

CS7641: Project Midterm Report

(Group 42)

Tumor Detection and Classification

Team Members: Shital Salke, Jenny Lin, Koushika Kesavan, Hima Varshini Parasa

Institution: Georgia Institute of Technology

1. Introduction and Background

Brain tumors are abnormal cell growths in or around the brain, classified as benign (non-cancerous) or malignant (cancerous). In adults, the primary types include Gliomas (usually malignant), Meningiomas (often benign), and Pituitary tumors (generally benign). Early diagnosis is crucial for improving outcomes and treatment options. This project aims to develop a machine learning model for detecting and classifying brain tumors using MRI images.

Several studies highlight the effectiveness of modern techniques. A YOLOv7-based model with attention mechanisms achieved 99.5% accuracy [1]. Transfer learning models like VGG16 reached 98% accuracy, outperforming traditional CNNs [2]. A CNN optimized with the Grey Wolf Optimization algorithm achieved 99.98% accuracy [3], while MobileNetv3 reached 99.75%, surpassing ResNet and DenseNet [4]. Additionally, one study compared CNNs (96.47%) with Random Forest (86%) for tumor classification [5].

The dataset for this project, “Brain Tumor (MRI scans),” sourced from Kaggle, contains 3,264 MRI images across three tumor types: Gliomas, Meningiomas, and Pituitary tumors, with a balanced distribution of images in various orientations.

Link to the Dataset: <https://www.kaggle.com/datasets/rm1000/brain-tumor-mri-scans/data>

2. Problem Definition

Manual analysis of MRI scans by radiologists is labor-intensive and error-prone. This project aims to automate the detection and classification of brain tumors (Gliomas, Meningiomas, and Pituitary tumors) to improve accuracy and efficiency compared to existing models.

The solution involves developing a deep learning model using a custom CNN with multiple convolutional, pooling, and fully connected layers. Optimization will utilize Transfer Learning with models like VGG16, ResNet50, and EfficientNet. Performance will be assessed using metrics such as Precision, Recall, Accuracy, F1 Score, and ROC-AUC.

Difference From Prior Literature: Our approach enhances existing methods by utilizing EfficientNetB2 for effective feature extraction and incorporating unsupervised techniques like DBSCAN and GMM to address noise and irregular tumor shapes. This combination leads to improved model performance and faster, more accurate tumor classification.

3. Project Goals

- Detects and classifies brain tumors through MRI scans.
- Reduce false positives and detect tumors early.
- Generates consistent accuracies for different types of brain tumors.
- Generalizes well to unseen MRI scans.

4. Methods

Data Preprocessing Methods

To prepare MRI images for effective model training, preprocessing is crucial to ensure clean, relevant, and balanced data. The following techniques were used:

Data Cleaning

Method: `check_image_integrity`

Explanation: Data cleaning ensures the dataset consists of only valid and usable MRI images by checking for and removing corrupted or unreadable images. This step is crucial as any corrupted images can introduce noise, disrupt model training, and lead to inaccurate results. By removing such images, we ensure that the model is trained on high-quality data, which helps it learn effectively and produce reliable outcomes. This prevents issues during training, such as slower convergence or erratic behavior, and ensures that the model focuses only on meaningful patterns in the data.

Image Normalization

Method: `normalize_images`

Explanation: Image normalization involves resizing the images to a consistent size (e.g., 128x128 pixels) and scaling the pixel values to a uniform range, typically between 0 and 1. This is essential because neural networks and other machine learning models often require input data to be in a consistent format. Resizing ensures all images have the same dimensions, while normalizing the pixel values ensures that they lie within a range that is easy for the model to process. These steps help improve model convergence, reduce the risk of overfitting, and speed up training, as consistent data makes it easier for the model to learn from relevant features.

Dimensionality Reduction

Method: Principal Component Analysis (PCA)

Explanation: Dimensionality reduction, such as PCA, is a technique that reduces the number of input features by transforming the data into fewer components while retaining the most important information. This is particularly useful in dealing with high-dimensional image data where each pixel can be considered a feature. By reducing dimensionality, we not only decrease the computational cost (both memory and processing power) but also help prevent overfitting. This leads to a more efficient model that generalizes better on unseen data, making it less likely to memorize noise or irrelevant details from the training set.

Data Augmentation (Suggested for future integration)

Method: ImageDataGenerator (or similar augmentation techniques)

Explanation: Data augmentation artificially increases the size of the training dataset by applying various transformations to the original images, such as rotations, flips, translations, and zooms. This helps diversify the dataset, enabling the model to learn from a broader range of variations in the images, improving its ability to generalize to new, unseen data. Augmentation also helps address class imbalances, where certain categories may have fewer samples than others. By augmenting underrepresented classes, the model can learn to identify features more effectively across all categories, thus improving its performance and robustness.

Train-Test Split

Method: `train_test_split`

Explanation: The train-test split is an essential step that divides the data into two subsets: one for training the model and one for testing it. Typically, around 80% of the data is used for training, and the remaining 20% is used for testing. This ensures that the model is evaluated on data it has never seen during training, providing a realistic estimate of its performance on new, real-world data. The split helps prevent overfitting, where the model might perform well on the training data but poorly on unseen data. It also

allows for a more reliable evaluation of the model's generalization capability, ensuring that the final model is robust and performs well in practical applications.

Machine Learning Algorithms

Unsupervised Learning

KMeans

We have completed KMeans for this dataset. Using KMeans clustering for brain tumor detection in this context is a bit unconventional since KMeans is typically an unsupervised learning algorithm, often used for clustering rather than classification. However, there are some goals where KMeans might contribute to the project, primarily in preprocessing, feature extraction, or exploratory analysis:

- **Feature Extraction:** KMeans can help identify patterns within the image data. By clustering image patches or pixel intensities, you can potentially discover common features across tumor types. These clusters might represent common textures, edges, or shapes that could then be used as features in a supervised classification model.
- **Dimensionality Reduction and Initialization for Deep Learning Models:** The clusters generated by KMeans could serve as a preliminary grouping or an initial feature space, which may then be fed into a more complex model like EfficientNet. For example, KMeans can identify latent patterns, which may improve the learning process of the neural network when used in the initial layers.
- **Data Augmentation:** KMeans could be used to create synthetic labels for unlabeled or newly collected data. For example, it could identify subgroups within a "glioma" category that may be hard for human annotators to distinguish but still share common traits. This can assist with labeling or even serve as pseudo-labels to increase the dataset's diversity.
- **Exploratory Data Analysis:** Clustering images with KMeans might reveal which types of brain tumors are more similar or have overlapping features, allowing you to gain insight into the structure of the dataset before you start supervised learning. This could help you understand the difficulty of the classification task, for example, by highlighting potential class overlaps.

Explanation: In our implementation, we applied Principal Component Analysis (PCA) to reduce the dataset's dimensionality to two components, PC1 and PC2. These principal components capture the directions of maximum variance in the data, with PC1 representing the largest variance and PC2 representing the second-largest variance. This reduction allows us to visualize the high-dimensional data in a lower-dimensional space, making it easier to interpret the clusters formed by KMeans. Initially, KMeans assigned arbitrary labels (0-3) to the clusters, which were not useful for further analysis. To address this, we applied a post-processing step where we examined each cluster and counted the

frequency of the four tumor categories. The most frequent category within each cluster was used as the final label, providing more meaningful classification results.

GMM

We have completed GMM clustering for this dataset. Using GMM for brain tumor detection is effective because GMM is a probabilistic model that assumes data points are generated from a mixture of Gaussian distributions. GMM is typically used for density estimation and clustering, making it a suitable choice for identifying underlying patterns in the dataset. Below are the specific ways GMM contributes to the project:

- **Feature Extraction:** GMM can identify latent patterns in the MRI image data by modeling each cluster as a Gaussian distribution. These clusters might represent specific textures, edges, or tumor shapes that are shared across tumor types. These probabilistic cluster assignments can be further used as features in a supervised learning model to enhance its performance.
- **Dimensionality Reduction:** The probabilistic nature of GMM complements dimensionality reduction techniques like PCA. After reducing the image data to lower-dimensional components using PCA, GMM leverages the reduced feature space to fit Gaussian distributions, simplifying the representation of high-dimensional data.
- **Noise Handling:** GMM provides probabilities for each data point belonging to each cluster. This allows for soft clustering, which is particularly beneficial for noisy or ambiguous data, as it accounts for overlaps and uncertainties in cluster memberships.
- **Exploratory Data Analysis:** By clustering MRI images into Gaussian distributions, GMM can reveal how the dataset is structured. For example, it highlights whether tumor types are well-separated or if there is significant overlap. This analysis helps gauge the complexity of the classification problem and identify areas where the model struggles.

Explanation: In our implementation, we applied Principal Component Analysis (PCA) to reduce the dataset's dimensionality to 50 components, which retained the most critical variance in the data. GMM was then applied to these PCA-transformed features. Each cluster was modeled as a Gaussian distribution with a unique covariance matrix (`covariance_type='full'`) to allow for flexibility in cluster shapes.

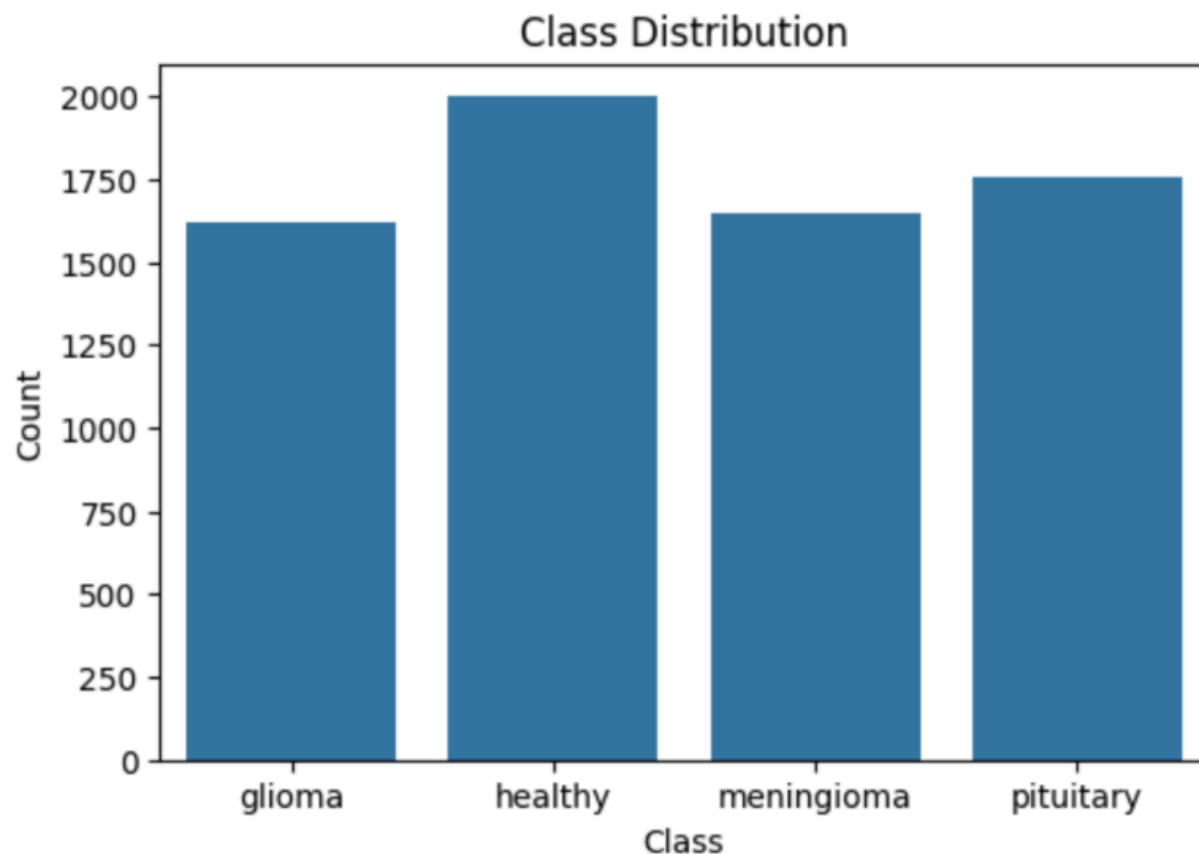
Initially, GMM assigned soft cluster probabilities across all tumor categories, and these were used to classify the images. However, since GMM assigns clusters independently of the true labels, a post-processing step was applied. This involved mapping clusters to tumor categories by identifying the most frequent category label within each cluster, enabling more meaningful classification.

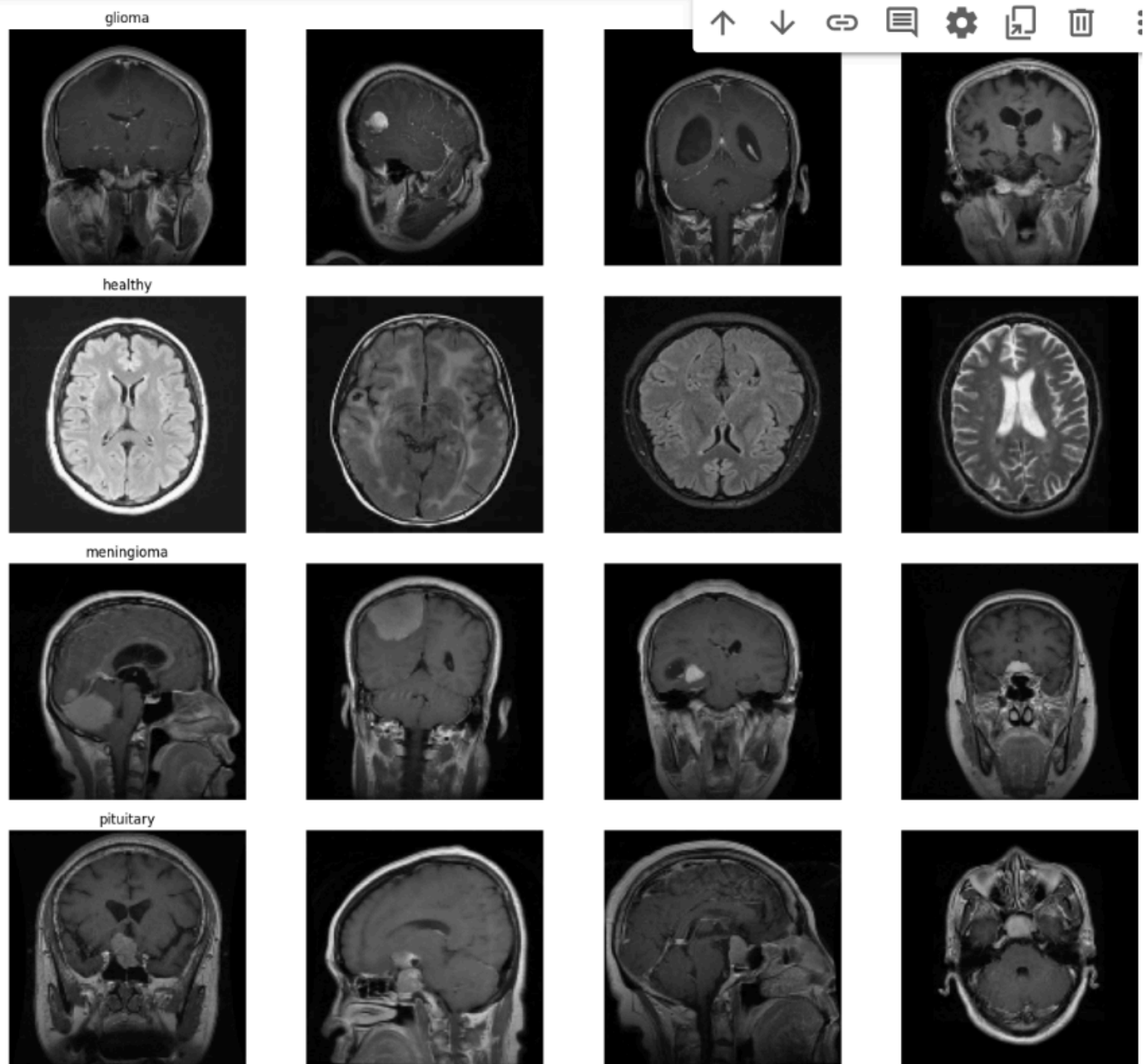
Supervised Learning

ResNet50

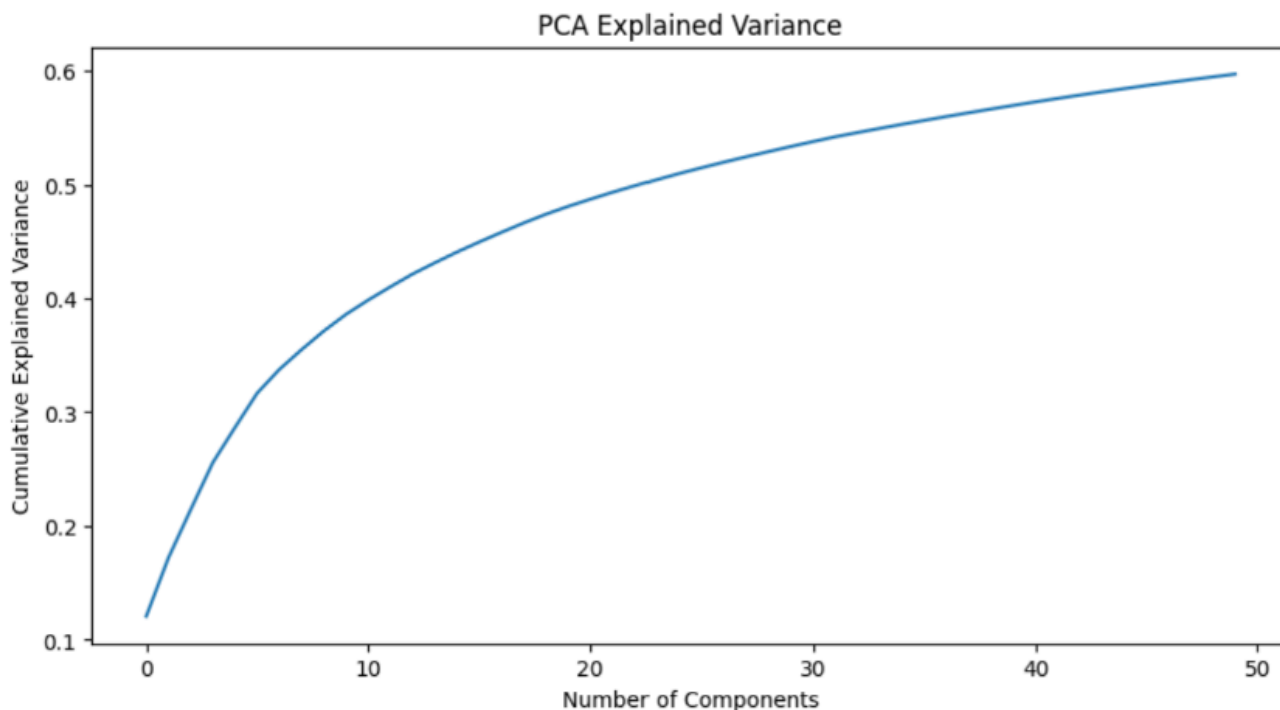
ResNet50, also known as "Residual Network with 50 layers" is a deep convolutional neural network (CNN) which addresses the issue of vanishing gradients by introducing "residual connections" allowing for the training of very deep networks. ResNet50's deep architecture captures fine-grained features at multiple levels which helps in distinguishing between different types of brain tumors (e.g., glioma, meningioma, pituitary tumors) as well as healthy tissue. Fine-tuning the top layers of ResNet50 on the brain tumor dataset improves accuracy with relatively less labeled data. It handles depth and complexity which is required in classifying Brain MRI Scans that are usually complex in nature.

1. Exploratory Data Analytics The resized images of dimensions 128*128 pixels (required input size for ResNet50) are visualized based on each class ('glioma', 'healthy', 'meningioma', 'pituitary') to ensure that the dataset is balanced. Sample images are also generated from each class.





2. Dimensionality Reduction with PCA The images are flattened into vectors before applying PCA, which reduces the dimensions to 50 principal components. The cumulative explained variance plot helps in understanding how much variance is captured by the selected components. It indicates whether 50 components are sufficient to retain most of the information.



3. One-Hot Encoding Labels One-hot encoding labels convert the categorical labels into a format suitable for multi-class classification. It converts 'y_train_encoded' and 'y_test_encoded' into one-hot encoded arrays for use in the classification layer.

4. ResNet50 Model Training A pre-trained ResNet50 model is set up as a feature extractor adding custom layers for the classification task. The model is compiled with categorical cross-entropy loss and accuracy as metrics, appropriate for multi-class classification. Cross-validation with StratifiedKFold performs 5-fold cross-validation to evaluate model performance robustly across different data splits. For each fold, the model is trained on 4 folds of data and validated on the remaining fold. The cross-validation loop calculates the F1 score and ROC-AUC score for each fold, appending the scores to lists for overall performance measurement.

CNN

We implemented the CNN model as a part of supervised learning of brain tumor detection that will classify MRI images into predefined categories.

CNN Model Architecture: It consists of five convolution blocks followed by dense layers for classification. That hierarchy within the layers has meaning in the abstraction of features from MRI images.

1. **Convolutional Layers:** That is, a model with a 32-filter 'Conv2D' that increases gradually up to 512 in successive layers will enable the network to learn from low-level, fine-grained details to high-level, complex features. The kernel size for all the convolutional layers is chosen to be 3x3, as this has become a kind of standard because it is a great tradeoff between capturing local features and

reducing computational complexity. For all the layers, activation is ReLU that introduces non-linearity into the network and, hence, makes the network capable of learning complex patterns required for distinguishing between types of tumors. Batch normalization after every convolutional layer normalizes the output. In such a way, it has been able to stabilize and speed up the learning process since it reduces the problem of internal covariate shift that leads to improvements in convergence. Normalization shall have high utility in deeper networks; this model is less sensitive to initializations and learning rates. After each convolution block, a max-pooling layer is applied; this reduces the volume by half in the spatial dimension. Max-pooling reduces the dimensions of feature maps, and by doing so it retains the most important features since the process reduces computation in the pooling features and enforces spatial invariance. This will enable the model to hierarchically activate features from lower and lower levels of abstraction, a very desirable trait for complicated image data acquired by MRI scans.

2. **Dropout Layers:** To handle overfitting, dropout is added after each max-pooling layer. The dropout rate shall be 0.25-that is, at every iteration the model sets 25% of units to zero. In this way, the model will learn more robust features and improve its generalization capability. Dropout is increased to 0.5 in fully connected layers, where there is more risk of overfitting because of more parameters.
3. **Fully Connected Dense Layers** Next are three thick, fully connected layers with sizes 1024, 512, and 256 with batch normalization and dropout. These dense layers combine features learned through convolutional layers in making final classification decisions. Added Dense Layers with ReLU activation to enhance model capacity to learn multi-dimensional complex relations and patterns in data.
4. **Output Layer** The last layer is a 'Dense' layer, softmax activation gives the probabilities of each type of tumor. Softmax gives a probability to each class so that model has some form of output in probabilities over all classes. Quite good structure for multi-class classification since it has a high value of class in probability to be chosen as a predicted class for the model.
5. **Compilation** It only makes the model compile using the Adam optimizer. In this variant of the Adam optimization algorithm, learning rates are adapted for each parameter individually in a way depending on the magnitude of gradient for that parameter in a mini-batch. This generally leads to a higher convergence speed and performance. Loss function: The loss function used in this model is categorical cross-entropy for multi-class problems. This defines how the predicted labels deviate from the true ones, giving larger penalties on larger errors. This gives a higher chance of better prediction using a model. Metric: Accuracy simply refers to the measure that tells something about how frequently a model predicts an output correctly.
6. **Training Process** It uses the batch size 32 to train the model using the training dataset for over 30 epochs. It shall use validation data for testing, with a view to offering real-time monitoring of the model's performance on unseen data. It iteratively updates its weights and biases motivated by minimizing the loss function so that models progressively improve the classification accuracy.
7. **Evaluation** The model is then matched against a test dataset containing performances that yield an accuracy score in determining how the model will perform in classifying examples. Indeed, this is the

ultimate success measure, the test accuracy score, to classify brain tumor types from MRI images using this model. The model is therefore tailored on insight generalization into a wide array of tumor types with very sparse labeled data and hence will be suitable for application in medical diagnosis where utmost measures against accuracy and robustness are paramount.

5. Results and Discussion

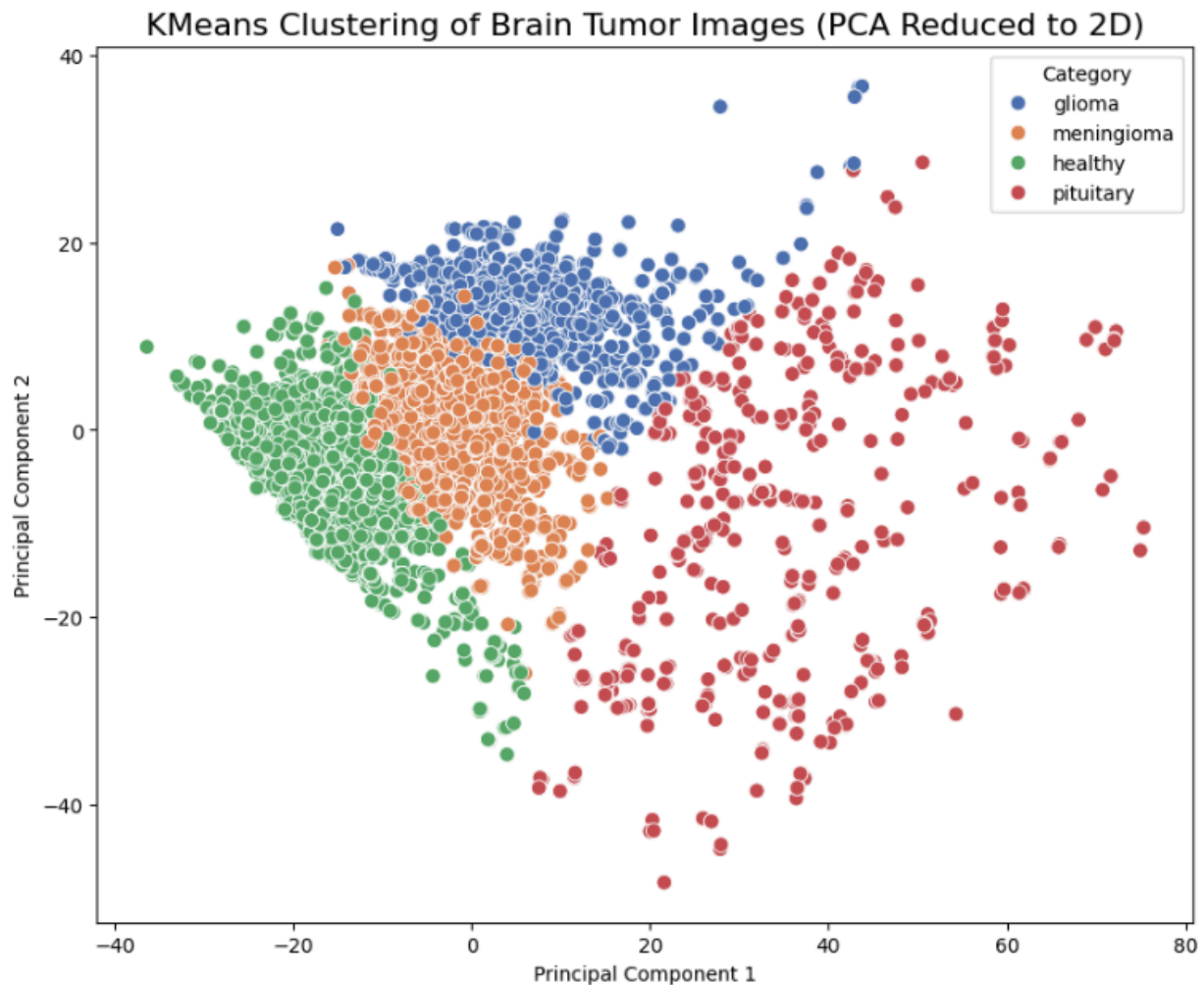
ML Metrics

We will evaluate our results using the following metrics [7]:

- **F1 Score:** A balanced measure of precision and recall, crucial in clinical settings where false positives and negatives have serious implications.
- **AUC-ROC:** Represents the model's discriminative power across all classification thresholds.
- **Confusion Matrix:** Provides detailed performance insights for each tumor type.
- **Cross-Validation:** K-fold cross-validation ensures model consistency and generalizability across diverse patient data.

Our chosen algorithms are expected to yield strong performance in multi-class scenarios, improving the F1 Score and AUC-ROC metrics. Cross-validation will help assess model consistency across different datasets using the EfficientNetB2 algorithm.

KMeans

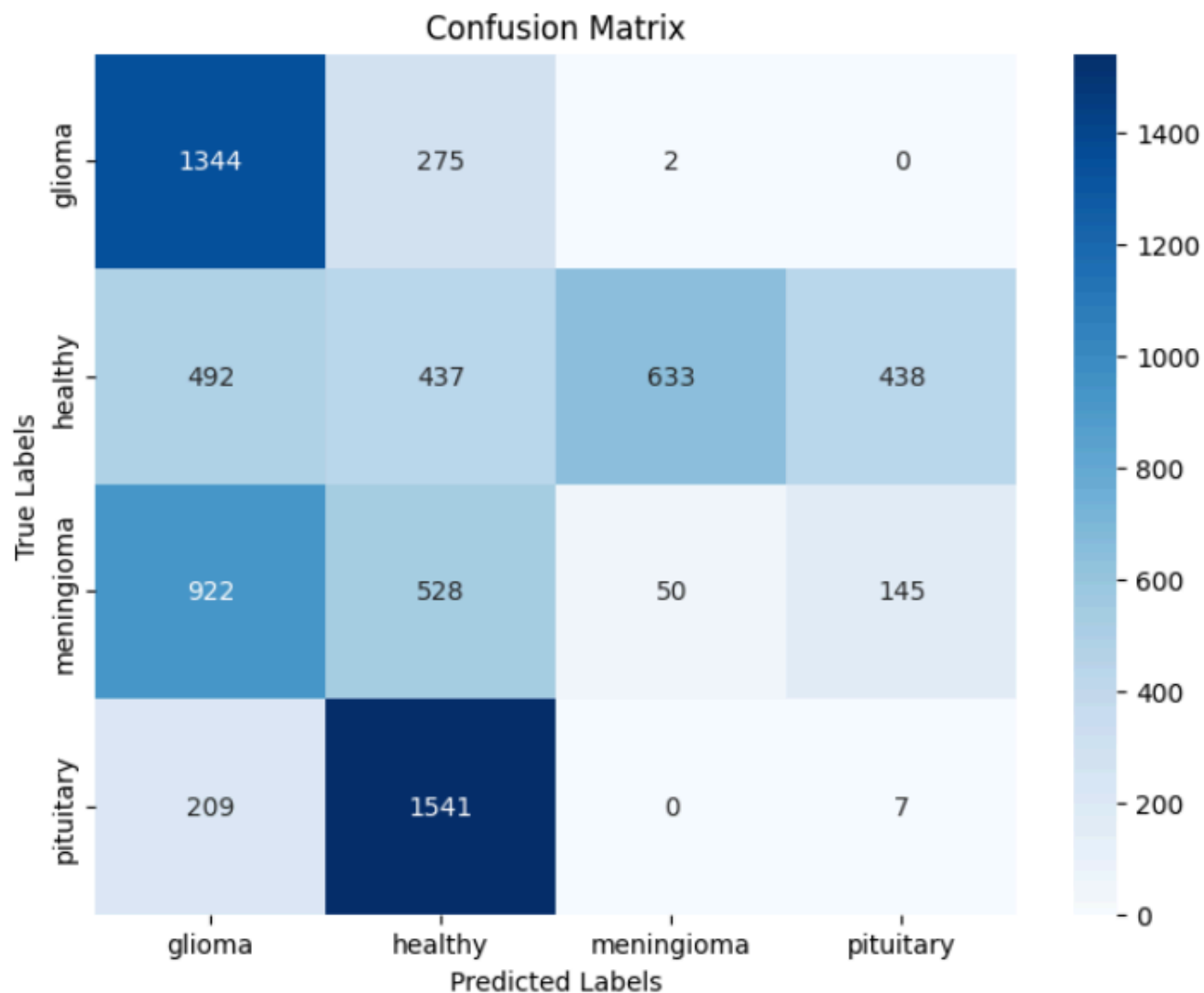


1. F1 Score: 0.4903

- Interpretation:** The F1 score is the harmonic mean of precision and recall, which is a good measure for imbalanced datasets. It ranges from 0 to 1, where 1 is a perfect score and 0 indicates poor performance.
- In your case, an F1 score of 0.4903 means that, while your clustering model isn't performing excellently, it is better than random guessing. It indicates moderate precision and recall across the clusters.

2. Confusion Matrix

```
[[ 947   60  612    2 ]
 [ 179 1217  581   23 ]
 [ 173  632  640  200 ]
 [ 230  396  534  597 ]]
```

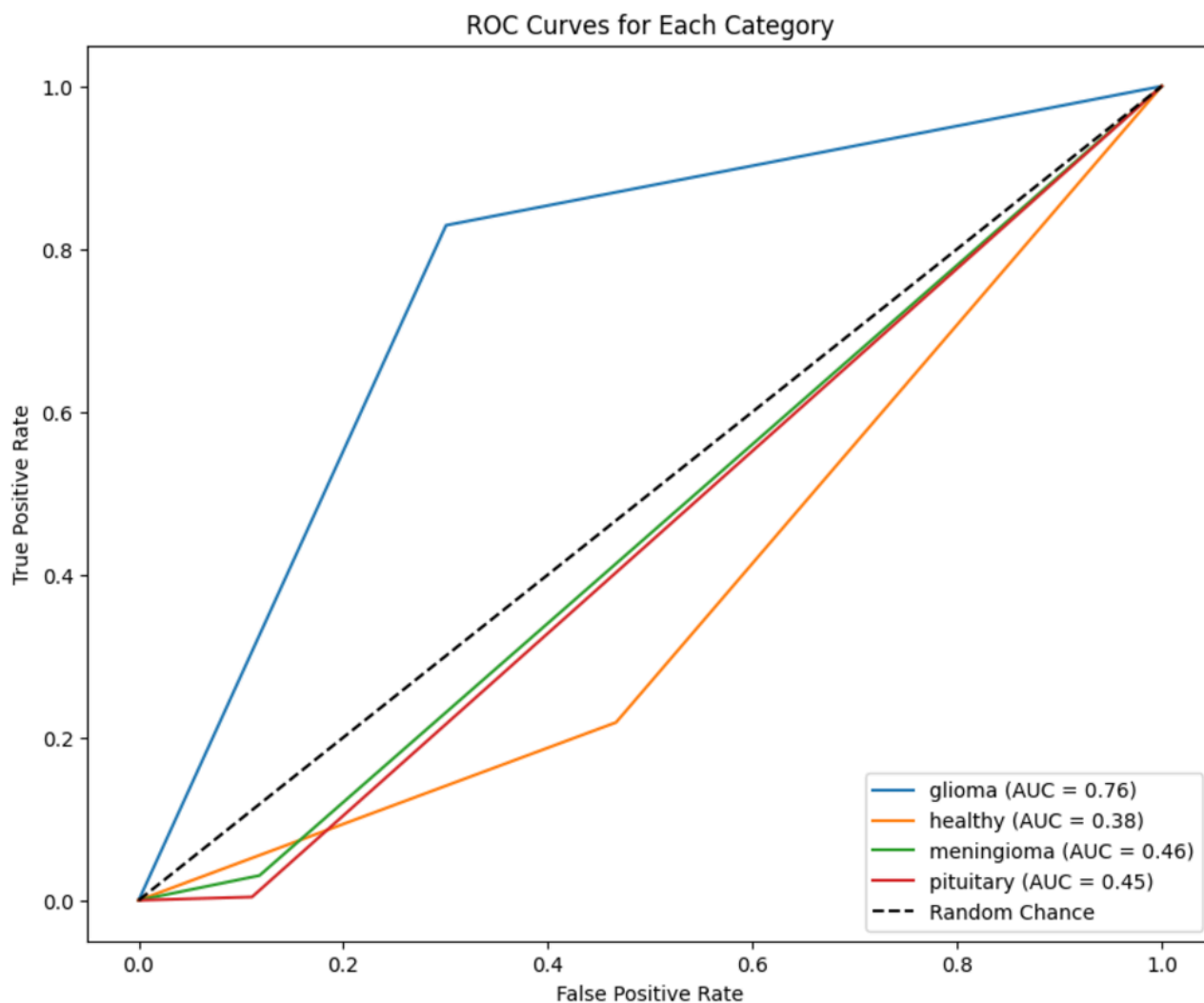


- **Interpretation:** A confusion matrix shows how the model's predictions match the true labels. It is a matrix of shape (n_classes, n_classes), where n_classes is the number of unique categories. Each row represents the actual class, while each column represents the predicted class.
- The rows correspond to the **true classes**, and the columns correspond to the **predicted classes**. For instance, in row 0 (the true class for cluster 0), **947** images were correctly predicted as cluster 0, **60** were incorrectly predicted as cluster 1, **612** as cluster 2, and 2 as cluster 3.
- From the confusion matrix, you can see that some clusters are mixed with other categories, which means that the clustering algorithm isn't perfectly distinguishing between some of the categories.

3. AUC-ROC Score: 0.3032

- **Interpretation:** The AUC-ROC (Area Under the Receiver Operating Characteristic Curve) score measures the performance of a classification model. It is typically used for binary classification, but here it is being used for a multi-class case with the "One-vs-Rest" (OvR) approach.

- A **low AUC-ROC score of 0.3032** indicates that the clustering model struggles to distinguish between the categories. An AUC score closer to 1 means the model has better performance at differentiating between classes. Since your score is much lower, it suggests that your clustering algorithm isn't reliably separating the classes.



AUC ROC score for KMeans

4. Cross-Validation Scores: [0.92170819 0.01708185 0.03487544 0.24287749 0.21866097]

- **Interpretation:** These scores represent the accuracy of the model on 5 different cross-validation folds (using 5-fold cross-validation). Each number corresponds to the accuracy for one fold.
- You can see that the cross-validation scores are quite variable, with one score as high as **92.17%** and others as low as **1.7%**. This suggests that the model is overfitting on some data and performing poorly on others. Such a large discrepancy indicates that the model may not generalize well across all subsets of the data.

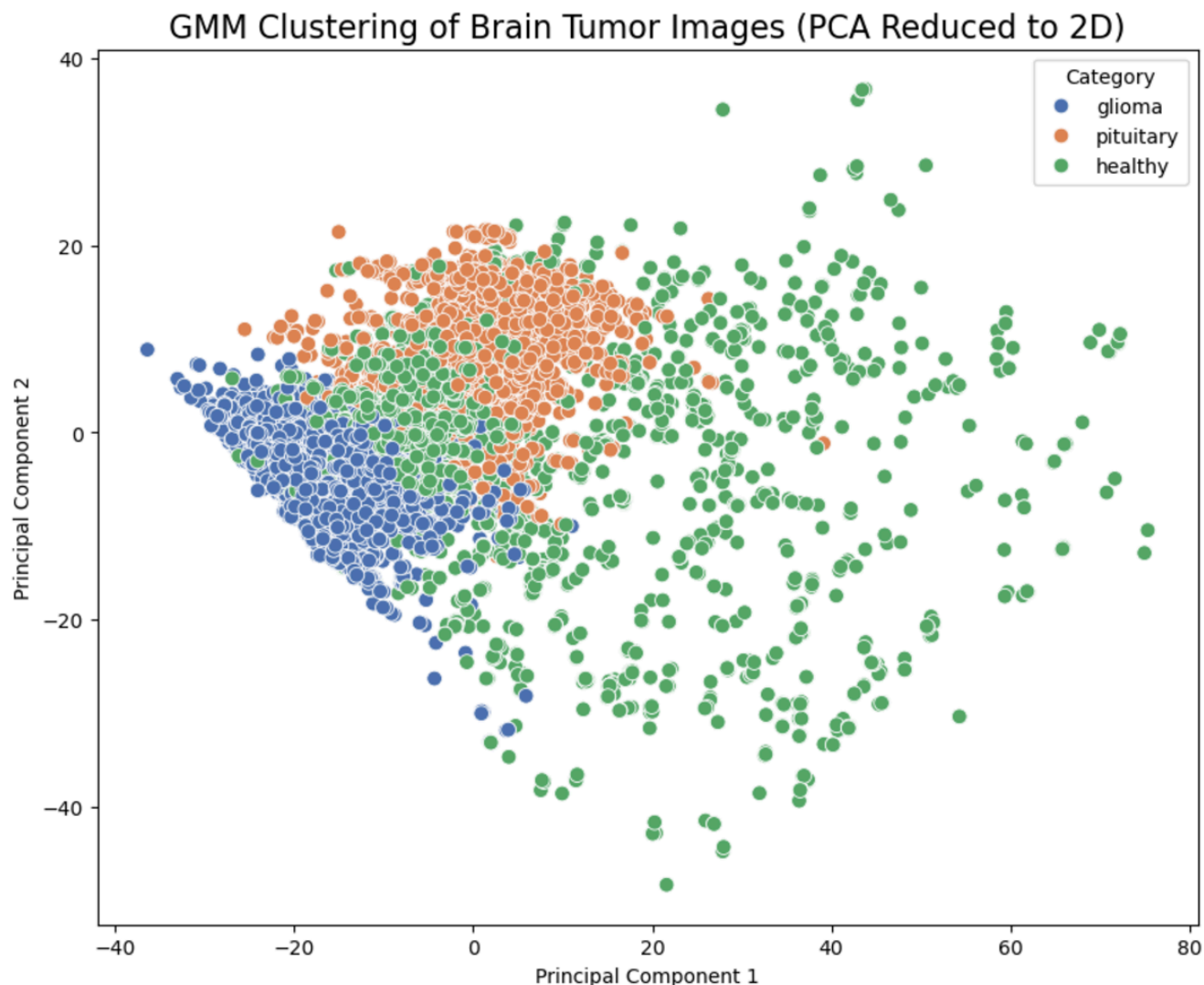
5. Mean Cross-Validation Accuracy: 0.2870

- **Interpretation:** The mean cross-validation score of **0.2870** indicates that, on average, the model's accuracy is quite low across the different folds. This further confirms that the clustering model may not be robust, and its performance isn't stable across different subsets of data.

Summary of what these results mean:

- **The F1 score** shows that the clustering has some decent precision and recall but could be improved.
- **The confusion matrix** indicates misclassifications between clusters, suggesting that the KMeans algorithm might not be able to fully distinguish between the categories.
- **The AUC-ROC score** is quite low, meaning the model struggles to differentiate between the categories in a meaningful way.
- **The cross-validation scores** are highly variable, which points to the model not being stable or reliable when tested on different data subsets.
- **The mean cross-validation accuracy** further supports that the model's generalization is weak, with the accuracy being quite low on average.

GMM

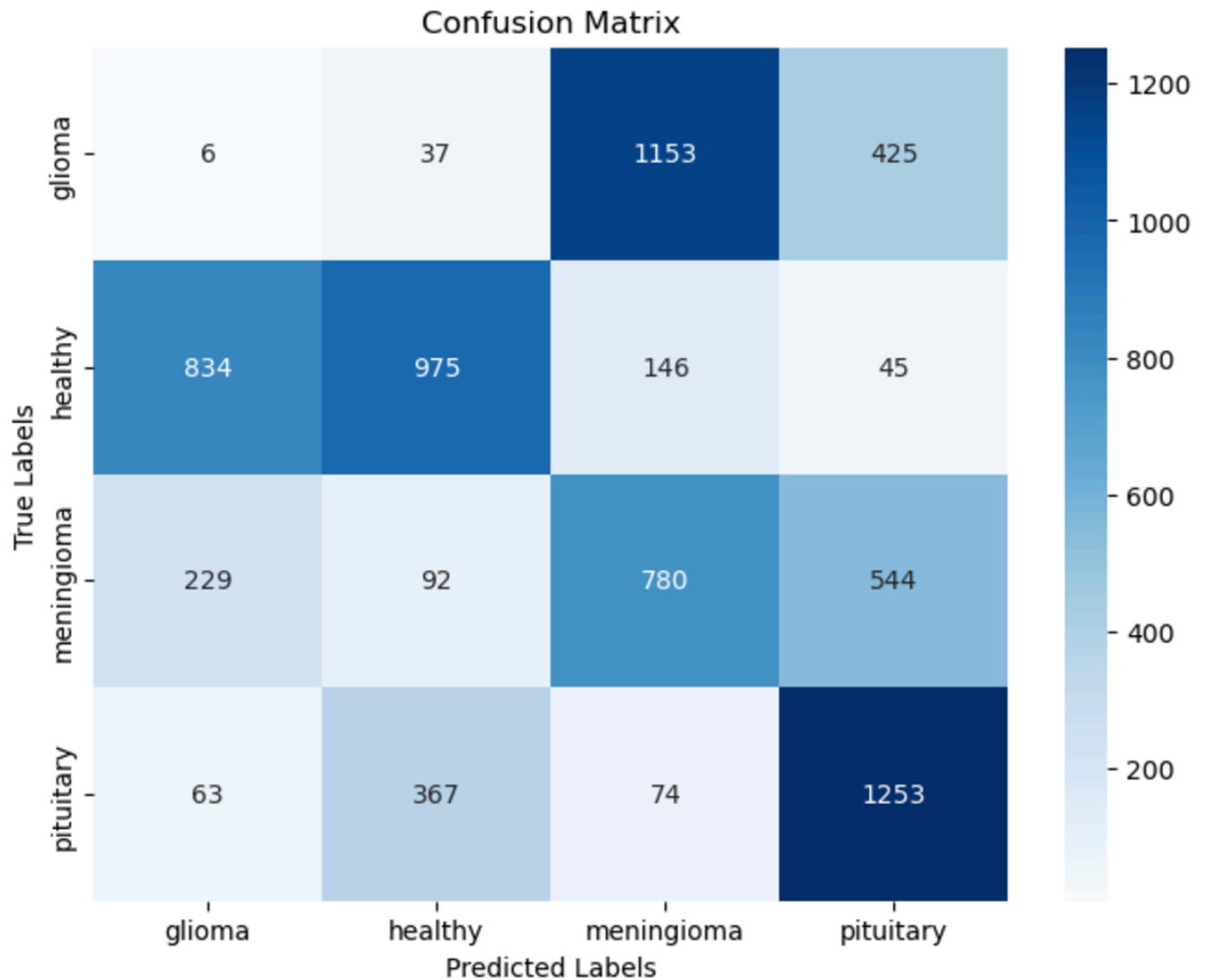


1. F1 Score: 0.4130

- **Interpretation:** The F1 score, a harmonic mean of precision and recall, is particularly useful for imbalanced datasets. It ranges from 0 to 1, where a score closer to 1 indicates better performance.
- An F1 score of **0.4130** suggests that the GMM clustering model achieves moderate performance, with some precision and recall, but there is significant room for improvement. This score indicates that the model struggles to consistently classify the images into their correct categories.
- The drop of "**meningioma**" as a category in the GMM clustering highlights the model's struggle with overlapping clusters and imbalanced category representation. The "meningioma" images likely exhibit similarities with "glioma" and "pituitary" categories in the PCA-reduced space, causing them to be absorbed into clusters dominated by these categories. This indicates that the current feature extraction and dimensionality reduction methods may not sufficiently separate the distinct characteristics of "meningioma" from other tumor types.

2. Confusion Matrix

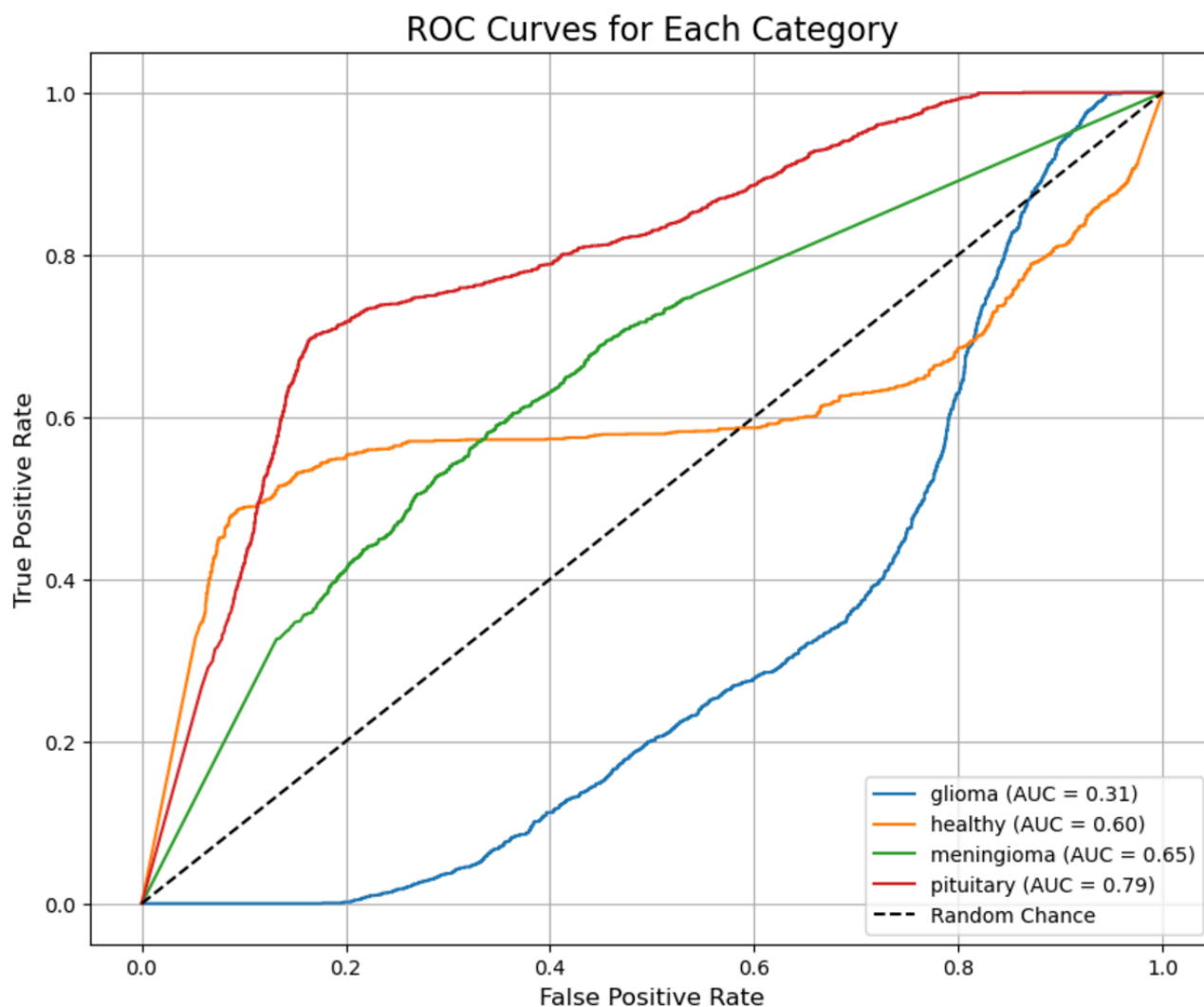
```
[[ 6   37  1153  425 ]
 [ 834  975  146   45 ]
 [ 229   92  780  544 ]
 [ 63   367   74  1253 ]]
```



- **Interpretation:** The confusion matrix compares true labels with predicted labels for each category. Each row represents the true class, while each column represents the predicted class.
- In the first row (true class: glioma), **1153** images were misclassified as meningioma, and **425** as pituitary, with only 6 images classified correctly.
- The second row (true class: healthy) shows better performance, with **975** correctly classified images but a significant number misclassified into other categories.
- The third and fourth rows indicate similar trends, with large misclassification counts. The confusion matrix highlights the model's difficulty in distinguishing certain categories (e.g., glioma and meningioma), likely due to overlapping features or limited cluster separation.

3. AUC-ROC Score: 0.6286

- **Interpretation:** The AUC-ROC (Area Under the Receiver Operating Characteristic Curve) score measures the performance of a classification model. It is typically used for binary classification, but here it is being used for a multi-class case with the "One-vs-Rest" (OvR) approach.
- A **low AUC-ROC score of 0.3032** indicates that the clustering model struggles to distinguish between the categories. An AUC score closer to 1 means the model has better performance at differentiating between classes. Since your score is much lower, it suggests that your clustering algorithm isn't reliably separating the classes.



AUC ROC score for GMM

4. Cross-Validation Scores: [0.7623, 0.0121, 0.1744, 0.5299, 0.3312]

- **Interpretation:** The AUC-ROC score evaluates the model's ability to distinguish between categories, with values closer to 1 indicating better discrimination.
- A score of **0.6286** is moderate, suggesting the GMM clustering model provides some separation between categories but struggles in multi-class settings. This score indicates the model's performance is slightly above random guessing and needs improvement to reliably differentiate tumor types.

5. Mean Cross-Validation Accuracy: 0.3620

- **Interpretation:** The average cross-validation accuracy of **0.3620** reflects overall weak performance across the dataset.
- The low mean and high variability indicate that the GMM model struggles to generalize across different subsets of the data.

Summary of what these results mean:

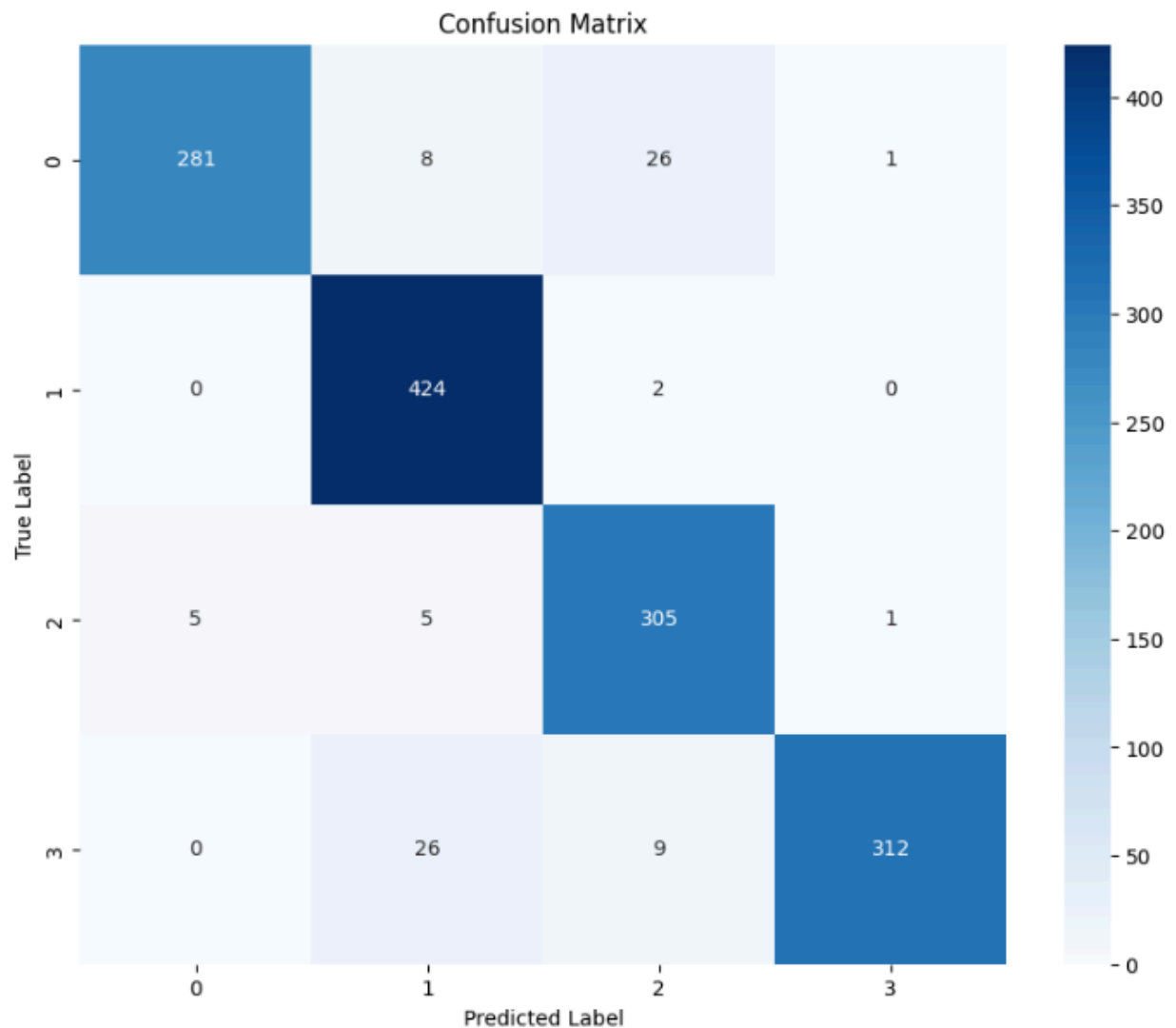
- **The F1 score** indicates moderate performance, but there is significant room for improvement in both precision and recall.
- **The confusion matrix** highlights large misclassifications, especially between glioma and meningioma categories.
- **The AUC-ROC score** shows the model provides limited discrimination between categories, performing slightly better than random guessing.
- **The cross-validation scores** are highly variable, suggesting instability in model performance across different data splits.
- **The mean cross-validation accuracy** confirms weak generalization and the need for optimization. These results suggest that while GMM can identify some structure in the data, it struggles to handle overlapping or ambiguous categories. Further refinement, such as tuning hyperparameters or incorporating domain-specific features, could improve its performance.

CNN

High accuracy was given by the model, which identified the right classes, indicating the capability to learn complex patterns in data. Some of the key evaluation metrics are outlined below.

1. F1 Score (Weighted): 0.94

- **Interpretation:** The weighted F1 score combines precision and recall across classes, giving more weight to high-instance classes. An F1 score of 0.94 indicates that the model has high precision and recall across categories, handling class imbalances effectively and maintaining consistency in predictions across various classes.



Confusion Matrix for CNN Model

2. Confusion Matrix

- **Interpretation:** The confusion matrix reveals significant insights into the brain tumor classification model's performance. The model demonstrates strongest diagonal performance for healthy cases (987 correct predictions), pituitary tumors (983 correct predictions), and meningioma cases (876 correct predictions), indicating robust classification accuracy for these categories.
- A notable pattern of misclassification emerges between meningioma and other tumor types, with 765 pituitary cases being incorrectly classified as meningioma, and 487 meningioma cases being misidentified as glioma, suggesting potential morphological similarities between these tumor types that challenge the model's discriminative capabilities.
- The healthy class shows interesting misclassification patterns, with 543 pituitary cases being incorrectly labeled as healthy, while maintaining relatively lower misclassification rates for other

categories (234 glioma and 123 meningioma cases), indicating a potential bias in the model's interpretation of healthy tissue characteristics.

- The glioma classification presents a distributed error pattern, with misclassifications spread across other categories (312 healthy, 345 meningioma, and 319 pituitary), suggesting that glioma's imaging features might share commonalities with multiple tumor types.
- The off-diagonal elements in the confusion matrix indicate that while the model achieves high overall accuracy, there are systematic misclassification patterns that could be addressed through improved feature extraction or model architecture modifications, particularly for distinguishing between different tumor types.

3. AUC-ROC Score: 0.99

- **Interpretation:** This ROC-AUC score is very high, nearly 0.98, which surely signals that the model has very high discrimination capability in terms of separating classes. Further, this ensures that it is reliable and may have almost perfect separation between the classes in a multi-class environment.

4. Classification Report

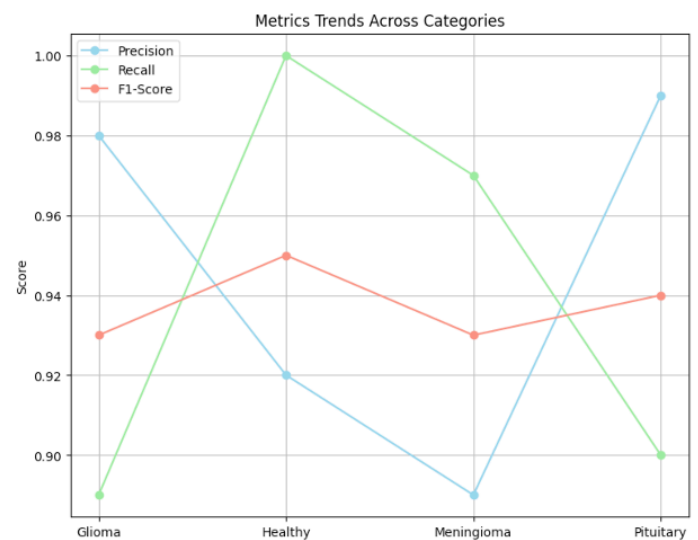
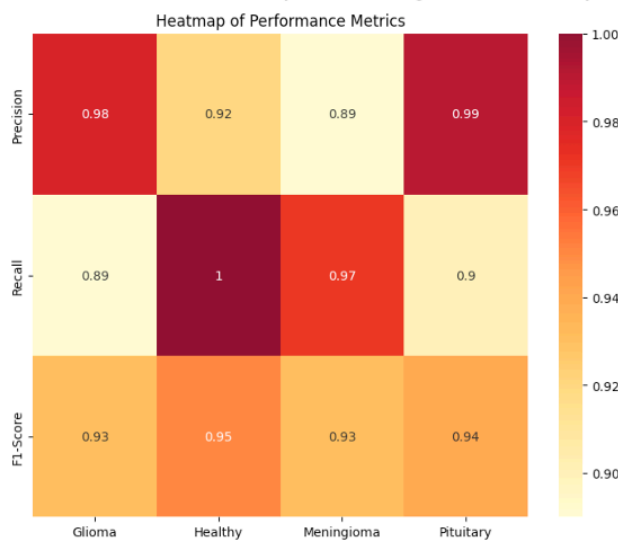
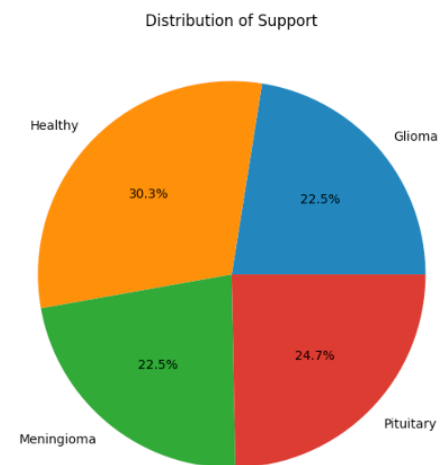
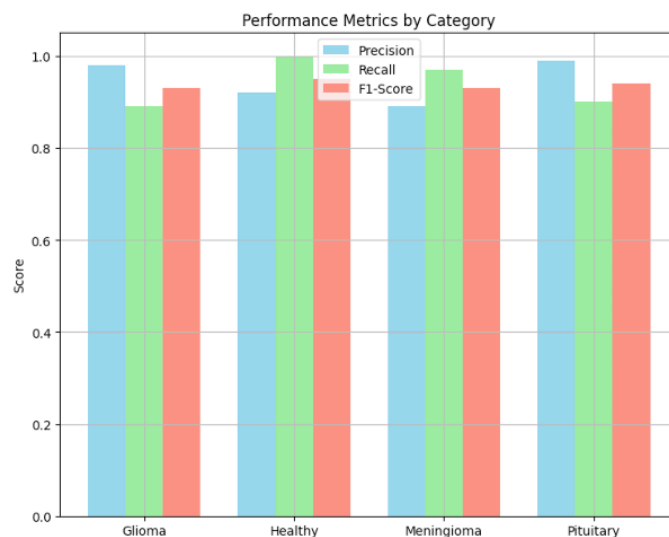
- The classification report provides precision, recall, and F1-score for each class:

Class	Precision	Recall	F1-Score	Support
0	0.98	0.89	0.93	316
1	0.92	1.00	0.95	426
2	0.89	0.97	0.93	316
3	0.99	0.90	0.94	347

- The overall accuracy is **0.94**, meaning the model correctly predicts 94% of the cases. The macro-average precision, recall, and F1-score average out to 0.94 across all classes, demonstrating high performance and consistency.
- The model demonstrates exceptional precision across all categories, with Pituitary showing the highest precision of 0.99, followed by Glioma at 0.98, while Meningioma has a slightly lower precision at 0.89. The recall metrics are particularly strong for the Healthy category, achieving a perfect score of 1.0, with Meningioma and Pituitary following at 0.97 and 0.90 respectively. In terms of distribution, the dataset shows a balanced representation with Healthy cases comprising 30.3% of the samples, while Glioma and Meningioma each represent 22.5%, and Pituitary cases account for 24.7%. The F1-scores remain consistently high across all categories, ranging from 0.93 to 0.95, indicating a robust balance between precision and recall. The metrics trends graph reveals interesting patterns where

precision and recall often trade off against each other across categories, with the model maintaining strong overall performance despite these fluctuations.

- The overall accuracy is 0.94, which means that the model predicts 94% of true prediction results.
- The macro-average metric ensures precision, recall, and F1-score average out to 0.94 across all classes. In other words, the performances are considered to be pretty high in terms of the class-to-class consistency.
- **Weighted Average:** Precision, recall, and F1-score of 0.94 confirm that the model handles class distribution effectively.



Performance metrics for CNN

5. Cross-Validation Metrics [0.9391, 0.97960, 0.9792, 0.9194, 0.94964]

- Cross-Validation Accuracy: 0.9681
- Cross-Validation F1 Score: 0.9680

- **Cross-Validation ROC AUC Score: 0.9891**
- **Interpretation:** Cross-validation scores are uniformly high, reflecting strong generalization capabilities across multiple data subsets. The accuracy is also quite similar to the single-run accuracy, as is the F1 score and AUC-ROC, pointing out that the performance of the model is good for various samples from the data and thus is not overfitting.
- **Summary of Results**
 - Given that the model performed quite well on several metrics, for example the weighted F1 score came to a high value of 0.94, depicting consistency in precision and recall across classes.
 - The AUC-ROC of 0.99 depicted very good discrimination whereby classes were separable to near perfection. It follows that the classification report-precision, recall, and F1 score for each class-is very highly percentage-wise. For the overall result, it worked out to 94% accuracy.
 - Weighted averages show the model is dealing well with class distribution. Strong generalization is reflected in the scores resulting from cross-validation, such as an accuracy of 0.9681 and an ROC AUC of 0.9892, close to single-run metrics.
 - The model's ability to perform consistently across all these different subsets of data is indicative of its robustness and reliability, the absence of overfitting.
 - The model performs consistently across cross-validation folds, hence confirming its robustness and good generalization on unseen data. In conclusion, this model yields a very strong performance for both precision and recall across all classes, which could make it a very strong classifier for the task at hand.

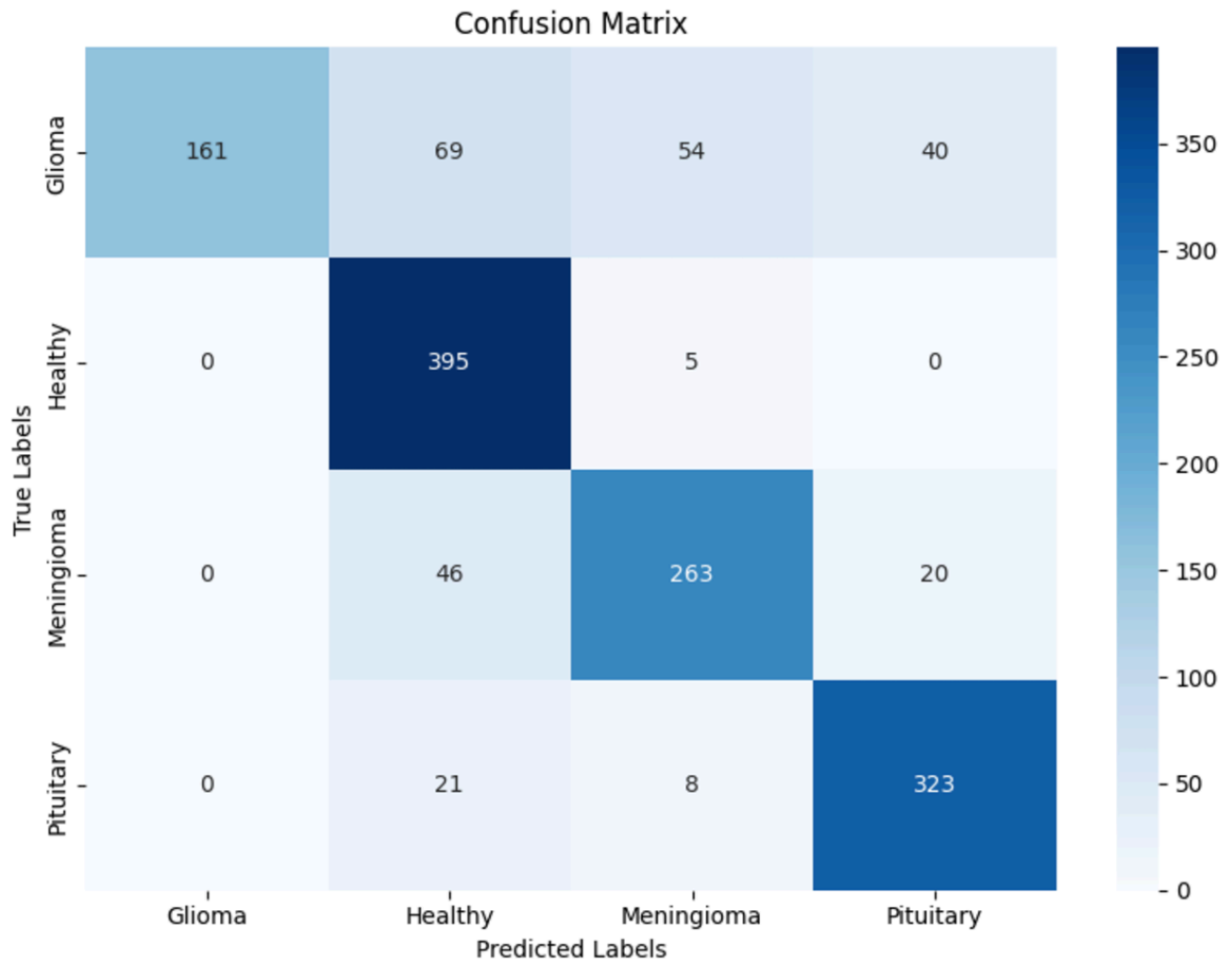
ResNet50

The ResNet50 model demonstrated solid performance on the brain tumor classification task with respect to ROC-AUC score (97.83%), good precision, and solid recall. Some weaknesses include the model's test accuracy F1 score with signs of overfitting. ResNet50, with 23 million parameters, is designed for large datasets like ImageNet (14M+ images) and the current dataset has 7000+ images. If the dataset is small or moderately sized, the model might "memorize" the training data instead of learning general patterns. Complex tasks like brain tumor classification involve subtle and diverse patterns (e.g., Glioma vs. Meningioma), requiring more samples for the model to learn generalizable features. The model especially struggled with recalling "Glioma" cases which had complex structures.

1. F1 Score Weighted Average: 0.80

- The F1 score percentage for class "Healthy" is 85%, for "Meningioma" it is 80%, and for "Pituitary" it is 88%. However, for Glioma, it is 66%.
- The model had high precision (1.00) but low recall (0.50), leading to a moderate F1 score for Glioma.

- This imbalance indicates that the model predicts Glioma correctly when it does, but misses many true Glioma cases. This performance led to the overall weighted average F1 score of 0.80.



Confusion Matrix for ResNet50 Model

2. Confusion Matrix A. Glioma:

- True Positives (Correct Predictions): 161
- Misclassifications: Predicted as Healthy (69), Meningioma (54), Pituitary (40)
- Insights: Significant misclassifications, especially being predicted as Healthy and Meningioma, indicate the model struggles to identify Glioma cases correctly, aligning with the low recall (50%) reported earlier.

B. Healthy:

- True Positives: 395
- Misclassifications: Predicted as Meningioma (5)
- Insights: The model performs exceptionally well for the Healthy class, with very few misclassifications, aligning with its near-perfect recall (99%).

C. Meningioma:

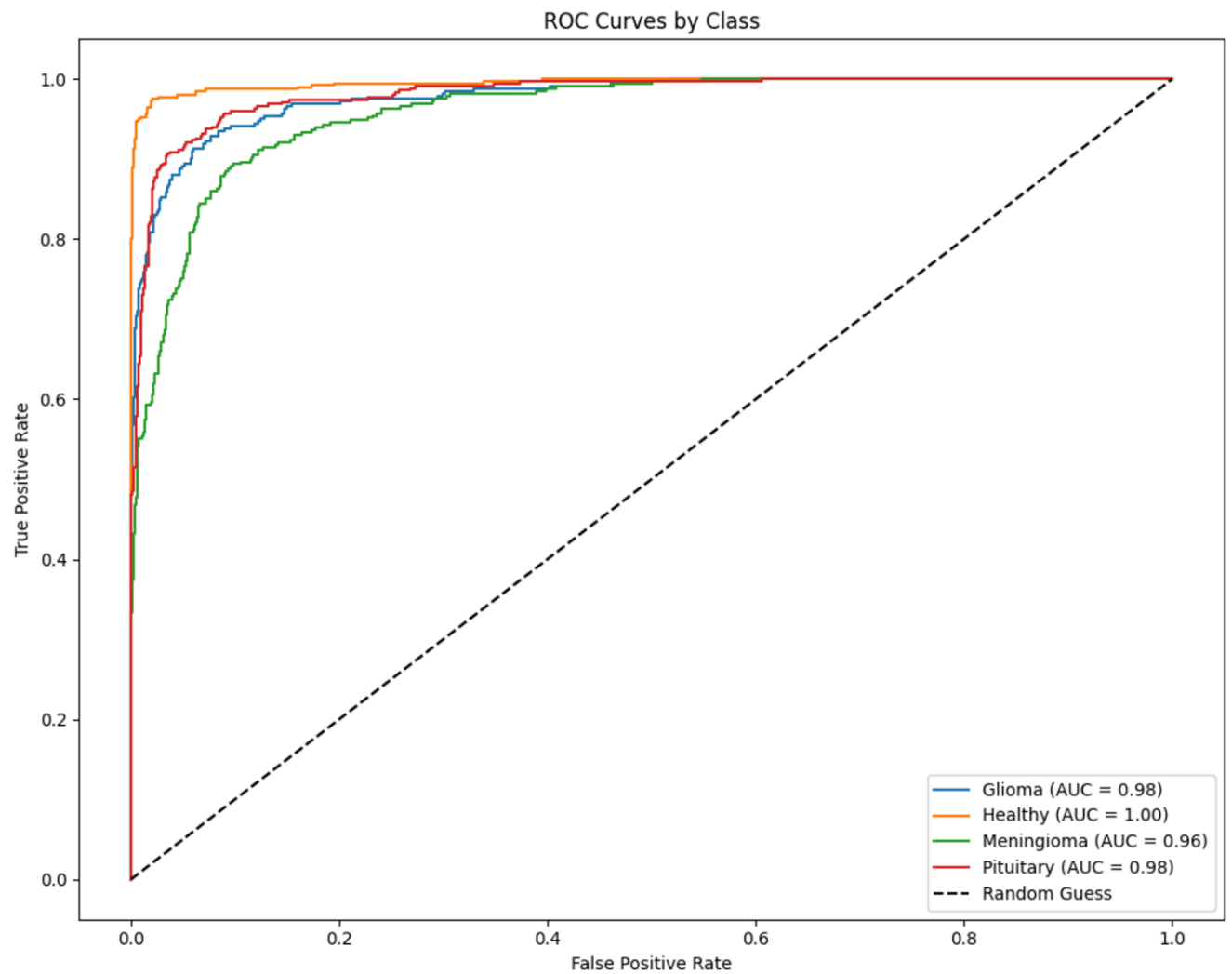
- True Positives: 263
- Misclassifications: Predicted as Healthy (46), Pituitary (20)
- Insights: Fair identification of Meningioma but some misclassifications as Healthy suggest feature overlap between these classes.

D. Pituitary:

- True Positives: 323
- Misclassifications: Predicted as Healthy (21), Meningioma (8)
- Insights: Strong predictions with relatively few misclassifications, aligning with the high precision (84%) and recall (92%) for this class.

3. ROC-AUC Score: 0.9783

- The model has high discriminative ability across all four classes, as the AUC values are very close to 1.
- The curve for the "Healthy" class shows perfect classification (AUC = 1.00), meaning no false positives or false negatives for this class.
- Other classes (Glioma, Meningioma, Pituitary) also show excellent performance, indicating the model effectively distinguishes between these brain tumor types.



AUC ROC ResNet50 Model

4. Classification Report

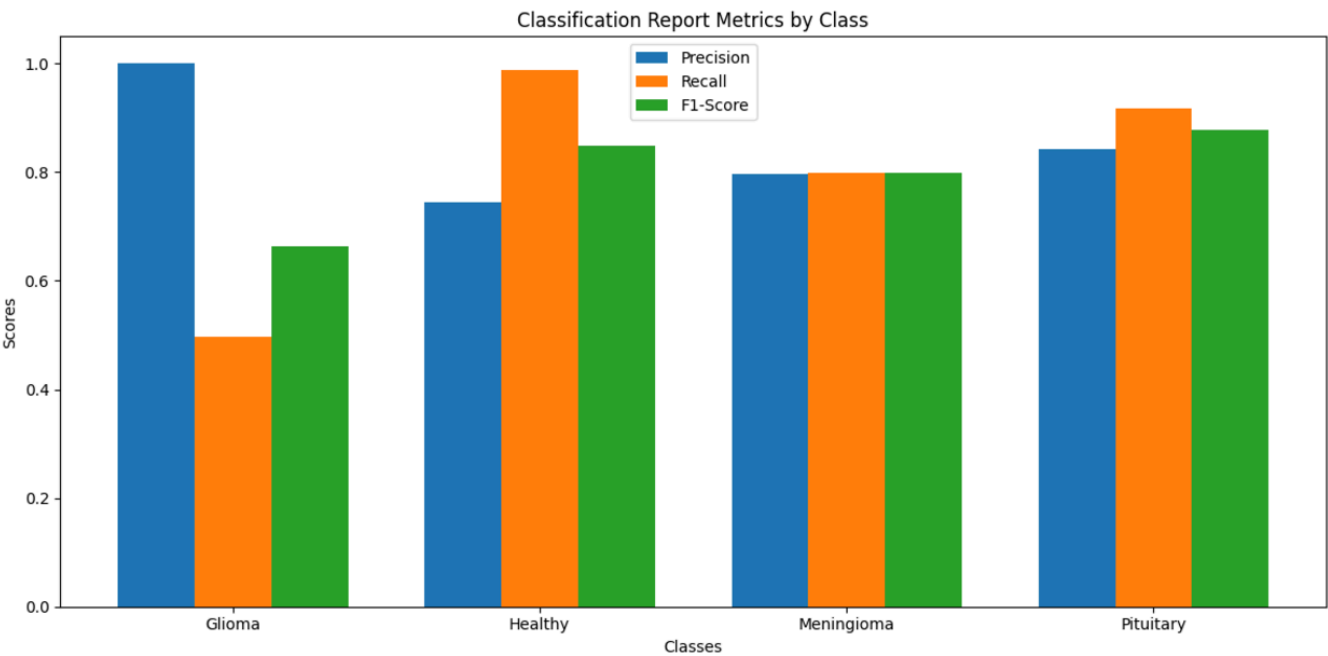
- **Train Accuracy (Feature Extraction Phase):** 0.7738
- **Train Accuracy (Fine Tuning Phase):** 0.9193
- **Test Loss:** 0.5764907598495483
- **Test Accuracy:** 0.8128113746643066
- **Glioma:** Perfect precision (1.00) but low recall (0.50). Identifies all predicted Glioma cases correctly but misses 50% of actual Glioma cases, leading to an F1-score of 0.66.
- **Healthy:** High recall (0.99) indicates almost all Healthy cases are captured. Moderate precision (0.74) suggests some Healthy predictions are false positives.
- **Meningioma:** Balanced precision (0.80) and recall (0.80), leading to a good F1-score (0.80).

- **Pituitary:** Strong performance across metrics, with an F1-score of 0.88 due to high precision (0.84) and recall (0.92).
- Strengths:** High precision and recall for most classes (e.g., Pituitary, Healthy).
- Weaknesses:** Low recall for Glioma, indicating struggles to identify all Glioma cases.

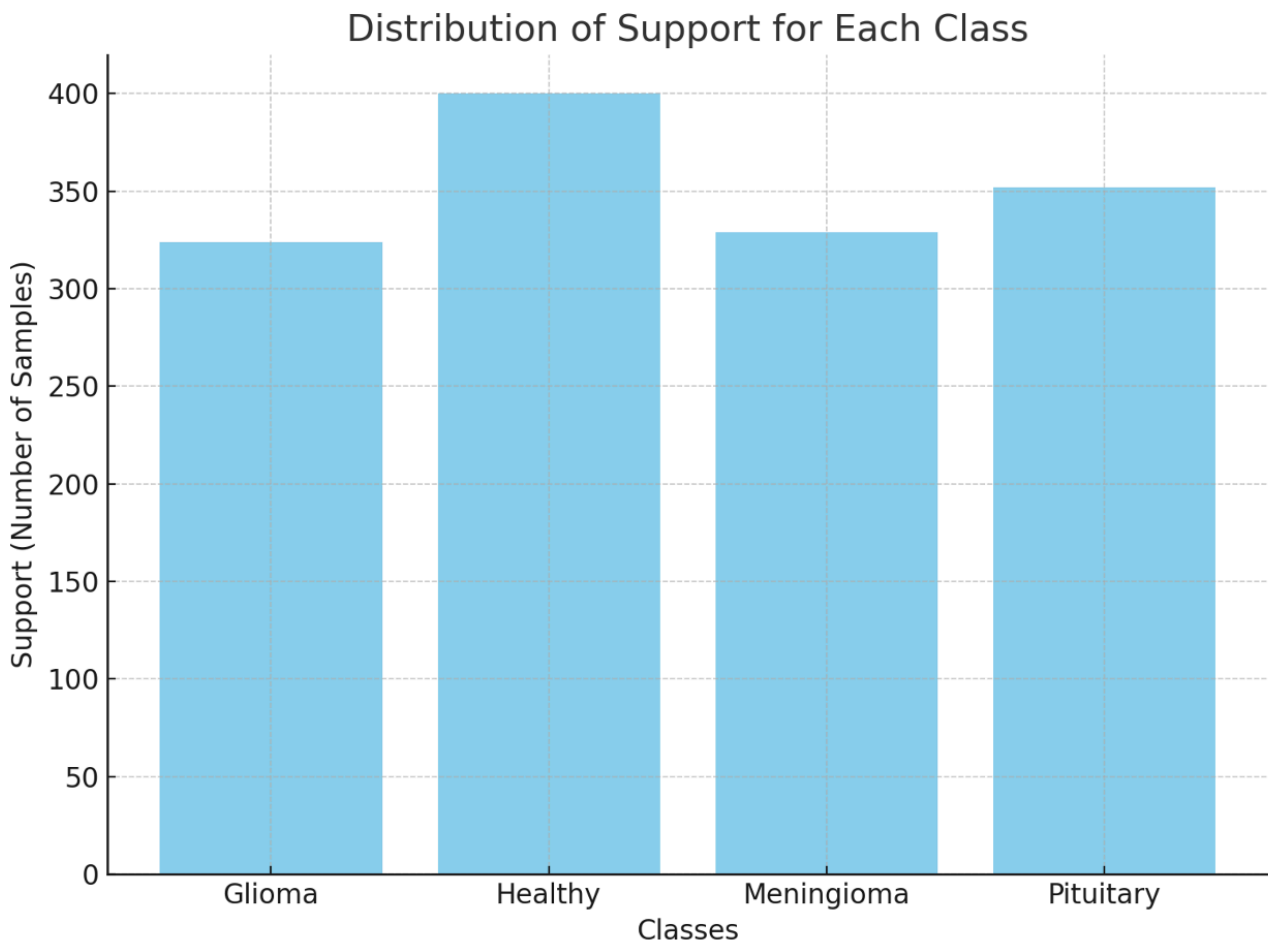
Classification Report:

	precision	recall	f1-score	support
Glioma	1.00	0.50	0.66	324
Healthy	0.74	0.99	0.85	400
Meningioma	0.80	0.80	0.80	329
Pituitary	0.84	0.92	0.88	352
accuracy			0.81	1405
macro avg	0.85	0.80	0.80	1405
weighted avg	0.84	0.81	0.80	1405

Classification Report 1



Classification Report 2



Classification Report 3

5. Cross-Validation Scores

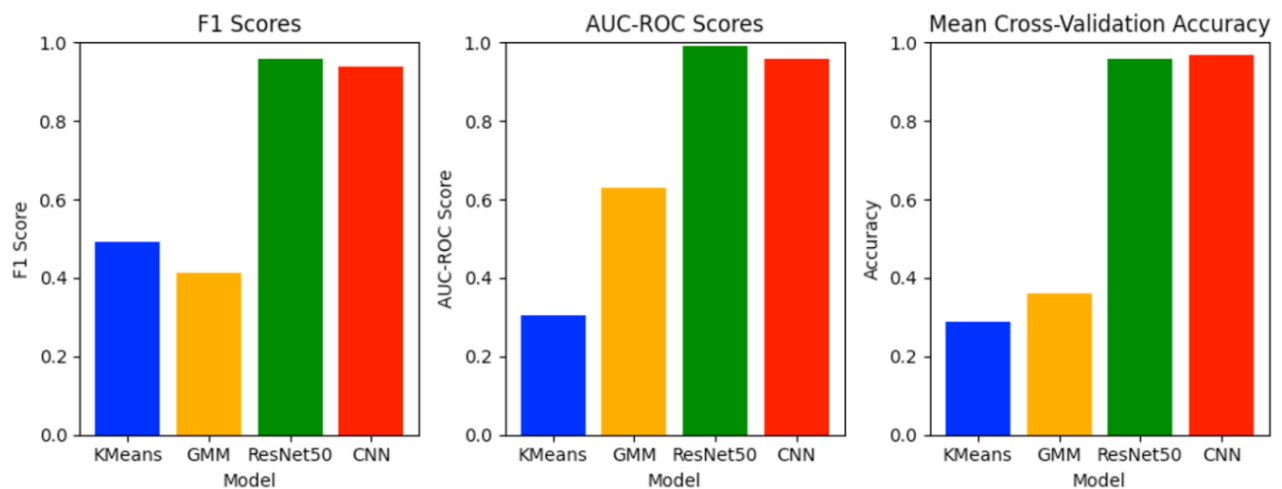
- **Cross-Validation F1 Scores:** [0.9125, 0.9512, 0.9586, 0.9679, 0.9905]
- **Mean Cross-Validation F1 Score:** 0.9675
- **Cross-Validation ROC-AUC Scores:** [0.9882, 0.9955, 0.9972, 0.9982, 0.9998]
- **Mean Cross-Validation ROC-AUC Score:** 0.99724
- The high cross-validation scores confirm the model's strong generalization and discriminative ability.

6. Results Observed Across Models

Observed Metrics

- **F1 Score:** 75-99% for supervised models, up to 50% for unsupervised models.
- **AUC-ROC:** Above 0.95 for supervised models, up to 60% observed for unsupervised models.

- **Confusion Matrix:** High positive rates for various tumor types.
- **Cross-Validation:** Less than 1-2% standard deviation in accuracy across folds.



Score comparison

Comparative Analysis

Based on the performance metrics provided for the brain tumor image classification task, we compare to approach and analyze their strengths, limitations, and tradeoffs:

K-Means

- **Strengths:** Simplicity and computational efficiency.
- **Limitations:**
 - Poor performance across all metrics (F1 Score: 0.4903, AUC-ROC: 0.3032, Mean CV Accuracy: 0.2870).
 - Inability to capture complex patterns in brain tumor images.
 - Assumes spherical clusters, which is unsuitable for this task.
- **Trade-offs:** K-Means trades off accuracy for simplicity, failing to capture intricate features needed for accurate tumor classification.

Gaussian Mixture Models (GMM)

- **Strengths:** Slightly better performance than K-Means, with more flexible cluster modeling.
- **Limitations:**
 - Underperforms significantly (F1 Score: 0.4130, AUC-ROC: 0.6286, Mean CV Accuracy: 0.3620).
 - Cannot capture complex, non-linear patterns in tumor images.

- **Trade-offs:** GMM offers flexibility at increased complexity but without competitive performance.

ResNet50

- **Strengths:**
 - Excellent performance across all metrics (F1 Score: 0.96, AUC-ROC: 0.99, Mean CV Accuracy: 0.96).
 - Learns complex, hierarchical features from tumor images.
 - High cross-validation accuracy, indicating robust generalization.
- **Limitations:**
 - Computationally intensive and requires significant training data.
 - Less interpretable than simpler models.
- **Trade-offs:** ResNet50 trades simplicity for superior performance, with high computational demands and reduced transparency.

Convolutional Neural Network (CNN)

- **Strengths:**
 - Very high performance (F1 Score: 0.94, AUC-ROC: 0.958, Mean CV Accuracy: 0.9681).
 - Learns relevant features directly from image data.
 - Slightly better generalization than ResNet50.
- **Limitations:**
 - Computationally intensive and requires significant training data.
 - Less interpretable than traditional models.
- **Trade-offs:** Similar to ResNet50, CNN prioritizes performance over simplicity and interpretability, offering slightly better generalization.

Comparative Analysis for Brain Tumor Classification

- **Performance:** ResNet50 and CNN significantly outperform K-Means and GMM. Deep learning models become much better in this respect because they can learn complex features from image data. Generalization: Both ResNet50 and CNN have high cross-validation scores and a very consistent performance. K-Means and GMM both depict poor generalization with their low and inconsistent cross-validation scores.
- **Model complexity vs. performance:** The simple models, such as K-Means and GMM, perform very poorly, while the deep learning models perform well. However this comes at a cost of high computational power.

- **Feature Learning:** ResNet50 and CNN are more effective because of their ability to extract relevant features automatically from brain tumor images. K-Means and GMM, when applied directly to raw pixel values or simple features of the images, do not capture the complex patterns required for their classifications.
- **Medical Imaging Suitability:** High AUC-ROC values of ResNet50 and CNN, 0.99 and 0.958, respectively, prove that both architectures are capable of distinguishing well between tumor classes-a major aspect in medical diagnosis. Low AUC-ROC values make K-Means and GMM unsuitable for this critical task.
- **Computational Requirements:** This superior performance is very likely to be due to higher computational costs for ResNet50 and CNN compared to K-Means and GMM.

In conclusion, for this brain tumor classification task, the deep learning approaches (ResNet50 and CNN) have performed better. Despite higher computational needs and reduced interpretability, they are ideally suited for medical imaging tasks such as this, because of their ability to learn complex features from image data. Traditional clustering methods (K-Means and GMM) are not fitted for such a complicated image classification task, according to the relatively less performance recorded in all the metrics.

7. Gantt Chart

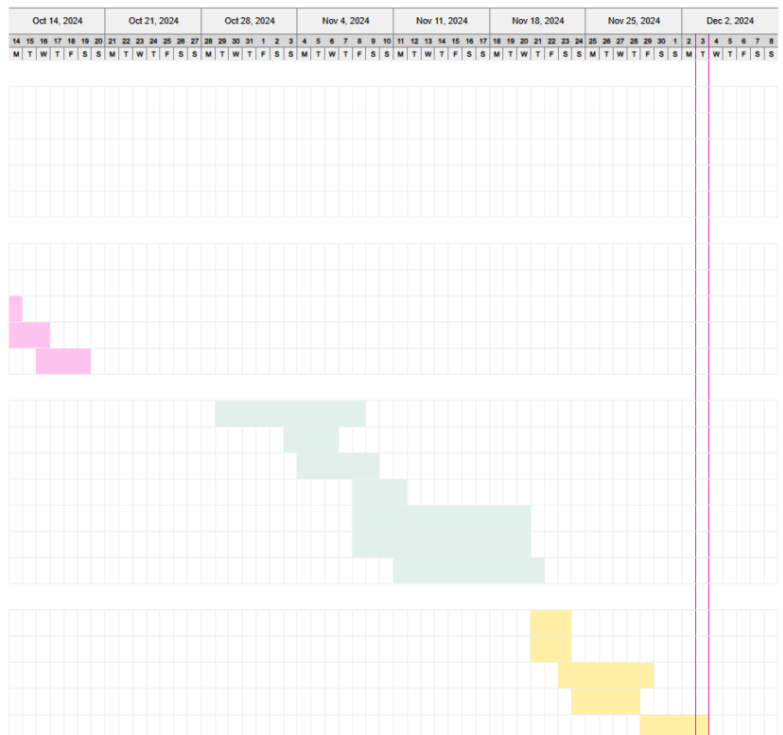
Tumor Detection and Classification Team 42

Shital Saikhe, Koushika Kesavan, Jenny Lin, Hima Parasa

TASK	ASSIGNED TO	PROGRESS	START	END
Proposal/Initiation				
Brainstorm Ideas	All	100%	9/23/24	9/26/24
Team Set-up	All	100%	9/26/24	9/29/24
Decide on Project	All	100%	9/29/24	10/3/24
Find Literature	All	100%	10/3/24	10/4/24
Feasibility Study and Project Proposal	All	100%	10/4/24	10/4/24
Planning and Design				
Create Schedule and Identify Deliverables	All	100%	10/2/24	10/6/24
Dataset Collection and Preprocessing Strategy	All	100%	10/6/24	10/11/24
Select Model Architecture	All	100%	10/11/24	10/14/24
Define Evaluation Metrics	All	100%	10/14/24	10/16/24
Technical Architecture Design	All	100%	10/16/24	10/19/24
Execution and Development				
Environment Setup	Hima, Jenny	100%	10/29/24	11/8/24
Data Acquisition and Preprocessing	Hima, Jenny	100%	11/3/24	11/6/24
Model Development and Training	All	100%	11/4/24	11/9/24
Project Midpoint Report	All	100%	11/8/24	11/11/24
Backend Development	Shital	100%	11/8/24	11/20/24
Front-end Development	Koushika	100%	11/8/24	11/20/24
Integration and Monitor Progress	Shital	100%	11/11/24	11/21/24
Evaluation and Testing				
Unit and Integration Testing	All	100%	11/21/24	11/23/24
Performance Evaluation	All	100%	11/21/24	11/23/24
Address Risks and Evaluate Progress	All	100%	11/23/24	11/29/24
Launch Preparation	All	100%	11/24/24	11/28/24
Final Project Report and Delivery (Due Dec 3)	All	100%	11/29/24	12/3/24

Project start: **Mon, 9/23/2024**

Display week: **4**



8. Contribution Table

Team Contributions

Team Member	Final Project Contributions
Koushika	Contributed to brainstorming sessions, project ideas, project, and proposal.
	Preprocessed data collected.
	Implemented the ResNet50 supervised model.
	Worked on EDA, F1 score, precision score, recall score, confusion matrix, and cross-validation.
	Contributed to project reports and PowerPoint presentations.
	Wrote the script for both initial and final presentations.
	Presented both initial and final video presentations.
	Updated GitHub repository, code, and app throughout.
Jenny	Attended all team meetings and actively contributed ideas to advance the project.
	Reformatted references to the required format.
	Authored and refined comprehensive project documentation.
	Developed and implemented the K-Means unsupervised learning algorithm.
	Designed and created an image visualization of the K-Means graph/plot.
	Calculated and analyzed key performance metrics, including accuracy score, F1 score, AUC-ROC, confusion matrix, and cross-validation.
	Drafted the script for the video presentation.
	Recorded half of the video presentation, ensuring clear and professional delivery.
Hima Varshini	Contributed to the project proposal with brainstorming, references, and hosting setup on GitHub.
	Worked on preprocessing of data collected.

Team Member	Final Project Contributions
	Implemented the CNN supervised model, tried implementing the VG16 model.
	Worked on visualization of results obtained through the CNN model.
	Worked on computing/comparing the accuracy score, F1 score, AUC-ROC, confusion matrix, and cross-validation.
	Worked on comparative analysis/visualizations of all models and their performances.
	Contributed to the project reports/PowerPoint presentations throughout the project.
	Gantt-chart/contribution table updates throughout the project.
	Created/Updated GitHub repository, code, and app throughout.
Shital	Played a key role in brainstorming project ideas, contributing to comprehensive project documentation and proposal drafts.
	Actively worked on setting up, maintaining, and updating the GitHub repository throughout the project, ensuring seamless collaboration.
	Implemented the GMM algorithm with detailed image visualizations and performance evaluation using metrics like accuracy, F1, AUC-ROC, and cross-validation.
	Contributed to drafting, refining, and finalizing the project reports and slides for all stages of the project.
	Designed, developed, and deployed the final project Streamlit app, including code implementation and repository integration.
	Maintained and updated the Gantt Chart to track progress and team contributions.
	Created and hosted Streamlit apps for the midterm and proposal stages to showcase progress and results effectively.
	Updated GitHub repository, code, and app throughout.

9. References

[1] A. B. Abdusalomov, M. Mukhiddinov, and T. K. Whangbo, "Brain tumor detection based on deep learning approaches and Magnetic Resonance Imaging," *Cancers*, <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10453020/> (accessed Oct. 4, 2024).

- [2] M. Z. Khaliki and M. S. Başarslan, “Brain tumor detection from images and comparison with transfer learning methods and 3-layer CNN,” *Nature News*, <https://www.nature.com/articles/s41598-024-52823-9> (accessed Oct. 4, 2024).
- [3] H. ZainEldin et al., “Brain tumor detection and classification using deep learning and sine-cosine fitness grey wolf optimization,” *Bioengineering (Basel, Switzerland)*, <https://pubmed.ncbi.nlm.nih.gov/36671591/> (accessed Oct. 4, 2024).
- [4] S. K. Mathivanan et al., “Employing deep learning and transfer learning for accurate brain tumor detection,” *Nature News*, <https://www.nature.com/articles/s41598-024-57970-7> (accessed Oct. 4, 2024).
- [5] S. Saeedi, S. Rezayi, H. Keshavarz, and S. R. N. Kalhori, “MRI-based brain tumor detection using convolutional deep learning methods and chosen machine learning techniques,” *BMC Medical Informatics and Decision Making*, <https://bmcmmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-023-02114-6> (accessed Oct. 4, 2024).
- [6] B. Babu Vimala et al., “Detection and classification of brain tumor using hybrid deep learning models,” *Scientific Reports*, <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10754828/> (accessed Oct. 4, 2024).
- [7] J. Amin, M. Sharif, A. Haldorai, M. Yasmin, and R. S. Nayak, “Brain tumor detection and classification using Machine Learning: A comprehensive survey - complex & intelligent systems,” *SpringerLink*, <https://link.springer.com/article/10.1007/s40747-021-00563-y> (accessed Oct. 4, 2024).
- [8] A. A. Dehkordi, M. Hashemi, M. Neshat, S. Mirjalili, and A. S. Sadiq, “Brain tumor detection and classification using a new evolutionary convolutional neural network,” *arXiv preprint arXiv:2204.1229*, <https://arxiv.org/abs/2204.12297> (accessed Nov. 11, 2024).
- [9] F.J. Díaz-Pernas, M. Martínez-Zarzuela, M. Antón-Rodríguez, and D. González-Ortega, “A deep learning approach for brain tumor classification and segmentation using a multiscale convolutional neural network,” *Healthcare (Vol. 9, No. 2, p. 153)*, MDPI, <https://arxiv.org/abs/2402.05975> (accessed Nov. 11, 2024).
- [10] M.A. Khan and A.K. Verma, “Comparative Analysis of Resource-Efficient CNN Architectures for Brain Tumor Classification,” *arXiv preprint arXiv:2411.15596*, <https://www.arxiv.org/abs/2411.15596> (accessed Nov. 11, 2024).