ELSEVIER

Contents lists available at ScienceDirect

Artificial Intelligence In Medicine

journal homepage: www.elsevier.com/locate/artmed



An explainable three dimensional framework to uncover learning patterns: A unified look in variable sulci recognition

Michail Mamalakis ^{a,b}, Héloïse de Vareilles ^a, Atheer Al-Manea ^a, Samantha C. Mitchell ^c, Ingrid Agartz ^f, Lynn Egeland Mørch-Johnsen ^{d,e}, Jane Garrison ^c, Jon Simons ^c, Pietro Lio ^b, John Suckling ^{a,1}, Graham K. Murray ^{a,1}

- ^a Department of Psychiatry, University of Cambridge, Cambridge, UK
- b Department of Computer Science and Technology, Computer Laboratory, University of Cambridge, Cambridge, UK
- ^c Department of Psychology, University of Cambridge, Cambridge, UK
- d Norment, Division of Mental Health and Addiction, Oslo University Hospital, Institute of Clinical Medicine, University of Oslo, Oslo, Norway
- e Department of Psychiatry and Department of Clinical Research, Østfold Hospital, Grålum, Norway
- f Department of Psychiatric Research, Diakonhjemmet Hospital, Oslo, Norway

ARTICLE INFO

Keywords: XAI Sulcal pattern Paracingulate Deep learning Brain classification Brain pattern

ABSTRACT

The significant features identified in a representative subset of the dataset during the learning process of an artificial intelligence model are referred to as a 'global' explanation. Three-dimensional (3D) global explanations are crucial in neuroimaging, where a complex representational space demands more than basic two-dimensional interpretations. However, current studies in the literature often lack the accuracy, comprehensibility, and 3D global explanations needed in neuroimaging and beyond. To address this gap, we developed an explainable artificial intelligence (XAI) 3D-Framework capable of providing accurate, lowcomplexity global explanations. We evaluated the framework using various 3D deep learning models trained on a well-annotated cohort of 596 structural MRIs. The binary classification task focused on detecting the presence or absence of the paracingulate sulcus (PCS), a highly variable brain structure associated with psychosis. Our framework integrates statistical features (Shape) and XAI methods (GradCam and SHAP) with dimensionality reduction, ensuring that explanations reflect both model learning and cohort-specific variability. By combining Shape, GradCam, and SHAP, our framework reduces inter-method variability, enhancing the faithfulness and reliability of global explanations. These robust explanations facilitated the identification of critical sub-regions, including the posterior temporal and internal parietal regions, as well as the cingulate region and thalamus, suggesting potential genetic or developmental influences. For the first time, this XAI 3D-Framework leverages global explanations to uncover the broader developmental context of specific cortical features. This approach advances the fields of deep learning and neuroscience by offering insights into normative brain development and atypical trajectories linked to mental illness, paving the way for more reliable and interpretable AI applications in neuroimaging.

1. Introduction

In both medical imaging and neuroscience, explainability holds paramount importance. Recently, the study by Mamalakis et al. [1] introduced the necessity of explanations in artificial intelligence (AI) healthcare applications, categorizing them into four types: self-explainable, semi-explainable, non-explainable applications, and new-pattern discovery. This categorization is based on the variability of expert opinions, the stability of the evaluation protocol, and the dimensionality of the problem.

In neuroscience AI applications, there commonly exists significant variability in evaluation protocols and decision-making among experts. The inherent high uncertainty creates a greater need for explainability (as the framework falls in the non-explainable category; [1]). To this end, applications usually need both "local" methods, which provide explanations for each AI prediction separately, and "global" methods, which derive explanations for the decision-making of the AI across the entire dataset. Such variability underscores the importance of thorough

^{*} Correspondence to: University of Cambridge, Department of Psychiatry, Herchel Smith Building, Forvie Site, Robinson Way, Cambridge, CB2 0SZ, UK. E-mail address: mm2703@cam.ac.uk (M. Mamalakis).

¹ Equal contribution

investigation and validation of AI model behavior. Even among experienced professionals, knowledge gaps can persist, and this is where AI has the potential to offer insights and stabilize the validity of key aspects of the evaluation protocols [2]. This is particularly true for classification tasks where the key features of a disease are not yet firmly established (matching the new-patterns discovery category; [1]).

The folding of the human cortical surface occurs during the perinatal period and remains constant for an individual for the rest of their life, much like a fingerprint, preserving early neurodevelopmental information [3]. Broadly, human brains share many features of cortical folding (sulci), however strong inter-individual variability creates inherent divergence in experts' opinions when it comes to labeling the more variable sulci, impeding the effort to automatically label sulci in the regions showing most variability.

While, to date, automated methods excel in the detection of most sulci, the variable shape and presence/absence of some sulci present a more complex computational hurdle [4]. Successful automation through generalized, unbiased annotation would greatly aid studies focusing on brain folding variations, which are proxies for a critical developmental period with information that may relate to cognitive, behavioral, and developmental outcomes, along with psychiatric and neurological disorders. Brain folding is linked to brain function, and specific folding patterns correlate with susceptibility to neurological issues [5]. In particular, the morphology of the paracingulate sulcus (PCS), a highly variable sulcus, is associated with cognitive performance and hallucinations in schizophrenia [6-9]. To this end, developing networks that can identify the presence or absence of PCS (or multiple PCS elements) in brain magnetic resonance imaging (MRI) and creating frameworks that can provide both local and global explanations would allow for a more systematic characterization and a more comprehensive understanding of whole-brain correlates related to its presence. This, in turn, would allow for both more precise and holistic assessments of its functional relevance and impact on brain functions $\lceil 10 \rceil$.

However, a significant challenge arises in neuroimaging due to the high representational dimensionality of AI applications, the anatomical complexity of the brain, the use of unstable evaluation protocols, and the intricate yet uncertain nature of expert annotations. These factors make tasks in this domain particularly difficult to explain and evaluate. Conventional 2D local or global XAI methods often fall short and, in some cases, may even produce misleading explanations [1]. Additionally, as highlighted in [11], the issue of inter-method variability—where different XAI methods emphasize different features as important—can significantly erode trust in AI within the scientific community. These challenges underscore the urgent need for 3D explanation frameworks that can surpass the limitations of traditional XAI methods and provide more reliable global explanations.

1.1. Aim and contribution of the study

This study introduces, for the first time, a comprehensive framework to validate and explain deep learning models, providing transformative insights into pattern learning while establishing a new standard of credibility and reliability for such networks. The proposed XAI 3D-Framework tackles the intricate challenges of explainability in neuroimaging, enabling the detection and interpretation of complex patterns related to the presence or absence of the paracingulate sulcus (PCS), a highly variable cortical structure associated with reality monitoring and psychotic conditions.

To achieve this, we developed an innovative methodology that employs two complementary local XAI methods, GradCam and SHAP, extended into 3D space to analyze the entire dataset used in the binary classification task. These methods are integrated with statistical patterns derived from the 3D brain inputs (Shape) via dimensionality reduction, producing global explanations that outperform traditional XAI methods in both interpretability and faithfulness. By combining

these complementary approaches, our framework not only reveals the underlying learning patterns within the network but also significantly enhances the accuracy and clarity of the results. For the PCS classification task, we implemented two distinct 3D deep learning architectures: a simple 3D convolutional neural network (simple-3D-CNN) and a two-headed attention layer network with diverse backbone options (2CNN-3D-MHL and simple-3D-MHL). These networks utilized high-resolution 3D brain inputs, including grey-white surface boundaries and sulcal skeletons from both hemispheres. Leveraging a well-annotated cohort of 596 subjects from the TOP-OSLO study [12], we trained, validated, and tested these networks, employing a 70%-20%–10% data split.

Our framework introduces several pivotal innovations to the domain of explainable AI in neuroimaging: (i) By integrating statistical features (Shape) that are correlated with reduced dimensionality information, the framework ensures that the discovered patterns are not only grounded in the AI model's learning but also reflect cohort-specific variability. This dual-layered approach bridges statistical data and model-derived insights, enabling a deeper and more contextually relevant understanding of the results. (ii) The extension of established XAI methods, such as GradCam and SHAP, into the 3D domain addresses the critical need for 3D explanations in neuroimaging applications. Conventional 2D local or global XAI methods often fall short and, in some cases, may even produce misleading explanations. (iii) The use of GradCam and SHAP in tandem reduces inter-method variability and bolsters the reliability of the explanations, setting a new benchmark for trustworthy AI applications. The proposed multi-method framework delivers robust and actionable insights, particularly in complex neuroimaging tasks such as cortical morphology studies.

Notably, the XAI 3D-Framework demonstrated superior performance compared to traditional XAI methods in terms of faithfulness for global explanations, successfully identifying significant sub-regions of an atlas brain (the ICBM 2009a Nonlinear Asymmetric atlas, [13,14]). This capability provides a transparent and reliable mechanism to trace the patterns driving network decisions, enhancing trust and enabling deeper exploration of deep learning model outputs in neuroscience. By combining methodological rigor with practical innovation, this framework opens new avenues for understanding and interpreting brain structure-function relationships, making it a foundational tool for advancing both research and clinical applications in neuroimaging.

2. Related work

2.1. The application: Sulcal pattern studies

Cortical folding, which develops during the perinatal period (i.e., in the last few months of gestation and the first few months after birth), results in significant inter-individual variability often overlooked in population studies. Understanding the variability of sulcal patterns is critical for multiple reasons: strict descriptive anatomy, refinement of inter-subject registration, investigation of neurodevelopmental mechanisms, and the search for anatomo-functional correlates. These patterns hold potential for investigating healthy functional variability (e.g., the relationship between cingulate folding patterns and functional connectivity [15]) and pathological outcomes (e.g., paracingulate folding linked to hallucinations in schizophrenia [16]).

While the study of cortical folding variability has been approached through global methods—considering whole-brain or regional sulcal parameters such as gyrification index or sulcal pits—finer investigations often require a focus on specific sulci. This underscores the need for automated sulcal recognition methods. Although several techniques have been developed for general sulcal labeling (reviewed in [17]), to the best of our knowledge, no current method can automatically label the PCS in a 3D approach; most rely on 2D analyses of specific MRI slices [18]. The omission of PCS in whole-brain labeling stems from its complex anatomical specification, defined not only by its location but

also by its orientation (parallel to the cingulate sulcus). Consequently, even newer labeling frameworks fail to identify the PCS within the medial frontal cortex [4].

This gap is particularly critical given the potential importance of investigating whole-brain anatomical correlates of PCS variability in understanding severe symptoms of psychosis. Automatic detection of the PCS is valuable for large dataset exploration, especially given its links to functional variability in reality monitoring [10]. Moreover, incorporating an explainability component through advanced XAI frameworks is unprecedented and transformative. It allows not only for the identification of the PCS but also for uncovering new patterns in its anatomical covariates, offering insights into the functional mechanisms underlying its role. The integration of XAI with AI classification thus provides a robust platform for both enhancing interpretability and advancing our understanding of the broader neurodevelopmental and pathological contexts associated with the PCS.

2.2. Explainable methods

Recent studies have seen a significant surge in the exploration of XAI within medical image analysis and neuroimaging domains [19–22]. XAI methodologies are broadly categorized into interpretable and posthoc approaches. Interpretable methods focus on models that possess inherent properties such as simulatability, decomposability, and transparency, often linked to linear techniques like Bayesian classifiers, support vector machines, decision trees, and K-nearest neighbor algorithms [23]. On the other hand, post-hoc methods are typically used with AI techniques to reveal nonlinear mappings within complex datasets [19,23].

A widely used post-hoc technique is Local Interpretable Model-Agnostic Explanations (LIME), which explains the network's predictions by building simple interpretable models that approximate the deep network locally, i.e. in the close neighborhood of the detected structure [24]. Post-hoc techniques include both model-specific approaches that address specific nonlinear behaviors and model-agnostic approaches that explore data complexity [19,23]. In computer vision, model-agnostic methods such as LIME and perturbation techniques are widely used, while model-specific methods encompass feature relevance, condition-based explanations, and rule-based learning [23,25, 26].

In medical imaging, explainable methods often focus on attribution and perturbation techniques [27]. Attribution techniques like LIME as well as Layer-wise Relevance Propagation (LRP), Gradient-weighted Class Activation Mapping (GradCAM), and Shapley Additive Explanations (SHAP) identify important features for a given prediction by assigning relevance scores to the input features. Perturbation techniques assess the sensitivity of an AI prediction to specific input features by systematically altering sub-groups of the input data [24,27]. GradCAM is notably prevalent among explainable methods in medical imaging due to its ease of application and understanding, as well as its ability to map significant features in the imaging space using the activations of the last convolutional layers [24].

By advancing these explainable methodologies we can better interpret complex models, enhancing their transparency and trustworthiness, particularly in applications like medical imaging, and neuroimaging. An important gap in the literature is the absence of three-dimensional representation frameworks that can explain complex models, such as those used in neuroscience, by providing faithful global explanations. Such frameworks could offer more accurate interpretations than established approaches, potentially improving AI transparency and uncovering new patterns in significant sub-regions involved in classification and prediction.

3. Materials and methods

In this study, we utilized different 3D classifiers to address the binary classification problem of determining the presence or absence of PCS. Furthermore, we developed a novel XAI 3D framework for non-explainable and new-pattern discovery tasks [1]. We employed two distinct explanation methods from the post-hoc family, SHAP and GradCam, which were expanded into three-dimensional space. The outcomes from these two explainable methods (SHAP and GradCam) were concatenated with results derived from a statistical feature extraction model (Shape). The statistical model is a transparent dimensionality reduction algorithm applied to the input data of the validation cohort. This comprehensive approach aims to mitigate potential biases and enhance the robustness of pattern learning discovery.

3.1. Architectural design of deep learning networks for the classification task

We used two different deep learning models for the binary classification task (Fig. 1b.,c.). The first network was a 3D Convolutional Neural Network (CNN) with five levels (see Fig. 1b.). Each level incorporated a CNN block with a 3D convolution layer, a 3D max-pooling layer, and a batch normalization layer. In the first three levels, the 3D convolution layer employed 64, 128, and 256 filters, respectively (as shown in Fig. 1b.). The last level connected to a multi-layer perceptron (MLP) for the final prediction (as illustrated in Fig. 1a.). The MLP comprised three distinct perceptrons and two dropout layers to estimate epistemic uncertainty. For the second network (Fig. 1c.), we used a combination of multi-head attention layers (MHL) to focus on the global diversity and variation of a backbone output. To reduce the biased choice of only one backbone selection, we opted for two distinct backbone networks: (i) a 3D Convolution layer block with two levels of 64 and 128 filters a Global Average Pooling, and a perceptron of 32400 hidden layers (2CNN-3D-MHL), and (ii) the straightforward 3D CNN network outlined previously (Fig. 1b., simple-3D-MHL). We used the multi-head attention mechanism described in [28] and as we used a two-head attention of the same input ("self-attention"). The two heads (heads) are given by:

$$head_c = AT(QW_c^Q, KW_c^K, VW_c^V)$$
 (1)

where c is the backbone output. There are two outputs corresponding to the two uses of the backbone, and each output is connected to a perceptron with N hidden layers (where N is equal to the product of the weight and height of the input image). The AT is the attention layer and it computed by:

$$AT(Q, K, V) = softmax(\frac{QK^{T}}{\sqrt{d_{k}}})V$$
 (2)

where the input matrix are combinations of queries and keys of dimension d_k , and values of dimension d_v . Queries are packed together into a matrix Q. The keys and values are also packed together into matrices K and V. The output of the MHL network is given by:

$$MHL = Concat(head_c, head_c)$$
 (3)

where c defines by the $backbone_{output}$. The output of the MHL was passed again from the MLP presented above to make the final prediction.

3.2. Extending explainability methods in 3D space

We used two distinct explainable techniques: the widely-utilized sensitivity local explainability technique in medical imaging applications known as the GradCam method [29], and a robust attribution explainability technique called SHAP [30].

The significance of the 3D space in neuroscience applications emphasizes a necessity to extend 2D XAI methods into three dimensions. In the pursuit of computing the class-discriminative localization map

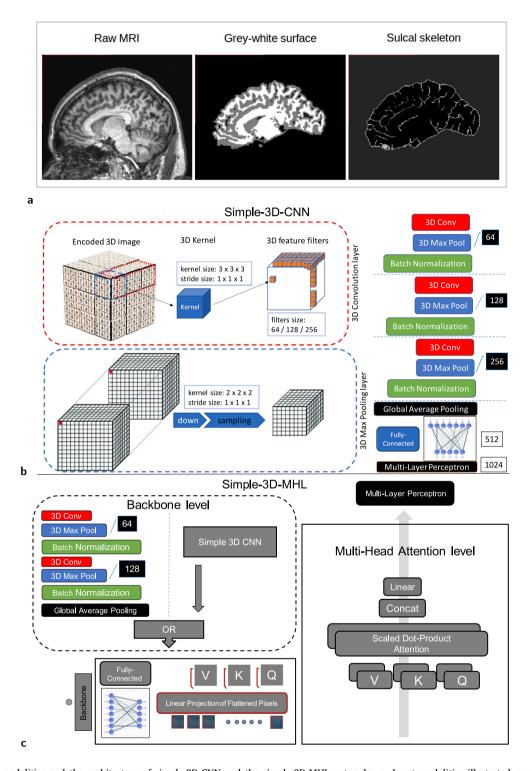


Fig. 1. The input modalities and the architecture of simple-3D-CNN and the simple-3D-MHL networks. a, Input modalities illustrated on a right hemisphere coronal slice. In more detail, the raw MRI of a given subject, the corresponding grey-white surface, and the corresponding sulcal skeleton. b, The architecture of simple-3D-CNN network with explanation of 3D Convolution layer (3D Conv) and 3D Max Pooling layer on the left. c, The three dimension MHL model with two different backbone choices, the full simple-3D-CNN (simple-3D-MHL) and the two level simple-3D-CNN layer (2CNN-3D-MHL).

encompassing width w, height h, and depth d within a specific 3D brain MRI corresponding to a class c (PCS or noPCS), the computation involves determining the gradient of the score for class c, denoted as y^c , in relation to the nth feature activation map (A^n) of the final convolution layer in each deep network. To determine the importance weights (α_n^c) for each n feature activation map, global average pooling

is employed over the width (i), height (j), and depth (k) of each feature.

$$\alpha_n^c = \frac{1}{Z} \sum_i \sum_j \sum_k \frac{dy^c}{dA_{ijk}^n} \tag{4}$$

where Z the summation of i, j and k. Moreover, we used a weighted combination of forward activation maps and a ReLU to deliver the final

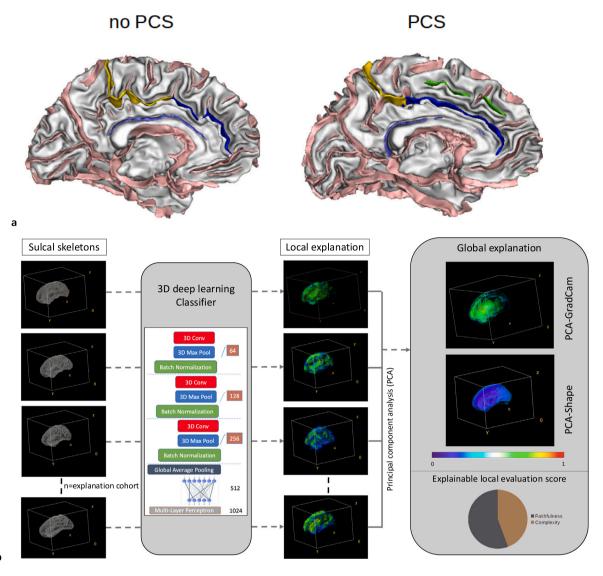


Fig. 2. Classification task determination and the transition from local to global 3D explaination. a, Illustration of the no PCS condition (no PCS), and the PCS condition (PCS, in green line) on two left hemisphere 3D white matter reconstructions obtained with BrainVISA. The cingulate sulcus is colored in yellow and blue and the callosal sulcus is colored in purple. b, The 3D explainable framework that provides both local and global interpretations and explanations of our deep learning 3D classification network's results. The ratio of the faithfulness and complexity metrics were computed at that stage. In this example we include only the GradCam explainability method for simplicity. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

GradCam activation map.

$$GradCam = ReLU(\sum_{n} \alpha_{n}^{c} A_{ijk}^{n})$$
 (5)

SHAP computes the attribution of each pixel of an input image for a specific prediction of a computer vision task. Attribution explainability methods follow the definition of additive feature attribution mainly as a linear function of:

$$g(f, x) = \phi_0 + \sum_{i=1}^{M} \phi_i x_i$$
 (6)

where f is the prediction network, g(f,x) is the explanation model, ϕ_i is the importance of each feature attribution ($\phi_i \in \mathbb{R}$), and M is the number of simplified input features (pixels). Shapley value estimation is one of the main mathematical formulations that the SHAP algorithm uses to assign an importance value to each feature, representing the effect on the model prediction of including that feature (attribution). If we define a subset S of the total feature space (F) of an input 3D image (i = 1...N), where S is the number of samples in the dataset),

and x_i is a 3D matrix of width w and height h and depth d for the ith sample, and x_S is the subset of chosen features in the 3D space, then:

$$\phi_i = \sum_{S/(i)} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f_{S(i)}(x_{S(i)}) - f_S(x_S)]$$
 (7)

Here, $f_{S(i)}$ is a model trained with the presented x_s features, and f_S is another model trained with the features withheld. For our study, we used Deep SHAP [30] to describe our deep learning network models. This approach uses a chain rule and linear approximation as described in [30].

3.3. The XAI 3D-framework

The XAI 3D-Framework introduces a groundbreaking approach to generating global explanations that facilitate the discovery of new patterns in neuroimaging studies. Our proposed framework uniquely integrates statistical features derived from cohort data (Shape) with insights from two complementary explainability methods, GradCam

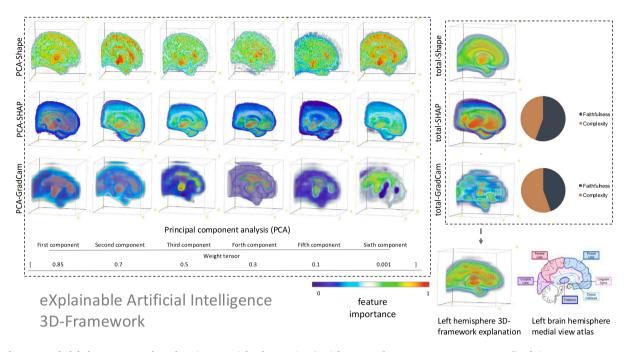


Fig. 3. The proposed global 3D-Framework explanation. A weighted averaging (Weight tensor: [0.85, 0.7, 0.5, 0.3, 0.1, 0.001]) of six PCA components produces the average PCA image for PCA-Shape, PCA-GradCam, and PCA-SHAP. Following this, a weighted averaging (Weight tensor: [0.85, 0.5, 0.1]) of the three Global PCA overlapping images extracted the total overlapping image. This total overlapping image was then registered on a sulcal probabilistic atlas (the ICBM 2009a Nonlinear Asymmetric atlas, [13,14]) to unveil the model's pattern for determining the presence or absence of the PCS.

and SHAP, in a three-dimensional space. By combining statistical information with model-driven learning, the framework provides a dual-layer understanding: one rooted in the cohort's inherent variability and the other reflecting the AI model's decision-making processes based on the classification task. Such integration not only enhances robustness but also minimizes inter-method variability and potential biases [11], ensuring reliable and interpretable outcomes. The use of multiple methods, as opposed to relying on a single explainability approach, represents a significant step forward in delivering more comprehensive and trustworthy explanations, especially in challenging contexts like cortical morphology.

To achieve its goals, the framework employs GradCam and SHAP, two widely used explainability methods, expanded into three-dimensional space to extract local explanations from deep learning classifiers (simple-3D-CNN and simple-3D-MHL; Fig. 2 b.). Faithfulness and complexity scores are assigned to evaluate the quality of these local explanations, ensuring they meet predefined thresholds (see Supplementary Material Table 2). To generate global explanations, we apply Principal Component Analysis (PCA; [31]) to both the sulcal skeleton and grey-white surface 3D brain inputs, capturing a variety of global feature importance patterns (PCA-Shape; see Fig. 3). After testing different configurations, the optimal solution—capturing over 80% of the cohort variance—was achieved with six PCA components.

The global explanations derived from GradCam and SHAP yielded superior faithfulness and complexity results when computed as a weighted average of their six PCA components using the tensor [0.85, 0.7, 0.5, 0.3, 0.1, 0.001], as shown in Supplementary Table 2 (section ii, iii), which compares this full aggregation approach (section ii) to using only the first PCA component (section iii). These weighted averages formed the final global explanations (total-GradCam and total-SHAP). In parallel, the same approach was applied in the Shape domain to determine the global statistical feature importance of sulcal skeleton and grey-white surface inputs (total-Shape; see Fig. 3).

By integrating statistical and model-driven insights, the framework provides a robust platform for uncovering meaningful patterns while maintaining the reliability and interpretability essential for applications in neuroimaging. The weighted average formulation is defined as follows:

$$G(X, W) = \frac{\sum w_i x_i}{\sum w_i} \tag{8}$$

where W is the weight tensor and the X is the pixel images tensor.

3.4. The global explanation of the 3D-framework

To compute the global explanation of the 3D framework, we first manually aligned and rescaled the two global XAI explanations (total-SHAP, total-GradCam) to the average of the six components of the PCA-Shape (total-Shape). Then, we used Eq. (8), incorporating a threecomponent weight tensor of [0.85, 0.5, 0.1]. The weight combination [0.85, 0.5, 0.1] was chosen after empirical testing of multiple combinations, to optimize faithfulness while maintaining reasonable complexity scores. Preliminary results showed that higher weights for total-Shape (e.g., 0.85) yielded superior faithfulness scores, reflecting the importance of feature importance in global explanations. This selection balances the trade-offs between faithfulness, complexity, and redundancy across the metrics. After we define the best combination, we conducted an ablation study to identify the optimal combination of the weight tensor. The different cases we examined were denoted as 851, 815, 185, 158, 518, and 581, which are fixed-order representations of the methods in the following sequence: total-Shape, total-SHAP, and total-GradCam explanations. The weights correspond to 0.85 as '8', 0.5 as '5', and 0.1 as '1'. For example, 3D-Framework-851 refers to the proposed 3D-Framework with weight values of 0.85 for total-Shape, 0.5 for total-SHAP, and 0.1 for total-GradCam. These components were derived from both the sulcal skeleton and the grey-white surface inputs of the total dataset, considering both the right and left hemispheres. In our specific study, the best combinations identified were 851 and 815 (see Supplementary Material Table 2).

To minimize potential bias that might occur by relying on only one deep learning network, we applied this approach to both the simple-3D-CNN and simple-3D-MHL networks. Identifying the significant features of the networks explanations provides insights into the mechanisms driving the network's decision-making process. For enhanced clarity,

the key brain sub-regions of interest, corresponding to the sulcal skeleton and grey-white surface, are visually depicted in Fig. 1a. Finally, the global explanation from the 3D framework was registered on a sulcal probabilistic atlas to illustrate the model's pattern associated with determining the presence or absence of the PCS. The registration process involved affine transformations, including translation, rotation, scaling, and geometry adjustments, to align with the probabilistic atlas (the ICBM 2009a Nonlinear Asymmetric atlas, [13,14]).

3.5. Cohort's description and pre-processing image analysis

We used the structural MRI of 596 participants from the TOP-OSLO study [12] for a binary classification task. The participants encompassed individuals with a diagnosis on the schizophrenia spectrum (183), on the bipolar disorder spectrum (151), and unaffected control participants (262). T1-weighted images were acquired using a 1.5 T Siemens Magnetom Sonata scanner (Siemens Medical Solutions, Erlangen, Germany). The mean age of the schizophrenia spectrum cohort was 31.0 years (SD = 9.1), with a sex distribution of 108 men and 75 women. For the bipolar spectrum cohort, the mean age was 34.3 years (SD = 11.6), with 65 men and 86 women. The healthy control group had a mean age of 34.7 years (SD = 9.7), with 138 men and 124 women. Diagnostic groups were balanced across demographic variables to minimize potential confounding effects in downstream analyses.

Two experts, A.A. and H.V., performed image annotations, categorizing them into two classes: 'no paracingulate sulcus' (noPCS) and 'paracingulate sulcus' (PCS). This annotation was based on the protocol described in [2] with further details available in the supplementary materials under '1.1 PCS classification of TOP-OSLO' and is illustrated in Fig. 2a. A more analytical description of the TOP-OSLO cohort details is available in the study [12].

3.6. Processing of two distinct images inputs

We initially processed the brain structural MRIs using the BrainVISA software [32] extracting two images as inputs for the classifier: the grey-white surface and the sulcal skeleton. These were extracted from the raw MRI with an established protocol consisting of bias correction, histogram analysis, brain segmentation, hemisphere separation, dichotomization of the white matter from the union of grey matter and cerebrospinal fluid, and skeletonization of the result, as detailed in [33]. Specifically, the grey-white surface was obtained by minimizing a Markov field and the segmentation used homotopic deformations of the hemisphere bounding box, resulting in the grey-white surface, where voxels are dichotomized into either grey or white. The skeleton was then derived from this object by applying a homotopic erosion embedding a watershed alogrithm that preserves the initial topology resulting in the sulcal skeleton. These two modalities were then used to train and evaluate our networks as well as for our explainability methods. Fig. 1a, b, c. show the structural MRI and the corresponding grey-white surface and sulcal skeleton outputs from BrainVISA.

3.7. Hyper-parameter initialization

After randomly shuffling the data, each dataset was split into training, validation, and testing sets containing 70%, 20%, and 10% of the total number of images, respectively. To account for diagnostic heterogeneity (schizophrenia, bipolar disorder, and controls), the dataset splits were stratified to preserve diagnostic proportions across training, validation, and testing subsets. We additionally ensured that age and sex distributions were balanced across splits. This procedure reduced the risk of demographic or diagnostic bias influencing the classification results. Sparse categorical cross-entropy was used as the cost function and the loss function were optimized using the Adam algorithm [34]. After manual hyper-parameter searching the best learning rate was a

value of 0.0001. We utilized a strategy of exponentially decreased during the first 50 epochs and then fixed at 0.0001 for the last 50 epochs. To train the networks, an early stopping criterion of 10 consecutive epochs was employed and a maximum of 100 epochs was used for both input modalities (sulcal skeleton and grey-white surface) for both the left (L) and right (R) hemispheres. Finally, we use data augmentation techniques including rotation (around the center of the image by a random angle from the range [-15°,15°]), width shift (up to 20 pixels), height shift (up to 20 pixels), and Zero phase Component Analysis (ZCA; [35]) whitening (add noise in each image) to avoid overfitting.

3.8. Evaluation metrics for the explanation

A crucial aspect of this study lies in evaluating how accurate and comprehensive were the local and global explanations. To derive a useful explanation, two primary scores play a pivotal role: **faithfulness** and **complexity**.

An intuitive way to assess the quality of an explanation is by measuring how accurately it reflects the model's actual decision-making behavior in response to input perturbations [36,37]. For a deep neural network f and input x, we define the **faithfulness** of an explanation method g as the Pearson correlation between the sum of attribution scores for a subset of features and the change in model output when those features are replaced by a baseline. This is formally expressed as:

$$M_{\text{faith}}(f, g; \mathbf{x}) = \text{corr}_{S} \left(\sum_{i \in S} g(f, \mathbf{x})_{i}, \ f(\mathbf{x}) - f(\mathbf{x}[\mathbf{x}_{s} = \mathbf{x}_{s}^{f}]) \right)$$
(9)

Here, S is a randomly selected subset of feature indices ($S \subseteq \{1, 2, \ldots, d\}$), x_s are the perturbed features replaced by baseline values x_s^f , and x_f are the remaining, unperturbed features.

To complement faithfulness, we also compute the **complexity** of an explanation, which captures how concentrated or diffuse the attribution is across input features. A lower complexity implies a more focused, and thus more interpretable, explanation. Complexity is measured using the entropy of the normalized attribution scores:

$$M_{\text{compx}}(f, g; \mathbf{x}) = \sum_{i=1}^{d} P_g(i) \cdot \log\left(\frac{1}{P_g(i)}\right)$$
(10)

where the attribution distribution $P_g(i)$ is defined as:

$$P_g(i) = \frac{|g(f, \mathbf{x})_i|}{\sum_{j=1}^d |g(f, \mathbf{x})_j|}$$
(11)

These two metrics offer complementary insights: faithfulness ensures that the explanation aligns with the model's behavior, while complexity encourages sparsity and interpretability.

In order to evaluate these two explainability metrics, we used the software developed by Hedström et al. [38]. This software package is a comprehensive toolkit that collects, organizes, and evaluates a wide range of performance metrics, proposed for explanation methods. Note that we used a zero baseline ('black'; $\boldsymbol{x}_s^f = \boldsymbol{0}$) and 70 random perturbations to calculate the faithfulness score.

Finally, to extract the faithfulness and complexity scores of the global explanations for the total-SHAP, total-GradCam, and the proposed 3D-Framework, we utilized again the software developed by Hedström et al. [38]. As a first step we manually aligned and rescaled the two global XAI explanations (total-SHAP and total-GradCam) to the total-Shape. The input image consisted of the total-Shape results, while the total-GradCam, total-SHAP, and 3D-Framework global explanation served as reference explanations, respectively. In this context, the score of a global explanation makes sense, as the individual input brains of the cohort were aligned in the same template and the most significant variability of each class was assigned in the total-Shape. Consequently, we anticipated the classifier to classify correctly whether a MRI has a PCS or not.

Table 1Performance metrics of global explanation scores of faithfulness, and complexity for the global state of the art explanation methods; total-SHAP, total-GradCam and the propose 3D-Framework of the simple-3D-MHL network.

Explainable evaluation score of used explainable methods								
XAI method	Left sulcal skeleton PCS/noPCS	Left white/grey PCS/noPCS	Right sulcal skeleton PCS/noPCS	Right white/grey PCS/noPCS				
total-SHAP	0.105/0.200	0.092/0.113	0.103/0.102	0.073/0.088				
total-IntGrad	0.108/0.202	0.140/0.150	0.112/0.115	0.097/0.109				
total-GradCAM	0.044/0.090	0.143/0.164	0.067/0.085	0.119/0.129				
3D-Framework	0.223/0.274	0.207/0.222	0.188/0.195	0.192/0.214				
total-SHAP	14.593/14.586	14.572/14.544	14.571/14.572	14.583/14.582				
total-IntGrad	14.594/14.590	14.583/14.589	14.581/14.585	14.589/14.590				
total-GradCAM	14.473/14.555	14.373/14.403	14.586/14.582	14.596/14.547				
3D-Framework	14.584/14.563	14.582/14.587	14.579/14.573	14.587/14.587				
	XAI method total-SHAP total-IntGrad total-GradCAM 3D-Framework total-SHAP total-IntGrad total-GradCAM	XAI method Left sulcal skeleton PCS/noPCS total-SHAP 0.105/0.200 total-IntGrad 0.108/0.202 total-GradCAM 0.044/0.090 3D-Framework 0.223/0.274 total-SHAP 14.593/14.586 total-IntGrad 14.594/14.590 total-GradCAM 14.473/14.555	XAI method Left sulcal skeleton PCS/noPCS PCS/noPCS total-SHAP 0.105/0.200 0.092/0.113 total-IntGrad 0.108/0.202 0.140/0.150 total-GradCAM 0.044/0.090 0.143/0.164 3D-Framework 0.223/0.274 0.207/0.222 total-SHAP 14.593/14.586 14.572/14.544 total-IntGrad 14.594/14.590 14.583/14.589 total-GradCAM 14.473/14.555 14.373/14.403	XAI method Left sulcal skeleton PCS/noPCS PCS/noPCS PCS/noPCS PCS/noPCS PCS/noPCS total-SHAP 0.105/0.200 0.092/0.113 0.103/0.102 0.140/0.150 0.112/0.115 0.112/0.115 0.140/0.150 0.143/0.164 0.067/0.085 0.143/0.164 0.067/0.085 0.143/0.164 0.067/0.085 0.143/0.164 0.105/0.105 0.188/0.195 0.105/0.223/0.274 0.207/0.222 0.188/0.195 0.143/0.164 0.105/0.164 0.105/0.164 0.105/0.165 0.165/0.165 0.165/0.165 0.165/0.165 0.165/0.165 0.165/0.165 0.165/0.165/0.165 0.165/0.				

4. Results

4.1. Classifiers performance for presence or absence of paracingulate sulcus

For brevity, we discuss in the main manuscript only the simple-3D-MHL results which outperformed the two-level CNN backbone network (2CNN-3D-MHL). The analytical tables and results for 2CNN-3D-MHL can be found in the supplementary material subsection 2.1 ('Classification results of the 2CNN-3D-MHL'; Table 1).

The performance of the simple-3D-MHL network in the left hemisphere was higher (around 73.00% in all testing metrics and 74.10% in all validation metrics) than that in the right hemisphere (around 58.00% in all testing metrics and 63.10% in the validation metrics). For the simple-3D-CNN, the performance of the network in the left hemisphere was higher (around 72.90% in all testing metrics and 74.00% in all validation metrics) than that in the right hemisphere (around 56.00% in all testing metrics and 63.00% in the validation metrics). The analytical Figures and outcomes for the simple-3D-MHL and simple-3D-CNN networks are presented in supplementary material subsection 2.2 ('Additional global explainability methods and different components PCA results'; Supplementary Fig. 1 a.,b.).

The discrepancy in performance between the left and right hemispheres was to be expected for two reasons. First, the PCS is more prominent in the left than in the right hemisphere [39,40], including in psychopathological conditions such as schizophrenia [7]. Furthermore, the left PCS has a greater number of associations with regional cortical thickness and sulcal depth than the right PCS [6], implying more covariability of anatomical features contained in either of our input modalities with the presence of the PCS in the left hemisphere than the right.

4.2. Global explanations and their PCA component results

We extracted and evaluated the explainability results from both networks to avoid biased observations and to investigate whether there was a clear cause-and-effect relationship between the quality of explanation and prediction performance. We present the results of the simple-MHL network in the main manuscript as it had slightly better performance compared to the simple-3D-CNN. The results for the simple-3D-CNN are thoroughly detailed in the supplementary material subsection 2.2 (refer to 'Additional global explainability methods and different components PCA results'; Supplementary Fig. 2 and 3).

The first component analysis of PCA explainability results of simple-3D-MHL networks on the left and right hemisphere of the grey-white surface inputs mainly focus on the frontal lobe (mostly inferior lateral and inferior medial), the cingulate gyrus, the temporal lobe, and occasionally the thalamus for detecting the presence or absence of the PCS. More specifically, for the detection of the presence of PCS (paracingulate sulcus; Fig. 4) in the left and right hemisphere the simple-MHL network focuses more in the frontal lobe (medial and inferior lateral), cingulate gyrus (mostly anterior), temporal lobe and

sometimes thalamus. Conversley, for the absence of PCS (No paracingulate sulcus; Fig. 4) in the left hemisphere the network focuses in the frontal lobe (mostly inferior medial and inferior frontal), the temporal lobe, the cingulate gyrus, occasionally the thalamus, and specifically for the PCA-GradCam, the corpus callosum. On the right hemisphere the simple-MHL network focuses in frontal lobe (medial and lateral), the temporal lobe, and the cingulate gyrus.

Fig. 5.a,b displays the comprehensive explanations of PCA-GradCam, PCA-SHAP, and PCA-Shape for the simple-3D-MHL network when using the sulcal skeleton inputs. The neural network shows distributed attention but still emphasizes some key regions in both hemispheres. In the left hemisphere, for the PCS class, there is focus on the superior temporal sulcus and its branches, the posterior sylvian fissure, and parts of the central and precentral sulci, with some attention on the cingulate sulcus and the medial frontal sulcus (containing the PCS) (Fig. 5a.). Conversely, in the right hemisphere for the same class, emphasis is on the superior and inferior temporal sulci, cingulate sulcus, medial frontal sulcus (containing the PCS), and the sub-parietal sulcus (Fig. 5b.). For the noPCS class, the left hemisphere shows similar yet less specific attention than for the PCS class, with an additional focus on the sylvian fissure and insula, and the anterior cingulate sulcus (Fig. 5a.). In the right hemisphere, additional focus is shifted to the anterior cingulate sulcus, the sub-parietal sulcus, and elements of the ventricles (Fig. 5b.).

To validate the pattern indicating the presence or absence of the PCS, we assessed the variability across all six components of the PCA. Notably, we observed differences in the six PCA components between the explanations and the global feature importance (refer to Supplementary Fig. 4). There were differences in the intensity and extent of regions highlighted between the sulcal skeleton and grey-white surface inputs, although the primary regions of focus remained consistent. We found that the faithfulness and complexity score of the PCA's first component for SHAP and GradCam methods performed poorly compared to the weighted average output of all six components as described in subsection 3.4 (total-SHAP, total-GradCam). Consequently, we used total-SHAP, total-GradCam, and total-Shape to compute the global explanation for the 3D framework.

4.3. Global explanations from the 3D-framework and the pattern learning results using grey-white surface inputs

For the simple-3D-MHL on grey-white surface inputs (see Fig. 6a.c), in the right hemisphere the focus was on the frontal lobe, the insula, and parietal lobe in the PCS existence and, in the PCS absence (noPCS) on the temporal lobe, frontal lobe, cingulate gyrus and parietal lobe. In the left hemisphere, the PCS decision mostly relied on the cingulate gyrus, frontal lobe parietal lobe, and corpus collosum. The noPCS condition predominantly focused on the frontal lobe, and lateral temporal lobe and cingulate gyrus (see Fig. 6a.c).

For the simple-3D-CNN on grey-white surface inputs (refer to Supplementary material Fig. 6 a.) in the right hemisphere, highlighted

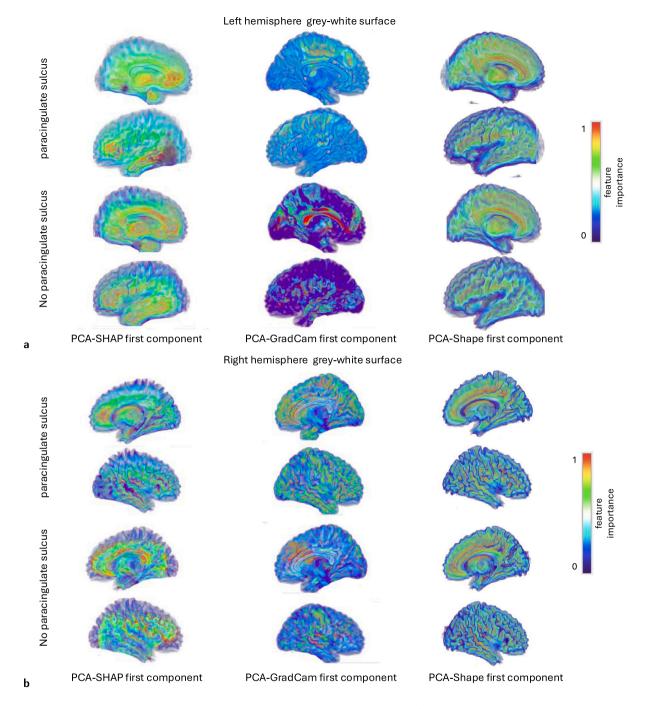


Fig. 4. Simple-3D-MHL results on the left and right hemisphere of grey-white surface brain inputs. a,b, show the explainability results for the PCS class images of the first component among the six components of PCA for the total input modality (PCA-Shape), the total corresponding GradCam results (PCA-GradCam), and the total corresponding SHAP results (PCA-SHAP). The feature's importance (pixel attribution) varies from 0 (blue color) to 1 (red color), with high importance being 1 for the PCA-GradCam and PCA-Shape results. The orientation of the results are based on the medial anatomical views. All the presented results are align and mapping in the ICBM 2009a Nonlinear Asymmetric atlas [13,14]. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

the lateral inferior frontal lobe, and inferior temporal lobe in the PCS condition and the thalamus, cingulate gyrus, lateral anterior occipital lobe, and posterior temporal lobe. In the noPCS condition, the focus was on the thalamus, cingulate gyrus, the medial frontal lobe, and the posterior temporal lobe. In the left hemisphere, both conditions globally focused on the same regions: the lateral middle frontal lobe, and inferior and superior parietal lobe, and the thalamus and in the cingulate, frontal and medial parietal lobes, with a special focus on the

anterior cingulate gyrus. The PCS condition additionally focused on the lateral view of the posterior temporal lobe.

For both networks, in the right hemisphere the focus was primarily on the medial aspect of the brain (except for the PCS condition in the right hemisphere) with the main contributions in the frontal lobe and cingulate gyrus, suggesting a rather constrained explainability of the PCS presence. Conversely, in the left hemisphere, the focus was much more broadly distributed, with strong contributions stemming from the

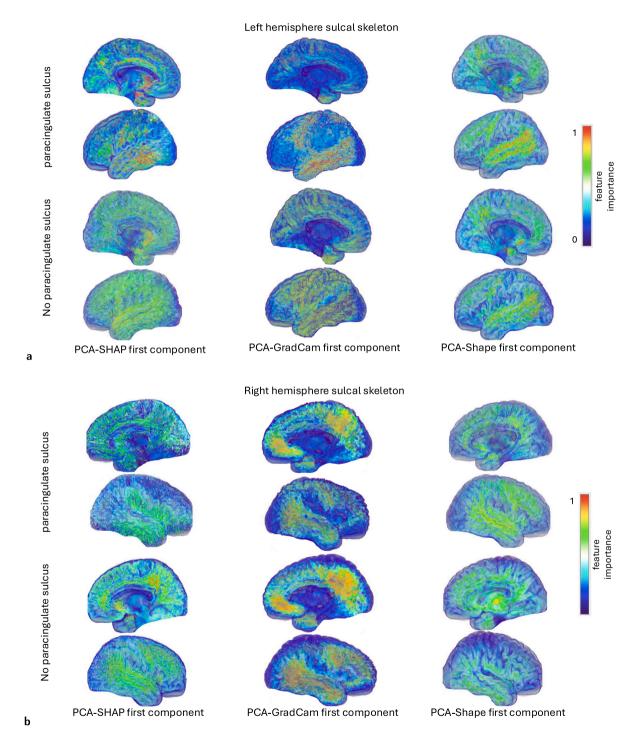


Fig. 5. Simple-3D-MHL results on the left and right hemisphere of sulcal skeleton brain inputs. a-b, show the explainability results for the noPCS class images of the first component among the six components of PCA for the total input modality (PCA-Shape), the total corresponding GradCam (PCA-GradCam), and the total corresponding SHAP results (PCA-SHAP). The feature's importance (pixel attribution) varies from 0 (blue color) to 1 (red color), with high importance being 1 for the PCA-GradCam and PCA-Shape results. The orientation of the results are based on the medial anatomical views. All the presented results are align and mapping in the ICBM 2009a Nonlinear Asymmetric atlas [13,14]. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

lateral cortex, suggesting that the developmental mechanisms leading to the presence of the PCS are related to the wider development of the brain [6].

4.4. Global explanation from the 3D-framework and the pattern learning results using sulcal skeleton inputs

A consistent method to overlay the outputs of the sulcal skeleton, similar to the process described for the grey-white surface outputs, was applied. For the simple-3D-MHL on sulcal skeleton inputs (Fig. 6b, c), in the right hemisphere, the presence of PCS focuses on the lateral superior temporal sulcus, sylvian fissure, inferior frontal sulcus, and the anterior cingulate sulcus. The absence of PCS (noPCS; Fig. 6c) focuses on the inferior temporal sulcus, superior temporal sulcus, ventricule, callosal sulcus, and inferior frontal sulcus. In the left hemisphere, the PCS and noPCS conditions had very different highlights. In the PCS condition, the main focuses were on the lateral posterior part of the lateral superior temporal sulcus, posterior inferior temporal sulcus, sylvian fissure, and the inferior parietal sulcus. In the noPCS condition, the main contributions were in the sylvian fissure, lateral superior temporal sulcus, superior frontal sulcus, internal parietal sulcus, anterior cingulate sulcus, and ventricle.

For the simple-3D-CNN on sulcal skeleton inputs, aimed at the accurate detection of PCS within the sulcal hemisphere inputs, the pivotal sub-regions encompassed the superior temporal sulcus, inferior precentral sulcus, sylvian fissure and sub-parietal sulcus (see Supplementary Material Fig. 6b). Conversely, when PCS was absent, the critical sub-regions within the right hemisphere sulcal skeleton inputs encompassed the ventricule, superior temporal sulcus, internal parietal sulcus and rostral sulcus. Transitioning to the left hemisphere, the sulcal skeleton inputs underscore the significance of the superior temporal sulcus the, sylvian fissure, and the internal parietal sulcus. When PCS was not present, the important left hemisphere sulcal skeleton inputs comprised the superior temporal sulcus, ventricle, inferior precentral sulcus, internal parietal sulcus (see Supplementary Material Fig. 6b).

We thereafter identified the common regions between the two networks' outputs. The overlap results of the two networks, simple-3D-CNN and simple-3D-MHL, for the presence and absence of PCS (noPCS) reveal several common regions of interest in both the right and left hemispheres. For the presence of PCS in the right hemisphere, both networks highlight the lateral superior temporal sulcus and sylvian fissure. In the left hemisphere under the PCS condition, both models emphasize the superior temporal sulcus, the sylvian fissure, and the internal parietal sulcus. For the absence of PCS (noPCS) in the left and right hemisphere, the networks overlap in highlighting the ventricle, part of superior temporal sulcus and the internal parietal sulcus.

4.5. Ablation study of our 3D explainability framework

An ablation study in the simple-3D-MHL model was conducted to evaluate various combinations of global explanations (total-GradCam and total-SHAP) and global feature importance (total-Shape) (see Table 2). The best results were achieved by assigning the highest weight (0.85) to total-Shape. This aligns with the idea that the feature importance of the inputs plays a crucial role in the explanations. For both hemispheres in the sulcal skeleton, a weight of 0.5 given to total-SHAP produced the highest faithfulness scores. For the grey-white surface inputs, total-GradCam with the same weight yielded the best results (a weight of 0.5). There were no significant differences observed in the complexity scores among the different combinations. The same patterns were observed for the simple-3D-CNN (Supplementary material; Table 2).

Table 2

Performance metrics assessing the global explanation scores of faithfulness and complexity were computed across various weight combinations assigned to the global XAI methods (total-SHAP, total-GradCam) and global feature extraction (total-Shape) within the proposed 3D-Framework for evaluating the global explanations of the simple-3D-MHL network. These combinations, denoted as 851, 815, 185, 158, 518, and 581, represent fixed-order representations of the assigned weights in the following sequence: total-Shape, total-SHAP, and total-GradCam explanations, with weight values of 0.85 represented as '8', 0.5 as '5', and 0.1 as '1'. For instance, 3D-Framework-851 refers to the proposed 3D framework with weight values of 0.85 in total-Shape, 0.5 in total-SHAP, and 0.1 in total-GradCam. The results shows the Right and Left white-grey (a) and sulcal skeleton (b) images.

XAI metrics	XAI method	Left white/grey PCS/noPCS	Right white/grey PCS/noPCS
Faithfulness	3D-Framework-851	0.166/0.197	0.172/0.182
	3D-Framework-815	0.207/0.222	0.192/0.214
	3D-Framework-185	0.118/0.124	0.113/0.133
	3D-Framework-158	0.143/0.173	0.155/0.162
	3D-Framework-518	0.125/0.152	0.136/0.142
	3D-Framework-581	0.087/0.102	0.103/0.112
Complexity	3D-Framework-851	14.582/14.582	14.595/14.592
	3D-Framework-815	14.582/14.587	14.587/14.587
	3D-Framework-185	14.585/14.582	14.584/14.587
	3D-Framework-158	14.587/14.583	14.593/14.592
	3D-Framework-518	14.592/14.584	14.594/14.594
	3D-Framework-581	14.593/14.592	14.597/14.592
(b) Ablation st	udy of the proposed 3D-F	ramework	
Faithfulness	3D-Framework-851	0.223/0.274	0.188/0.195
	3D-Framework-815	0.204/0.233	0.165/0.174
	3D-Framework-185	0.105/0.122	0.053/0.133
	3D-Framework-158	0.074/0.095	0.122/0.076
	3D-Framework-518	0.115/0.106	0.144/0.134
	3D-Framework-581	0.126/0.138	0.135/0.142
Complexity	3D-Framework-851	14.584/14.563	14.579/14.573
Complexity			
Complexity	3D-Framework-815	14.594/14.595	14.582/14.576
Complexity	3D-Framework-815 3D-Framework-185	14.594/14.595 14.595/14.584	14.582/14.576 14.572 /14.595
Complexity		,	
Complexity	3D-Framework-185	14.595/14.584	14.572 /14.595

4.6. Evaluation of the global explanation from the 3D-framework and the XAI methods and the pattern learning results

To evaluate whether the global explanation from the 3D-Framework was superior to those provided by SHAP or GradCam, we scored the explanations with respect to faithfulness and complexity (see Table 1). For the simple-3D-MHL network, our proposed 3D framework outperformed total-GradCam and total-SHAP in terms of faithfulness score in the left hemisphere with values exceeding 0.21 compared to scores of less than 0.16 for total-GradCam, and less than 0.11 for total-SHAP. In the right hemisphere, our proposed 3D framework again outperformed total-GradCam and total-SHAP, achieving faithfulness scores over 0.18 compared to scores of less than 0.13 for total-GradCam, and less than 0.10 for total-SHAP. The 3D framework achieved the second-best result in complexity scores with total-GradCam having the lowest score in the left hemisphere and total-SHAP in the right hemisphere. To ensure the robustness of our findings and to mitigate potential biases arising from relying solely on linear combination-based explanation methods such as SHAP and GradCAM, we further compared our framework against Integrated Gradients (total-IntGrad), a wellestablished attribution technique [41]. Our framework demonstrated superior performance on both explanation metrics, faithfulness and complexity, consistently across all hemispheres and input modalities (see Table 1).

To strengthen the evaluation of our framework, we performed a detailed statistical analysis comparing explanation quality across

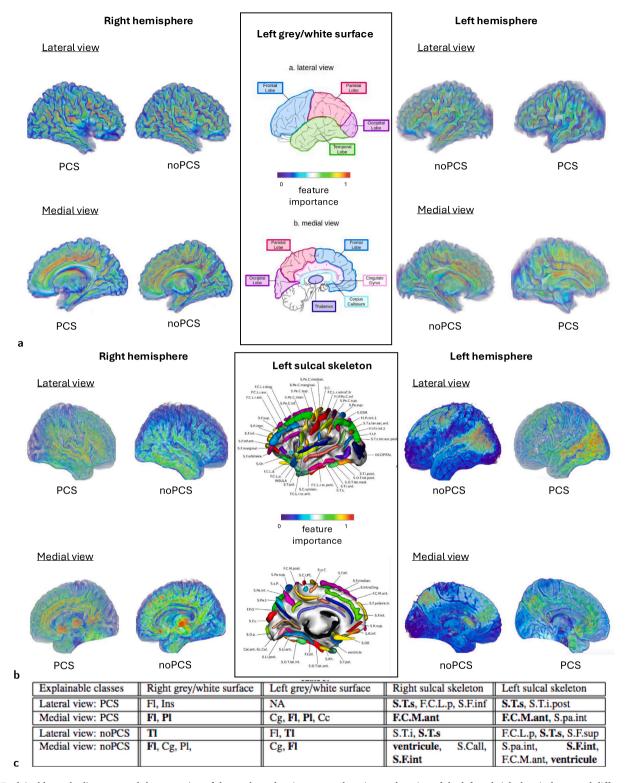


Fig. 6. Explainable method's scores and the extraction of the total overlapping pattern learning explanation of the left and right hemispheres and different inputs with expert's observation. a shows the total overlapping pattern learning results for the right and left hemisphere of the brain for grey-white surface images of the simple-3D-MHL network. b shows the total overlapping pattern learning results for the right and left hemispheres of the brain for sulcal skeleton inputs of the simple-3D-MHL network. All the presented results are align and mapping in the ICBM 2009a Nonlinear Asymmetric atlas [13,14]. c shows the pattern learning results from the overlapping of the simple-3D-MHL deep learning networks on the total overlapping pattern learning results of the lateral and medial views based on experts' observation. For the grey-white surface input we used the acronyms, T: thalamus, H: hypothalamus, Fl: Frontal lobe, Ol: occipital lobe, Tl: temporal lobe, Pl: parietal lobe, Cc: corpus callosum, Cg: cingulate gyrus, NA: none. For the skeleton sulcal input we used the acronyms, S.T.s: superior temporal sulcus, S.T.i: inferior temporal sulcus, F.C.M.ant: anterior cingulate sulcus, S.pa.int:internal parietal sulcus, F.C.L.p: sylvian fissure, S.F.sup: superior frontal sulcus, S.F.int: internal frontal sulcus, S.F.inf: inferior frontal sulcus, S.Call: callosal sulcus.

Table 3

Comparison of explanation metrics (Faithfulness and Complexity) across PCS and noPCS classes and anatomical regions. Values represent the mean and 95% confidence interval (CI) of our method, and p-values from paired t-tests comparing SHAP and GradCAM with our method.

Region	Class	Metric	Ours (Mean [95% CI])	SHAP p-value	GradCAM p-value
Left sulcal skeleton	PCS	Faithfulness	0.287 [0.219, 0.355]	0.040	0.940
Left grey/white surface	PCS	Faithfulness	0.266 [0.155, 0.376]	0.038	0.749
Right sulcal skeleton	PCS	Faithfulness	0.234 [0.222, 0.246]	0.049	0.034
Right grey/white surface	PCS	Faithfulness	0.317 [0.168, 0.466]	0.135	0.046
Left sulcal skeleton	noPCS	Faithfulness	0.138 [0.038, 0.238]	3.85×10^{-3}	1.43×10^{-2}
Left grey/white surface	noPCS	Faithfulness	0.329 [0.168, 0.427]	3.68×10^{-3}	0.920
Right sulcal skeleton	noPCS	Faithfulness	0.201 [0.036, 0.365]	0.151	2.27×10^{-4}
Right grey/white surface	noPCS	Faithfulness	0.189 [0.108, 0.270]	6.39×10^{-4}	7.85×10^{-4}
Left sulcal skeleton	PCS	Complexity	13.33 [13.29, 13.37]	8.48×10^{-9}	5.25×10^{-9}
Left grey/white surface	PCS	Complexity	12.28 [12.23, 12.33]	4.76×10^{-8}	3.74×10^{-5}
Right sulcal skeleton	PCS	Complexity	13.22 [13.17, 13.27]	1.68×10^{-10}	8.96×10^{-3}
Right grey/white surface	PCS	Complexity	12.27 [12.19, 12.34]	3.30×10^{-5}	4.51×10^{-10}
Left sulcal skeleton	noPCS	Complexity	13.29 [13.25, 13.33]	6.70×10^{-8}	9.95×10^{-1}
Left grey/white surface	noPCS	Complexity	12.25 [12.21, 12.30]	5.38×10^{-6}	1.10×10^{-4}
Right sulcal skeleton	noPCS	Complexity	13.28 [13.23, 13.34]	2.51×10^{-10}	2.80×10^{-1}
Right grey/white surface	noPCS	Complexity	12.25 [12.19, 12.32]	2.14×10^{-4}	2.06×10^{-6}

methods. Specifically, we conducted two-sided paired t-tests for each anatomical region to assess significant differences between our framework and the SHAP and GradCAM baselines. Results are reported alongside 95% confidence intervals, enhancing the interpretability and transparency of our findings. Each comparison was based on n=10 participants per group (PCS and noPCS). The analysis demonstrates that our method significantly outperforms the baselines in both faithfulness and complexity metrics across most regions and input modalities (see Tables 3), thereby underscoring the statistical robustness and consistency of our explanation approach.

Up to this point, we have mainly explored the explanation results visually. However, we aimed to automate the process to identify the most significant subregions of interest based on the hypothesis (the classification task). To this end, we applied an affine registration to the total overlapping results of sulcal skeleton inputs from each hemisphere onto a probabilistic atlas of sulci [42]. For this task, we explored the sulcal skeleton output of the simple-3D-MHL network as it slightly outperformed the simple-3D-CNN (see Supplementary Material Table 2 ii.) in the classification task and delivered better global explanations, faithfulness, and complexity scores than the simple-3D-CNN. Additionally, it follows patterns based on evidence from the literature [15,43,44].

Fig. 7 presents the distribution of the most relevant voxels for outcome decisions within sulcal probabilistic areas according to different thresholds. The decisions in the right hemisphere were highly focused on specific sulci, with up to three sulci contributing to the 20% threshold, which we retained as the lower threshold (blue). Conversely, the left hemisphere predictions were based on broader considerations, with a number of sulci already contributing to the decision at the 5% threshold, which we retained as the lower threshold for the left hemisphere (blue).

For both conditions and both hemispheres, a specific focus was on the superior temporal sulcus and its posterior branches with a smaller but consistent contribution of the internal parietal and sub-parietal sulci. The noPCS condition on the left hemisphere additionally focused on the precentral sulcus and the Sylvian fissure. Interestingly, no specific focus was oriented towards the internal frontal sulcus (S.F.int), the probabilistic region in which the PCS is located when present.

5. Discussion

Explainability is essential in medical imaging. Healthcare applications of AI need to be able to explain their decision-making to build trust and ensure that their predictions align with other symptoms and signs that affect health. Neuroimaging, the combination of brain images and computational methods, is a research application of medical

imaging. Here, explainability for AI predictions supports the assessment of the validity of results, but can also identify key contributors to decisions that in themselves reveal new patterns and directions for future investigation.

Our prior study [1] categorized the need for AI explanations into self-explainable, semi-explainable, non-explainable applications, and new-pattern discovery, based on the variability of expert opinions, the stability of the evaluation protocol and the representation dimensionality of the application. We applied the proposed guidance in a binary classification task related to a sub-region of the medial surface of the brain where secondary sulci are highly reproducible related with symptoms of psychosis, specifically hallucinations. This application was of the new-pattern discovery class. The output of the explainability indicated a wide distribution of brain regions on which the predictions depend suggesting covariant development of these regions during the perinatal period [45].

Automatic classification of psychotic and control patients based on structural MRI is a challenging task [46]. In most cases, an acceptable detection rate is around 80.0% in clinical applications. However, in highly heterogeneous and variable cohorts, a lesser performance can be acceptable; approximately 60.0-70.0% acceptance accuracy; [47, 48]. The TOP-OSLO cohort was particularly difficult for classification tasks of psychotic and control patients using structural MRI, with an accuracy of less than 60.0% [46]. The variability of the paracingulate region in the TOP-OSLO cohort and the heterogeneity of the dataset create a highly challenging context for the automatic binary classification task of PCS presence. In this study, an accuracy of more than 70.0% in the left hemisphere and more than 60.0% in the right hemisphere was achieved, delivering an acceptable automated 3D deep learning network [47,48] to apply global explainability methods for new-pattern discovery [1]. While these classification results are modest, especially in the right hemisphere, they are in line with performance reported in similarly complex neuropsychiatric tasks. We acknowledge this constraint on predictive accuracy and emphasize the study's primary contribution in explainability and novel pattern discovery.

For the binary classification task of PCS presence, we developed two different 3D deep learning networks: a simple 3D convolutional neural network and a two-headed attention layer network. These networks utilized 3D brain inputs derived from preprocessed structural MRI scans, which included grey-white surface boundaries and sulcal skeletons from both hemispheres of a well-annotated cohort of 596 subjects. The performance of all networks was higher in the left hemisphere than in the right hemisphere. This discrepancy in performance was expected as the PCS is more prominent in the left hemisphere, including in psychopathological situations such as schizophrenia [7]. Moreover,

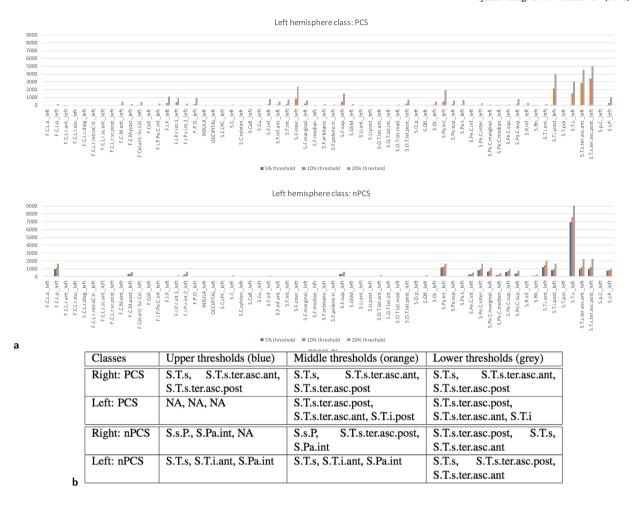


Fig. 7. A representation depicting the presence and absence of the paracingulate sulcus on the left hemisphere sulcal skeleton input **a**, is the histogram of the number of voxels per sulcus based on the probabilistic mapping of sulci for both conditions (PCS and noPCS) and both hemispheres, using sulcal skeleton brain input images with the simple-3D-MHL network. The voxels are extracted after thresholding for highest explainability values, with thresholds of the highest 5, 10, and 20% intensity on the left hemisphere (high threshold: blue, medium threshold: orange, low threshold: grey). **b**, the total overlapping pattern learning results of the three most significant sulcal sub-region for the three different level of intensity thresholding. The acronyms for all sulci are defined in Supplementary Fig. 5., and NA: undefined. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

the left PCS has more associations with regional cortical thickness and sulcal depth than the right PCS, implying greater covariability of anatomical features in our input modalities with the presence of the PCS in the left hemisphere [6].

We developed an innovative XAI 3D-Framework to address the need for accurate, low-complexity global explanations in neuroimaging, where traditional 2D methods fall short in capturing the intricacies of 3D representational spaces. Designed for the binary classification of PCS presence, our framework provides robust, faithful global explanations that outperform GradCam and SHAP in faithfulness. The shift from traditional 2D explainability to a 3D volumetric approach represents a key strength of our framework, providing anatomically grounded insight that 2D methods cannot capture. Key novelties include the integration of statistical features (Shape) with reduced dimensionality information, ensuring explanations reflect both model learning and cohort-specific variability, and the combined use of GradCam and SHAP to reduce inter-method variability and enhance reliability. Additionally, the shift from local, subject-level explanations to global, population-wide interpretations enables identification of stable patterns of anatomical relevance, thereby enhancing model trustworthiness and biological plausibility. This multi-method framework sets a new standard for explainable AI in neuroimaging, offering actionable insights for complex tasks like cortical morphology analysis. Our global explanations surpassed those produced by GradCam and SHAP in terms of faithfulness, providing a reliable interpretation of the deep networks for this classification task.

The overall explainability outputs cover wide regions of the brain, but we can notice some repetitive patterns through the different pipelines and modalities. In particular, there is a repeat of the cingulate region, the posterior temporal region, and both the medial and inferior frontal cortices. This may reflect some neurodevelopmental intertwining of the macroscopical development of the PCS and these regions. Regarding the cingulate and medial frontal regions, this intertwining is self-explanatory as the PCS is located in the medial frontal region, directly adjacent to the cingulate region, and as such the developmental events leading to the formation of a PCS are very likely to affect these regions. The two other notable regions are the inferior frontal region and the posterior temporal region. In the fetus, the sulci matching these regions (namely the inferior frontal sulcus and the posterior superior temporal sulcus) have both been reported to start appearing at 26 weeks of gestational age (w GA), while the cingulate sulcus appears earlier (around 23w GA) and the "secondary cingular sulci", which encompass the PCS, appear later, at 31w GA [49]. This may point towards a time-window which is decisive to the development of the PCS, prior to its actual apparition.

In terms of functional interpretation, it is interesting to notice a striking similarity between the regions on which the AI mostly focuses and the regions which have been reported to show the most functional connectivity with the paracingulate region [15,43], and the anterior cingulate region [44]. These functional overlaps support the robustness and neurobiological validity of the regions highlighted by our global XAI framework. These studies report relevant functional connectivity between the medial frontal lobe (including the PCS) and the temporal region (including a focus on the posterior superior temporal region), the inferior frontal region, and the medial parietal region, which are all regions showing particular interest in the present work. Both [15,43] investigate the relation between the presence of a PCS and the related functional connectivity in these regions, and obtain more focused results (respectively in the cerebellum and superior anterior temporal region, or in the medial frontal region), but the important functional relationships between these regions support the relevance of the regions highlighted by our results.

The effectiveness of extracting generalized patterns using our proposed framework underscores the importance of incorporating data from multiple cohorts. This generalizability is a further strength, allowing our framework to be adapted for large-scale studies across diverse populations and clinical contexts. To this end, we plan to apply the framework to additional cohorts, such as BeneMin [50] and Biobank [51], to identify shared patterns in the classification of PCS presence and absence. Combining XAI techniques with dimensionality reduction methods may further reveal overlapping aspects of the data. Advanced approaches, such as t-SNE for non-linear dimensionality reduction, could also provide deeper insights into these relationships. Additionally, we aim to extend the framework's application to other neurological conditions and classification tasks, including schizophrenia and bipolar disorder, by leveraging external datasets and improving interpretability techniques. Future work will also investigate uncertainty quantification techniques, such as test-time dropout and ensemble modeling, to assess the robustness and confidence of both model predictions and explanation maps. This would further strengthen the clinical interpretability and reliability of our proposed framework.

Despite its strengths, the present study has a few limitations. First, all results are based on the TOP-OSLO dataset, and broader generalization requires external validation on independent cohorts, as the current findings may pose a risk of overfitting. Future validation using external and larger public datasets such as the Human Connectome Project [52] or UK Biobank [51] will be essential. Second, the framework currently assumes access to reliable manual annotations of PCS, which may not always be available. Inter-rater variability and the scalability of manual labeling protocols present practical challenges that merit further investigation. Third, while we use SHAP, GradCAM and Integrated Gradients due to their complementarity and strong prior use in neuroimaging explainability, we acknowledge that other explainability methods, such as Layer-wise Relevance Propagation (LRP) or Local Interpretable Model-agnostic Explanations (LIME) are not included in our framework. Future extensions could explore their integration for additional robustness.

While our current implementation is voxel-based, we acknowledge that surface-based learning may provide a more natural and anatomically meaningful representation of sulcal morphology, including the PCS. Sulci are inherently defined on the cortical surface, and leveraging surface-based methods, such as mesh-based CNNs or spherical mapping techniques, could lead to finer and more topology-consistent explanations. Although volumetric modeling allows us to maintain compatibility with many existing clinical pipelines and explainability tools, future extensions of our framework could explore the incorporation of surface-based architectures to improve both classification accuracy and interpretability. We also acknowledge that while we evaluated three CNN/MHL models, extending the framework to additional architectures (e.g., transformers, graph neural networks; [53,54]) remains an important direction. Moreover, the computational complexity of MHL grows exponentially with voxel input size; therefore, efficient

attention-based alternatives such as Performer [55] should be considered. Although GradCAM provides only indirect feature localization, we mitigate this limitation by combining GradCAM with SHAP and PCA, which produces more robust and interpretable explanations. In addition, AI-based segmentation methods could be leveraged in future extensions by building on the existing TOP-OSLO annotations, while surface-based learning approaches [56,57] also represent a promising direction for further development.

In summary, this study introduces a novel, integrated, and scalable 3D explainability framework that bridges methodological gaps in neuroimaging AI and lays a foundation for a systematic, biologically grounded exploration of sulcal variability. This study establishes a foundation for systematic exploration of sulcal variability through deep learning, with the potential to advance our understanding of cognitive and functional variability as well as pathological changes.

6. Conclusion

In this study, we classified the presence or absence of the paracingulate sulcus (PCS) in a diverse cohort of 596 structural MRIs using 3D deep convolutional neural networks and attention mechanisms. To address the lack of robust global explainability methods for 3D neuroimaging data, we developed an innovative XAI 3D-Framework. This framework provides accurate and low-complexity global explanations for PCS detection by integrating statistical features (Shape) with XAI methods (GradCam and SHAP) alongside reduced dimensionality information, ensuring that the explanations capture both model learning and cohort-specific variability. Furthermore, the combined application of GradCam and SHAP mitigates inter-method variability, thereby enhancing the reliability and robustness of the explanations.

Our framework outperformed established methods like GradCam and SHAP in faithfulness, enabling the robust identification of subregions critical for decision-making through a fusion of global explanations and statistical features. Key patterns identified include a focus on the posterior temporal and internal parietal regions on the sulcal skeleton, and on the cingulate region and thalamus when analyzing the grey-white surface. These findings indicate potential co-variation between these structures, likely underpinned by shared genetic or developmental mechanisms. Such insights hold significant implications for both neurodevelopmental and pathological research, providing a foundational framework for guiding future investigative trajectories.

Our work advances both deep learning and neuroscience by enabling automated, unbiased annotations and delivering unprecedented insights into sulcal variability and its functional or pathological relevance. The XAI 3D-Framework sets the stage for broader applications in medical imaging and other complex computer vision tasks, providing a foundation for comprehensive exploration of neuroanatomy and developmental mechanisms.

CRediT authorship contribution statement

Michail Mamalakis: Writing – review & editing, Writing – original draft, Visualization, Validation, Supervision, Software, Methodology, Investigation, Formal analysis, Conceptualization. Héloïse de Vareilles: Writing – original draft, Visualization, Data curation. Atheer Al-Manea: Writing – review & editing, Data curation. Samantha C. Mitchell: Writing – review & editing, Data curation. Ingrid Agartz: Writing – review & editing, Validation, Data curation. Lynn Egeland Mørch-Johnsen: Writing – review & editing, Validation, Data curation. Jane Garrison: Writing – review & editing, Validation. Jon Simons: Writing – review & editing, Validation. Pietro Lio: Writing – review & editing, Validation, Resources, Investigation. Graham K. Murray: Writing – review & editing, Validation, Resources, Funding acquisition, Conceptualization.

Code availability

The code developed in this study is written in the Python programming language using pytorch, Keras, tensorflow (Python) libraries. For training and testing of deep learning networks, we have used an NVIDIA cluster with 4 GPUs and 64 GB RAM memory. The code is publicly available in https://github.com/ece7048/3Dsulci.

Funding

This study was supported by funding from the Medical Research Council, grant number: MR/W020025/1.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Graham K Murray reports financial support was provided by UK Research and Innovation Medical Research Council. Graham K Murray reports a relationship with ieso Digital health that includes: consulting or advisory. GKM consults for ieso digital health. All other authors declare that they have no competing interests. If there are other authors, they declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

All research at the Department of Psychiatry in the University of Cambridge is supported by the NIHR Cambridge Biomedical Research Centre (NIHR203312) and the NIHR Applied Research Collaboration East of England. The views expressed are those of the author(s) and not necessarily those of the NIHR or the Department of Health and Social Care. This study was supported by funding from the Medical Research Council, grant number: MR/W020025/1. We acknowledge the use of the facilities of the Research Computing Services (RCS) of University of Cambridge, UK. GKM consults for ieso digital health. All other authors declare that they have no competing interests.

Appendix A. Supplementary data

Supplementary material related to this article can be found online at https://doi.org/10.1016/j.artmed.2025.103286.

Data availability

This study used the dataset of the TOP-OSLO cohort [12] which can be obtained from University of Oslo upon request, subject to a data transfer agreement.

References

- [1] Mamalakis M, de Vareilles H, Murray GK, Lio P, Suckling J. The explanation necessity for healthcare ai. In: 2025 IEEE Symposium on Trustworthy, Explainable and Responsible Computational Intelligence (CITREX Companion). 2025, p. 1–5. http://dx.doi.org/10.1109/CITREXCompanion65208.2025.10981502.
- [2] Mitchell SC, De Vareilles H, Garrison JR, Al-Manea A, Suckling J, Murray GK, Simons JS. Paracingulate sulcus measurement protocol V2. Apollo - Univ Camb Repos 2023. http://dx.doi.org/10.17863/CAM.102040, URL: https://www.repository.cam.ac.uk/handle/1810/358381.
- [3] Cachia A, Borst G, Tissier C, Fisher C, Plaze M, Gay O, Rivière D, Gogtay N, Giedd J, Mangin J-F, Houdé O, Raznahan A. Longitudinal stability of the folding pattern of the anterior cingulate cortex during development. Dev Cogn Neurosci 2016;19:122–7. http://dx.doi.org/10.1016/j.dcn.2016.02.011, URL: https://linkinghub.elsevier.com/retrieve/pii/S1878929315300943.

- [4] Borne L, Rivière D, Mancip M, Mangin J-F. Automatic labeling of cortical sulci using patch- or CNN-based segmentation techniques combined with bottom-up geometric constraints. Med Image Anal 2020-05;62:101651. http://dx.doi.org/10.1016/j.media.2020.101651, URL: https://linkinghub.elsevier.com/retrieve/nii/S1361841520300189
- [5] Jiang X, Zhang T, Zhang S, Kendrick KM, Liu T. Fundamental functional differences between gyri and sulci: implications for brain function, cognition, and behavior. Psychoradiology 2021;1(1):23–41. http://dx.doi.org/10.1093/psyrad/ kkab002, URL: https://academic.oup.com/psyrad/article/1/1/23/6187507.
- [6] Fornito A, Yücel M, Wood SJ, Proffitt T, McGorry PD, Velakoulis D, Pantelis C. Morphology of the paracingulate sulcus and executive cognition in schizophrenia. Schizophr Res 2006;88(1):192–7. http://dx.doi.org/10.1016/j.schres.2006.034, URL: https://linkinghub.elsevier.com/retrieve/pii/S0920996406003021.
- [7] Garrison JR, Fernyhough C, McCarthy-Jones S, Haggard M, The Australian Schizophrenia Research Bank, Simons JS. Paracingulate sulcus morphology is associated with hallucinations in the human brain. Nat Commun 2015;6(1):8956. http://dx.doi.org/10.1038/ncomms9956, URL: http://www.nature.com/articles/ ncomms9956.
- [8] Gay O, Plaze M, Oppenheim C, Gaillard R, Olié J-P, Krebs M-O, Cachia A. Cognitive control deficit in patients with first-episode schizophrenia is associated with complex deviations of early brain development. J Psychiatry Neurosci 2021;42(2):87–94. http://dx.doi.org/10.1503/jpn.150267, URL: http://www.jpn.ca/lookup/doi/10.1503/jpn.150267.
- [9] Ćurčić-Blake B, de Vries A, Renken RJ, Marsman JBC, Garrison J, Hugdahl K, Aleman A. Paracingulate sulcus length and cortical thickness in schizophrenia patients with and without a lifetime history of auditory hallucinations. Schizophr Bull 2023;49(Supplement_1):S48–57. http://dx.doi.org/10.1093/schbul/sbac072.
- [10] Simons JS, Garrison JR, Johnson MK. Brain mechanisms of reality monitoring. Trends Cogn Sci 2017-06;21(6):462–73. http://dx.doi.org/10.1016/j.tics.2017. 03.012, URL: https://linkinghub.elsevier.com/retrieve/pii/S1364661317300554.
- [11] Mamalakis M, Mamalakis A, Agartz I, Mørch-Johnsen LE, Murray GK, Suckling J, Lio P. Solving the enigma: enhancing faithfulness and comprehensibility in explanations of deep networks. AI Open 2025;6:70–81. http://dx.doi.org/10. 1016/j.ajopen.2025.02.001.
- [12] Mørch-Johnsen L, Nesvåg R, Jørgensen KN, Lange EH, Hartberg CB, Haukvik UK, Kompus K, Westerhausen R, Osnes K, Andreassen OA, Melle I, Hugdahl K, Agartz I. Auditory cortex characteristics in schizophrenia: Associations with auditory hallucinations. Schizophr Bull 2017;43(1):75–83. http://dx.doi.org/10.1093/schbul/sbw130, URL: https://academic.oup.com/schizophreniabulletin/article/43/1/75/2503785.
- [13] Fonov V, Evans AC, Botteron K, Almli CR, McKinstry RC, Collins DL. Unbiased average age-appropriate atlases for pediatric studies. NeuroImage 2011;54(1):313–27. http://dx.doi.org/10.1016/j.neuroimage.2010.07.033, URL: https://www.sciencedirect.com/science/article/pii/S1053811910010062.
- [14] Fonov V, Evans A, McKinstry R, Almli C, Collins D. Unbiased nonlinear average age-appropriate brain templates from birth to adulthood. NeuroImage 2009;47:S102. http://dx.doi.org/10.1016/S1053-8119(09)70884-5, URL: https://www.sciencedirect.com/science/article/pii/S1053811909708845. Organization for Human Brain Mapping 2009 Annual Meeting.
- [15] Fedeli D, Del Maschio N, Caprioglio C, Sulpizio S, Abutalebi J. Sulcal pattern variability and dorsal anterior cingulate cortex functional connectivity across adult age. Brain Connect 2022;10(6):267–78. http://dx.doi.org/10.1089/brain. 2020.0751, URL: https://www.liebertpub.com/doi/10.1089/brain.2020.0751.
- [16] Rollins CPE, Garrison JR, Arribas M, Seyedsalehi A, Li Z, Chan RCK, Yang J, Wang D, Liò P, Yan C, Yi Z-h, Cachia A, Upthegrove R, Deakin B, Simons JS, Murray GK, Suckling J. Evidence in cortical folding patterns for prenatal predispositions to hallucinations in schizophrenia. Transl Psychiatry 2020;10(1):387. http://dx.doi.org/10.1038/s41398-020-01075-y, URL: http://www.nature.com/articles/s41398-020-01075-y.
- [17] Mangin J-F, Perrot M, Operto G, Cachia A, Fischer C, Lefèvre J, Rivière D. Sulcus identification and labeling. In: Brain mapping. Elsevier; 2015, p. 365–71. http://dx.doi.org/10.1016/B978-0-12-397025-1.00307-9, URL: https://linkinghub.elsevier.com/retrieve/pii/B9780123970251003079.
- [18] Yang J, Wang D, Rollins C, Leming M, Liò P, Suckling J, Murray G, Garrison J, Cachia A. Volumetric segmentation and characterisation of the paracingulate sulcus on MRI scans. 2019-11-29, http://dx.doi.org/10.1101/859496, URL: http://biorxiv.org/lookup/doi/10.1101/859496.
- [19] Samek W, Montavon G, Lapuschkin S, Anders CJ, Müller K-R. Explaining deep neural networks and beyond: A review of methods and applications. Proc IEEE 2021;109(3):247–78. http://dx.doi.org/10.1109/JPROC.2021.3060483.
- [20] van der Velden BH, Kuijf HJ, Gilhuijs KG, Viergever MA. Explainable artificial intelligence (XAI) in deep learning-based medical image analysis. Med Image Anal 2022;79:102470. http://dx.doi.org/10.1016/j.media.2022.102470, URL: https://www.sciencedirect.com/science/article/pii/S1361841522001177.
- [21] Quellec G, Al Hajj H, Lamard M, Conze P-H, Massin P, Cochener B. Explain: Explanatory artificial intelligence for diabetic retinopathy diagnosis. Med Image Anal 2021;72:102118. http://dx.doi.org/10.1016/j.media.2021.102118, URL: https://www.sciencedirect.com/science/article/pii/S136184152100164X.

- [22] Mamalakis M, Garg P, Nelson T, Lee J, Wild JM, Clayton RH. MA-SOCRATIS: An automatic pipeline for robust segmentation of the left ventricle and scar. Comput Med Imaging Graph 2021;93:101982. http://dx.doi.org/10.1016/j.compmedimag.2021.101982, URL: https://www.sciencedirect.com/science/article/pii/S0895611121001312.
- [23] van der Velden BHM. Explainable AI: current status and future potential. Eur Radiol 2023/08/17. http://dx.doi.org/10.1007/s00330-023-10121-4.
- [24] Singh A, Sengupta S, Lakshminarayanan V. Explainable deep learning models in medical image analysis. 2020, arXiv:2005.13799.
- [25] Bach S, Binder A, Montavon G, Klauschen F, Müller K-R, Samek W. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS One 2015;10(7):1–46. http://dx.doi.org/10.1371/journal.pone. 0130140.
- [26] Rajani NF, McCann B, Xiong C, Socher R. Explain yourself! leveraging language models for commonsense reasoning. 2019, arXiv:1906.02361.
- [27] Tjoa E, Guan C. A survey on explainable artificial intelligence (XAI): Toward medical XAI. IEEE Trans Neural Netw. Learn Syst 2020;1–21. http://dx.doi.org/ 10.1109/tnnls.2020.3027314.
- [28] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, Polosukhin I. Attention is all you need. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, Garnett R, editors. In: Advances in neural information processing systems, vol. 30, Curran Associates, Inc.; 2017.
- [29] Selvaraju RR, Cogswell M, Das A, Vedantam R, Parikh D, Batra D. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE international conference on computer vision. 2017.
- [30] Lundberg S, Lee S-I. A unified approach to interpreting model predictions. 2017, http://dx.doi.org/10.48550/ARXIV.1705.07874, URL: https://arxiv.org/abs/1705.07874.
- [31] Pearson K. LIII. On lines and planes of closest fit to systems of points in space. Lond Edinb Dublin Philos Mag J. Sci 1901;2(11):559–72. http://dx.doi.org/10. 1080/14786440109462720.
- [32] Cointepas Y, Mangin J-F, Garnero L, Poline J-B, Benali H. BrainVISA: Software platform for visualization and analysis of multi-modality brain data. NeuroImage 2001;13(6, Supplement):98. http://dx.doi.org/10.1016/S1053-8119(01)91441-7, URL: https://www.sciencedirect.com/science/article/pii/S1053811901914417. Originally published as Volume 13, Number 6, Part 2.
- [33] Rivière D, Mangin J-F, Papadopoulos-Orfanos D, Martinez J-M, Frouin V, Régis J. Automatic recognition of cortical sulci of the human brain using a congregation of neural networks. Med Image Anal 2002-06;6(2):77–92. http: //dx.doi.org/10.1016/S1361-8415(02)00052-X, URL: https://linkinghub.elsevier. com/retrieve/pii/S136184150200052X.
- [34] Kingma DP, Ba J. Adam: A method for stochastic optimization. 2014, http://dx. doi.org/10.48550/ARXIV.1412.6980, URL: https://arxiv.org/abs/1412.6980.
 [35] Bell AJ, Sejnowski TJ. The "independent components" of natural scenes
- [35] Bell AJ, Sejnowski TJ. The "independent components" of natural scenes are edge filters. Vis Res 1997;37(23):3327–38. http://dx.doi.org/10.1016/ S0042-6989(97)00121-1, URL: https://www.sciencedirect.com/science/article/ pii/S0042698997001211.
- [36] Yeh CK, Hsieh CY, Suggala AS, Inouye DI, Ravikumar P. On the (in)fidelity and sensitivity for explanations. 2019, http://dx.doi.org/10.48550/ARXIV.1901. 09392, URL: https://arxiv.org/abs/1901.09392.
- [37] Bhatt U, Weller A, Moura JMF. Evaluating and aggregating feature-based model explanations. 2020, http://dx.doi.org/10.48550/ARXIV.2005.00631, URL: https://arxiv.org/abs/2005.00631.
- [38] Hedström A, Weber L, Krakowczyk D, Bareeva D, Motzkus F, Samek W, Lapuschkin S, Höhne MMM. Quantus: An explainable AI toolkit for responsible evaluation of neural network explanations and beyond. J Mach Learn Res 2023;24(34):1–11, URL: http://jmlr.org/papers/v24/22-0142.html.
- [39] Paus T, Tomaiuolo F, Otaky N, MacDonald D, Petrides M, Atlas J, Morris R, Evans AC. Human cingulate and paracingulate sulci: Pattern, variability, asymmetry, and probabilistic map. Cerebral Cortex 1996;6(2):207–14. http://dx.doi.org/10.1093/cercor/6.2.207, URL: https://academic.oup.com/cercor/article-lookup/doi/10.1093/cercor/6.2.207.
- [40] Yücel M, Stuart GW, Maruff P, Velakoulis D, Crowe SF, Savage G, Pantelis C. Hemispheric and gender-related differences in the gross morphology of the anterior cingulate/paracingulate cortex in normal volunteers: An MRI morphometric study. Cerebral Cortex 2001;9.
- [41] Sundararajan M, Taly A, Yan Q. Axiomatic attribution for deep networks. In: International conference on machine learning. PMLR; 2017, p. 3319–28.
- [42] Perrot M, Rivière D, Mangin J-F. Cortical sulci recognition and spatial normalization. Med Image Anal 2011;15(4):529–50. http://dx.doi.org/10.1016/j.media.2011.02.008, URL: https://linkinghub.elsevier.com/retrieve/pii/S1361841511000302.
- [43] Harper L, Strandberg O, Spotorno N, Nilsson M, Lindberg O, Hansson O, Santillo AF. Structural and functional connectivity associations with anterior cingulate sulcal variability. Brain Struct Funct 2024-06-20. http://dx.doi.org/10.1007/s00429-024-02812-5, URL: https://link.springer.com/10.1007/s00429-024-02812-5.

- [44] Lopez-Persem A, Verhagen L, Amiez C, Petrides M, Sallet J. The human ventromedial prefrontal cortex: Sulcal morphology and its influence on functional organization. J Neurosci 2019-05-08;39(19):3627-39. http://dx.doi.org/10.1523/JNEUROSCI.2060-18.2019, URL: https://www.jneurosci.org/lookiny/doi/10.1523/JNEUROSCI.2060-18.2019
- [45] De Vareilles H, Rivière D, Pascucci M, Sun Z-Y, Fischer C, Leroy F, Tataranno M-L, Benders MJ, Dubois J, Mangin J-F. Exploring the emergence of morphological asymmetries around the brain's Sylvian fissure: a longitudinal study of shape variability in preterm infants. Cerebral Cortex 2023;bhac533. http://dx.doi.org/10.1093/cercor/bhac533, URL: https://academic.oup.com/cercor/advance-article/doi/10.1093/cercor/bhac533/7005629.
- Nunes A, Schnack HG, Ching CRK, Agartz I, Akudjedu TN, Alda M, Alnæs D, Alonso-Lana S. Bauer J. Baune BT. Bøen E. Bonnin CdM. Busatto GF. Canales-Rodríguez EJ, Cannon DM, Caseras X, Chaim-Avancini TM, Dannlowski U, Díaz-Zuluaga AM, Dietsche B, Doan NT, Duchesnay E, Elvsåshagen T, Emden D, Eyler LT, Fatjó-Vilas M, Favre P, Foley SF, Fullerton JM, Glahn DC, Goikolea JM, Grotegerd D, Hahn T, Henry C, Hibar DP, Houenou J, Howells FM, Jahanshad N, Kaufmann T, Kenney J, Kircher TTJ, Krug A, Lagerberg TV, Lenroot RK, López-Jaramillo C, Machado-Vieira R, Malt UF, McDonald C, Mitchell PB, Mwangi B, Nabulsi L, Opel N, Overs BJ, Pineda-Zapata JA, Pomarol-Clotet E, Redlich R, Roberts G, Rosa PG, Salvador R, Satterthwaite TD, Soares JC, Stein DJ, Temmingh HS, Trappenberg T, Uhlmann A, van Haren NEM, Vieta E, Westlye LT, Wolf DH, Yüksel D, Zanetti MV, Andreassen OA, Thompson PM, Hajek T, for the ENIGMA Bipolar Disorders Working Group. Using structural MRI to identify bipolar disorders -13 site machine learning study in 3020 individuals from the ENIGMA bipolar disorders working group. Mol Psychiatry 2020/09/01;25(9):2130-43. http://dx.doi.org/10.1038/s41380-018-0228-9.
- [47] Iniesta R, Hodgson K, Stahl D, Malki K, Maier W, Rietschel M, Mors O, Hauser J, Henigsberg N, Dernovsek MZ, Souery D, Dobson R, Aitchison KJ, Farmer A, McGuffin P, Lewis CM, Uher R. Antidepressant drug-specific prediction of depression treatment outcomes from genetic and clinical variables. Sci Rep 2018/04/03;8(1):5530. http://dx.doi.org/10.1038/s41598-018-23584-z.
- [48] Hosmer DW, Lemeshow S. Applied logistic regression (Wiley Series in probability and statistics). 2nd ed. Wiley-Interscience Publication; 2000, URL: http://www.amazon.com/Applied-logistic-regression-probability-statistics/dp/0471356328%3 FSubscriptionId%3D192BW6DQ43CK9FN0ZGG2%26tag%3Dws%26linkCode%3Dxm2%26camp%3D2025%26creative%3D165953%26creativeASIN%3D04713563
- [49] Garel C, Chantrel E, Brisse H, Elmaleh M, Luton D, Oury J-F, Sebag G, Hassan M. Fetal cerebral cortex: Normal gestational landmarks identified using prenatal MR imaging. AJNR Am J Neuroradiol 2001;6.
- [50] Deakin B, Suckling J, Barnes TRE, Byrne K, Chaudhry IB, Dazzan P, Drake RJ, Giordano A, Husain N, Jones PB, Joyce E, Knox E, Krynicki C, Lawrie SM, Lewis S, Lisiecka-Ford DM, Nikkheslat N, Pariante CM, Smallman R, Watson A, Williams SCR, Upthegrove R, Dunn G. The benefit of minocycline on negative symptoms of schizophrenia in patients with recent-onset psychosis (BeneMin): a randomised, double-blind, placebo-controlled trial. Lancet Psychiatry 2018;5(11):885–94. http://dx.doi.org/10.1016/S2215-0366(18)30345-6, URL: https://www.sciencedirect.com/science/article/pii/S2215036618303456.
- [51] Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, Downey P, Elliott P, Green J, Landray M, Liu B, Matthews P, Ong G, Pell J, Silman A, Young A, Sprosen T, Peakman T, Collins R. UK biobank: An open access resource for identifying the causes of a wide range of complex diseases of middle and old age. PLOS Med 2015;12(3):1–10. http://dx.doi.org/10.1371/journal.pmed.1001779.
- [52] Van Essen DC, Smith SM, Barch DM, Behrens TE, Yacoub E, Ugurbil K. The WU-Minn human connectome project: an overview. NeuroImage 2013;80:62–79. http://dx.doi.org/10.1016/j.neuroimage.2013.05.041.
- [53] Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In: International conference on learning representations. 2021.
- [54] Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks. In: International conference on learning representations. 2017.
- [55] Choromanski K, Likhosherstov V, Dohan D, Song X, Gane A, Sarlos T, et al. Rethinking attention with performers. In: International conference on learning representations. 2021.
- [56] Fischl B. Freesurfer. NeuroImage 2012;62(2):774–81. http://dx.doi.org/10.1016/j.neuroimage.2012.01.021.
- [57] Glasser MF, Sotiropoulos SN, Wilson JA, Coalson TS, Fischl B, Andersson JL, Xu J, Jbabdi S, Webster M, Polimeni JR, Van Essen DC, Jenkinson M. The minimal preprocessing pipelines for the human connectome project. NeuroImage 2013;80:105–24. http://dx.doi.org/10.1016/j.neuroimage.2013.04.127.