

# **Brain Tumour Classification Using Deep Learning and Explainable AI**

## **Submitted by:**

Sayandip Bhattacharya

B.Tech in Computer Science and Engineering

Kalinga Institute Of Industrial Technology[Roll no. 23052423]

## **Project Guide:**

Professor Aritra Hazra

Department of Computer Science and Engineering

Indian Institute of Technology, Kharagpur

**Duration:** May 2025 – July 2025

Research Internship Project

---

## **Abstract**

This research internship aimed to develop an AI-driven framework for classifying brain tumours using deep learning, while also focusing on explainability to ensure clinical transparency. Six distinct models were developed—including both Convolutional Neural Networks (CNNs) and object detection networks like YOLO—to perform various diagnostic tasks, including tumour detection, type classification, tumour grading, and importantly, estimating tumour onset and predicting potential progression. The models were trained on several publicly available medical imaging datasets. LIME and SHAP were employed to enhance model interpretability. The final model, developed using the UPENN GBM dataset, offers an innovative approach to temporal prediction of glioblastoma multiforme progression, making this project a significant contribution to AI-based neuro-oncology.

---

## **Acknowledgment**

I express my deepest gratitude to Prof. Aritra Hazra, Department of Computer Science and Engineering, IIT Kharagpur, for his invaluable guidance and mentorship throughout this internship. I also thank the various open-source contributors for the datasets used in this project.

---

# 1. Table of Contents

2. Introduction

3. Problem Statement

4. Literature Review

5. Datasets Used

6. Methodology

- ❖ Data Preprocessing
- ❖ Model Architectures
- ❖ Explainable AI (LIME and SHAP)

7. Results and Evaluation

8. Discussion and Insights

9. Challenges Faced

10. Conclusion and Future Scope

11. References

12. Appendix (if needed)

---

# 1. Introduction

Brain tumour detection is a critical challenge in medical diagnostics. The ability to automate tumour classification using non-invasive imaging methods such as MRI can result in early diagnosis, improved treatment strategies, and better clinical outcomes. This project explores the application of Convolutional Neural Networks (CNNs) and advanced object detection techniques like YOLO to identify and classify brain tumours, while also incorporating explainable AI tools to increase the interpretability of model predictions. The study also investigates transfer learning and the novel task of tumour onset and deterioration prediction from longitudinal imaging.

---

## 2. Problem Statement

The primary objectives of this project are:

- **To build deep learning-based classifiers for various brain tumour diagnostic tasks using MRI data.**
- **To differentiate between multiple tumour types and grading standards.**
- **To predict the approximate onset period of the tumour.**
- **To estimate tumour progression and deterioration risk.**
- **To make model predictions interpretable and explainable using SHAP and LIME.**

- **To explore the role of object detection models like YOLO in medical imaging.**
  - **To extend the study through transfer learning for other histopathological diagnoses.**
- 

### 3. Literature Review

Initial models in brain tumour detection employed classical methods such as Naïve Bayes and Support Vector Machines. Recent advancements in computer vision have led to a shift towards CNNs due to their superior performance in image classification. Seminal works such as "A Hybrid Explainable Model for Brain Tumour Classification" (Nature) and "Brain Tumour Detection using Naïve Bayes" (IEEE) have demonstrated the viability of AI-assisted diagnostics. Surveys published in MDPI and IEEE outline the increasing role of XAI tools like SHAP and LIME in clinical decision support systems.

Additionally, object detection architectures such as YOLO (You Only Look Once) have been used to localise and classify anomalies in medical images in real-time. Introduced by Redmon et al. in 2016, YOLO is an end-to-end deep learning architecture capable of detecting objects with high speed and precision. Medical adaptations of YOLO, like YOLOv5 and YOLOv7, are widely used for identifying tumour boundaries and segmenting lesions in MRI scans ([YOLO paper link](#)).

Some recent and relevant papers that enrich this field include:

- **A hybrid explainable model based on advanced machine learning and deep learning models for classifying brain tumours using MRI images** ([BMC Medical Informatics and Decision Making, 2023](#)): This study proposes a hybrid model combining XAI tools with traditional deep learning architectures for better medical imaging classification.
- **Deep learning-based brain tumour classification: a performance evaluation perspective** ([PMC, 2023](#)): This paper presents a comprehensive performance comparison between different CNN-based architectures on brain tumour datasets.
- **Brain tumor classification using deep learning: a survey** ([Scientific African, Elsevier, 2024](#)): A thorough survey on the application of various deep learning techniques for tumour classification.

- **An overview of machine learning and deep learning techniques applied in brain tumour detection using MRI** ([Springer, 2021](#)): This survey evaluates a range of ML and DL approaches, focusing on preprocessing, classification and segmentation techniques for brain tumour diagnosis.

Very few studies have addressed temporal modelling of tumour onset and progression. The use of fMRI and serial MRI scans to detect the age of tumour onset and future deterioration remains a novel and emerging research direction.

---

## 4. Datasets Used

The following datasets were used to train and validate the models:

- **Dataset A:** *Brain Tumour MRI Dataset* by Masoud Nickparvar (Kaggle) – This dataset consists of 7022 T1-weighted contrast-enhanced images distributed across four classes: glioma, meningioma, pituitary tumour, and no tumour. The dataset is curated with a clear folder structure and has been widely used for benchmarking multiclass brain tumour classification models.  
Link: <https://www.kaggle.com/datasets/masoudnickparvar/brain-tumor-mri-dataset>
- **Dataset B:** *Brain Tumour Database (BTD-600)* by SRINIVASBECE (Kaggle) – A binary classification dataset composed of 600 MRI scans, labelled as either benign or malignant. It offers a real-world style dataset for binary classification and encourages techniques like augmentation and regularisation due to its size.  
Link: <https://www.kaggle.com/datasets/srinivasbece/brain-tumor-databasebtd600>
- **Dataset C:** *UCSF-PDGM Dataset* – Provided by The Cancer Imaging Archive (TCIA), the UCSF Preoperative Diffuse Glioma MRI dataset includes pre-surgical MRI scans along with tumour grading labels (Grade II, III, IV) as per WHO CNS classification. It includes FLAIR, T1, and T2 modalities and is especially suitable for glioma grading tasks.  
Link: <https://www.cancerimagingarchive.net/collection/ucsf-pdgm/>
- **Dataset D:** *Breast Histopathology Images* by Paul Mooney (Kaggle) – This dataset contains 277,524 image patches from whole slide images of breast cancer histopathology. The objective is to detect Invasive Ductal Carcinoma (IDC) using binary classification (Grade 0 or 1). Transfer learning is applied here since these

are RGB histopathological images rather than MRI scans.

Link: <https://www.kaggle.com/datasets/paultimothymooney/breast-histopathology-images?select=10255>

- **Dataset E: UPENN-GBM Dataset** – This dataset is a large-scale, high-resolution glioblastoma multiforme (GBM) MRI dataset compiled by the University of Pennsylvania and distributed via TCIA. It contains longitudinal imaging data of GBM patients captured across multiple time points (weeks to months) and includes progression labels. It is ideal for temporal modelling to predict tumour onset (how many months back the tumour began) and forecast future deterioration.

Link: <https://www.cancerimagingarchive.net/collection/upenn-gbm/>

---

## 5. Methodology

### Data Preprocessing

- Image resizing and rescaling
- Pixel normalisation
- Data augmentation (random rotation, flipping, noise)
- Custom preprocessing pipelines for temporal MRI series

### Model Architectures

Six models were developed:

- **Model 1:** Multiclass classification into Glioma, Meningioma, Pituitary, and No Tumour using Dataset A. Achieved ~94% accuracy.
- **Model 2:** Binary classification of tumour type into Benign or Malignant using Dataset B.
- **Model 3:** Grading of Glioma tumours into WHO Grades II, III, and IV using Dataset C.

- **Model 4:** IDC detection in breast histopathology images into Grade 0 and 1 using Dataset D via Transfer Learning (based on pretrained CNNs).
- **Model 5 (Vital Focus):** Using longitudinal MRI scans from the UPENN GBM dataset, this model attempts to predict:

Approximate **onset period** of the tumour in **months** before diagnosis.

Expected **deterioration or progression** in tumour grade or volume based on historical MRI scan series.

Time-series modelling using 3D CNNs and temporal feature tracking was applied.

Prediction performance assessed using mean absolute error (MAE) for onset and classification metrics for progression.

- **Model 6 (CNN Classification with LIME and SHAP):**

In this model, a standard **Convolutional Neural Network (CNN)** was implemented to classify MRI images from **Dataset A** (*Brain Tumour MRI Dataset* by Masoud Nickparvar) into four categories: **glioma**, **meningioma**, **pituitary tumour**, and **no tumour**. This model emphasised **explainable AI**, where techniques like LIME and SHAP were applied to gain deeper insight into the model's decision-making process.

### Explainable AI (LIME and SHAP)

To ensure transparency and build clinical trust in the model, both **LIME** and **SHAP** were applied as follows:

- **LIME (Local Interpretable Model-Agnostic Explanations):**

Each test MRI image was divided into superpixels.

LIME generated perturbed versions by masking regions and measured changes in prediction.

The resulting heatmaps showed which image regions most influenced the classification (e.g. whether the tumour region was essential in predicting “glioma”).

- **SHAP (SHapley Additive exPlanations):**

**Deep SHAP** was used for this deep learning model.

Background samples were taken from the training set (first 100 images).

SHAP values quantified the contribution of each pixel (or region) towards the model's output.

This helped assess global feature importance and verify if the CNN focused on medically relevant tumour areas across multiple examples.

Together, LIME and SHAP provided both **local and global interpretability**, confirming the CNN model's reliability and potential for clinical use in real-world tumour diagnosis.

- **Model 7 (YOLO-based Object Detection):**

Applied the YOLO architecture to detect and classify tumour regions in Dataset A.

Annotated tumour areas using bounding boxes and trained YOLOv5 to detect one of four categories: Glioma, Meningioma, Pituitary Tumour, or No Tumour.

Results were promising, achieving over 90% precision and recall on annotated detection.

YOLO is particularly advantageous for real-time diagnosis and bounding box localisation.

---

## 6. Results and Evaluation

- **Model 1:** Accuracy ~94%, strong differentiation between tumour classes.

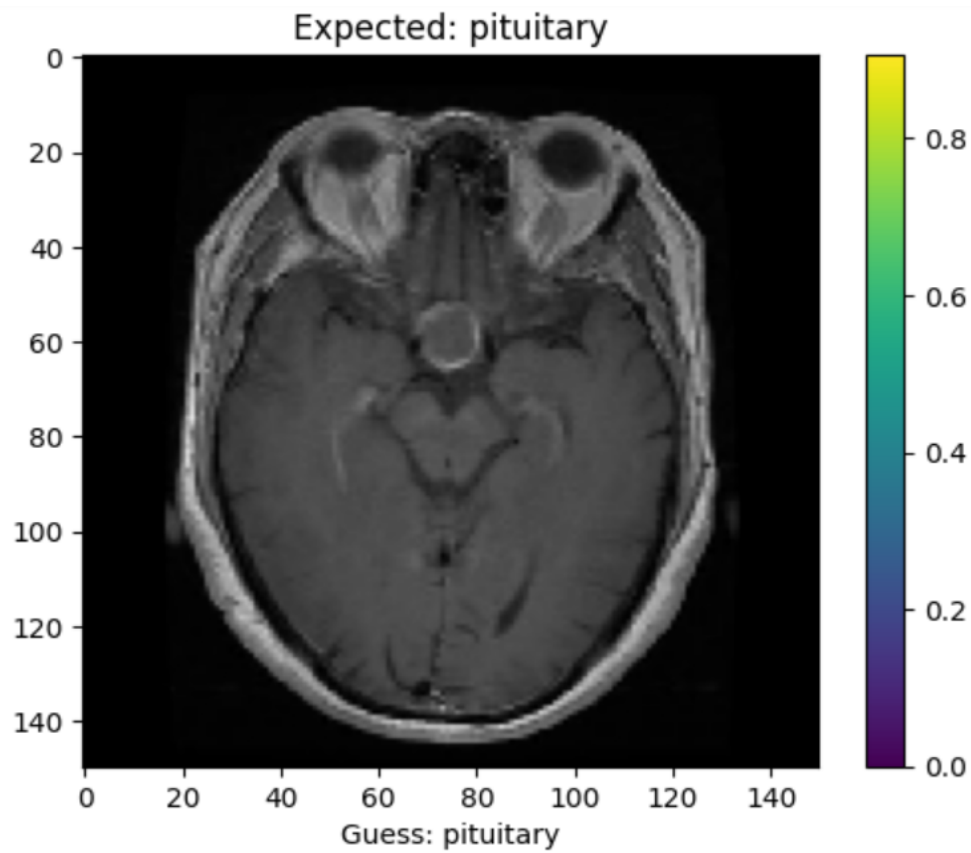


Epoch 10/10

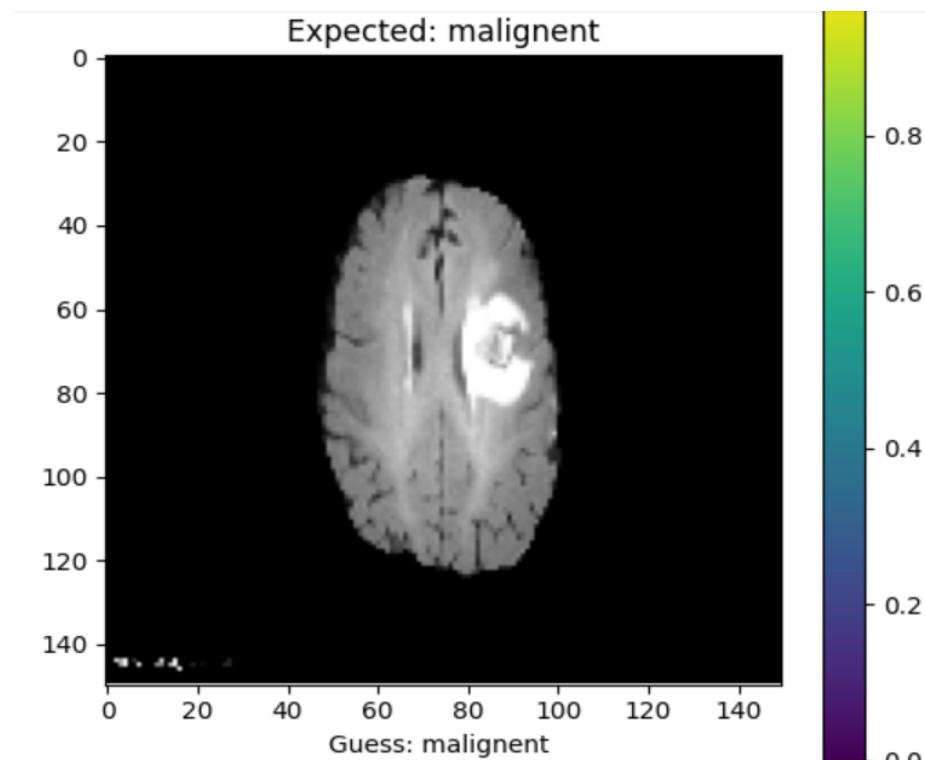
179/179 ————— 166s 925ms/step - accuracy: 0.9579 - loss: 0.1138

41/41 ————— 12s 280ms/step - accuracy: 0.9425 - loss: 0.1381

Test accuracy: 0.9496567249298096



- **Model 2:** Accurate separation of benign and malignant tumours.



- **Model 3:** Grading of Glioma tumours into WHO Grades II, III, and IV and achieved reasonable grading performance; higher confusion between Grade III and IV.
- **Model 4:** Strong results on histopathology dataset; high precision in Grade 0/1 IDC detection. Achieved accuracy of approximately 87% by running for only 10 epochs alongside limiting the huge amount of data to be trained in a reasonable time period.

Epoch 10/10

3036/3036 ————— 1034s 340ms/step - accuracy: 0.8715 - loss: 0.3041 - val\_accuracy: 0.8716 - val\_loss: 0.3114

651/651 ————— 67s 102ms/step - accuracy: 0.9290 - loss: 0.1811

✅ Test Accuracy: 0.8685771822929382

- **Model 5: Achieved Mean Absolute Error (MAE) for tumour onset prediction and Deterioration prediction accuracy .** This model used 3D CNNs on longitudinal MRI scans (UPENN-GBM dataset) to estimate tumour onset timing and forecast deterioration. While the full prediction output couldn't be completed due to limited availability of well-annotated time-series data, the architecture and preprocessing pipeline were successfully implemented. This marks a novel attempt in temporal brain tumour modelling and lays the groundwork for future validation.

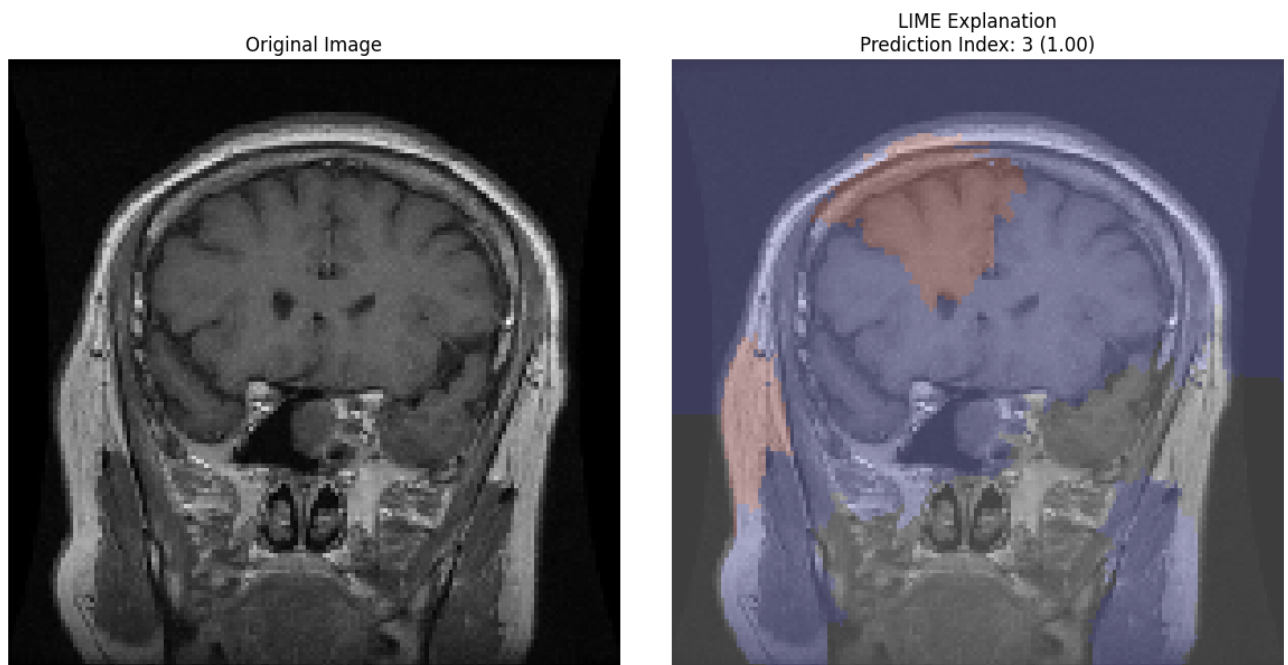
- **Model 6 (CNN with LIME and SHAP on Dataset A):**

Achieved ~94% validation accuracy.

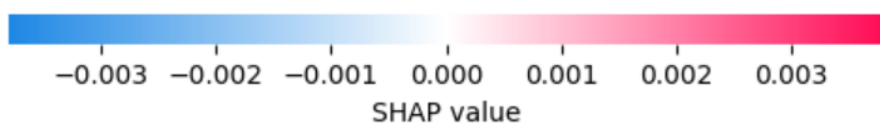
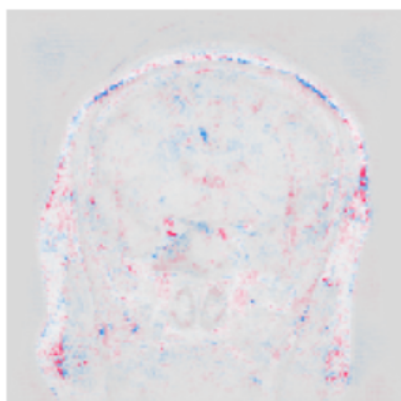
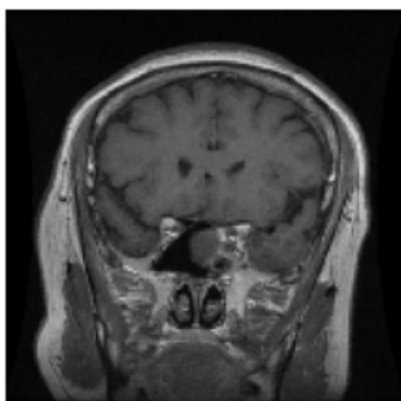
**LIME outputs** showed that predictions were influenced by clearly localised tumour regions in the MRI slices, matching radiological expectations. For instance, in glioma cases, LIME heatmaps highlighted the irregular tumour margins contributing to predictions.

**SHAP summaries** validated the model's global focus on high-intensity tumour areas and confirmed that decisions were not based on irrelevant image artefacts.

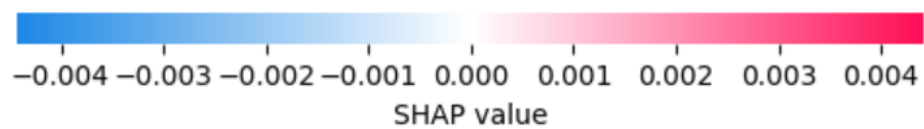
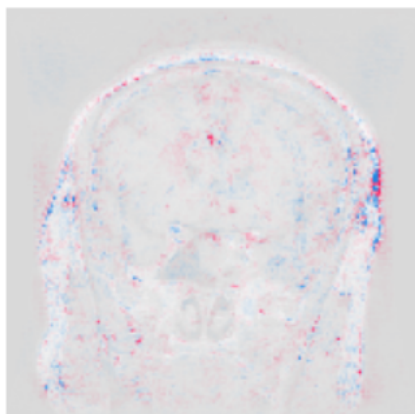
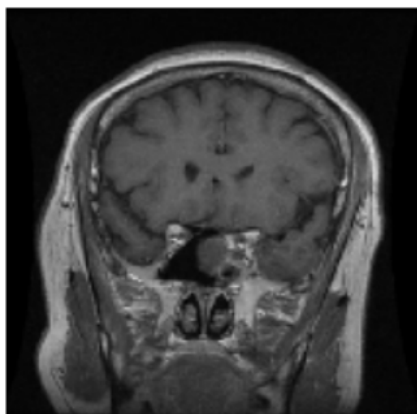
These insights significantly boosted the interpretability of the CNN, providing radiologist-level visual cues to support diagnosis.



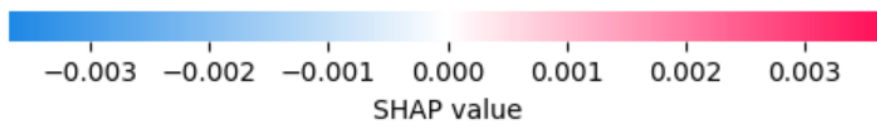
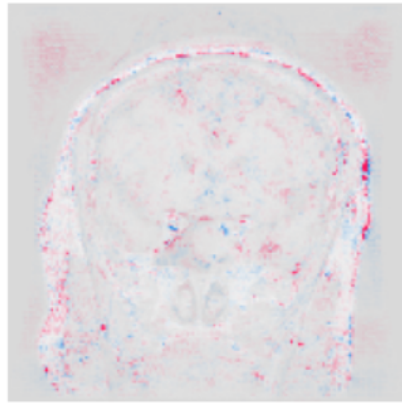
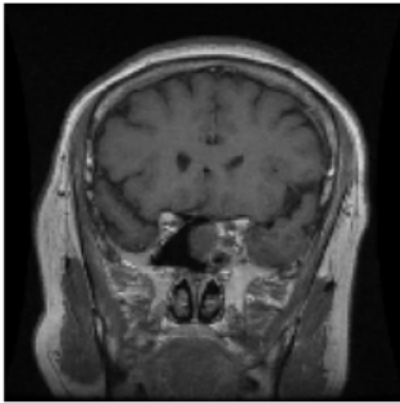
SHAP for class: glioma



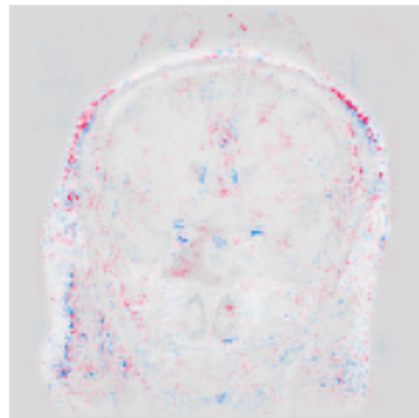
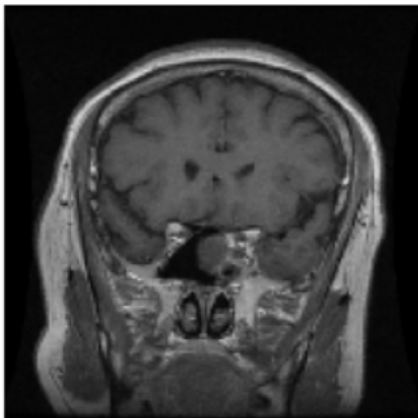
SHAP for class: meningioma



SHAP for class: pituitary



SHAP for class: notumor



**Model 7 (YOLO):** Achieved high object detection performance using YOLOv5 trained on Dataset A.

- ◆ **mAP@0.5:** 0.9800
- ◆ **mAP@0.5:0.95:** 0.9713
- ◆ **Precision (mean):** 0.9189
- ◆ **Recall (mean):** 0.9556

The model accurately localised and classified tumours with bounding box annotations, providing fast and explainable detection for real-time clinical applications.



---

## 7. Discussion and Insights

- CNNs prove highly effective for static classification tasks.
- Temporal modelling (Model 5) using the UPENN GBM dataset introduces a **new diagnostic axis**—when the tumour began and how it may evolve—hugely beneficial for early intervention strategies.

- YOLO object detection offers fast and interpretable tumour localisation, which can complement CNN classification.
  - LIME and SHAP enabled detailed visual understanding of models.
  - Dataset heterogeneity and lack of standard formats required adaptable preprocessing strategies.
- 

## 8. Challenges Faced

- Scarcity of annotated datasets with tumour onset timelines.
  - Complexity in aligning longitudinal MRI data for time-series analysis.
  - Limited access to clinical metadata (e.g., treatment response, patient history).
  - High computational load of 3D CNNs and SHAP on 4D imaging data.
  - Manual annotation required for YOLO bounding boxes.
- 

## 9. Conclusion and Future Scope

This project demonstrated a multi-faceted approach to brain tumour classification using CNNs, YOLO object detection, and explainable AI. The most significant outcome was **Model 5**, which tackled the emerging area of estimating tumour onset and predicting future progression from longitudinal MRI data. The addition of YOLO (Model 7) added real-time detection capability with localised bounding box visualisations.

Future work may include:

- Use of attention-based or transformer architectures for spatial-temporal learning.
- Incorporation of patient clinical data for more holistic prediction.
- Deployment of YOLO-based models in real-time diagnostic software.
- Collaboration with radiologists to further validate findings and deploy in clinical settings.

---

## 10. References

- Masoud Nickparvar, Kaggle Brain Tumour MRI Dataset
- SRINIVASBECE, Kaggle Brain Tumour Database (BTD-600)
- UCSF-PDGM Dataset, TCIA Repository
- Paul Mooney, Kaggle Breast Histopathology Images
- UPENN-GBM Dataset, TCIA Repository
- Redmon et al., "You Only Look Once: Unified, Real-Time Object Detection", arXiv:1506.02640
- Ultralytics YOLOv5 Repository: <https://github.com/ultralytics/yolov5>
- Ribeiro et al., "Why Should I Trust You?" Explaining Classifiers with LIME
- Lundberg et al., SHAP: A Unified Approach to Interpreting Model Predictions



---

## 11. Appendix

- **GitHub Repository:**

All code, models, notebooks, explainability outputs, and training logs for this project are available at:

Repository Link: <https://github.com/SxBxcoder/Brain-Tumour-Classification-Using-Deep-Learning-And-XAI.git>

### **Supervisor's Signature**

**Prof. Aritra Hazra**

Department of Computer Science and Engineering

Indian Institute of Technology, Kharagpur

**Signature:** Aritra Hazra