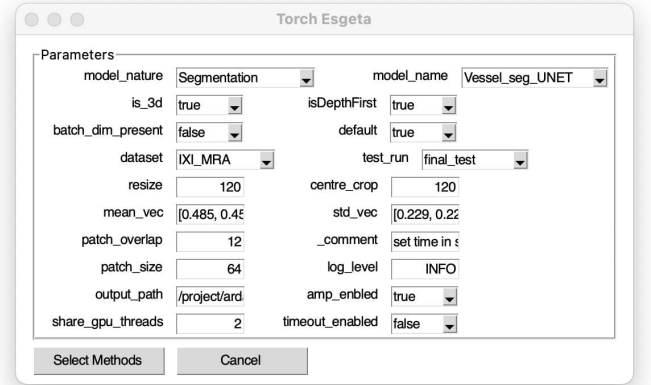


Figure 2: GUI: parameter selection

### 3.4. Evaluation Methods

In order to evaluate and compare the methods, both qualitative and quantitative aspects have been considered.

### 3.4.1. Qualitative evaluation

For qualitative or visual evaluation, the cascading randomisation (Adebayo et al., 2018) technique was implemented, in which the model weights are randomised successively, from the top to the bottom layers. The learned weights of the layers are destroyed from top to bottom by this technique. While doing so, the interpretability and explainability techniques are applied to each state of randomisation. If the interpretability-explainability results are dependent upon the model's weights - how it is supposed to be - then the quality would decrease as the amount of randomisation increases. If they are not dependent, then the results will be unaffected by the randomisation - hence, the technique is not accurate or unsuitable for the current model.
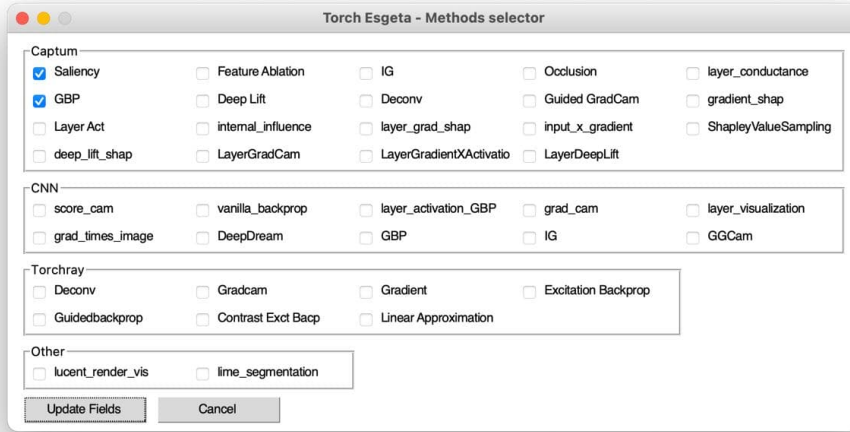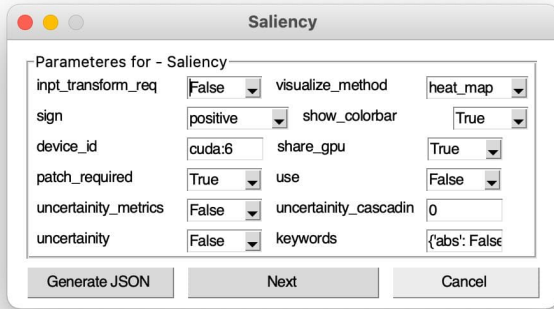
8

Figure 3: GUI: method selection



Figure 4: GUI: Method-specific parameter selection a the chosen method

### 3.4.2. Quantitative evaluation

For the quantitative evaluation, the uncertainty method was chosen due to the unavailability of ground truth data for the attribution images. Two metrics have been used for computing the uncertainty values of the attribution methods: Infidelity (Yeh et al., 2019)] and Sensitivity (Yeh et al., 2019). These metrics can be used only on model attribution methods as of now. It is to be noted that the qualitative metrics can be used on top of cascading randomisation as an additional level of evaluation.

Infidelity[11] represents the expected mean-squared error between the explanation multiplied by a meaningful input perturbation and the differences between the predictor function at its input and perturbed input. Sensitivity[11] measures the extent of explanation change when the input is slightly perturbed.

---

[11] https://captum.ai/api/metrics.html

## 4. Use case: Results and Analysis

To evaluate the TorchEsegeta pipeline for segmentation models, a use-case of vessel segmentation was chosen.

### 4.1. Models

The segmentation network models chosen for this use case are from the DS6 paper (Chatterjee et al., 2020a). The models are UNet, UNet MSS and UNet MSS with Deformation. The difference between the UNetMSS and UNetMSS with Deformation is a change in the number of up-sampling and down-sampling layers and a modified activation function. UNetMSS performs downsampling five times using convolution striding and uses transported convolution for upsampling, and it includes instance normalisation and Leaky ReLU in the convolution blocks. On the contrary, the modified version applies four down-sampling layers using max-pool and performs up-sampling using interpolation combined with Batch normalisation and ReLU in its convolution block.

For UNetMSS with Deformation, a small amount of variable elastic deformation is added at the time of training, along with each volume input. The authors have given a comprehensive overview of the method, and based on their results, one can see that UNetMSS with Deformation is the best performing model.

### 4.2. Use Case Experiment

As a use case study, the previously mentioned interpretability and explainability techniques were explored while analysing the results of a vessel segmentation model trained on Time-of-fight (TOF) Magnetic Resonance Angiogram (MRA) images of the human brain called DS6 (Chatterjee et al., 2020a). The model automatically segments vessels from 3D 7 Tesla TOF-MRA images.

A study about the attribution outputs of the different methods across the three models was conducted. The zoomed-in

portions of the attribution images show even more closely the model's focus areas.

Figure 5 shows examples of interpretability techniques compared across the three models.

Figure 6 portrays similar interpretability results from U-Net MSS Deformation model.

In all of these attribution maps, whichever pixel is highlighted by red represents high activation, indicating the most important pixels in the input according to the respective method, and the other region represents lower to no activation.

### 4.3. Notable observations

In Figure 5, the CNN Visualization Layer Activation Guided Backpropagation attribution map shows a lot of attributions, most of which are on parts other than the brain vessels, i.e. the region of interest. The respective zoomed images also confirm the observation.

On the contrary, methods such as Captum Deeplift and TorchRay Gradient provide fewer attributions and miss out on major portions of the brain vessels.

Torchray Deconvolution provides more attributions comparatively; however, most of the attributions are concentrated on certain areas of the brain.

CNN Visualization Vanilla Backpropagation, Captum Deconvolution and Captum Saliency provide better results compared to the others - they provided more human-interpretable results. These methods mainly attribute on the region of interest. However, a look at the respective zoomed images would reveal that for Captum Saliency, there are attributions even in the areas surrounding the brain vessels. For the other two, the attributions are mainly on the brain vessels only.

Figure 6 shows a comparison of some of the better performing interpretability methods for model U-Net MSS Deformation model.

In Figure 7 and Figure 8, the layer-wise attributions for the three models are shown, for the methods: Excitation Backpropagation, Layer Conductance, Layer DeepLift and Layer Gradient X Activation. For all the methods, the maximum attributions are shown in the Conv3 layer. Among the three models, U-Net MSS Deformation seems to focus better than the other two, as it attributions can be seen all over the brain contrary to the other two (it is to be noted the vessels are present all over the brain, and not focused in a specific region). The network focus shifts while going from the first to the last layer of the models. In the initial layers, the focus is more on selective regions of the brain. However, towards the deeper layers like Conv3, the focus is almost on the entire brain.

### 4.4. Evaluation

A comparative study of the cascading randomisation technique is shown for the outputs of 4 interpretability methods for all the models in Figure 9.

Moreover, to show the functionality of the quantitative evaluation part of the pipeline, a few methods were compared quantitatively using infidelity and sensitivity, the scores are shown in Table 1 and 2.

|  | UNet | UNetMSS | UNetMSS_Deform |
|---|---|---|---|
| Guided Backprop | 5.87e−17 | 2.08e−17 | 2.59e−18 |
| Deconvolution | 3.20e−16 | 2.62e−16 | 1.48e−16 |
| Saliency | 1.95e−15 | 1.54e−15 | 1.23e−16 |

Table 1: Infidelity Scores

|  | UNet | UNetMSS | UNetMSS_Deform |
|---|---|---|---|
| Guided Backprop | 1.156 | 0.917 | 0.831 |
| Deconvolution | 1.210 | 1.188 | 1.140 |
| Saliency | 1.171 | 1.197 | 1.153 |

Table 2: Sensitivity Scores

## 5. Conclusion

In this work, various interpretability and explainability methods were adopted for segmentation models and were used to interpret the network. A pipeline has been developed and made public on GitHub[10], TorchEsegeta, comprised of those interpretability and explainability methods that can be applied on 2D or 3D image-based deep learning models for classification and segmentation. The pipeline was experimented with using three models – UNet, UNetMSS and the UNetMSS-Deform, for the task of vessel segmentation from MRAs. The evaluation of the methods have been done qualitatively using cascading randomisation and quantitatively using evaluation metrics. This pipeline can be used by data scientists to improve their models by tweaking their models based on the shown interpretability and explainability - to improve the reasoning of the models, in turn improving the performance. Moreover, they can use this pipeline to show it to the model's users, for them to have trust in the model while using them in high-risk situations. On the other hand, this pipeline can be used by expert decision-makers, like clinicians, as a decision support system - by understanding the reasoning done by the models, they can get assistance in decision making.

It is noteworthy that the current pipeline can be extended for reconstruction purposes, and new interpretability and explainability methods can be added to the existing pipeline. Ground truth based pixel-wise interpretability for segmentation models can be implemented, which will add a new dimension to the existing work. Nevertheless, the real evaluation of the interpretations and explanations can only be done by the domain experts - in this case of vessel segmentation, by the clinicians - who can judge whether these results are actually useful or not. This step of the evaluation was not performed under the scope of the current work and will be performed in the near future.
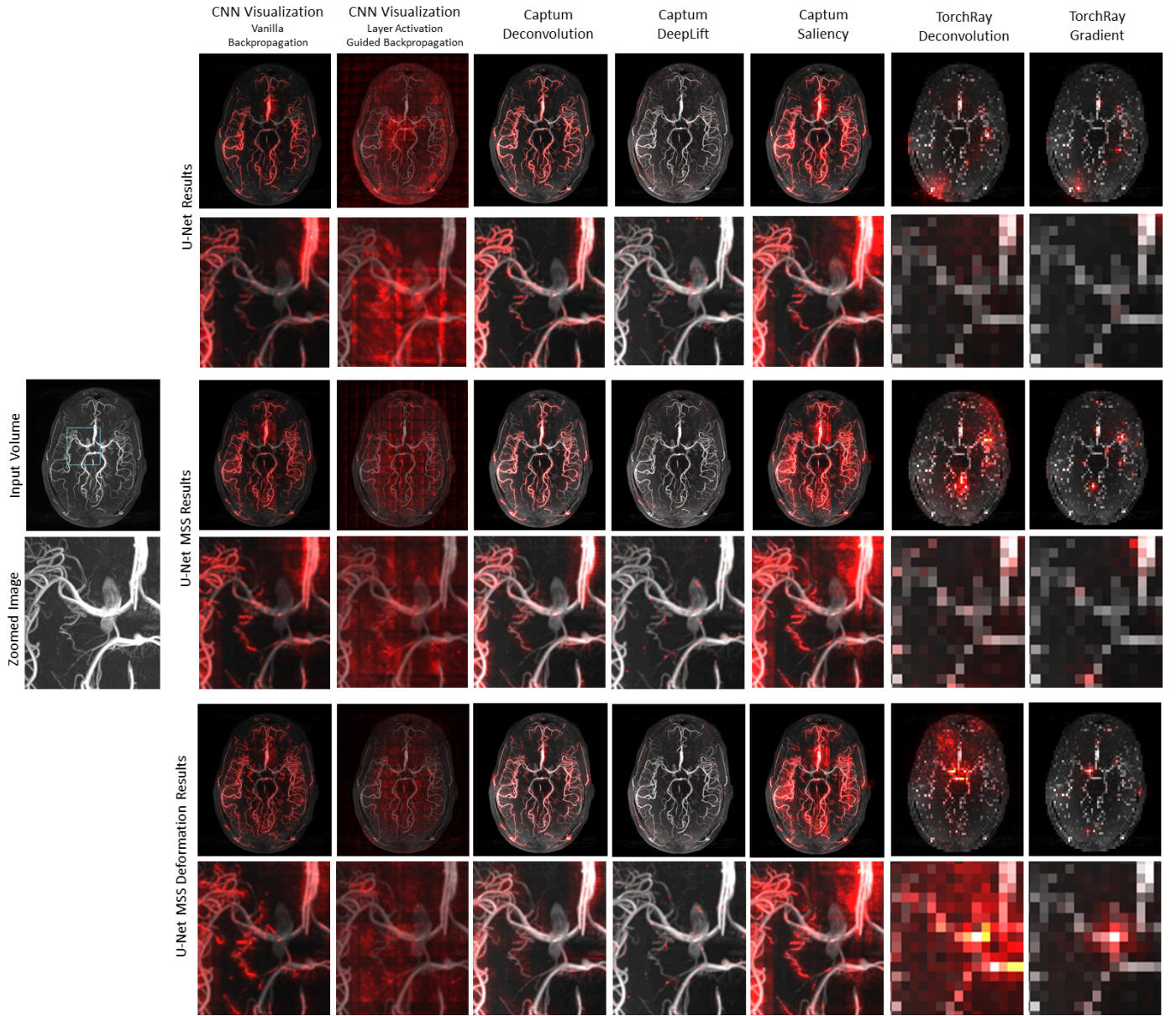
Figure 5: Example of interpretability techniques compared across three models: U-Net, U-Net MSS and U-Net MSS with deformation, along with its corresponding zoomed image.
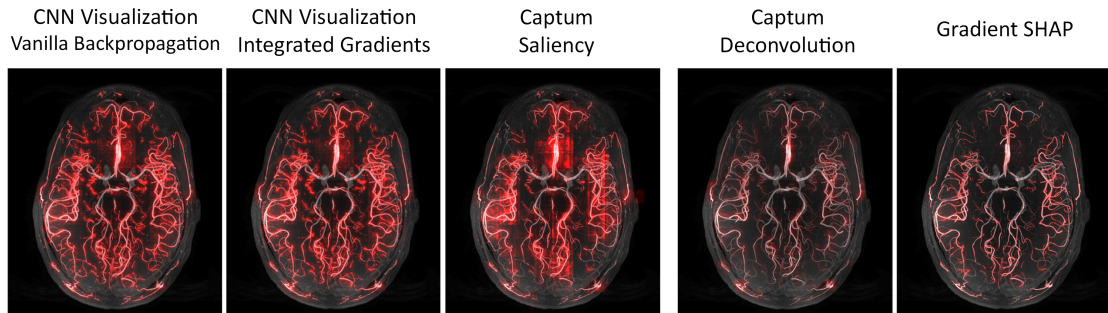


Figure 6: Some of similar interpretability results of U-Net MSS Deformation model.
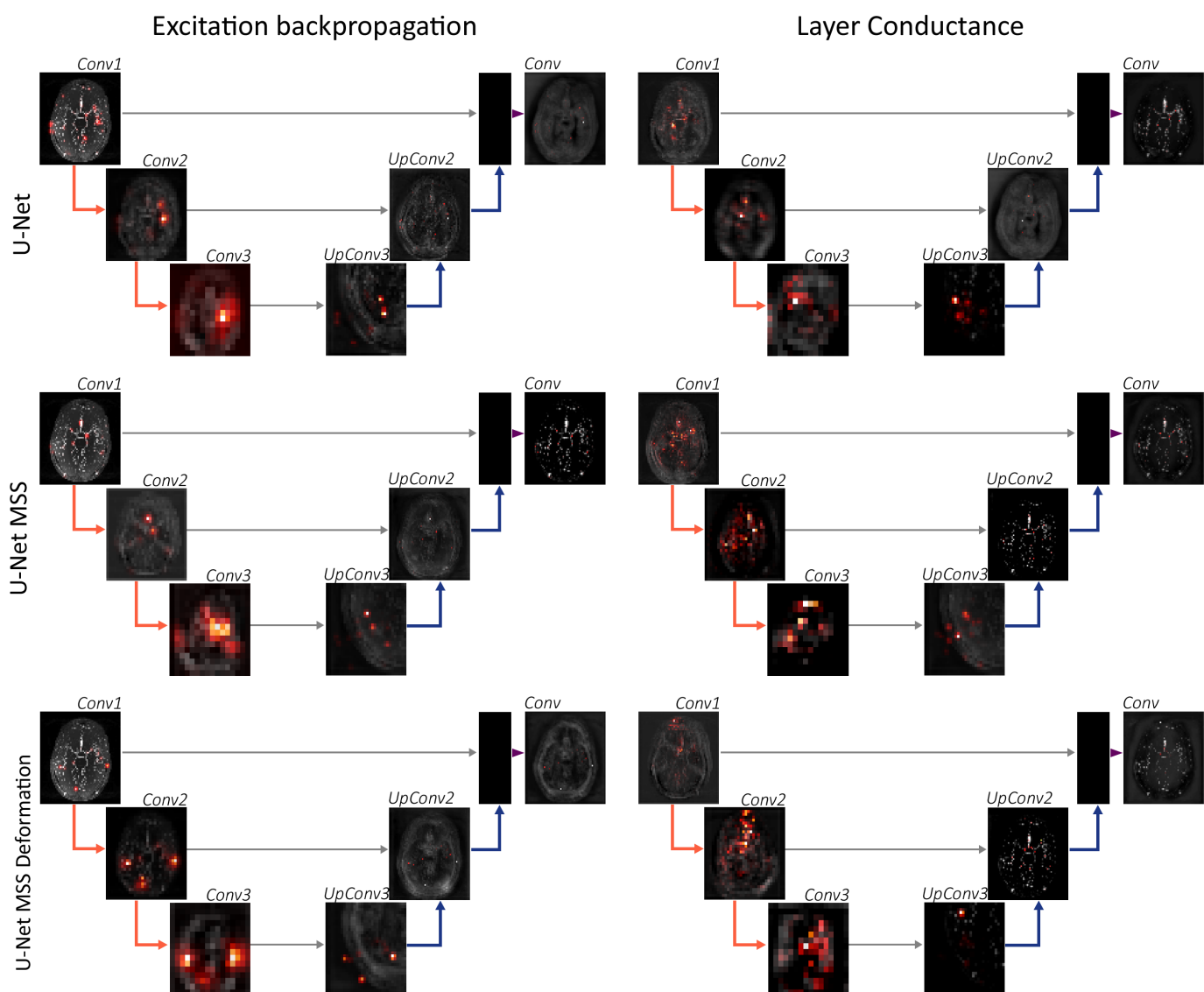
Figure 7: Layer-based interpretability methods: Excitation Backpropagation and Layer Conductance, shown using representative UNets
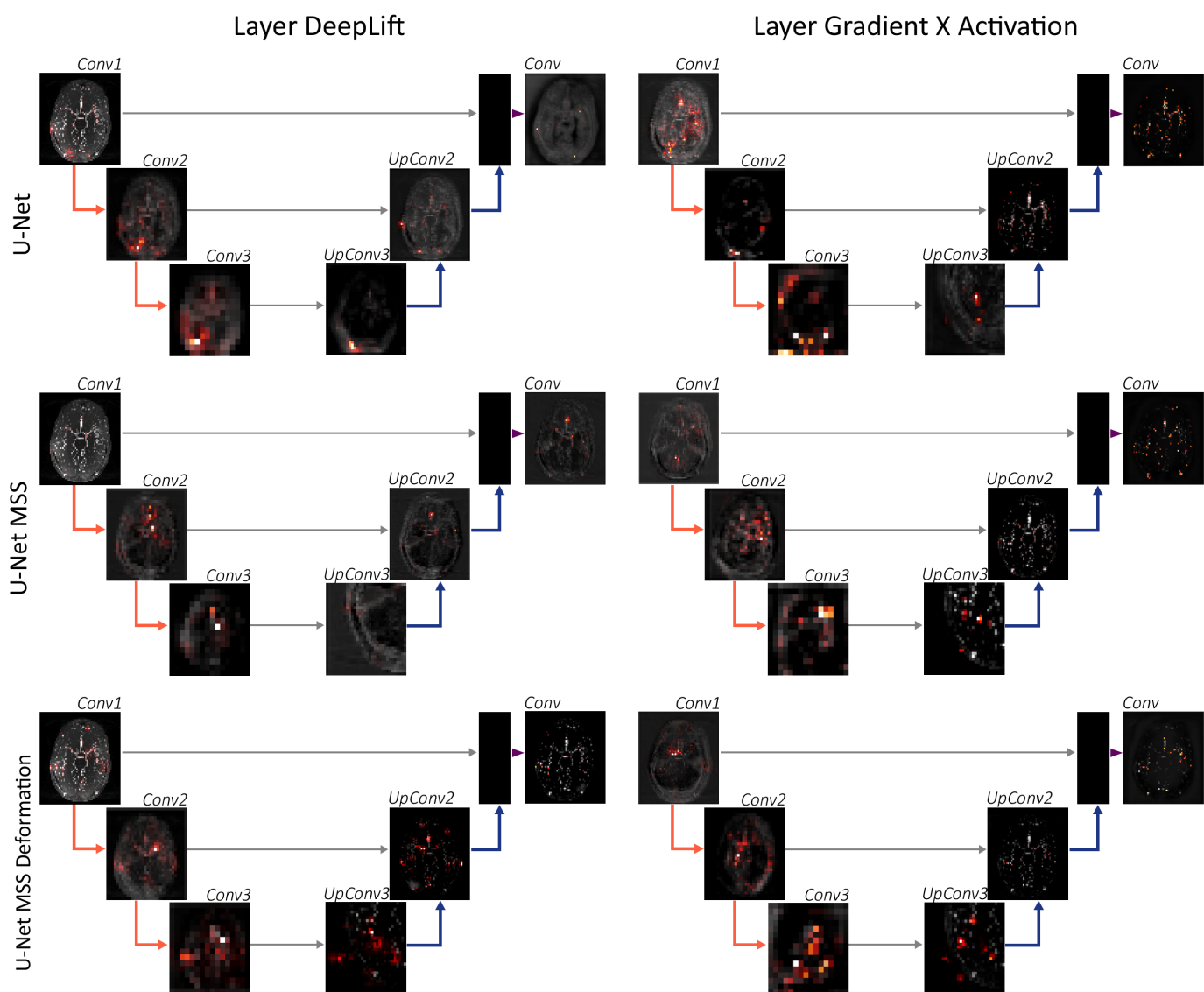
Figure 8: Layer-based interpretability methods: Layer Deeplift and Layer Gradient X Activation, shown using representative UNets
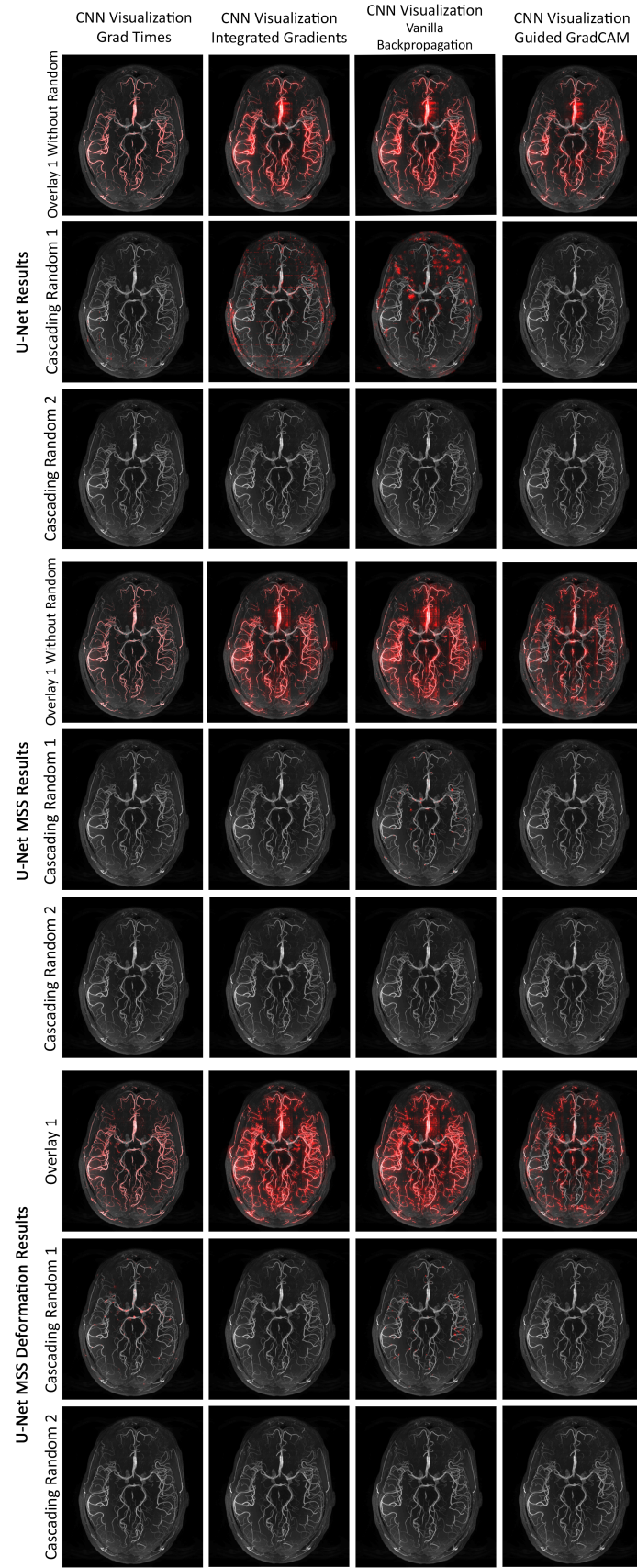
Figure 9: Example of cascading radomisation outputs of the three models: U-Net, U-Net MSS and U-Net MSS with deformation.

# References

Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., Kim, B., 2018. Sanity checks for saliency maps. arXiv preprint arXiv:1810.03292 .

Ancona, M., Ceolini, E., Öztireli, C., Gross, M., 2017. Towards better understanding of gradient-based attribution methods for deep neural networks. arXiv preprint arXiv:1711.06104 .

Arrieta, A.B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., García, S., Gil-López, S., Molina, D., Benjamins, R., et al., 2020. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. Information Fusion 58, 82–115.

Belle, V., Papantonis, I., 2020. Principles and practice of explainable machine learning. arXiv preprint arXiv:2009.11698 .

Breiman, L., 2001. Random forests. Machine learning 45, 5–32.

Castro, J., Gómez, D., Tejada, J., 2009. Polynomial calculation of the shapley value based on sampling. Computers & Operations Research 36, 1726–1730.

Chakraborty, S., Tomsett, R., Raghavendra, R., Harborne, D., Alzantot, M., Cerutti, F., Srivastava, M., Preece, A., Julier, S., Rao, R.M., et al., 2017. Interpretability of deep learning models: a survey of results, in: 2017 IEEE smartworld, ubiquitous intelligence & computing, advanced & trusted computed, scalable computing & communications, cloud & big data computing, Internet of people and smart city innovation (smartworld/SCALCOM/UIC/ATC/CBDcom/IOP/SCI), IEEE. pp. 1–6.

Chatterjee, S., Prabhu, K., Pattadkal, M., Bortsova, G., Sarasaen, C., Dubost, F., Mattern, H., de Bruijne, M., Speck, O., Nürnberger, A., 2020a. Ds6, deformation-aware semi-supervised learning: Application to small vessel segmentation with noisy training data. arXiv preprint arXiv:2006.10802 .

Chatterjee, S., Saad, F., Sarasaen, C., Ghosh, S., Khatun, R., Radeva, P., Rose, G., Stober, S., Speck, O., Nürnberger, A., 2020b. Exploration of interpretability techniques for deep covid-19 classification using chest x-ray images. arXiv preprint arXiv:2006.02570 .

Chattopadhay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N., 2018. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks, in: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV), IEEE. pp. 839–847.

Choo, J., Liu, S., 2018. Visual analytics for explainable deep learning. IEEE computer graphics and applications 38, 84–92.

Dhamdhere, K., Sundararajan, M., Yan, Q., 2018. How important is a neuron? arXiv preprint arXiv:1805.12233 .

Emmert-Streib, F., Yli-Harja, O., Dehmer, M., 2020. Explainable artificial intelligence and machine learning: A reality rooted perspective. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery 10, e1368.

Fisher, A., Rudin, C., Dominici, F., 2019. All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously. J. Mach. Learn. Res. 20, 1–81.

Fong, R., Patrick, M., Vedaldi, A., 2019. Understanding deep networks via extremal perturbations and smooth masks, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 2950–2958.

Guo, M., Zhang, Q., Liao, X., Zeng, D.D., 2020. An interpretable neural network model through piecewise linear approximation. arXiv:2001.07119 .

Leino, K., Sen, S., Datta, A., Fredrikson, M., Li, L., 2018. Influence-directed explanations for deep convolutional networks, in: 2018 IEEE International Test Conference (ITC), IEEE. pp. 1–8.

Liu, H., Brock, A., Simonyan, K., Le, Q.V., 2020. Evolving normalization-activation layers. arXiv preprint arXiv:2004.02967 .

Lundberg, S., Lee, S.I., 2017. A unified approach to interpreting model predictions. arXiv preprint arXiv:1705.07874 .

Mahendran, A., Vedaldi, A., 2014. Understanding deep image representations by inverting them. arXiv:1412.0035.

Marcinkevičs, R., Vogt, J.E., 2020. Interpretability and explainability: A machine learning zoo mini-tour. arXiv preprint arXiv:2012.01805 .

Micikevicius, P., Narang, S., Alben, J., Diamos, G., Elsen, E., Garcia, D., Ginsburg, B., Houston, M., Kuchaiev, O., Venkatesh, G., et al., 2017. Mixed precision training. arXiv preprint arXiv:1710.03740 .

Ozbulak, U., 2019. Pytorch cnn visualizations. https://github.com/utkuozbulak/pytorch-cnn-visualizations.

Petsiuk, V., Das, A., Saenko, K., 2018. Rise: Randomized input sampling for explanation of black-box models. arXiv:1806.07421.

Ribeiro, M.T., Singh, S., Guestrin, C., 2016. " 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, pp. 1135–1144.

Samek, W., Montavon, G., Lapuschkin, S., Anders, C.J., Müller, K.R., 2020. Toward interpretable machine learning: Transparent deep neural networks and beyond. arXiv preprint arXiv:2003.07631 .

Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D., 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization, in: Proceedings of the IEEE international conference on computer vision, pp. 618–626.

Shrikumar, A., Greenside, P., Kundaje, A., 2017a. Learning important features through propagating activation differences, in: International Conference on Machine Learning, PMLR. pp. 3145–3153.

Shrikumar, A., Greenside, P., Shcherbina, A., Kundaje, A., 2016. Not just a black box: Learning important features through propagating activation differences. arXiv preprint arXiv:1605.01713 .

Shrikumar, A., Greenside, P., Shcherbina, A., Kundaje, A., 2017b. Not just a black box: Learning important features through propagating activation differences. arXiv:1605.01713.

Simonyan, K., Vedaldi, A., Zisserman, A., 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. arXiv preprint arXiv:1312.6034 .

Smilkov, D., Thorat, N., Kim, B., Viégas, F., Wattenberg, M., 2017. Smoothgrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825 .

Springenberg, J.T., Dosovitskiy, A., Brox, T., Riedmiller, M., 2014. Striving for simplicity: The all convolutional net. arXiv preprint arXiv:1412.6806 .

Springenberg, J.T., Dosovitskiy, A., Brox, T., Riedmiller, M., 2015. Striving for simplicity: The all convolutional net. arXiv:1412.6806.

Sundararajan, M., Taly, A., Yan, Q., 2017. Axiomatic attribution for deep networks, in: International Conference on Machine Learning, PMLR. pp. 3319–3328.

Wang, H., Wang, Z., Du, M., Yang, F., Zhang, Z., Ding, S., Mardziel, P., Hu, X., 2020. Score-cam: Score-weighted visual explanations for convolutional neural networks, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 24–25.

Yeh, C.K., Hsieh, C.Y., Suggala, A.S., Inouye, D.I., Ravikumar, P., 2019. On the (in) fidelity and sensitivity for explanations. arXiv preprint arXiv:1901.09392 .

Zeiler, M.D., Fergus, R., 2014. Visualizing and understanding convolutional networks, in: European conference on computer vision, Springer. pp. 818–833.

Zhang, J., Bargal, S.A., Lin, Z., Brandt, J., Shen, X., Sclaroff, S., 2018. Top-down neural attention by excitation backprop. International Journal of Computer Vision 126, 1084–1102.