

---

# UNCOVERING NEUROIMAGING BIOMARKERS OF BRAIN TUMOR SURGERY WITH AI-DRIVEN METHODS

---

**Carmen Jiménez-Mesa**

Department of Communication Engineering  
University of Málaga  
Spain

**Yizhou Wan**

Department of Clinical Neurosciences  
University of Cambridge  
United Kingdom

**Guilio Sansone**

Department of Neuroscience  
University of Padova  
Italy

**Francisco J. Martínez-Murcia**

Department of Signal Theory, Telematics  
and Communications, University of Granada  
Spain

**Javier Ramirez**

Department of Signal Theory, Telematics  
and Communications, University of Granada  
Spain

**Pietro Lio**

Department of Computer Science and Technology  
University of Cambridge  
United Kingdom

**Juan M. Gorriz**

Department of Signal Theory, Telematics  
and Communications, University of Granada  
Spain

**Stephen J. Price**

Department of Clinical Neurosciences  
University of Cambridge  
United Kingdom

**John Suckling**

Department of Psychiatry  
University of Cambridge  
Cambridge and Peterborough NHS Foundation Trust  
United Kingdom

**Michail Mamalakis**

Department of Psychiatry  
Department of Computer Science and Technology  
University of Cambridge  
United Kingdom  
mm2703@cam.ac.uk

July 8, 2025

## ABSTRACT

Brain tumor resection is a complex procedure with significant implications for patient survival and quality of life. Predictions of patient outcomes provide clinicians and patients the opportunity to select the most suitable onco-functional balance. In this study, global features derived from structural magnetic resonance imaging in a clinical dataset of 49 pre- and post-surgery patients identified potential biomarkers associated with survival outcomes. We propose a framework that integrates Explainable AI (XAI) with neuroimaging-based feature engineering for survival assessment, offering guidance for surgical decision-making. In this study, we introduce a global explanation optimizer that refines survival-related feature attribution in deep learning models, enhancing interpretability and reliability. Our findings suggest that survival is influenced by alterations in regions associated with cognitive and sensory functions, indicating the importance of preserving areas involved in decision-making and emotional regulation during surgery to improve outcomes. The global explanation optimizer improves both fidelity and comprehensibility of explanations compared to state-of-the-art XAI methods. It effectively identifies survival-related variability, underscoring its relevance in precision medicine for brain tumor treatment.

**Keywords** Brain Tumor · explainable AI · feature engineering · Machine Learning · PCA

## 1 Introduction

Gliomas, the most frequent primary brain tumors, vary in aggressiveness, prognosis, and histopathology. Their treatment often involves surgical resection, followed by radiotherapy and chemotherapy. The extent of resection significantly affects survival, with surgery needing to balance tumor removal and brain function preservation [1]; known as onco-functional balance. Post-surgical brain reorganization plays a key role in recovery, but the mechanisms behind these changes are not fully understood [2]. Accurate assessment of these changes is crucial for improving patient outcomes and rehabilitation strategies.

Structural Magnetic Resonance Imaging (sMRI) may provide insights into brain reorganization, but its high dimensionality poses challenges for analysis. Machine learning (ML) techniques have demonstrated superior performance over traditional methods in predicting postoperative complications and inpatient length of stay for patients with brain tumours [3–5]. They have been effectively applied to tumor segmentation, classification, and pre- and post-operative MRI analysis [6–9], as well as in predicting neurosurgical outcomes, including survival and recurrence [10, 11]. The application of ML-based latent space representations, such as Principal Component Analysis (PCA) [12], has proven valuable in neuroimaging analysis by aiding in the interpretation of structural variations. However, no studies have employed latent spaces to assess structural changes in brain tumors before and after surgery, which this work addresses.

EXplainable Artificial Intelligence (XAI) techniques have effectively identified patterns in hypothesis-driven research [13–15]. Local XAI methods provide interpretations of individual model predictions, whereas global methods offer cohort-level insights into the model’s overall decision-making process, thereby enhancing our understanding of its behavior across populations. By combining these methodologies, we can improve our understanding of brain tumor-related structural reorganization and facilitate personalized treatment approaches. However, different local techniques lead to varying explanations, increasing uncertainty rather than trust [16] in AI-driven solutions. The same problem is observed in global explanations. Part of this work aims to address the problem of inter-method variability in global explanations by proposing a global XAI optimizer.

This study presents a novel framework combining latent-PCA spaces with XAI to analyze variability in brain tumors among patients who survived and those who did not, following the surgery. We introduce a global explanation optimizer to improve the accuracy and clarity of survival-related biomarkers, eliminating the inter-method variability in explanations. Additionally, we investigate the relationship between significant brain regions and survival outcomes to provide guidance that may reduce surgical fatality rates. Given the limited availability of real-world longitudinal data in neuro-oncology, particularly covering both pre- and post-surgical stages, data used in this study offers a rare opportunity to investigate the impact of surgery on brain function and survival. We focus on two factors: time (pre- and post-surgery) and survival (longer-term and shorter-term). The key contributions of our study are: (i) a global explanation optimizer to identify survival-related biomarkers, enhancing both optimal accuracy and comprehensibility of the explanations, (ii) a framework that combines XAI with neuroimaging-based feature engineering to offer novel insights into the relationship between sMRI features and survival risk, and (iii) providing suggestions and potential guidance for surgical decision-making to improve survival outcomes for patients with brain tumors.

## 2 The proposed framework

This work combines latent space feature engineering and XAI methods to identify biomarkers related to surgical outcomes. A summary of the implemented framework is presented in Figure 1. First, sMRI images are processed, including spatial alignment, skull stripping, and masking of tumor regions to ensure consistency across subjects and facilitate subsequent analyses. The main part consists of two phases: feature engineering through dimensionality reduction and a global explanation optimizer integrated with DL networks and XAI methods.

### 2.1 Phase I: Feature engineering based on dimensionality reduction

We utilized PCA to extract the most relevant patterns of variation across the brain and tumor cohorts, based on four groups derived from the time/survival factors. Two main approaches were used: first, analyzing PCA component variability across groups, and second, quantifying variability across PCA components (see Figure 1):

#### 2.1.1 First PCA component variability across groups

The first PCA component, representing the highest variance, was compared across groups (shorter-term and longer-term survivals) to identify dominant differentiation patterns. To do so, spatial variability between the two conditions (pre- and

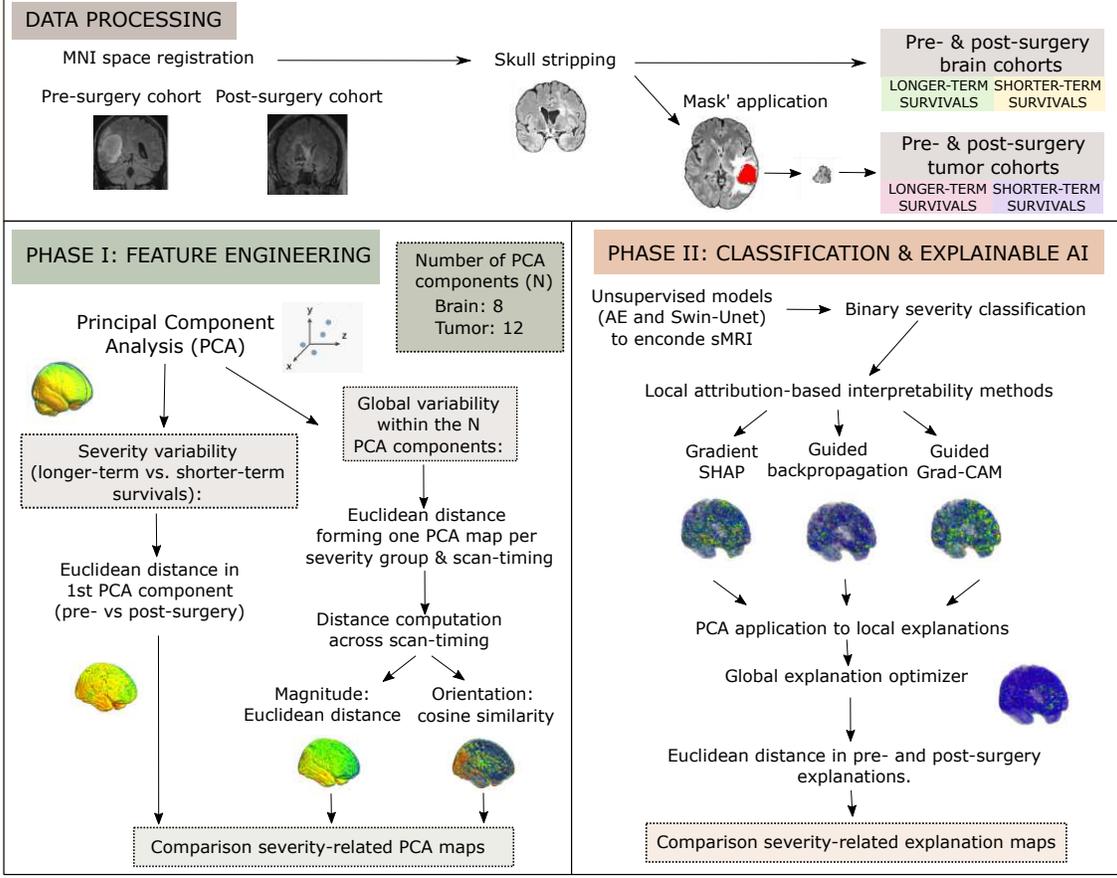


Figure 1: Framework proposed. Imaging processing is followed by two different analysis. Phase I: in the feature engineering study PCA components are extracted from the different cohorts analyzing variability between and within groups. Phase II: to enhance interpretability and robustness of the outcomes, an analysis is conducted by means of a binary classifier where three different XAI techniques are applied: Gradient-SHAP, Guided-Backpropagation and Guided-GradCAM. Their outputs serve as input to a global explanation optimizer, generating a map of the most relevant global patterns for each severe condition.

post-surgery) was compute using voxel-wise Euclidean distance. The Euclidean distance between two PCA-transformed representations,  $\mathbf{p}_A$  and  $\mathbf{p}_B$ , can be mathematically described as:

$$d_E = \|\mathbf{p}_A - \mathbf{p}_B\|_2 = \sqrt{\sum_{i=1}^K (p_{A,i} - p_{B,i})^2} \quad (1)$$

where a larger distance indicates greater structural change. The comparison of these variability maps allow to assess whether the PCA has captured meaningful group distinctions.

### 2.1.2 Variability quantification across PCA Components

Local variability maps were generated by computing voxel-wise Euclidean distances across the  $k$  PCA components of each of the four subgroups, summarizing the total magnitude of variations captured by PCA and quantifying regional brain variability. This approach enabled the estimation of global variability within groups (pre- vs. post surgery) by analyzing both magnitude (Euclidean distance) and orientation (cosine similarity) from the local maps. The cosine similarity can be mathematically described as:

$$S_{\cos} = \frac{\mathbf{p}_A \cdot \mathbf{p}_B}{\|\mathbf{p}_A\| \|\mathbf{p}_B\|} \quad (2)$$

Once this is done, the global maps of shorter-term and longer-term survivals can be compared to assess spatial variability.

## 2.2 Phase II: Feature identification based on cohort-level explanations, integrated with deep learning networks and tailored to survival classification

Figure 1 presents the explainable AI framework used to identify global (cohort-level) patterns in survival outcomes after brain tumor surgery. An unsupervised learning approach was used, in which an encoder-decoder architecture learns the distribution of sMRI data. A binary classification model was then trained and validated to distinguish between subjects with shorter and longer survival. The classification task explored various combinations of frozen and fine-tuned encoder layers and an ablation study of different networks for prediction from the unsupervised learning stage to enhance generalization and improve accuracy. Finally, cohort-level explanations were identified to the survival classification task, incorporating our proposed global explanation optimizer to improve the clarity and consistency of the global explanation patterns.

### 2.2.1 Unsupervised learning of structural MRI

Two deep learning architectures: a convolutional autoencoder (AE) comprising three encoder and three decoder blocks, and the Swin-Unet [17] were trained in an unsupervised setting to reconstruct entire pre- and post-operative 3D structural MRI scans. An ablation study was conducted to compare the effect of two different cohort training strategies aiming to determine which approach enhances the networks’ reconstruction performance and generalization capability. A detailed explanation is given in Section 3.2.

### 2.2.2 Survival classification of structural MRI

For the survival classification task, we used the encoder components of the previously trained unsupervised AE and Swin-Unet models. An ablation study was conducted to evaluate different output layer configurations: (i) a three-layer multilayer perceptron (MLP) for binary classification, and (ii) a cross-attention (Attention) mechanism applied to the four encoder stages of the Swin-Unet [17]. We explored three training strategies: (1) freezing the encoder (freeze) and training only the output layer, (2) fine-tuning the encoder by unfreezing its weights (unfreeze), and (3) re-initializing and jointly training both the encoder and the output layer (full training).

### 2.2.3 Global explanations of structural MRI

To enhance interpretability in the survival classification task, we used six local attribution-based methods: Guided Backpropagation [18], Guided GradCam [19], and Gradient Shap [20], Input  $\times$  Gradient [21], Integrated Gradients [21], and Kernel SHAP [20]. The goal was to uncover global patterns distinguishing between longer-term and shorter-term survivals outcomes by generating global explanations from pre- and post-surgery sMRI. To achieve this, we first estimated the global (cohort-level) pre-surgery and post-surgery explanations using the six different local explanation methods. We then applied PCA to the local explanations generated by each of these XAI methods to obtain a globalized representation. Finally, Euclidean distances were used to quantify differences between the global pre- and post-surgery explanations. To assess the accuracy of these explanations, we evaluated sparseness [22] and faithfulness [23]. These explainability metrics were computed using the software developed by [24], a comprehensive toolkit designed to collect, organize, and assess various performance metrics proposed for XAI methods. We note that a zero baseline (“black”) and 20 random perturbations were used to compute the faithfulness score.

### 2.2.4 The proposed global explanation optimizer of structural MRI

To identify potential biomarkers, reduce inter-method global explanation variability, and extract actionable insights for improving surgical outcomes, we aimed to generate a global explanation for the binary survival classification task. To this end, we proposed a global explanation optimizer, building on the methodology introduced by [16] for optimizing explanation representations. Our framework follows the foundational design of the original approach, including a non-linear encoder-decoder architecture (Swin-Unet) and a multi-objective cost function. A key distinction in our implementation lies in the evaluation strategy: we assess the optimized global explanation by comparing it to the first principal component extracted via PCA on the sMRI data. This comparison enables quantitative assessment of structural relevance using the Structural Similarity Index Measure (SSIM).

We extracted the first three principal components via PCA from the total cohort saliency maps, generated using three of the six widely used attribution methods employed in this study: Guided Backpropagation, Guided Grad-CAM, and Gradient SHAP. These components, along with their weighted average—computed according to the procedure described in [16]—were used as four inputs to the proposed global explanation optimizer.

The cost function guiding the optimization integrates three key components: sparseness, as defined in [22]; faithfulness [23], to ensure consistency with model predictions; and similarity, to align the optimized explanation with a structural

representation. This composite objective supports the generation of explanations that are both interpretable and clinically meaningful.

The resulting SSIM score between the optimized global explanation and the first PCA component of the structural MRI inputs is reported as follows:

$$loss_{sim}(\mathbf{x}, \mathbf{y}) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)} \quad (3)$$

where  $\mathbf{x}$  represents the derived explanation by the global optimizer,  $\mathbf{y}$  denotes the first component of PCA of the structural MRI,  $\mu_x$  indicates the average of  $\mathbf{x}$ ,  $\sigma_x^2$  signifies the variance of  $\mathbf{x}$ ,  $\sigma_{xy}$  represents the covariance of  $\mathbf{x}$  and  $\mathbf{y}$ , and  $c_1$  and  $c_2$  are two parameters utilized to stabilize the division with a weak denominator [25]. The total loss function was given by:

$$loss_{total}(\mathbf{x}, \mathbf{y}) = l_1 \frac{1}{M_{faith}(f, g; \mathbf{x})} + l_2 M_{sparse}(f, g; \mathbf{x}) + l_3 loss_{sim}(\mathbf{x}, \mathbf{y}) \quad (4)$$

where  $M_{sparse}$ ,  $M_{faith}$  are the metrics for sparseness [22] and faithfulness [23], respectively and the  $g$  global explanation for the network  $f$ .

### 3 Experimental setup

#### 3.1 Dataset

The main dataset was from Addenbrooke’s Hospital (Cambridge, UK). The UK’s Research Authority provided ethical approval (ref:19/WM/0152) and data were anonymized before analyses. The dataset consists of 49 MRI T2-weighted scans acquired both before and after surgical resection of the tumour. These scans were spatially normalized to MNI space using SPM12 ([fil.ion.ucl.ac.uk/spm/](http://fil.ion.ucl.ac.uk/spm/)) and resampled to a  $1 \times 1 \times 1$  mm<sup>3</sup> resolution resulting in final image dimensions of  $157 \times 189 \times 156$  mm. Skull-stripping was performed using HD-BET [26]. Patients were categorized into two outcome groups: longer-term (32) and shorter-term (17) survivals. Most patients (42, 85%) had a glioblastoma, but there were also cases of astrocytoma (1), gliosarcoma (3) and others (3). The shorter-term survival group comprised patients who had died within 10 months after the postoperative scan. In contrast, the longer-term group included those who survived for more than 10 months. The UK’s Research Authority provided ethical approval (ref:19/WM/0152); data were anonymized before analyses

In Phase II, an additional dataset from the 2025 Brain Tumor Segmentation (BraTS) Glioma Challenge [27] was employed. Hereafter, we refer to this dataset as BraTS2025. This dataset comprises pre- and post-treatment T2-weighted MRI scans. We used a total of 1453 images (1251 pre-treatment and 202 post-treatment). These scans were used to train the unsupervised learning models (see 2.2.1). Demographic information was not provided for this dataset.

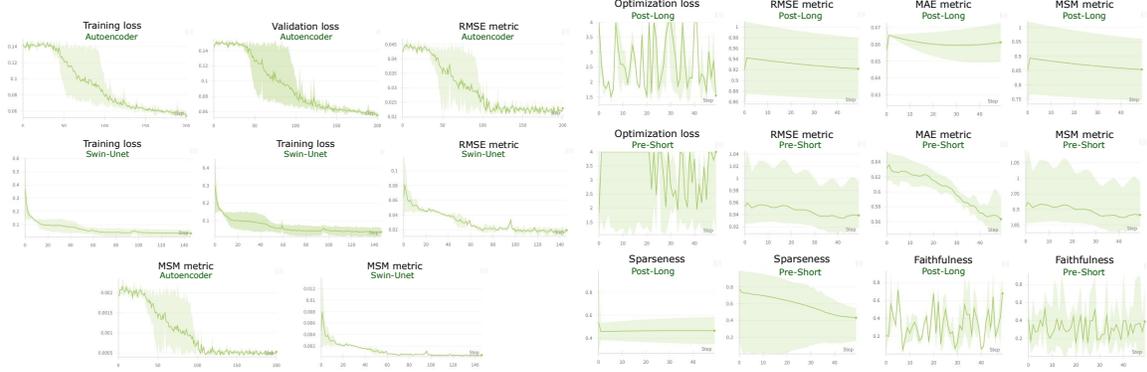
#### 3.2 Implementation details

For PCA computation, sMRI scans were vectorized and standardized with zero-mean, unit-variance scaling. No intensity normalisation was applied to the tumour masks due to their binary nature and spatial variability. PCA outcomes were normalized using min-max scaling to  $[0, 1]$ . Eight components were selected for brain images and 12 for tumor images based on cumulative variance with 8 components explaining over 80% of variance in both severity conditions. Tumor images required 12 components for similar variance.

To enhance reproducibility and facilitate result interpretation, the outcomes of these and subsequent analyses were mapped onto the Human Connectome Project (HCP) HCP-MMP1 atlas [28].

For the unsupervised learning task, a fixed-step learning rate ( $5 \times 10^{-4}$ ) and the Adam optimizer [29] were used to minimize a SSIM-based loss function [25] (see 3). The learning rate remained constant, with early stopping after 10 epochs of no improvement (max 200 epochs). Two cohort training strategies were evaluated: (i) using only the Addenbrooke’s Hospital dataset, and (ii) combining the Addenbrooke’s Hospital and BraTS2025 datasets. For the Addenbrooke’s Hospital dataset, the 96 available scans were randomly shuffled and divided into five folds for cross-validation (CV) across the entire cohort. In the combined dataset scenario, a 60/40 training/validation (1453 and 96 3D-MRI scans) split was employed.

For survival classification, sparse categorical cross-entropy was used as the loss function, optimized with Adam. The learning rate was constant for the first 100 epochs and then reduced by a factor of 0.1 every 100 epochs. Early stopping



(a) Hyperparameter learning rate tuning of unsupervised learning architectures. (b) Hyperparameter tuning of various combinations of cost functions ( $l_1$ ,  $l_2$ , and  $l_3$ ) for global explanation models in both post-surgery longer-term and pre-surgery shorter-term survival groups.

Figure 2: Examples of hyperparameter tuning results for (a) unsupervised learning, (b) global explanation models on structural MRI. Abbreviations: RMSE–Root Mean Squared Error, MSM–Mean Squared Magnitude, MAE–Mean Absolute Error.

was applied after 100 epochs of no improvement (max 400 epochs). A 5-fold CV was used. Both tasks employed data augmentation, including rotation ( $[-15^\circ, 15^\circ]$ ), width/height shift (up to 20 pixels), and intensity shift (up to 20%). Hyperparameter tuning tested learning rates:  $5 \times 10^{-2}$ ,  $5 \times 10^{-3}$ ,  $5 \times 10^{-4}$ , and  $5 \times 10^{-5}$  (see Figure 2a.). The XAI task used the Adam optimizer, but no data augmentation. The cost function was (4). For 3D tasks, training lasted up to 100 epochs, with early stopping after 10 epochs of no improvement beyond the first 50. Hyperparameter tuning tested the same learning rates as previously and various combinations of the  $l_1$ ,  $l_2$ , and  $l_3$  parameters in (4) with the best combination of parameters determined as  $l_1 = 0.4$ ,  $l_2 = 0.3$ ,  $l_3 = 0.3$  and a learning rate  $5 \times 10^{-5}$  (see Figure 2b.). Codes were implemented in Python using PyTorch and trained on one A100 GPU with 64 GB RAM. It will be publicly available on GitHub.

### 3.3 Explanation Quality Metrics

A critical component of this study is the evaluation of *how accurate and comprehensive an explanation is*. To this end, we focus on two essential metrics: faithfulness and complexity. One intuitive and widely adopted approach for assessing explanation quality is to examine how well it captures the behavior of a predictive model under input perturbations [30].

#### 3.3.1 Faithfulness Metric

Let  $f$  denote a deep neural network, and let  $\mathbf{x} \in \mathbb{R}^d$  represent an input with  $d$  features. We aim to assess whether the attribution scores—also known as feature importance scores—accurately reflect the impact of each feature on the model’s output.

Consider a subset  $S \subseteq \{1, 2, \dots, d\}$  of input features, and let  $\mathbf{x}_S$  denote the corresponding sub-vector of  $\mathbf{x}$ , with  $\mathbf{x}_S^f$  being the baseline (reference) values for those features. If  $g(f, \mathbf{x}) \in \mathbb{R}^d$  is the attribution vector provided by explanation method  $g$ , then the faithfulness is measured by the Pearson correlation between the sum of attributions for the features in  $S$  and the change in the model’s output when those features are set to baseline:

$$M_{\text{faith}}(f, g; \mathbf{x}) = \text{corr}_S \left( \sum_{i \in S} g(f, \mathbf{x})_i, f(\mathbf{x}) - f(\mathbf{x}[\mathbf{x}_S = \mathbf{x}_S^f]) \right) \quad (5)$$

where  $\mathbf{x}_F = \mathbf{x} \setminus \mathbf{x}_S$  denotes the unchanged features.

#### 3.3.2 Sparseness Metric

To quantify the complexity of an explanation, we evaluate the sparseness of the attribution vector. Sparseness indicates whether the explanation highlights only the most relevant features, which is desirable for interpretability.

We use the *Gini Index*, a well-established measure of inequality, to assess sparseness [31]. Given a non-negative vector  $\mathbf{v} \in \mathbb{R}_{\geq 0}^d$ , let  $v_{(k)}$  be the  $k$ -th smallest value after sorting. The Gini Index is defined as:

$$G(\mathbf{v}) = 1 - 2 \sum_{k=1}^d \frac{v_{(k)}}{\|\mathbf{v}\|_1} \cdot \left( \frac{d - k + 0.5}{d} \right), \quad (6)$$

where  $\|\mathbf{v}\|_1 = \sum_{i=1}^d v_i$  is the  $\ell_1$ -norm.

To measure the sparseness of an attribution vector  $\phi^{(k)}$ , we apply the Gini Index to the vector of its absolute values:

$$\text{Sparseness} \left( \phi^{(k)} \right) = G \left( \left| \phi^{(k)} \right| \right), \quad (7)$$

where  $\left| \phi^{(k)} \right| = \left( |\phi_1^{(k)}|, |\phi_2^{(k)}|, \dots, |\phi_d^{(k)}| \right)$ . Higher values indicate greater sparseness. A value of 1 implies that the attribution is entirely concentrated on a single feature, while 0 corresponds to equal attribution across all features.

## 4 Results

### 4.1 Structural patterns identified using feature engineering based on PCA

Once the PCA components (brain: 8, tumor: 12) were computed across groups, structural variability was quantified to explore spatial differences in tumor and brain patterns. The localization of the first PCA component in the tumor cohorts within the cerebral space is illustrated in the top right of Figure 3, revealing group-specific spatial distributions. To evaluate brain-wide structural changes, voxel-wise Euclidean distances were computed on the first PCA component, producing variability maps across groups (Figure 3, top left). The shorter-term survival group showed greater distances between pre- and post-surgery scans, suggesting more pronounced structural alterations. Moreover, this group exhibited higher spatial variability in the tumor PCA component, both before (grayscale) and after surgery (red), suggesting increased heterogeneity in tumor location and size.

To quantify global structural variability, we computed voxel-wise Euclidean distances across PCA components, generating variability maps that highlight key differences between groups. This approach allowed us to capture the overall magnitude of structural differences at each voxel, revealing patterns of brain alterations associated with disease progression. To ensure a robust characterization, we evaluated both the magnitude and orientation of variations in PCA space, comparing pre- and post-surgery subgroups to assess changes relative to disease severity. These maps are displayed in Figure 3 (*Global Variability Maps* section) and offered a global depiction of structural variability across the brain, highlighting areas where voxel-wise differences between pre- and post-surgery scans were most pronounced in each survival group.

To identify the most relevant brain regions, we used atlas-based segmentation and applied both intensity and volume criteria. For the Euclidean distance maps, a brain region was considered significant if it met two conditions: it contained at least one voxel above the 95th percentile (indicating a strong local effect), and at least 50% of its voxels exceeded the 80th percentile (reflecting a substantial spatial extent). For the cosine similarity analysis, we focused on regions with the lowest similarity values, as they reflect the greatest divergence in directionality of the PCA patterns. Specifically, we selected regions where the lowest voxel values fell below the 5th percentile, and applied a volume threshold of the 20th percentile to ensure spatial relevance.

Columns *Euclidean maps* and *Cosine maps* of Table 1 summarize the key brain regions identified through PCA-based feature engineering. These regions exhibit the greatest dissimilarity between pre- and post-surgery states in both the longer-term and shorter-term survival groups (*Brain regions* rows), as well as the largest changes observed within the tumour masks before and after surgery (*Surgical regions* row), reflecting differences between tumour locations and the surgical removal area.

### 4.2 Ablation Study of Unsupervised Pretraining and Fine-Tuning Strategies

We conducted an ablation study comparing two different cohort training strategies (see 2.2.1). Based on Table 2, the best validation results were achieved using the second cohort strategy, particularly with the Addenbrooke’s Hospital and BraTS2025 datasets. Specifically, the Swin-Unet model achieved the lowest error values across both strategies, with the best performance in the second strategy; an RMSE of 0.008 compared to 0.010, MSM of 0.001 in both cases, and MAE of 0.005 compared to 0.009. These results highlight the superiority of the strategy involving both the Addenbrooke’s Hospital and BraTS2025 datasets over the strategy using only the Addenbrooke’s Hospital dataset.

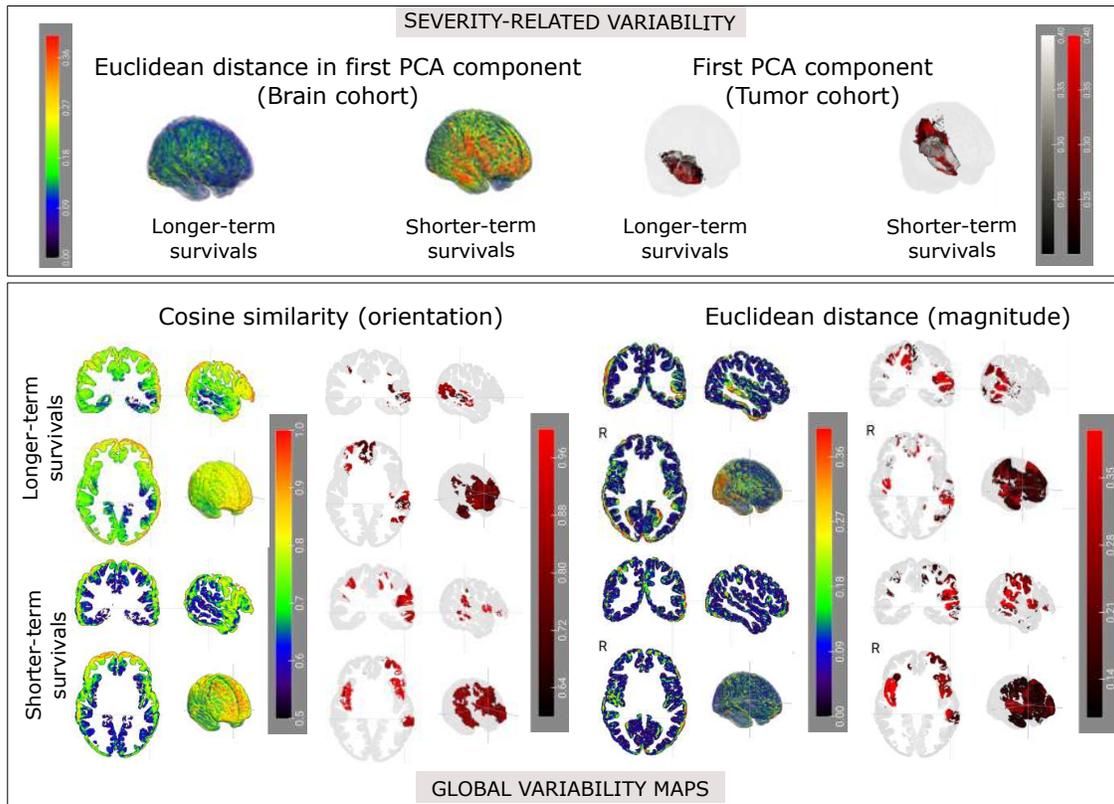


Figure 3: Variability across severity conditions. Top: First PCA components analysis showing atlas-based Euclidean distances in the new space between pre- and post-surgery subgroups in the brain cohort. The first PCA distribution is presented for the tumor cohort before surgery (grayscale) and after surgery (red). Bottom: Global variability in PCA components, displaying magnitude and orientation results for the comparison between pre- and post-surgery groups for both brain and tumor cohorts.

Table 1: Key brain regions with significant 3D volume differences pre- vs. post-surgery and surgical regions highlighting dissimilarities between tumor volumes and surgical removal areas across survival groups.

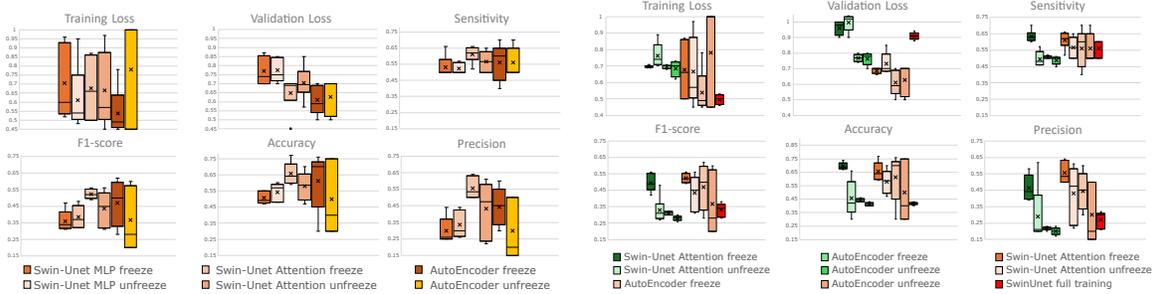
		PCA Euclidean maps	PCA cosine maps	First PCA from local explanations	Global optimizer explanation
Longer-term survivals	Surgery regions	DLP, EA, IFO, PLMC	AA, ACMP, OPF	N/A	N/A
	Brain regions	ACMP, EA	AA, DSV, EA, IFO, LT, MT, PC, VSV	AA, IF, LT, MT, MTV, OPF, Premotor	AA, ACMP, DLP, IFO, MT, OPF, PC, PLMC, PO, SP, VSV
Shorter-term survivals	Surgery regions	ACMP, EA, IFO, PO	DLP, OPF	N/A	N/A
	Brain regions	EA, IFO, OPF, PC	EA, IFO, MT, MTV, PC, PO, VSV	DSV, IF, IFO, MTV, VSV	AA, EA, MT, OPF, PC, PO, VSV

AA: Auditory Association, ACMP: Anterior Cingulate and Medial Prefrontal, DLP: Dorsolateral Prefrontal, DSV: Dorsal Stream Visual, EA: Early Auditory, IF: Inferior Frontal, IFO: Insular and Frontal Opercular, LT: Lateral Temporal, MT: Medial Temporal, MTV: MT+ Complex and Neighboring Visual Areas, OPF: Orbital and Polar Frontal, PC: Posterior Cingulate, PLMC: Paracentral Lobular and Mid Cingulate, PO: Posterior Opercular, SP: Superior Paretal, VSV: Ventral Stream Visual.

Table 2: Training and validation metrics from unsupervised learning of structural MRI, 5-fold cross-validation was used in the Addenbrooke’s Hospital case and 60% - 40% training validation split in the Addenbrooke’s Hospital and BraTS2025.

	Addenbrooke’s Hospital		Addenbrooke’s Hospital and BraTS2025	
	Swin-Unet	Autoencoder	Swin-Unet	Autoencoder
Training loss	0.020 ± 0.004	0.040 ± 0.003	0.003	0.017
Validation loss	0.040 ± 0.050	0.060 ± 0.030	0.004	0.018
RMSE metric	0.010 ± 0.005	0.020 ± 0.005	0.008	0.020
MSM metric	0.001 ± 0.002	0.001 ± 0.001	0.001	0.001
MAE metric	0.009 ± 0.007	0.019 ± 0.006	0.005	0.017

RMSE: Root Mean Squared Error, MSM: Mean Squared Magnitude, MAE: Mean Absolute Error.



(a) Ablation study evaluating incorporating MLP and atten- (b) Ablation study under different unsupervised training co-  
tion modules under encoder freeze and unfreeze strategies hort configurations, including encoder freezing and full train-  
during fine-tuning. Metrics highlight the impact of architec- ing scenarios. Metrics demonstrate the impact of training  
tural choices and training configurations. strategies on model performance.

Figure 4: Examples of training and validation results in the ablation study of Swin-Unet and AutoEncoder variants (a) during fine-tuning with frozen and unfrozen encoder settings; (b) under different unsupervised training configurations.

The performance outcomes for fine-tuning in the survival binary classification task using sMRI data (see 2.2.2) are illustrated in Figure 4. We conducted an ablation study across three encoder-decoder configurations: Swin-Unet with an MLP output layer, Swin-Unet with an attention-based output layer, and a baseline AutoEncoder. Each model was evaluated under two fine-tuning strategies: (i) freezing the pre-trained encoder, and (ii) unfreezing the encoder during downstream training. The variability reported reflects results obtained via 5-fold CV. As shown in Figure 4a, frozen encoders exhibited higher variability across most metrics compared to their unfrozen counterparts. Surprisingly, the frozen configurations also achieved higher average performance. Among all models, the Swin-Unet with an attention-based output layer and frozen encoder achieved the best overall results, with an average F1-score of 0.52, accuracy of 0.67, sensitivity of 0.64, and precision of 0.55. Its maximum values across folds reached an F1-score of 0.56, accuracy of 0.77, sensitivity of 0.66, and precision of 0.65. Although the AutoEncoder architecture achieved slightly higher maximum values in F1-score and sensitivity, it exhibited considerably higher CV variability across all metrics and consistently lower precision (below 0.57), indicating a higher false-positive rate compared to the Swin-Unet with the attention-based output layer. These findings highlight a trade-off between performance stability and sensitivity, and suggest that attention-based decoding in transformer-style architectures offers a more reliable and interpretable solution for domain-specific fine-tuning in neuroimaging applications.

Lastly, Figure 4b provides further evidence from the ablation study on different unsupervised training cohorts in the fine-tuning task, confirming the superiority of the strategy that leverages both the Addenbrooke’s Hospital and BraTS2025 datasets compared to the approach that relies solely on the Addenbrooke’s Hospital dataset. Training from scratch on the Addenbrooke’s Hospital dataset without any fine-tuning resulted in substantially poorer performance compared to either of the two unsupervised training cohorts strategies.

### 4.3 Interpretable Deep Learning for Survival Classification

As the Swin-Unet model with an attention-based output layer and a frozen, unsupervised pre-trained encoder trained on both the Addenbrooke’s Hospital and BraTS2025 datasets outperformed all other configurations, we applied the explanation framework exclusively to this model.

Table 3: Training and validation metrics from Global explanations of structural MRI.

Method	RMSE	MAE	MSM	Sparseness	Faithfulness
Global optimizer (proposed)	<b>0.964 ± 0.12</b>	<b>0.610 ± 0.11</b>	<b>0.967 ± 0.22</b>	0.537 ± 0.31	<b>0.913 ± 0.04</b>
Gradient SHAP	1.066 ± 0.20	0.665 ± 0.22	1.160 ± 0.47	0.441 ± 0.01	0.370 ± 0.38
Guided Backpropagation	1.061 ± 0.21	0.678 ± 0.26	1.175 ± 0.46	0.427 ± 0.01	0.380 ± 0.17
Guided GradCam	1.067 ± 0.20	0.643 ± 0.19	1.166 ± 0.47	<b>0.611 ± 0.05</b>	0.362 ± 0.31
Input X Gradient	1.095 ± 0.12	0.674 ± 0.26	1.189 ± 0.35	0.10 ± 0.02	0.273 ± 0.19
Integrated Gradient	1.095 ± 0.12	0.681 ± 0.25	1.189 ± 0.35	0.445 ± 0.01	0.386 ± 0.26
Kernel SHAP	1.095 ± 0.15	0.690 ± 0.23	1.189 ± 0.37	0.444 ± 0.01	0.35 ± 0.16

RMSE: Root Mean Squared Error, MSM: Mean Squared Magnitude, MAE: Mean Absolute Error.

### 4.3.1 Metrics and interpretations of XAI models

The proposed global explanation optimizer outperformed both the baseline explanation methods used during its training and testing; namely, Gradient SHAP, Guided Backpropagation, and Guided Grad-CAM, as well as established explanation techniques not involved in its training process, including Input  $\times$  Gradient [21], Integrated Gradients [21], and Kernel SHAP [20]. In terms of faithfulness, the optimizer achieved a score of 0.913 (see Table 3). It also had the lowest average RMSE (0.964), MAE (0.610), and MSM (0.967). Standard deviation was assessed across four global explanations: pre- and post-surgery as well as shorter-term and longer-term survivals. While Guided GradCam showed the highest sparseness (0.612), its faithfulness was below 0.362. The optimized method had the highest reliability aligning closely with the first PCA component of sMRI images and preserving key PCA-derived features. Figure 2b illustrates results for the post-surgery longer-term survivals and pre-surgery shorter-term survivals using different  $l_1$ ,  $l_2$ , and  $l_3$  parameter combinations from (4).

### 4.3.2 Patterns identified in XAI explanations

Figure 5 displays the Euclidean distances between pre- and post-surgery scans for both the longer-term and shorter-term survival cohorts. These distances are shown for the first PCA component derived from local explanations, the local explanations themselves (*Gradient SHAP*, *Guided Backpropagation*, and *Guided GradCAM*), as well as for the global explanation (*Global explanation*). The global explanation optimizer outperforms the other methods in terms of sparsity and faithfulness, offering better insights into the global patterns. Thus, we focus primarily on discussing the results from the first PCA component obtained from local explanations and the global explanation maps in Figure 5. By comparing with the atlas using the same thresholding criterion as applied in the PCA Euclidean distance maps, at least 50% of voxels exceeding the 80th percentile and at least one voxel above the 95th percentile, the significant regions are summarised in Table 1 (columns *First PCA from local explanations* and *Global optimizer explanation*).

Overlap patterns between both (first PCA and global optimizer) explanations show that for the longer-term survival cohort, the common regions are Orbital and Polar Frontal, while for the dead cohort, the overlapping regions are Auditory Association (AA), Medial Temporal (MT) and Orbital and Polar Frontal (OPF). The superiority of the global explanation optimizer results is evident, as it includes regions that are observed in the Euclidean distance to the entire PCA component space (see 4.1). For example, the Posterior Cingulate (PC) for the shorter-term survival cohort and the Anterior Cingulate and Medial Prefrontal (ACMP) regions for the longer-term survival cohort are identified. These regions are not included in the first PCA component highlighted in Figure 5.

## 5 Discussion

A suggested guidance based on Table 1 follows the pattern below: Feature engineering (PCA-based Euclidean and cosine maps) revealed that the *surgery regions*, i.e. those showing the largest pre- vs. post-surgery changes within tumour masks PCA-space, partially overlap with (and may help explain) the post-surgery alterations observed in *brain regions*. A key example is the Early Auditory (EA) cortex, which consistently appeared across both survival groups and map types in both the surgical and brain-level results, suggesting it is a core region affected by tumour resection and a hub of post-operative reorganisation [32]. Similarly, Insular and Frontal Opercular (IFO) areas and the OPF cortex were commonly involved, indicating that disruption to sensory and frontal integration areas may play a central role in shaping global connectivity changes [33]. In longer-term survivors, surgical effects were more confined to frontal and midline structures (e.g. ACMP), with downstream changes in executive and motor regions, possibly engaging compensatory networks such as the frontoparietal control system. In contrast, shorter-term survivors showed surgical involvement in posterior and multimodal sensory areas (e.g. Posterior Opercular, PO, or Ventral Stream Visual, VSV), paralleled

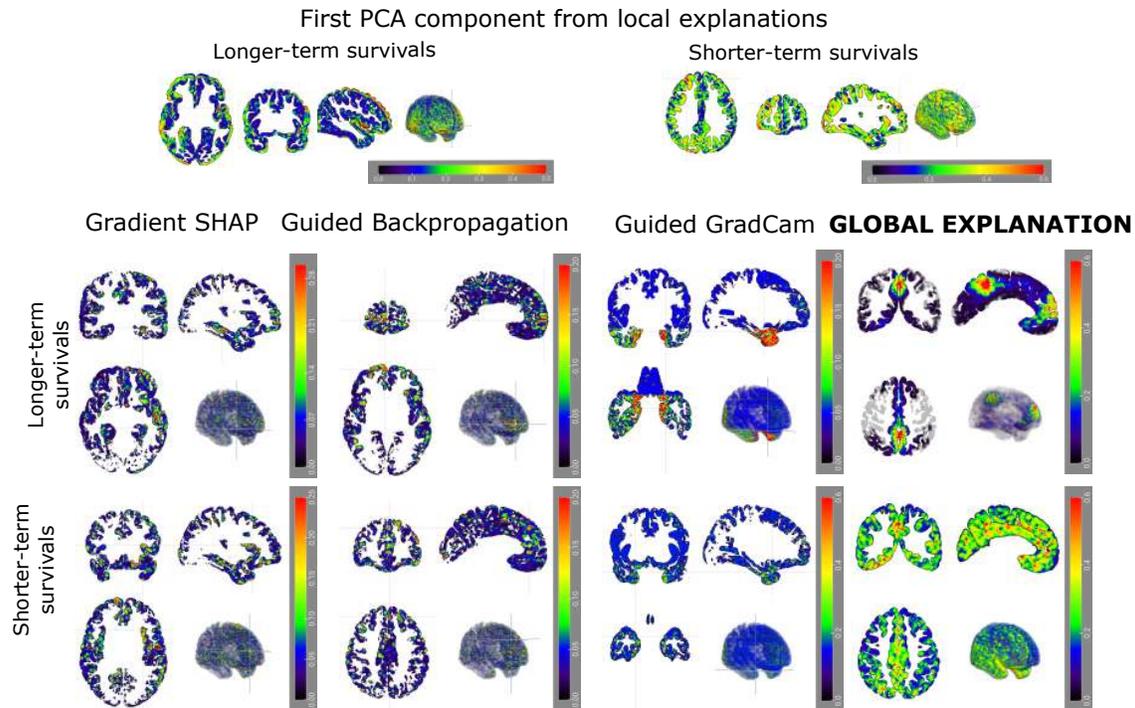


Figure 5: Explainability analysis results. For each severity condition (shorter-term and longer-term survivals), multi-planar slices of the representations obtained are displayed. Top: the first PCA component. Bottom: from left to right the outcomes of applying gradient SHAP, guided backpropagation, guided Grad-CAM, and the final result associated with the global explanation optimizer.

by more diffuse alterations in visual and perceptual cortices, which may reflect greater network fragility or reduced plasticity.

The last two columns of Table 1, representing local and global explanation methods, highlight key brain regions contributing to group distinctions in the binary severity classification. In longer-term survivals, local explanations highlight frontal and temporal areas such as Inferior Frontal (IF), Lateral Temporal (LT) or MT regions. These areas support language, memory, executive functions, and motor planning, indicative of preserved or adaptable networks facilitating recovery. Global explanations in this group further stress integrative hubs like the ACMP, Dorsolateral Prefrontal Cortex (DLP), MT, or OPF, which are implicated in emotional regulation, high-order cognition, and multisensory integration [34]. In shorter-term survivors, the first PCA component based on local explanations also include frontal regions and MTV but is more focused on posterior sensory and association cortices, such as the Dorsal Stream Visual (DSV) and VSV areas. This suggests a heavier impact on visual and interoceptive processing systems that may have less neuroplastic potential.

Overall, the patterns observed across PCA-derived feature maps and model explanation methods suggest that focal surgical changes, particularly in frontal and opercular regions such as EA, IFO, and OPF, are associated with broader post-operative alterations in structurally and functionally connected brain areas. Longer-term survivors consistently show more engagement of fronto-cingulate and temporal structures (e.g., DLP, ACMP, PLMC), linked to executive function and cognitive control, whereas shorter-term survivors exhibit more extensive involvement of posterior sensory and visual association regions (e.g., PO, DSV, VSV), suggesting less focal and potentially less compensable network disruptions [35]. Figure 6 illustrates the regions with the most significant differences in both severity-related groups. The recurrence of regions such as EA and OPF across all analysis methods and survival groups highlights their centrality in the brain’s structural and functional adaptation to tumour resection. These findings support the notion that the surgical impact on specific cortical hubs may influence the extent and efficacy of post-operative neuroplasticity, ultimately modulating clinical outcomes.

A major limitation of this study lies in the limited availability of paired pre- and post-surgical structural MRI data, which restricts the statistical power and the generalisability of the results. Longitudinal imaging of brain tumours remains scarce in clinical datasets, particularly when considering real-world cases with pre- and post-surgical scans. While large public datasets are increasingly available, they often lack longitudinal follow-up or rely on synthetic or

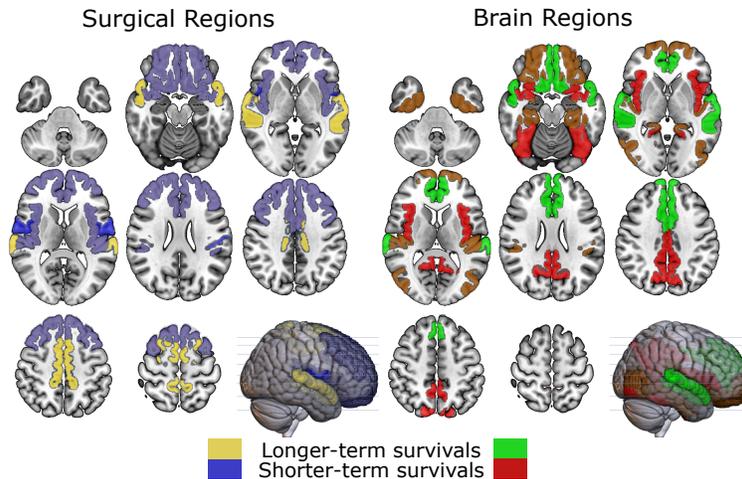


Figure 6: Surgical and brain regions with the most significant changes according to the combined framework for longer-term and shorter-term survival groups. Surgical regions include all areas identified across frameworks, while brain regions are limited to those consistently detected by at least two different frameworks within the same survival group.

preprocessed data that may not fully reflect clinical variability. In contrast, our dataset is composed of real clinical cases, collected under routine care conditions. We believe this provides added value, as it captures the heterogeneity, imaging artefacts, and surgical effects that are often absent in curated public repositories. This specificity allows us to more accurately model the structural impact of surgery on brain tissue and better capture individual variability in post-operative brain reorganisation.

To address this sample size limitation, we are actively collecting additional longitudinal cases to increase cohort size and improve model reliability. In parallel, future work will focus on applying the current framework to other neuroimaging modalities, such as functional MRI (fMRI) or diffusion MRI (dMRI). These modalities offer complementary information that could enhance both the predictive performance of survival models and the interpretability of regional alterations. Integrating functional and structural perspectives will allow a more comprehensive understanding of surgical impact and disease progression.

## 6 Conclusions

Our proposed framework integrates XAI with neuroimaging-based feature engineering to predict the survival of brain tumor patients, providing guidance for surgical decision-making to achieve the necessary onco-functional balance. Our findings correlate dissimilarities between tumor volumes and surgical removal areas with their structural impact on the brain post-operation. By extracting global explanations based on DL approaches for predicting short- and long-term survival, the framework serves as a predictive guideline. The results highlight the involvement of sensory and cognitive regions, with greater disruptions observed in shorter-term survivors. The proposed global explanation optimizer enhances biomarker identification by improving faithfulness and interpretability compared to alternative global XAI methods, while also mitigating inter-method global explanation variability that undermines the trustworthiness of explainable AI. It effectively distinguishes between groups, emphasizing its role in identifying survival-related variability.

## Acknowledgment

This project has received funding from the European Union’s Horizon Europe research and innovation program under grant agreement No 101147319 (EBRAINS 2.0). It was also supported in part by the PID2022-137629OA-I00 and PID2022-137451OB-I00 projects, funded by the MICIU/AEI/10.13039/and by “ERDF/EU”. C.J.M is supported by grant JDC2023-051807-I funded by MICIU/AEI/10.13039/501100011033 and by ESF+. 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.

## References

- [1] Nicholas B. Dadario, Bledi Brahimaj, Jacky Yeung, and Michael E. Sughrue. Reducing the cognitive footprint of brain tumor surgery. *Frontiers in Neurology*, 12, August 2021.
- [2] Christina Drewes, Lisa Millgård Sagberg, Asgeir Store Jakola, and Ole Solheim. Perioperative and postoperative quality of life in patients with glioma—a longitudinal cohort study. *World neurosurgery*, 117:e465–e474, 2018.
- [3] Christiaan HB Van Niftrik and et al. Machine learning algorithm identifies patients at high risk for early complications after intracranial tumor surgery: registry-based cohort study. *Neurosurgery*, 85(4):E756–E764, 2019.
- [4] Jan-Oliver Neumann, Stephanie Schmidt, Amin Nohman, Paul Naser, Martin Jakobs, and Andreas Unterberg. Routine ICU surveillance after brain tumor surgery: Patient selection using machine learning. *Journal of Clinical Medicine*, 13(19):5747, 2024.
- [5] Whitney E Muhlestein, Dallin S Akagi, Jason M Davies, and Lola B Chambless. Predicting inpatient length of stay after brain tumor surgery: developing machine learning ensembles to improve predictive performance. *Neurosurgery*, 85(3):384–393, 2019.
- [6] R Pugalenth, MP Rajakumar, J Ramya, and V Rajinikanth. Evaluation and classification of the brain tumor mri using machine learning technique. *Journal of Control Engineering and Applied Informatics*, 21(4):12–21, 2019.
- [7] Francisco Javier Díaz-Pernas, Mario Martínez-Zarzuela, Míriam Antón-Rodríguez, and David González-Ortega. A deep learning approach for brain tumor classification and segmentation using a multiscale convolutional neural network. *Healthcare*, 9(2):153, February 2021.
- [8] P Sobha Xavier, G Raju, and SU Asawthy. Pre and post operative brain tumor segmentation and classification for prolonged survival. In *International Conference on Soft Computing and Pattern Recognition*, pages 608–616. Springer, 2021.
- [9] Jakub Nalepa and et al. Deep learning automates bidimensional and volumetric tumor burden measurement from mri in pre-and post-operative glioblastoma patients. *Computers in biology and medicine*, 154:106603, 2023.
- [10] Joeky T Senders and et al. Machine learning and neurosurgical outcome prediction: a systematic review. *World neurosurgery*, 109:476–486, 2018.
- [11] Siraj Y Abualnaja, James S Morris, Hamza Rashid, William H Cook, and Adel E Helmy. Machine learning for predicting post-operative outcomes in meningiomas: a systematic review and meta-analysis. *Acta Neurochirurgica*, 166(1):1–14, 2024.
- [12] Karl Pearson. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901.
- [13] Michail Mamalakis, Héloïse de Vareilles, Graham Murray, Pietro Lio, and John Suckling. The explanation necessity for healthcare ai, 2024.
- [14] Luca Longo and et al. Explainable artificial intelligence (xai) 2.0: A manifesto of open challenges and interdisciplinary research directions. *Information Fusion*, 106:102301, 2024.
- [15] Juan M Górriz and et al. Computational approaches to explainable artificial intelligence: advances in theory, applications and trends. *Information Fusion*, 100:101945, 2023.
- [16] Michail Mamalakis et al. Solving the enigma: Enhancing faithfulness and comprehensibility in explanations of deep networks, 2025.
- [17] Ali Hatamizadeh, Vishwesh Nath, Yucheng Tang, Dong Yang, Holger R. Roth, and Daguang Xu. Swin UNETR: Swin transformers for semantic segmentation of brain tumors in MRI images. In Alessandro Crimi and Spyridon Bakas, editors, *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries*, pages 272–284, Cham, 2022. Springer International Publishing.
- [18] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net, 2015.
- [19] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2):336–359, October 2019.
- [20] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

- [21] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 3319–3328. PMLR, 2017.
- [22] Prasad Chalasani, Jiefeng Chen, Amrita Roy Chowdhury, Xi Wu, and Somesh Jha. Concise explanations of neural networks using adversarial training. In *International conference on machine learning*, pages 1383–1391. PMLR, 2020.
- [23] Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the (in) fidelity and sensitivity of explanations. *Advances in neural information processing systems*, 32, 2019.
- [24] Anna Hedström and et al. Quantus: An explainable ai toolkit for responsible evaluation of neural network explanations and beyond. *Journal of Machine Learning Research*, 24(34):1–11, 2023.
- [25] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- [26] Fabian Isensee and et al. Automated brain extraction of multisequence mri using artificial neural networks. *Human brain mapping*, 40(17):4952–4964, 2019.
- [27] Maria Correia de Verdier, Rachit Saluja, Louis Gagnon, Dominic LaBella, Ujjwal Baid, Nourel Hoda Tahon, Martha Foltyn-Dumitru, Jikai Zhang, Maram Alafif, Saif Baig, et al. The 2024 brain tumor segmentation (brats) challenge: glioma segmentation on post-treatment mri. *arXiv preprint arXiv:2405.18368*, 2024.
- [28] Matthew F Glasser and et al. A multi-modal parcellation of human cerebral cortex. *Nature*, 536(7615):171–178, 2016.
- [29] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014.
- [30] Umang Bhatt, Adrian Weller, and José M. F. Moura. Evaluating and aggregating feature-based model explanations. 2020.
- [31] Prasad Chalasani, Jiefeng Chen, Amrita Roy Chowdhury, Xi Wu, and Somesh Jha. Concise explanations of neural networks using adversarial training. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pages 1383–1391, Virtual Event, online, July 13–18 2020. PMLR. Originally released as arXiv:1810.06583 (2018).
- [32] Salla-Maarit Kokkonen, Vesa Kiviniemi, Minna Mäkiranta, Sanna Yrjänä, John Koivukangas, and Osmo Tervonen. Effect of brain surgery on auditory and motor cortex activation: a preliminary functional magnetic resonance imaging study. *Neurosurgery*, 57(2):249–256, 2005.
- [33] Shengyu Fang, Yinyan Wang, and Tao Jiang. The influence of frontal lobe tumors and surgical treatment on advanced cognitive functions. *World Neurosurgery*, 91:340–346, 2016.
- [34] J Hornak, J O’doherly, Jessica Bramham, Edmund T Rolls, Robin G Morris, PR Bullock, and CE Polkey. Reward-related reversal learning after surgical excisions in orbito-frontal or dorsolateral prefrontal cortex in humans. *Journal of cognitive neuroscience*, 16(3):463–478, 2004.
- [35] Riho Nakajima, Masashi Kinoshita, Hirokazu Okita, Tetsutaro Yahata, and Mitsutoshi Nakada. Glioma surgery under awake condition can lead to good independence and functional outcome excluding deep sensation and visuospatial cognition. *Neuro-Oncology Practice*, 6(5):354–363, 2019.