# The Explanation Necessity for Healthcare AI

Michail Mamalakis[1,2][*], Héloïse de Vareilles[1], Graham Murray[1], Pietro Lio[2], John Suckling[1]

[1]Department of Psychiatry, University of Cambridge, Hills Road, Cambridge, CB2 2QQ, Cambridgeshire, United Kingdom.
[2]Department of Computer Science and Technology, University of Cambridge, 15 JJ Thomson Ave, Cambridge, CB3 0FD, Cambridgeshire, United Kingdom.

*Corresponding author(s). E-mail(s): mm2703@cam.ac.uk;
Contributing authors: hd488@cam.ac.uk; gm285@cam.ac.uk;
pl219@cam.ac.uk; js369@cam.ac.uk;

**Abstract**

Explainability is often critical to the acceptable implementation of artificial intelligence (AI). Nowhere is this more important than healthcare where decision-making directly impacts patients and trust in AI systems is essential. This trust is often built on the explanations and interpretations the AI provides. Despite significant advancements in AI interpretability, there remains the need for clear guidelines on when and to what extent explanations are necessary in the medical context. We propose a novel categorization system with four distinct classes of explanation necessity, guiding the level of explanation required: patient or sample (local) level, cohort or dataset (global) level, or both levels. We introduce a mathematical formulation that distinguishes these categories and offers a practical framework for researchers to determine the necessity and depth of explanations required in medical AI applications. Three key factors are considered: the robustness of the evaluation protocol, the variability of expert observations, and the representation dimensionality of the application. In this perspective, we address the question: When does an AI medical application need to be explained, and at what level of detail?

# Main

Explainable artificial intelligence (XAI) has become a critical concern in digital devices and artificial intelligence (AI) affecting various fields like environmental science, climate studies, automotive technology, and medicine. In particular, the use of XAI in medical practice is crucial due to its significant role in the diagnosis of disease and the care patients receive. XAI plays a key role in fostering trust in algorithms as it aids in understanding risks and identifies therapeutic targets. Additionally, it offers insights into disease progression, treatment response, decision-making, and enables closed-loop control. To this end, a robust explanation of an AI framework can contribute to the design of safety parameters for regulatory consideration of potential therapies [1].

Although many studies have proposed methods to enhance the interpretability of AI systems, there remains a gap regarding when and at what level of explainability is truly required. Specifically, the literature lacks practical guidance on distinguishing between when the explanation necessity required is for predictions of individual patients or samples, the *local* level, and when it is required to decode the entire model for predictions of the whole cohort or dataset, the *global* level [1].

In this perspective, we address the question of explanation necessity in the field of AI with a focus on medical applications. We present a categorization that identifies different explanation needs across a range of AI tasks, providing clear algorithmic guidance on when to utilize none, local, global, or both types of explanations. We parameterize the classes of explanation necessity based on the robustness of the evaluation protocol, the degree of agreement among experts' observations, and the representational dimensionality of the application. We propose a mathematical representation of the different categories and discuss various frameworks for delivering explanations. Additionally, we explore different AI tasks and provide examples using our framework.

# Background of explainable methods

There has been a notable surge in recent publications concerning XAI and machine learning in the field of digital devices and especially in medical applications [2],[3],[4],[5]. XAI can be broadly classified into two methodological approaches: post-hoc and transparent.

Post-hoc methods are employed alongside AI techniques in a post-prediction setting to explain the (otherwise non-interpretable or 'black-box') AI predictions and unveil nonlinear mappings within complex datasets. They include both model-specific approaches that address particular nonlinear behaviors and model-agnostic approaches that explore data complexity [2, 6]. A widely used post-hoc technique is Local Interpretable Model-Agnostic Explanations (LIME), which clarifies the network's predictions by constructing simple, interpretable models that locally approximate the deep network; i.e., in the immediate vicinity of the prediction of interest [7]. LIME is an example of a perturbation technique that evaluates the sensitivity of an AI prediction to the input features by systematically altering sub-groups of the training data [7, 8]. In contrast, Gradient-weighted class activation mapping (Grad-CAM) [7] is a model-specific method that maps significant features in the imaging space using

the activations of the last convolutional layers of the AI architecture. This method is amongst the most common XAI methods used in medical imaging due to its ease of application and interpretation. Other post-hoc methods are the attribution explainability methods, such as Shapley Additive Explanations (SHAP; [9]) and Layer-wise Relevance Propagation (LRP; [10]), which identify important features for a given prediction by assigning relevance scores to the features of a given input. These attribution methods can be used for either local or global explanations, in contrast to majority of the post-hoc methods that are used mainly locally.

Transparent methods focus on AI models that exhibit inherent properties such as simulatability, decomposability, and transparency ('white-box') and are closely associated with linear techniques such as Bayesian classifiers, support vector machines, decision trees, and K nearest neighbor algorithms [6].

In many applications involving medical images, XAI faces the "curse of dimensionality" arising from the high-dimensional nature of the data. This challenge underscores the need for simpler models and variable selection techniques to deliver interpretability even if this sacrifices the accuracy and efficiency typically achieved by AI models like deep learning networks. Despite the trade-offs, XAI can foster trust in algorithms, aid in understanding risks and side effects, help identify therapeutic targets, offer insights into the progression of diseases and their response to treatments, support decision-making, enable closed-loop control, and contribute to designing safety parameters for regulated therapies [1].

Despite numerous studies proposing methods to improve the interpretability of AI systems, there remains a significant gap in understanding when and at what level the explainability is necessary. This absence can lead to confusion about the best approaches to apply in different scenarios. To address these issues, the AI research community is in need of a framework that outlines when and how to use local and global explainability techniques. This would in turn signpost the appropriate application of XAI in medicine, and other fields, ensuring that AI tools not only deliver accurate results, but are also transparent and trustworthy.

## The categories of explanation necessity

We propose a categorization of four distinct classes that define the need for explanation and indicate when to use local versus global explanations. This categorization is based on the robustness of the evaluation protocol, the degree of variability in expert opinions, and the representation dimensionality of the specific tasks. The classes are:

1. *Self-explainable applications* pertaining to tasks where interpretation of the internal mechanisms of AI are unnecessary due to very low variability in experts' opinion, a very robust evaluation protocol, a low representation dimensionality of AI application, and the direct understanding of AI predictions. No explanation is needed in such cases.
2. *Semi-explainable applications* that have a robust evaluation protocol with low variability in the opinion of experts and low to medium representation dimensionality of the AI application, requiring the provision of explanations within the AI learning process to ensure effective training. This category demands partial explanation

to confirm the accuracy of the AI's training process. There is a need for local explanations.

3. *Non-explainable* AI applications are characterized by the lack of a robust evaluation protocol, high variability in expert opinions, and medium to high representation dimensionality of the AI application. In such cases, there is a need for both local and global explanations.

4. *New-patterns discovery* AI applications are characterized by a lack of a robust evaluation protocol, significant variability in expert opinions, high representation dimensionality of the AI application, and a substantial gap in understanding the mechanisms and functions behind AI predictions. In such cases, there is a need for both local and global explanations, along with further evaluation to validate the new patterns captured through these explanations.

In the final section of this manuscript, we will present examples of these categories to better illustrate the terminology.

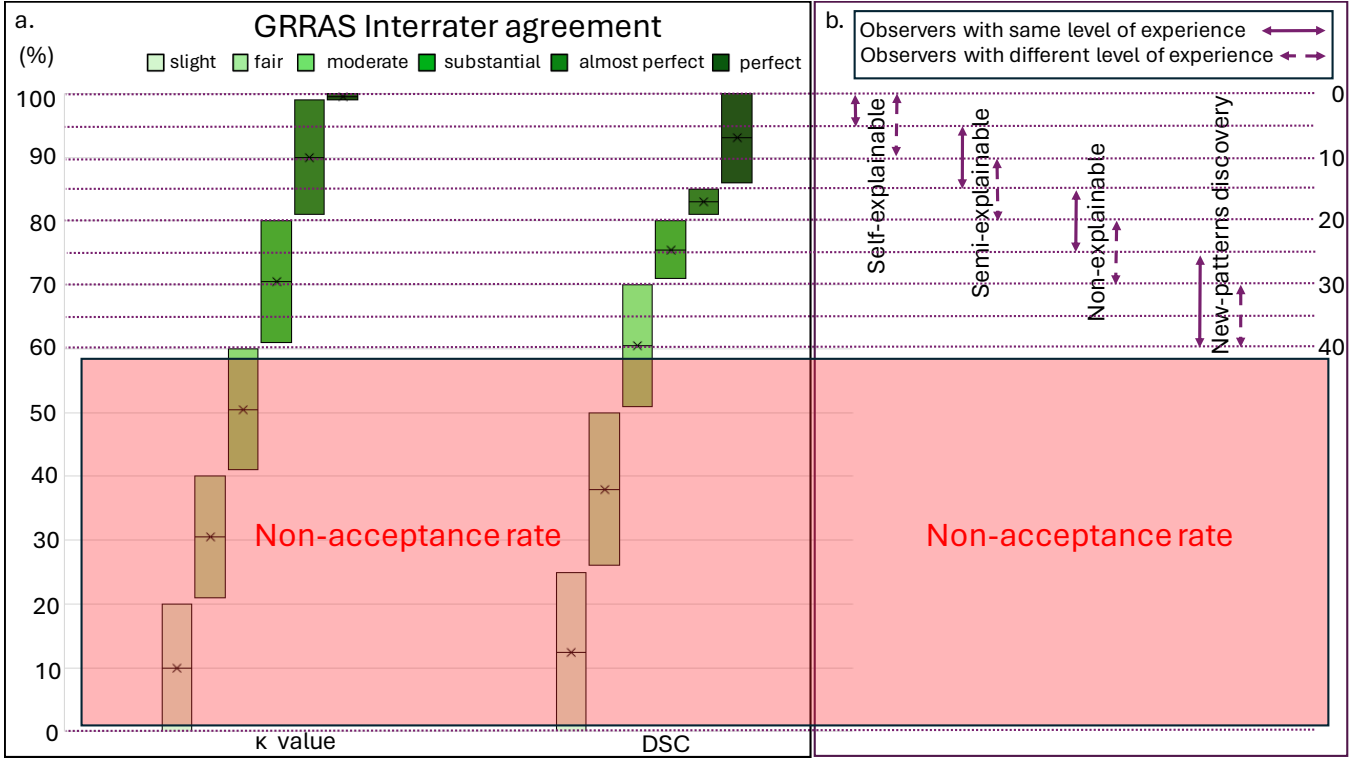# Parameters for classifying explanation necessity

To distinguish the different categories of XAI necessity, we used three key parameter: i) the variability in experts' observations (variability in observations for observers with same level of experience); ii) the robustness of the evaluation protocol (variability in observations for observers with different level of experience); and iii) the representational dimensional of the AI application.

## Variability of experts' observations and robustness of the evaluation protocol

To assess the variability in experts' observations, we propose capturing the diversity in annotations or answers provided by experts for each case. In this manuscript we adapt the Guidelines for Reporting Reliability and Agreement Studies (GRRAS) terminology. We primarily focus on "Agreement" which denotes the degree to which scores or observations are the same, and "Inter-rater (or inter-observer) agreement" which signifies the degree to which two or more observers achieve identical results under similar assessment conditions [11]. A common approach for scoring inter-observer agreement is the calculation of Kappa ($\kappa$) statistics and their variations, including Cohen's, Fleiss's, Light's, and weighted $\kappa$, as reported in thirty-one prior studies (0.39) [12]. Additionally, the utilization of the Landis and Koch interpretation of $\kappa$ is prevalent, found in forty-four prior studies (0.56) [12].

In medical applications, inter-observer variability (observers with same level of experience), a $\kappa$-value between 0.00 and 0.20 is classified as "slight," while values between 0.21 and 0.40 are deemed "fair." "Moderate" agreement falls between 0.41 and 0.60, while "substantial" agreement ranges from 0.61 to 0.80, and "almost perfect" agreement is between 0.81 and 1.00 [11, 12]. Generally, values of 0.60, 0.70, or 0.80 serve as the minimum standards for the labels for reliability coefficients, but higher values, like 0.90 or 0.95, are recommended for critical individual decisions [13–15]. To take as an example the segmentation of lesions or other pathologies from medical

**Fig. 1** Variability of experts' observations and robustness of the evaluation protocol for the classification of explanation necessity and the threshold regions.

**a.** Based on the Guidelines for Reporting Reliability and Agreement Studies (GRRAS) "Interrater (or inter-observer) agreement," scores of $\kappa$-value between 0.00 and 0.20 are classified as "slight," while values between 0.21 and 0.40 are deemed "fair." "Moderate" agreement falls between 0.41 and 0.60, while "substantial" agreement ranges from 0.61 to 0.80, and "almost perfect" agreement is between 0.81 and 0.99 "perfect" is between 0.99 and 1.00 [11]. For the agreement dice similarity coefficient (DSC), the values 0.00 to 0.25, are classified as "slight," while values between 0.26 and 0.50 are deemed "fair." "Moderate" agreement falls between 0.51 and 0.70, while "substantial" agreement ranges from 0.71 to 0.80, and "almost perfect" agreement is between 0.81 and 0.86 and "perfect" is between 0.86 and 1.00. The red transparent box highlights where the agreement rate is very low and not acceptable for medical applications. **b.** To classify the different categories of explanation necessity, we used as key parameters: (i) the robustness of the evaluation protocol; and (ii) the variability in experts' observations. The robustness of the evaluation protocol ('Observers with different level of experience') is measured as the variability across different observers with varying levels of experience (Inexperienced, Experienced, Expert). The threshold ratios are 0.00-0.10, 0.11-0.20, 0.21-0.30, and 0.31-0.40, for the self-explainable, semi-explainable, non-explainable, and new pattern discovery classes, respectively. The variability in experts' observations (observers with same level of experience) is measured directly with threshold ratios: 0.00-0.05, 0.06-0.15, 0.16-0.25, and 0.26-0.40, for the self-explainable, semi-explainable, non-explainable, and new pattern discovery classes, respectively.

images, the agreement DSC is utilized with thresholds: DSC $\geq 0.85$ considered "High Agreement," $0.85 >$ DSC $\geq 0.70$ as "Medium Agreement," $0.7 >$ DSC $\geq 0.5$ as "Low Agreement," and DSC $< 0.5$ as "Very Low Agreement" [16, 17] (see Fig. 1a.). The proposed number of experts is two to four with the same level of experience in the topic of interest.

The second key parameter of the explanation necessity is the robustness of the evaluation protocol. We suggest measuring the variability among observers with varying levels of experience (Inexperienced, Experienced, Expert). A robust evaluation protocol is defined by low variability in responses, indicating a clear, well-defined explainable protocol that can be adapted to different experience levels. To this end, we modify the proposed boundaries of the GRRAS Inter-rater agreement discussed above by $\pm$ 5% to account for the uncertainty arising from varying levels of experience. Typically, a suitable sample size for obtaining robust results consists of two to four observers selected across differing levels of experience (see Fig. 1b.).

Fig. 1b. presents the thresholds that categorize the explanation needs of an AI application based on the robustness of the evaluation protocol (variability in observer experience, 'purple dashed line') and the variability in expert opinions ('purple line'). The thresholds are set according to the level of uncertainty given by the probability $1 - \kappa$ (for classification, regression, etc.) or $1 - DSC$ (for segmentation, registration, overlapping regions etc.) value (see Fig. 1b.), that can be tolerated for a specific task, helping to classify the explanation requirements for different AI applications. These thresholds may vary depending on the application, as in some cases diversity of opinion amongst experts can be significant (such as in survival protocols or critical individual decisions [13–15]).

With these parameters, the proposed framework on explanation necessity is as follows: the self-explainable AI applications cover tasks where the protocols are established (0.00-0.10 observer with different level of experience; Fig. 1b.) and the variability in experts is low (0.00-0.05 observer with same level of experience; Fig. 1b;). The semi-explainable category is applied for AI applications where the protocols are established (0.11-0.20 observer with different level of experience; Fig. 1b.) and the variability in experts is low to middle (0.06-0.15 observer with same level of experience; Fig. 1b;). The non-explainable category is applied in AI applications where the protocols are not established (0.21-0.30 observer with different level of experience; Fig. 1b.) and the variability in experts is middle to high (0.16-0.25 observer with same level of experience; Fig. 1b;). Lastly, the new-patterns discovery category is applied in AI applications where the protocols are very unstable (0.31-0.40 observer with different level of experience; Fig. 1b.) and the variability in experts is high (0.26-0.40 observer with same level of experience; Fig. 1b;). For thresholds lower than 0.60 agreement we assume that the acceptance rate is invalid due to an error in the annotation process by the observers or a flaw in the protocol.

In studies involving high risks and critical individual decisions [13–15], it becomes imperative to adapt our proposed thresholds accordingly. In such cases, the acceptable 'Inter-rater agreement' values should ideally surpass the standard thresholds of 0.70, 0.80, or even 0.95, serving as the minimum benchmarks for reliability coefficients (the standard threshold for acceptance rate for 'Inter-rater agreement' typically

exceeds 0.60; see Fig. 1.). This adjustment ensures a heightened level of reliability and robustness in decision-making processes, crucial for maintaining safety and minimizing potential risks.

Aside from medical applications, the parameterization of explanation necessity in ecological studies utilizes similar categorical variables to gauge agreement between citizen scientists and experts. This approach enhances the reliability of environmental assessments [18, 19] following similar thresholds as proposed in GRRAS. The (driverless) automotive vision sector presents unique challenges in inter-observer variability assessment, primarily focused on navigation accuracy and safety [20]. While XAI plays a crucial role in optimizing AI models for safety purposes, the focus shifts from explainability to robustness and transparency to prevent potential traffic accidents [20]. This highlights the importance of maintaining reliability and robustness in automated, safety-critical systems, emphasizing the need for standardized evaluation protocols in diverse domains. However, further investigations are needed to identify similar thresholds (like Fig. 1) for other computer vision fields.

### Representational dimensionality of AI applications

Explanations in applications typically fall in two or three dimensions. The complexity of the representation dimensionality of the application can be further correlated with the level of explanation needed. Thus, low representation dimensionality applications in medical applications, common in semi-explainable and self-explainable applications (see Fig. 2,a; [21, 22]). However, more complex dimensions and applications, for example involving time-series data or multiple modalities inputs, are generally classified as non-explainable and new pattern discovery (Fig. 2,a). These tasks are more variable, involving complex protocols and intricate expert annotations, making them challenging to explain and evaluate (see Fig. 2,b). Furthermore, in these highly challenging fields, ordinary 2D explainability and validation of only local explanations will fall short and may even prove erroneous (Fig. 2,b). Thus, when assessing the need for explanations, it is important to consider the dimensionality representation of the AI application. This should be considered alongside the two key parameters and the mathematical formulation already provided.
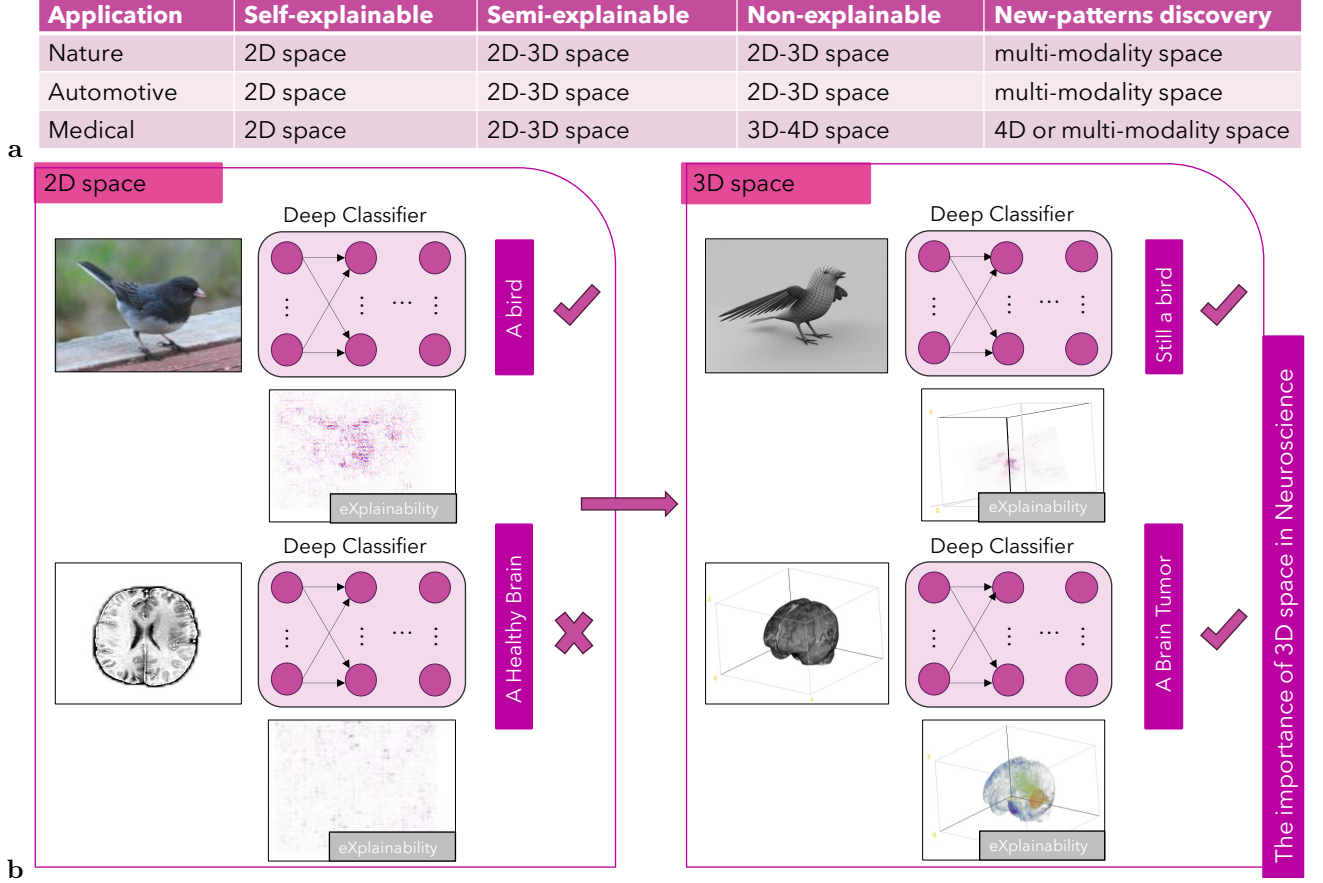
Remote sensing from satellites ([23, 24]), automotive vision, and other similar AI applications ([25, 26]) typically deal with 2D to 3D inputs (images and video) and can be categorized into self-explanatory, semi-explanatory, and non-explanatory (see Fig. 2,a). The concept of discovering new patterns is not well-defined in these fields. We suggest that multi-modality applications, where the complexity of AI increases significantly, might be suitable for this purpose.

## Mathematical formulation of the problem

For an input $\mathbf{x} \in R^d$, we define a deep learning model function $f(\mathbf{x}; \theta)$ with $f : R^d \to R^c$, where $c$ is the output dimension, $d$ is the input dimension and $\theta$ consists of the parameters of the model in a computer vision problem (e.g. in computer vision, for a classification task $c$ is the number of classes, and for a segmentation task $c \le d$). The inference of the model is denoted as $y = f(\mathbf{x}; \theta)$ where $\mathbf{y} \in R^c$ is the predicted

probability of the corresponding dimension $c$ (Fig. 3c). An explanation method $g$ from a post-hoc family of explanation methods $G$ takes as inputs the $f$, $\mathbf{x}$ and $\mathbf{y}$ and returns an explanation map $\mathbf{z}$.

**Fig. 2** Representational dimensionality of AI applications and the need of explanation.

| Application | Self-explainable | Semi-explainable | Non-explainable | New-patterns discovery |
|---|---|---|---|---|
| Nature | 2D space | 2D-3D space | 2D-3D space | multi-modality space |
| Automotive | 2D space | 2D-3D space | 2D-3D space | multi-modality space |
| Medical | 2D space | 2D-3D space | 3D-4D space | 4D or multi-modality space |



**a,** The table presents the representational dimensionality for each AI application and the corresponding level of explanation necessity, categorized as self-explainable, semi-explainable, non-explainable, and new-patterns discovery. **b,** The figure illustrates the importance of three-dimensional space in computer vision tasks in medical imaging compared to animals as a traditional computer vision classification task.

The $G(\mathbf{x}, \mathbf{y}, f)$ contains multiple explanation methods, denoted as $g^j$ where $j$ is the index of each post-hoc method. The explanation map of each input $\mathbf{x}$ can be given by $z = g(\mathbf{x}, \mathbf{y}, f)$ where $\mathbf{z} \in R^d$ and has the same dimension space as the input $\mathbf{x}$. For computer vision applications, we define $D : \mathbb{R}^d \times \mathbb{R}^c \times \mathbb{R}^d \mapsto \mathbb{R} \geq 0$ as a tuple of datasets $D = (\mathbf{x}_i, \mathbf{y}_i, \mathbf{z}_i)_{i=1}^n$, that denotes input-output-explanations triples, where $n$ is the total number of the validation examples. $D_x$ denotes all $\mathbf{x}_i$, $D_y$ denotes all $\mathbf{y}_i$

and $D_z$ denotes all $\mathbf{z}_i$ in $D$. In this study, we define as *"global explanation"* a map $g_z$ that represents all the local explanations for the whole validation dataset $D_z$. The global explanation can be viewed as a dimensionality reduction function $g_z = R(\mathbf{z})$ of the whole local explanations $D_z$ (for example, in a setting of a principal component analysis, $g_z$ it would correspond to the eigenvectors). We define four distinct applications of explainability needs: self-explainable, semi-explainable, non-explainable, and new-pattern learning (Fig. 3c).
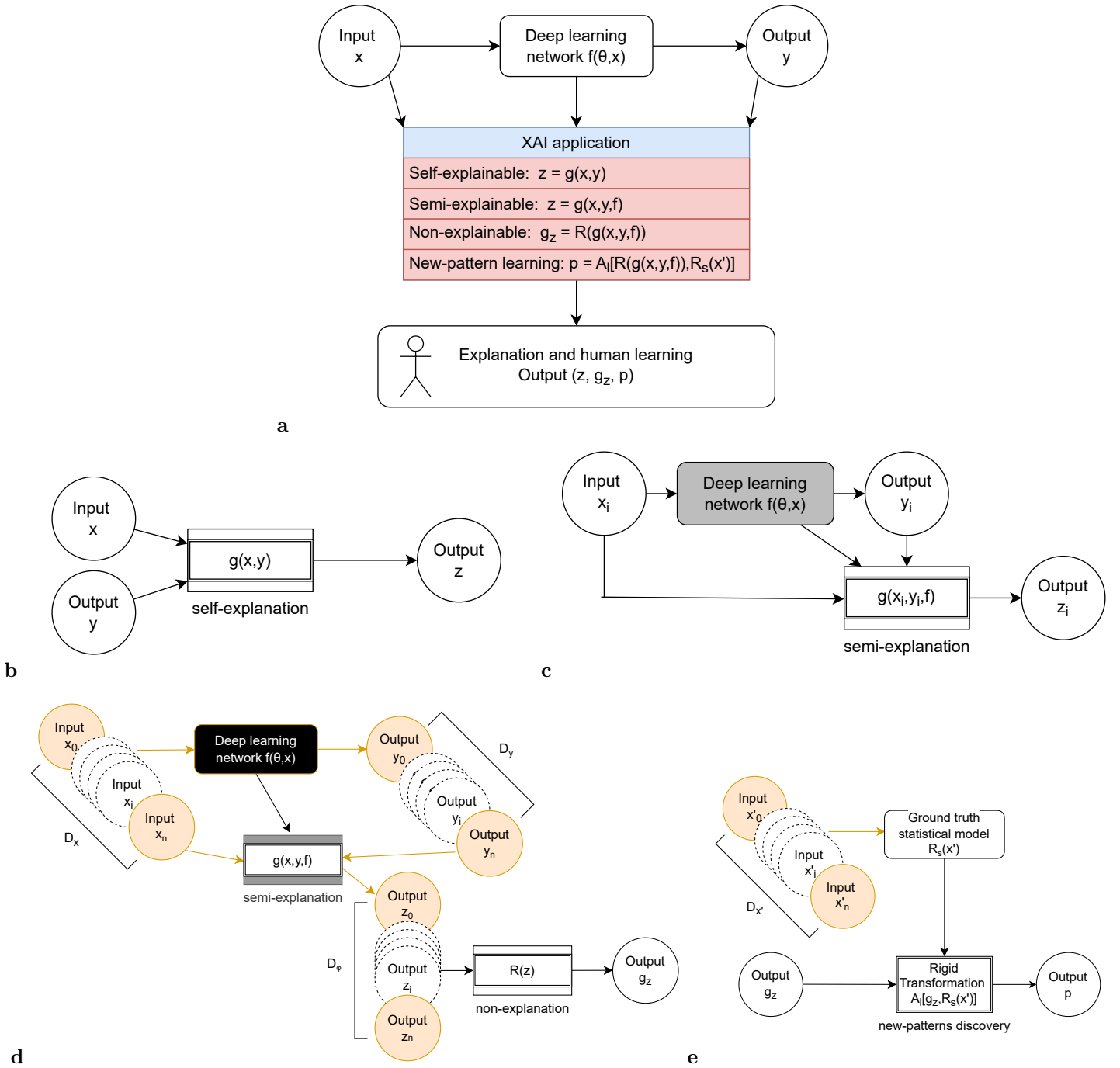
The mathematical representation of the self-explainable application involves an interpretable function $g(\mathbf{x}, \mathbf{y})$ that solely employs the inputs $\mathbf{x}$ and outputs $\mathbf{y}$ of the deep learning model $f$. Given an established, clear evaluation protocol, low variability in experts' observations, low representation dimensionality of the AI application, and sufficient correlation between the inputs and outputs of the deep learning network, there is no necessity for explaining the hidden parameters of the network. Therefore, this application in literature is named as a *'white box application'* (Fig. 3a).

The mathematical description of the semi-explainable application comprises a set of explanation methods $g(\mathbf{x}, \mathbf{y}, f)$ that utilize both the inputs $\mathbf{x}$ and outputs $\mathbf{y}$, along with the hidden parameters of the deep learning model $f(\mathbf{x}; \theta)$. Due to the established evaluation protocol, low to medium variability in experts' observations, low to medium representation dimensionality of the AI application, and insufficient correlation between the inputs and outputs of the deep learning network, a sub-group $S_D = (\mathbf{x}_i, \mathbf{y}_i, \mathbf{z}_i)_{i=1}^{o}$ of local explanations $\mathbf{z}$ of the hidden parameters are needed (where $o$ represents a fixed number of samples from the validation dataset D, $S_D \subset D$ and $o < n$). Consequently, this application in literature is named as a *'grey box application'* (Fig. 3b).

The mathematical representation of non-explainable applications involves a set of explanation methods $g(\mathbf{x}, \mathbf{y}, f)$ that relies on the inputs $\mathbf{x}$ and outputs $\mathbf{y}$ in addition to the hidden parameters of the deep learning model $f(\mathbf{x}; \theta)$.

With a non established and clear evaluation protocol, medium to high variability in experts' observations and medium to high representation dimensionality of the AI application, the correlation between the inputs and outputs of the deep learning network is inadequate. Thus, there is a requirement for local and additionally global explanations, $g_z$, by using the whole validation dataset $D = (\mathbf{x}_i, \mathbf{y}_i, \mathbf{z}_i)_{i=1}^{n}$ (where $n$ encompasses the total number of samples from the validation dataset $D$). This application in the literature is named a *'black box application'* (Fig. 3d). In the new-patterns discovery application, a collection of potentially significant markers can be computed for the initial AI task by aligning the global explanations $g_z$ with a ground truth statistical model $R_S(\mathbf{x}')$. This statistical model is generated by inputs $\mathbf{x}'$ of a super dataset $D_{x'}$ of the validation dataset $D_x$ ($D_x \subseteq D_{x'}$). The ground truth statistical model tries to capture the generalized features of interest and the generalized shape related with the AI task. The model's inference is given by $r_{\mathbf{x}'} = R_S(\mathbf{x}')$ where $r_{x'}$ is the output of the ground truth statistical model that best describes the $D_{x'}$ dataset. The alignment can be done by a rigid transformation function $A_l(g_z, r_{\mathbf{x}'})$ of the $r_{\mathbf{x}'}$ and the $g_z$ (Fig. 3e). The new-patterns discovery category is used in cases with non-established and clear evaluation protocols, high variability in experts' observations and high representation dimensionality of the AI application.

**Fig. 3** A mathematical formulation of the explainability need across deep learning applications



**a,** The mathematical representation presents the overall XAI framework for a specific method $g$ across different explainability applications. **b,** Self-explanations utilize methods $g(x, y)$ focusing solely on the inputs $x$ and outputs $y$ of the deep learning model $f$. With a clear evaluation protocol and strong input-output correlation, explaining the hidden parameters becomes unnecessary. This application categorized as a 'white box application'. **c,** Semi-explanations employ methods $g(x, y, f)$ integrating inputs $x$, outputs $y'$, and the hidden parameters of $f(\theta, x)$. Due to varied annotation and insufficient input-output correlation, local explanations for specific samples in the validation dataset $D$ become necessary. This application is categorized as a 'grey box application'. **d,** Non-explanations involve methods $g(x, y, f)$ relying on inputs $x$, outputs $y$, and the hidden parameters of $f(\theta, x)$. Due to significant variance in annotations and discrepancies in ground truth extractions, a necessity arises for global explanations $g_z$ using the entire dataset $D = (x_i, y_i, z_i)_{i=1}^{n}$. This application is categorized as a 'black box application'. **e,** In the new-patterns discovery application, crucial markers for the initial classification task are derived by aligning global explanations $g_z$ with a ground truth statistical model $R(x')$ generated from a supergroup $D_{x'}$ describing the classification application's overall shape. A rigid transformation function $A_l(g_z, r_{x'})$ is used to align the output of the statistical shape model with the global explanations.

10

# The proposed framework

In this study, we propose a framework to classify the explanation necessity of AI medical applications. The framework consists of two main flows: one for determining inter-observer variability (First Flow; see Fig. 2) and one for representing dimensionality (Second Flow; see Fig. 2).

Firstly, the user needs to compute the average values of inter-observer variability from a group of observers with the same and different experience levels to justify the variability of experts' observations ('Same level of experience'; Fig. 2), and the robustness of the evaluation protocol ('Different level of experience'; Fig. 2). Using the average thresholds ('Calculate Average Threshhold'; Fig. 2) of these two inter-observer variabilities, the user identifies the two 'Initial Explanation necessity classification' (see Fig. 2; Table 1. ; Table 2.). If the two categories are not the same, an adjudicating expert identifies which class is more appropriate for the case and they continue to the next step.

The second flow of the framework involves the representation dimensionality of the application ('Representation dimentionality'; Fig. 2). Lastly, the results of the two flows are passed through an 'XAI Need Decision' statement where if the results are the same the final class of XAI need is determined ('Category decision'). If the explanation class is not the same, an adjudicating expert identifies which class is more appropriate for the application (see Fig. 2)

# Examples and applications

By utilizing our proposed framework, anyone can determine when and at the explanation necessity of an application. In this section, we present some examples and applications of our framework in a variety of medical imaging applications.

Some applications in medicine require minimal understanding of the inner mechanisms of AI due to low variability in the evaluation protocols (0.00-0.10; Fig. 1), low variability between experts' observations (0.00-0.05; Fig. 1), the two-dimensionality representation of the application, and the straightforward nature of AI predictions. Examples include human organ segmentation from abdominal computed tomography (CT) and registration of multi-modal images from the same individual [5, 27]. These types of applications might benefit from XAI methods for optimization purposes rather than for enhancing trust. Consequently, the performance of AI models can be reliably assessed without the need for additional explanation.

Other applications involve greater variability in evaluation protocols, requiring local explanations to ensure proper training, such as in classification tasks with established disease evaluation protocols (0.05-0.15; Fig. 1) and two-dimensional representation of the applications [21].

With an aging global population, neurodegenerative diseases are likely to become increasingly prevalent. Binary AI classification from MRI scans of the brain for diagnosing Alzheimer's disease or healthy aging is a task with low inter-observer variability (0.05-0.15; Fig. 1) as brain atrophy is clearly visible when present, and does not require multi-modality datasets for high performance (low representation dimensionality). This kind of application is self-explainable or semi-explainable based on the
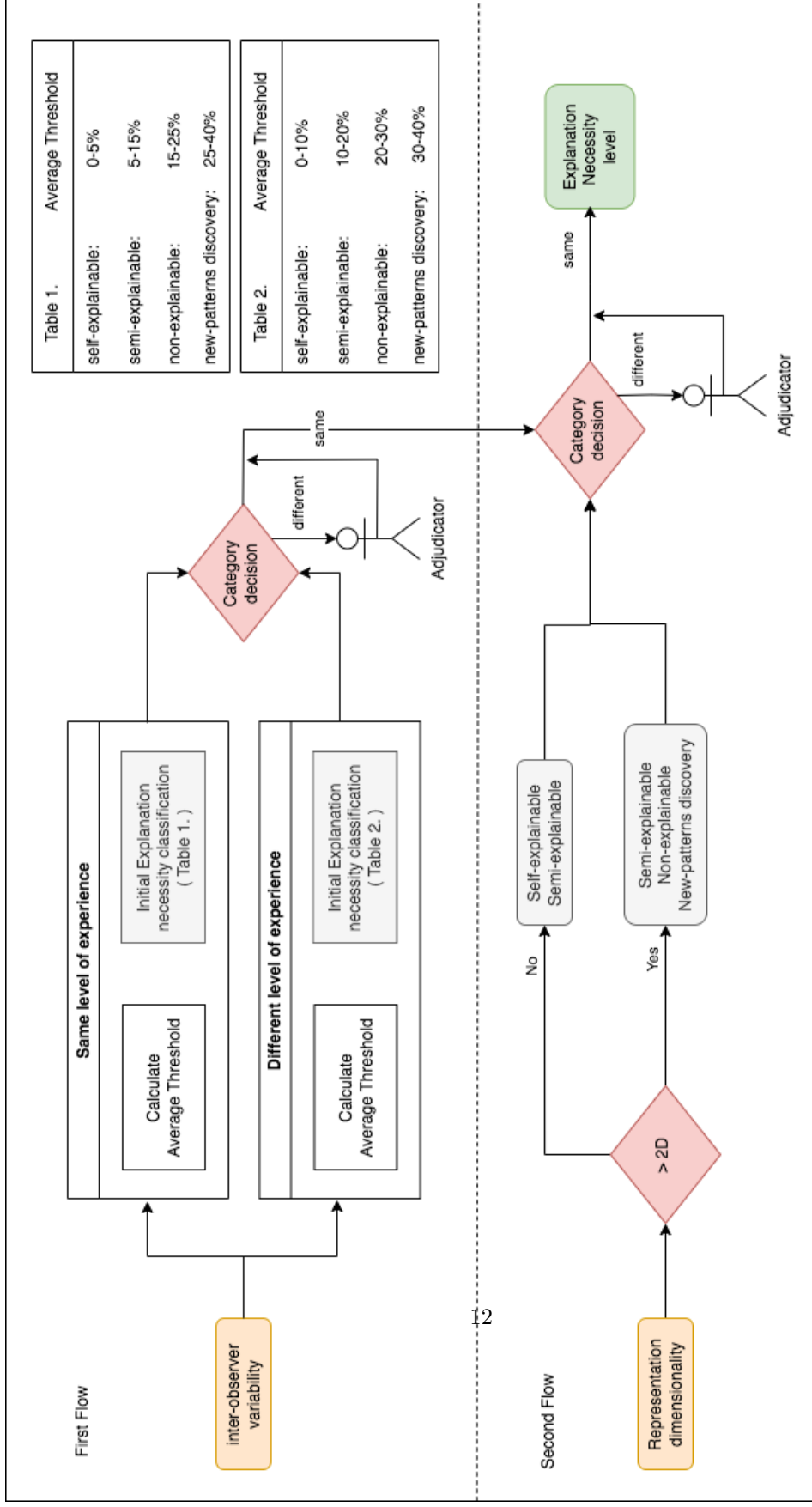
11

**Fig. 4** Proposed framework for explanation necessity. The framework consists of two main flows: one for assessing inter-observer variability and the other for representing dimensionality. Initially, users compute the average values of inter-observer variability of observers with 'Same level of experience' and with 'Different level of experience'. The thresholds (Fig. 1) are then applied to identify the two 'Initial Explanation necessity classification' from Table 1. and Table 2. If these categories differ ('different'), an adjudicating expert determines the most suitable class for the case. The second flow focuses on the representation dimensionality of the application, outlined in Fig. 2. Finally, the results undergo an 'Category decision' statement. If they align ('same'), the final XAI need class is determined ('Explanation Necessity Level'). Otherwise, an adjudicating expert identifies the most appropriate class for the application.

proposed framework (see Fig. 2). Detecting early stages many years before diagnosis is significantly more challenging [28].

Even among experienced professionals, knowledge gaps can persist, where AI has the potential to offer insights and stabilize the validity and key aspects of the protocols (0.25-0.40; Fig. 1; [29]). This is particularly true for classification tasks where disease evaluation protocols are not yet firmly established (new-patterns discovery) [30].

Ovarian cancer is one of the most common cancers in women, with an uncertain prognosis (0.20-0.40; Fig. 1), and is difficult to detect at an early stage even with multi-modal imaging (MRI, ultrasound, and computed tomography; [31]). This AI application is categorized as non-explainable or even as a new-patterns discovery applications.

Sepsis is a life-threatening acute immune response to infection that causes organ damage. Early-stage diagnosis, when treatments are effective, is complex. Prognosis is particularly poor where access to healthcare is limited. Along with clinical and laboratory assessments, chest X-rays [32] and whole-body computed tomography [33] are helpful in diagnosis and disease management. Thus, in the majority of AI applications in this medical topic, inter-observer variability among experts is high (0.25-0.40; Fig. 1), the evaluation protocol robustness is low (0.30-0.40; Fig. 1), and the representation dimensionality of the application needs to be multi-modality. These applications are new-patterns discovery applications.

The proposed framework can be applied in various computer vision fields such as natural or automotive. For the sake of generalization, we provide examples of applications in these fields. However, accurate determination of protocols and thresholds (Fig. 1) are necessary in each field. In natural computer vision applications, like animal classification in images and climate regression, typically local explanation (semi-explainable) is required. This need arises because expert knowledge varies minimally (0.05-0.10; Fig. 1), the robustness of the evaluation protocols is straightforward [23] and the dimensionality representation of the application is usually two-dimensions. In contrast, automotive computer vision, which encompasses tasks like vehicle and object classification and semantic segmentation for self-driving cars, generally does not require explanation. This is because these tasks are relatively straightforward with clear evaluation metrics, minimal expert involvement and mainly two- to three-dimensionality representation of the application (self-explainable); [25, 34].

# Outlook

Explainability, together with accuracy and consistency, are important aspects of AI systems to gain the trust of scientists and healthcare professionals, even without them fully understanding how the algorithms work. While the use of XAI is generally important, it becomes crucial in clinical contexts as decisions taken relying on AI-driven tools may directly impact a patient's health. While many studies focus on enhancing the interpretability of AI systems, we highlight the lack of user-directed recommendations on when to utilize explainability techniques and to what extent (global, local, or both).

In this perspective, we address this important gap in the literature by categorizing the necessity for AI explanations into four distinct groups: self-explainable, semi-explainable, non-explainable applications, and new-patterns discovery. These classifications are informed by the variability of experts' observations, the stability of the evaluation protocol and the representation dimensionality of the application.

By accessing the average variability of experts' observations across different experience levels and comparing them with the average variability of observations from experts of the same level, we can establish an initial categorization. If the clinical application is identified as high-risk, it becomes imperative to adjust the proposed thresholds to align with the risk level of the application. For instance, rather than the less than 0.60 non-acceptance rate initially proposed, the application may necessitate higher thresholds, such as 0.80, 0.90, or even 0.95.

We also take into account the dimensionality representation of the AI application to revise the explanation necessity category according to our recommendations. Consistent with the proposed framework, we present a mathematical formulation of these classes to cover a broad range of explanation requirements. This mathematical formulation and suggested framework can be employed to provide the essential explanations required for the AI application.

We have developed a comprehensive framework that researchers can readily customize for their AI applications. Our framework assists in determining the most suitable explanation necessity for their specific medical application. This enables them to furnish the requisite explanations, supporting the delivery of a transparent, safe, and trustworthy AI framework while also strengthening safety parameters for regulated therapies.

## Code and dataset availability

N/A

**Supplementary information.** The manuscript does not contains supplementary material.

*

*

Author information

\*

Authors and Affiliations

**Department of Psychiatry, University of Cambridge, Cambridge, UK.**

Michail Mamalakis, Héloïse de Vareilles, Graham Murray, John Suckling

**Department of Computer Science and Technology, Computer Laboratory, University of Cambridge, Cambridge, UK.**

Michail Mamalakis, Pietro Lio

\*

Contributions M.M, conceived and supervise the study. M.M, collect and analyzed the data and results. M.M. contributed to pulling deep learning and XAI methods and conducted chart reviews. M.M contributed to the framework design and evaluation protocol. M.M. J.S. and P.L. was in charge of overall direction and planning. All authors contributed to the interpretation of the results. M.M., and J.S. visualized the study and extracted figures, and M.M and H.V. and J.S. drafted the manuscript, which was reviewed, revised and approved by all authors.

\*

Corresponding author

Correspondence to: Michail Mamalakis

\*

Ethics declarations N/A

## Competing interests

The authors declare no competing interests.

# References

[1] Longo, L. *et al.* Explainable artificial intelligence (xai) 2.0: A manifesto of open challenges and interdisciplinary research directions. *Information Fusion* **106**, 102301 (2024). URL https://www.sciencedirect.com/science/article/pii/S1566253524000794.

[2] Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J. & Müller, K.-R. Explaining deep neural networks and beyond: A review of methods and applications. *Proceedings of the IEEE* **109**, 247–278 (2021).

[3] van der Velden, B. H., Kuijf, H. J., Gilhuijs, K. G. & Viergever, M. A. Explainable artificial intelligence (xai) in deep learning-based medical image analysis. *Medical Image Analysis* **79**, 102470 (2022). URL https://www.sciencedirect.com/science/article/pii/S1361841522001177.

[4] Quellec, G. *et al.* Explain: Explanatory artificial intelligence for diabetic retinopathy diagnosis. *Medical Image Analysis* **72**, 102118 (2021). URL https://www.sciencedirect.com/science/article/pii/S136184152100164X.

[5] Mamalakis, M. *et al.* Ma-socratis: An automatic pipeline for robust segmentation of the left ventricle and scar. *Computerized Medical Imaging and Graphics* **93**, 101982 (2021). URL https://www.sciencedirect.com/science/article/pii/S0895611121001312.

[6] van der Velden, B. H. M. Explainable ai: current status and future potential. *European Radiology* (2023). URL https://doi.org/10.1007/s00330-023-10121-4.

[7] Singh, A., Sengupta, S. & Lakshminarayanan, V. Explainable deep learning models in medical image analysis (2020). 2005.13799.

[8] Tjoa, E. & Guan, C. A survey on explainable artificial intelligence (xai): Toward medical xai. *IEEE Transactions on Neural Networks and Learning Systems* 1–21 (2020).

[9] Lundberg, S. & Lee, S.-I. A unified approach to interpreting model predictions (2017). URL https://arxiv.org/abs/1705.07874.

[10] Bach, S. *et al.* On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE* **10**, 1–46 (2015). URL https://doi.org/10.1371/journal.pone.0130140.

[11] Kottner, J. *et al.* Guidelines for reporting reliability and agreement studies (grras) were proposed. *International Journal of Nursing Studies* **48**, 661–671 (2011). URL https://www.sciencedirect.com/science/article/pii/S0020748911000368.

[12] Quinn, L. *et al.* Interobserver variability studies in diagnostic imaging: a methodological systematic review. *British Journal of Radiology* **96**, 20220972 (2023). URL https://doi.org/10.1259/bjr.20220972.

[13] Kottner, J. & Dassen, T. Interpreting interrater reliability coefficients of the braden scale: A discussion paper. *International Journal of Nursing Studies* **45**, 1238–1246 (2008). URL https://www.sciencedirect.com/science/article/pii/S002074890700199X.

[14] Polit-O'Hara, D. & Beck, C. T. *Nursing research : generating and assessing evidence for nursing practice / Denise F. Polit, Cheryl Tatano Beck.* Eighth edition. edn (Wolters Kluwer, Philadelphia, 2008).

[15] Thorndike, R. M. Book review : Psychometric theory (3rd ed.) by jum nunnally and ira bernstein new york: Mcgraw-hill, 1994, xxiv + 752 pp. *Applied Psychological Measurement* **19**, 303–305 (1995). URL https://doi.org/10.1177/014662169501900308.

[16] Wong, J. *et al.* Effects of interobserver and interdisciplinary segmentation variabilities on ct-based radiomics for pancreatic cancer. *Scientific Reports* **11**, 16328 (2021). URL https://doi.org/10.1038/s41598-021-95152-x.

[17] Kothari, G. *et al.* The impact of inter-observer variation in delineation on robustness of radiomics features in non-small cell lung cancer. *Scientific Reports* **12**, 12822 (2022). URL https://doi.org/10.1038/s41598-022-16520-9.

[18] Earp, H. *et al.* Do you see what i see? quantifying inter-observer variability in an intertidal marine citizen science experiment. *Citizen Science: Theory and Practice* (2022).

[19] Landis, J. R. & Koch, G. G. The measurement of observer agreement for categorical data. *Biometrics* **33**, 159–174 (1977). URL http://www.jstor.org/stable/2529310.

[20] Kuznietsov, A., Gyevnar, B., Wang, C., Peters, S. & Albrecht, S. V. Explainable ai for safe and trustworthy autonomous driving: A systematic review (2024). 2402.10086.

[21] Mamalakis, M., Macfarlane, S. C., Notley, S. V., Gad, A. K. B. & Panoutsos, G. A novel framework employing deep multi-attention channels network for the autonomous detection of metastasizing cells through fluorescence microscopy (2023). 2309.00911.

[22] Ma, J. *et al.* Segment anything in medical images. *Nature Communications* **15**, 654 (2024). URL https://doi.org/10.1038/s41467-024-44824-z.

[23] Mamalakis, A., Ebert-Uphoff, I. & Barnes, E. A. Neural network attribution methods for problems in geoscience: A novel synthetic benchmark dataset. *Environmental Data Science* **1**, e8 (2022).

[24] Bi, K. *et al.* Accurate medium-range global weather forecasting with 3d neural networks. *Nature* **619**, 533–538 (2023). URL https://doi.org/10.1038/s41586-023-06185-3.

[25] Ghosh, S. *et al.* Segfast-v2: Semantic image segmentation with less parameters in deep learning for autonomous driving. *International Journal of Machine Learning and Cybernetics* **10**, 3145–3154 (2019). URL https://doi.org/10.1007/s13042-019-01005-5.

[26] Assidiq, A., Khalifa, O. O., Islam, M. R. & Khan, S. *Real time lane detection for autonomous vehicles*, 82–88 (2008).

[27] Mamalakis, M. *et al.* Artificial intelligence framework with traditional computer vision and deep learning approaches for optimal automatic segmentation of left ventricle with scar. *Artificial Intelligence in Medicine* **143**, 102610 (2023). URL

https://www.sciencedirect.com/science/article/pii/S0933365723001240.

[28] Bron, E. E. *et al.* Ten years of image analysis and machine learning competitions in dementia. *NeuroImage* **253**, 119083 (2022). URL https://www.sciencedirect.com/science/article/pii/S1053811922002129.

[29] Mitchell, S. C. *et al.* Paracingulate sulcus measurement protocol v2 (2023). URL https://www.repository.cam.ac.uk/handle/1810/358381.

[30] Mamalakis, M. *et al.* A transparent artificial intelligence framework to assess lung disease in pulmonary hypertension. *Scientific Reports* **13**, 3812 (2023). URL https://doi.org/10.1038/s41598-023-30503-4.

[31] Mathieu, K. B., Bedi, D. G., Thrower, S. L., Qayyum, A. & Bast Jr, R. C. Screening for ovarian cancer: imaging challenges and opportunities for improvement. *Ultrasound in Obstetrics & Gynecology* **51**, 293–303 (2018). URL https://obgyn.onlinelibrary.wiley.com/doi/abs/10.1002/uog.17557.

[32] Carrillo-Bayona, J. A. & Arias-Alvarez, L. *Diagnostic Imaging in Sepsis of Pulmonary Origin*, 51–65 (Springer New York, New York, NY, 2018). URL https://doi.org/10.1007/978-1-4939-7334-7_5.

[33] Pohlan, J. *et al.* Relevance of ct for the detection of septic foci: diagnostic performance in a retrospective cohort of medical intensive care patients. *Clinical Radiology* **77**, 203–209 (2022). URL https://www.sciencedirect.com/science/article/pii/S0009926021005262.

[34] Yu, X. *et al.* The robust semantic segmentation uncv2023 challenge results (2023). 2309.15478.