# BrainGAT

Graph Attention Network Approach for Multi-Class Brain Tumor Classification

**Team Members:**
- Ghandouz Amina
- Benghenima Hafsa
- Abdelmalek Lamia
- Yehyaoui Aya

**Professor:**
- Mr Khaldi Belkacem

# Contents

# List of Equations

# List of Figures

# 1 Introduction:

## 1.1 Context:

Brain tumors are among the most life-threatening neurological disorders, requiring precise diagnosis and treatment planning. Magnetic Resonance Imaging (MRI) has become the primary imaging technique for detecting and analyzing brain tumors due to its high spatial resolution and contrast. However, manual detection of tumor regions from MRI scans is time consuming, subjective and prone to variability between radiologists. Automated brain tumor detection and classification has therefore emerged as a critical task in medical image analysis, enabling faster and more consistent delineation of tumor boundaries. Achieving accurate classification directly impacts clinical decision-making, surgical planning and therapy assessment.

Traditional deep learning models, particularly Convolutional Neural Networks (CNNs), have achieved remarkable success in image classification. Yet, they struggle to model non-Euclidean relationships between image regions, a limitation when representing the complex spatial and structural dependencies of brain tissue. Graph Neural Networks (GNNs) overcome this challenge by representing an image as a graph, where nodes correspond to regions or superpixels and edges capture spatial or semantic relationships. Among GNN variants, the Graph Attention Network (GAT) introduces an attention mechanism that allows the model to focus on the most relevant neighboring nodes during feature aggregation. This adaptive capability makes GATs particularly suitable for medical imaging, where anatomical structures exhibit irregular geometries and context-dependent interactions.

## 1.2 Objectives:

This project aims to explore the use of Graph Attention Networks (GATs) for multi-class brain tumor classification using MRI data. Specifically, the objectives are:

- To analyze and understand the architecture and methodology proposed in the paper *Multi-class Brain Tumor Segmentation using Graph Attention Network*.

- To implement and evaluate a GAT-based classification pipeline using the Brain MRI Dataset from Kaggle.

- To assess the performance of the model in identifying and classifying different tumor types and analyze its advantages over conventional CNN approaches.

- To document findings and insights for potential extension toward future research or professional applications, such as clinical diagnostic support systems.

# 2 Literature Review:

## 2.1 Graph Neural Networks (GNNs):

### 2.1.1 Basics of Graphs:

Graphs are non-euclidean data structures (non-structure data representations) used to simulate data from complex real-world scenarios where regular structures such as images, audio, or sequential text might fail allowing unfixed sizes and forms. Some of these complex scenarios include brain networks, chemical components and the classic example of social networks. A **graph G** as **G = {N, E}**, this means **G** contains a set of nodes **N** that represent the elements of the graph and a set of edges **E** which determine the relationships between the nodes, nodes and edges can also be associated with features to indicate the characteristics of the entities and the relationships they represent. For example, a molecule of water can be modeled as a graph with three nodes: one for the oxygen and two for the hydrogen. Each node can include information related to the electric charge or the diameter of the atom. On the other hand, edge features can be used to distinguish between strong or weak bonds of the atoms they connect. Also, treating the molecule as a graph allows us to exploit the natural structure of the molecule without having to transform the sample into a different data type like an image or text. This is useful to better capture fine-grained relationships within the nodes and the different graph samples.

Another example is the brain network, which can be produced by associating certain regions of the brain with nodes and linking these nodes through the electrical signals they share. In reality, we can use any relevant information to determine when an edge should be added between the nodes. So, in this example, we could also connect nodes if they belong to brain regions that are relevant for certain processes like recognizing the face of a friend, learning a new word, etc.



Figure 1: Graph example

## 2.1.2 GNNs:

GNNs are a special type of neural network architecture specifically designed to handle ML problems related to graph data. It takes a whole graph as input where each node in the graph has a set of attributes generate the original embedding for that node. During training, we multiply these node embeddings with the (randomly initialized) neural layer's weights and apply an activation function to get the updated embedding. We further update this embedding by looking at the connected nodes and doing some sort of aggregation. We repeat this many times.

Training a model using back-propagation to adjust its weights to get the best embedding for each node is relatable to anyone who has ever trained any form of a neural network. The novelty and pecularity of a GNN is how we handle the topology of the graph. In other words, if a particular node **X** is connected to 4 neighbouring nodes, we somehow want the information from the 4 connected nodes to be reflected in the embedding of **X** node, those 4 nodes may be connected to several other nodes and information from those nodes should flow to them and reflect in their embeddings. Moreover, now that the 4 connected nodes of node **X** have updated embeddings, we need to consider updating the embedding of node **X** to reflect this new information and this is what we called *message passing* and happens via the aggregation operations like sum, mean etc where for each node, the embeddings from all the connected nodes are summed (assuming a sum operation) and this sum is made to reflect in the node's own embedding, we do it for all the nodes in the graph then recalculate each node's embedding again (since neighbour's embeddings would have got updated). We repeat a few times till there is stability. This is the core of a GNN, exchange information between neighbours via message passing until equilibrium is reached. There are many architectures that implement this core concept.



Figure 2: Single Node Aggregates Messages Example

$$h_v^{(k)} = \sigma \left( \sum_{u \in \mathcal{N}(v)} f\left(h_v^{(k-1)}, h_u^{(k-1)}, e_{vu}\right) \right) \tag{1}$$

### 2.1.3 Types of GNNs:

- **Graph Convolutional Network(GCN):**

  - **Idea:** Treats each node as a signal and performs a graph convolution using the normalized adjacency matrix.

  - **Use Case:** Semi-supervised learning on citation or social networks.

  - **Limitation:** Treats all neighbors equally and limited to shallow architectures.

- **Graph Attention Network(GAT):**

  - **Idea:** Learns attention weights between a node and its neighbors allowing more focus on important ones.

  - **Use Case:** Social graphs, influence modeling or noisy heterogeneous networks.

  - **Strength:** Adaptive and expressive without needing prior graph statistics.

- **Sample and Aggregate(GraphSAGE):**

  - **Idea:** Samples a fixed-size neighborhood and learns aggregation functions (mean, LSTM, pooling).

  - **Use Case:** Large-scale inductive learning where new nodes are added frequently.

  - **Strength:** Efficient, scalable, works for unseen data.

- **Graph Isomorphism Network(GIN):**

  - **Idea:** Proven to be as powerful as the Weisfeiler-Lehman graph isomorphism test using sum aggregators and MLPs

  - **Use Case:** Molecular property prediction or graph classification.

  - **Strength:** High expressive power.

## 2.2 Graph Attention Networks (GATs):

### 2.2.1 Overview:

GAT is designed to update each node's representation by aggregating information from its neighbor. But unlike earlier models, it learns how much attention to give to each neighbor. Let's assume each graph has three main ingredients: **Nodes (V)** which are the entities in the graph, **Edges (E)** the relationships or connections between nodes and **Features (X)** where each node has a feature vector (could be raw attributes like age or an embedding like a 512-dimensional text vector). At a high level, a GAT layer works like this:

- **Input:** Start with node features

- **Linear Transformation:** Project features into a new space.

- **Attention Scores:** Compute how important each neighbor is.

- **Softmax Normalization:** Convert scores into probabilities.

- **Aggregation:** Update each node by combining neighbors' features, weighted by attention.

- **Output:** A new representation for every node.

As an example, we'll use a graph with 4 nodes: A, B, C and D.



Figure 3: GAT: 4 Nodes Graph Example

Each node has a feature vector. For a feature dimension equals to 3 and the new dimension (hyperparameter) to 2, we get:

$$Node\ A:\ [1.0,\ 0.5,\ 0.2]$$

$$Node\ B:\ [0.9,\ 0.1,\ 0.3]$$
$$Node\ C:\ [0.4,\ 0.7,\ 0.8]$$
$$Node\ D:\ [0.2,\ 0.3,\ 0.9]$$

So the Feature Matrix **X** will be :

$$X = \begin{bmatrix} 1.0 & 0.5 & 0.2 \\ 0.9 & 0.1 & 0.3 \\ 0.4 & 0.7 & 0.8 \\ 0.2 & 0.3 & 0.9 \end{bmatrix}, \quad X \in \mathbb{R}^{4\times3}$$

And the Initialization Weight Matrix **W** will be:

$$W = \begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \end{bmatrix}, \quad W \in \mathbb{R}^{2\times3}$$

In the Linear Transformation step and using the weight matrix **W**, GAT projects node features into a new space to allowsthe model to learn a better representation of the features and reduce or expand the feature dimension depending on our task, we used:

$$h_i' \in \mathbb{R}^3 = W \in \mathbb{R}^{2\times3} * h_i \ \in \mathbb{R}^2 \tag{2}$$

The Transformed Feature Matrix **X'** will be:

$$X' = \begin{bmatrix} 0.26 & 0.77 \\ 0.20 & 0.59 \\ 0.42 & 0.99 \\ 0.35 & 0.77 \end{bmatrix}, \quad X' \in \mathbb{R}^{4\times2}$$

For the Attention Mechanism step, we compute attention scores which decide how important one node is to another when aggregating information. For each edge (i,j), we have:

$$e_{ij} = \text{LeakyReLU}\left(a^\top [h_i' \,\|\, h_j']\right) \tag{3}$$

Where
$$a \in \mathbb{R}^{2F'} = [0.5, 0.6, 0.7, 0.8]$$

is a learnable attention vector. We get:

$$e_{AB} = 1.204$$
$$e_{AC} = 1.678$$
$$e_{BA} = 1.252$$
$$e_{BD} = 1.315$$
$$e_{CA} = 1.602$$

$$e_{DB} = 1.249$$

These values are unnormalized attention scores. They can be any real number. therefore in the next step which is Softmax Normalization we normalize them across each node's neighbors using softmax so that they become comparable (like probabilities).

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}(i)} \exp(e_{ik})} \tag{4}$$

So we get:

| Operation | $\alpha_{ij}$ |
|-----------|---------------|
| $A \rightarrow B$ | 0.383 |
| $A \rightarrow C$ | 0.616 |
| $B \rightarrow A$ | 0.484 |
| $B \rightarrow D$ | 0.515 |
| $C \rightarrow A$ | 1.000 |
| $D \rightarrow B$ | 1.000 |

Now that we have attention coefficients $\alpha_{ij}, each node updates its representation by aggregating the features of its neighbors, weighted by the attention scores, then passing through a nonlinearity as following$ :

$$h_i'' = \sigma\left( \sum_{j \in \mathbb{N}(i)} \alpha_{ij} * h_j' \right) \tag{5}$$

So we get:

*Node A: [0.335,0.836]*
*Node B: [0.306,0.769]*
*Node C: [0.260,0.770]*
*Node D: [0.200,0.590]*

That was a single attention head (transform node features, compute attention scores, normalize with softmax and aggregate neighbor features) but in practice GATs don't just use one attention mechanism, we use multiple heads in parallel because A single head might put too much weight on one neighbor, multiple heads smooth things out and each head learns its own weight matrix and attention vector, this means different heads can focus on different aspects of the neighborhood. For **K** heads, the updated representation of node **i** is:

$$h_i'' = \Big\|_{k=1}^{K} \sigma\left( \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^{(k)} h_j' \right) \tag{6}$$

### 2.2.2 Advantages of GATs over GCNs:

Graph Attention Networks offer several improvements over traditional Graph Convolutional Networks. The primary advantage is the attention mechanism which allows GAT to assign different importance weights to neighboring nodes during aggregation, rather than treating all neighbors equally or using fixed normalized weights as GCN does.

A significant strength of GAT is its ability to handle variable neighborhood sizes naturally. Unlike GCN which relies on normalization factors dependent on node degrees, GAT computes attention coefficients that automatically adapt to different numbers of neighbors without requiring manual adjustment.

GAT also learns adaptive relationships without needing explicit graph normalization. The attention mechanism computes edge weights based on node features, eliminating the need for preprocessing steps like symmetric normalization of the adjacency matrix that GCN requires. This makes GAT more flexible when working with diverse graph structures.

Additionally, GAT is more interpretable due to attention coefficients. These learned weights provide direct insights into which neighboring nodes are most influential for each node's representation, offering a level of explainability that is difficult to achieve with GCN's fixed aggregation scheme. This interpretability is valuable for understanding model decisions and debugging in practical applications.

### 2.2.3 Limitations:

Despite their effectiveness and adaptability, GATs present several limitations that affect their scalability and computational efficiency like:

- **High Computational Cost for Large Graphs:** The attention mechanism requires computing pairwise attention coefficients between each node and its neighbors and for large-scale graphs with millions of nodes or high average degrees, this results in significant computational and memory overhead where the complexity grows approximately linearly with the number of edges

- **Limited Scalability for Very Dense Graphs:** In dense graphs, each node connects to many others, leading to an explosion in attention operations which doesn't only increase computation time but also makes it challenging to store intermediate attention matrices in memory.

- **Over-smoothing and Over-fitting:** When stacking multiple GAT layers, node embeddings tend to become overly similar (over-smoothing) and due to the limited amount of labeled graph data in many domains, GATs are prone to over-fitting especially when trained on small datasets.

- **Sensitivity to Hyperparameters:** The performance of GATs heavily depends on careful tuning of hyperparameters such as the number of attention heads, learning rate and dropout rate. Improper choices can lead to unstable training or poor convergence.

Recent research has proposed various improvements, such as Sparse GATs, Hierarchical GATs and Graph Transformers, to mitigate these challenges by reducing computational complexity and enhancing scalability on large or dense graphs.

# 3 Related Works on Image Classification Using GNNs:

## 3.1 Traditional approaches:

### 3.1.1 Early CNN-based Architectures:

Traditional medical image classification tasks including brain tumor analysis have primarily relied on Convolutional Neural Networks (CNNs). Architectures such as U-Net, Fully Convolutional Network (FCN) and DeepLab have been widely used for segmentation and classification due to their ability to automatically extract hierarchical spatial features. U-Net introduced an encoder–decoder structure enabling precise localization by combining low-level spatial details with high-level semantic information while FCN extended CNNs to dense prediction tasks by replacing fully connected layers with convolutional layers, allowing pixel-wise classification. On the other hand, DeepLab incorporated dilated convolutions and Conditional Random Fields (CRFs) to enhance multi-scale context aggregation and boundary refinement. These architectures achieved significant success on MRI and CT imaging datasets, enabling automated detection and classification of tumors with high accuracy.

### 3.1.2 Limitations of CNN-based Methods:

Despite their dominance, CNNs are inherently limited by their grid-based convolutional operations:

- They assume Euclidean structure (regular grids), making it difficult to model irregular or non-local relationships between image regions or brain structures.

- Contextual dependencies between distant pixels or anatomically related regions are often poorly captured due to local receptive fields.

- CNNs require large labeled datasets and their feature aggregation is constrained to spatial proximity, limiting generalization to complex topologies or small data regimes.

These shortcomings highlight the need for graph-based models that can flexibly represent non-Euclidean structures such as relationships between superpixels, tissue regions or features across slices, leading to the adoption of Graph Neural Networks (GNNs) in medical image analysis.

## 3.2 GNN-based classification models:

With the rise of Graph Neural Networks (GNNs), medical image analysis has shifted toward non-Euclidean representations, allowing models to reason about complex spatial and structural relationships beyond grid-based convolutional constraints.

### 3.2.1 Superpixel-based Graph Representations:

A prominent line of research constructs graphs from medical images by representing superpixels or image patches as nodes, while edges encode spatial or feature similarities, this approach captures local topology and semantic connectivity between anatomical regions. For instance, in brain MRI classification tasks, each superpixel can represent a homogeneous tissue region and the graph edges can model intensity correlations or spatial adjacency. Such representations enable compact, topology-aware learning, often leading to better tumor boundary delineation and robustness to noise compared to pixel-level CNNs.

### 3.2.2   Advanced GNN Architectures: Graph U-Net and Dynamic GCNs:

Recent GNN architectures such as Graph U-Net introduce hierarchical graph pooling and unpooling mechanisms analogous to CNN encoder–decoder designs, allowing multi-scale reasoning on graph-structured data. Similarly, Dynamic Graph Convolutional Networks (DGCNs) adaptively update graph connections during training, enabling the model to learn task-driven relationships among brain regions or organ structures rather than relying on static adjacency matrices. These architectures have shown strong performance in classifying diseases and anatomical regions by learning structural dependencies not captured by traditional convolutional filters.

### 3.2.3   Hybrid CNN–GNN Frameworks:

Another emerging paradigm combines the feature extraction capability of CNNs with the relational reasoning power of GNNs. Typically, CNNs first extract local features from MRI or CT scans, which are then projected onto graph nodes representing spatial regions or feature clusters. The GNN layer then models inter-region dependencies, global context and spatial coherence, producing more interpretable and biologically consistent predictions. Such CNN–GNN hybrid systems have been successfully applied to brain tumor classification, Alzheimer's detection and organ segmentation, demonstrating superior generalization across datasets with limited training samples.

## 3.3   Medical imaging focus:

GNNs have demonstrated strong potential across diverse medical imaging tasks where spatial dependencies and anatomical connectivity play a critical role in diagnosis and classification.

### 3.3.1   Brain Tumor Detection and Classification

In brain imaging, GNNs are increasingly used to capture structural dependencies between brain regions or tissue segments. Superpixel-based and region-adjacency graph representations allow models to reason over the global context of tumor spread and its interaction with surrounding tissues.

Compared to CNNs, which process images locally, GNNs can integrate long-range spatial relationships improving the differentiation of tumor types such as glioma, meningioma and pituitary tumors. This structural modeling has proven especially useful in MRI-based classification frameworks, where graph connectivity reflects both geometry and tissue heterogeneity.

### 3.3.2   Retinal Vessel and Ophthalmic Image Analysis:

In ophthalmology, GNNs have been applied to retinal vessel segmentation and classification, modeling the vascular network as a graph structure. Nodes represent vessel junctions or segments, while edges encode their connectivity patterns and by propagating information across the vessel graph, GNNs can better capture hierarchical vascular organization, leading to more accurate identification of pathological changes such as diabetic retinopathy or glaucoma-related anomalies.

### 3.3.3   Histopathology and Cellular-Level Analysis:

At the microscopic scale, GNNs have been employed to analyze histopathology images by treating cells or nuclei as graph nodes connected through spatial proximity or morphological similarity. This graph-based modeling enables the network to learn cellular interactions and tissue architecture, which

are crucial for distinguishing between benign and malignant tissue patterns. Such relational learning surpasses pixel-based convolutional methods by directly exploiting tissue topology and microstructural organization.

### 3.3.4 Why GNNs Excel in Medical Imaging:

Across these domains, GNNs outperform conventional methods when the underlying spatial or anatomical relationships matter. Their ability to encode non-Euclidean structures, propagate contextual information and reason about inter-regional dependencies makes them particularly suited for medical imaging tasks involving complex spatial hierarchies and biologically structured data.

## 3.4 Research trends:

Recent research in graph-based medical image analysis has increasingly focused on enhancing the representational power and interpretability of GNNs through advanced architectural designs. A dominant trend involves attention-based and hierarchical GNNs that integrate pixel-level and region-level contextual information within a unified framework. Attention mechanisms, as introduced in Graph Attention Networks (GATs), enable models to adaptively weight node relationships, allowing the network to focus on the most diagnostically relevant regions in brain MRIs or histopathology images. This selective aggregation improves both classification accuracy and explainability, providing insights into which structures contribute most to a clinical prediction.

Hierarchical approaches, on the other hand, combine multi-scale reasoning where low-level pixel features capture local texture patterns while high-level regional graphs represent broader anatomical dependencies, by merging these complementary levels of information, hierarchical GNNs achieve superior spatial consistency and robustness to variations in imaging modalities or tumor morphology.

Overall, the field is moving toward hybrid, multi-level graph frameworks that unify local feature extraction, global reasoning and interpretability, setting the foundation for next-generation models such as BrainGAT, which leverage attention mechanisms to enhance relational learning in complex brain tumor classification tasks.

# 4 Analysis of the Chosen Paper:

## 4.1 Paper Overview:

**Title:** Multi-class Brain Tumor Segmentation using Graph Attention Network

**Authors:** Dhrumil Patel, Dhaivat Patel, Rudra Saxena, and T. Akilan

**Publication Venue:** 2023 8th International Conference on Signal and Image Processing (ICSIP 2023)

**Year:** 2023 (Published on arXiv in February 2023)

### 4.1.1 Research Problem and Motivation:

Brain tumor segmentation from MRI is essential for diagnostic radiology and there is significant demand for automatic segmentation algorithms to overcome practical limitations of manual approaches. The manual segmentation of brain tumors by radiologists is time-consuming, subjective and prone to inter-observer variability. The complexity and heterogeneity of brain tumors, combined with the need for precise delineation of tumor boundaries for treatment planning, motivated the development of automated solutions.

The authors addressed the challenge of leveraging both spatial relationships and contextual information in volumetric MRI data, which traditional convolutional approaches often struggle to capture effectively. Their work sought to exploit recent advancements in graph neural networks to model the structural relationships between different brain regions more explicitly.

### 4.1.2 Novelty and Contribution:

The paper's primary novelty lies in its representation of volumetric MRI data as a Region Adjacency Graph (RAG) and the application of Graph Attention Networks for tumor segmentation. Key contributions include:

1. **Graph-Based MRI Representation:** Converting 3D MRI volumes into graph structures where nodes represent brain regions and edges encode spatial adjacencies and relationships.

2. **Attention Mechanism for Medical Imaging:** Applying GAT's attention mechanism to dynamically weight the importance of neighboring regions, allowing the model to focus on tumor-relevant features while suppressing irrelevant background information.

3. **Multi-class Segmentation:** Extending the approach to handle multiple tumor sub-regions including whole tumor, tumor core and enhancing tumor components.

4. **Competitive Performance:** The model achieved mean dice scores of 0.91, 0.86 and 0.79, and mean Hausdorff distances in the 95th percentile of 5.91, 6.08 and 9.52 mm for whole tumor, core tumor and enhancing tumor segmentation respectively on the BraTS2021 validation dataset.

## 4.2 Methodology:

### 4.2.1 Data Preprocessing and Graph Construction:

The methodology transforms traditional volumetric MRI data into a graph-based representation through the following steps:

- **MRI Preparation:**

  – Multi-modal MRI sequences (T1, T1-contrast enhanced, T2 and FLAIR) are preprocessed through standard steps including skull stripping, bias field correction and intensity normalization.

  – Images are co-registered to ensure alignment across different modalities.

  – The volumetric data is standardized to a consistent resolution.

- **Graph Construction:**

  – The volumetric MRI is converted into a Region Adjacency Graph (RAG) representation.

  – Node Definition where each node in the graph represents a super-pixel or region in the MRI volume, extracted using segmentation techniques like SLIC (Simple Linear Iterative Clustering) or watershed algorithms.

  – Node Features where each node is characterized by multi-modal intensity features, texture descriptors and spatial coordinates extracted from the corresponding region.

  – Edge Definition where edges connect spatially adjacent regions, with edge weights potentially encoding similarity measures between neighboring regions based on intensity, texture or spatial proximity.

This graph representation enables the model to capture both local neighborhood relationships and long-range dependencies across the brain volume.

### 4.2.2 GAT Architecture:

1. **Network Structure:** The Graph Attention Network architecture consists of multiple attention layers stacked sequentially:

   - **Input Layer:** Takes the node feature matrix and adjacency matrix as inputs.

   - **Graph Attention Layers:** Multiple GAT layers (2-4 layers) where each layer applies attention mechanisms to aggregate information from neighboring nodes.

   - **Attention Mechanism:** For each node, the attention coefficients are computed using a learnable weight matrix transforms node features, the attention scores are calculated between node pairs using a shared attention mechanism, softmax normalization ensures attention weights sum to 1 then multi-head attention is employed to capture different aspects of node relationships.

   -

2. **Key Hyperparameters:**

   - **Number of Attention Heads:** 4-8 heads per layer to capture diverse relational patterns.

   - **Hidden Dimensions:** 64-256 units depending on layer depth.

   - **Dropout Rate:** Applied to attention coefficients and features (commonly 0.3-0.5) to prevent overfitting.

   - **Activation Functions:** LeakyReLU or ELU activations between layers.

---

3. **Output Layer:**

   - A final classification layer maps the learned node representations to tumor class probabilities.

   - Softmax activation produces probability distributions over tumor classes (background, whole tumor, core, enhancing regions)

4. **Training Details:** The model is trained using a combination of loss functions to handle class imbalance and improve segmentation quality:

   - **Dice Loss:** Focuses on maximizing overlap between predictions and ground truth.

   - **Cross-Entropy Loss:** Provides pixel-wise classification supervision.

   - A weighted combination of both losses balances region-level and pixel-level optimization.

5. **Optimizer:**

   - Adam optimizer with an initial learning rate typically in the range of 1e-3 to 1e-4.

   - Learning rate scheduling with decay or cyclic strategies to improve convergence.

   - Weight decay (L2 regularization) to prevent overfitting.

6. **Training Configuration:**

   - **Epochs:** Training conducted for 100-300 epochs depending on dataset size.

   - **Batch Processing:** Due to graph representation, batching strategies are adapted to handle variable graph sizes.

   - **Data Augmentation:** Random rotations, flips, scaling and elastic deformations applied to training samples.

7. **Evaluation Metrics:**

   - **Dice Similarity Coefficient (DSC):** Measures overlap between predicted and ground truth segmentations.

   - **Hausdorff Distance (HD95):** Quantifies boundary accuracy and worst-case segmentation errors.

   - **Sensitivity and Specificity:** Assess true positive and true negative rates.

   - **Intersection over Union (IoU):** Alternative overlap metric complementing Dice score.

## 4.3   Dataset Description:

- **Dataset**: BraTS2021 (Brain Tumor Segmentation Challenge 2021).

- **Source:** The dataset is publicly available through the Medical Image Computing and Computer Assisted Intervention (MICCAI) BraTS challenge, containing multi-institutional pre-operative MRI scans of glioblastoma and lower-grade glioma patients.

- **Dataset Composition:**

    - **Training Set:** 1,251 cases with expert-annotated ground truth labels.

    - **Validation Set** 219 cases for model evaluation during challenge participation.

    - The dataset includes both high-grade gliomas (HGG) and low-grade gliomas (LGG).

- **Classes and Annotations:** The dataset provides voxel-wise annotations for four main categories:

    1. **Necrotic and Non-Enhancing Tumor Core (NCR/NET):** Hypointense regions representing dead tissue.

    2. **Peritumoral Edema (ED):** Surrounding swelling visible in FLAIR sequences.

    3. **Enhancing Tumor (ET):** Active tumor regions enhanced in T1-contrast images.

    4. **Background:** Healthy brain tissue.

- **MRI Modalities:** Each case includes four co-registered MRI sequences:

    - T1-weighted (T1).

    - T1 contrast-enhanced (T1ce).

    - T2-weighted (T2).

    - Fluid-Attenuated Inversion Recovery (FLAIR).

- **Preprocessing:**

    - **Skull Stripping:** Brain extraction to remove non-brain tissues.

    - **Co-registration:** All modalities aligned to the same anatomical space.

    - **Resampling:** Volumes standardized to 1mm$^3$ isotropic resolution.

    - **Intensity Normalization:** Z-score normalization or histogram matching applied to each modality independently.

    - **Patch Extraction:** For graph construction, the volume may be divided into overlapping or non-overlapping patches.

- **Dataset Split:** While the paper used the official BraTS2021 split, typical training strategies involve:

    - **Training:** 70-80% of available annotated data.

    - **Validation:** 10-15% for hyperparameter tuning.

    - **Testing:** 10-15% or using the official challenge validation/test sets.

– Cross-validation strategies (e.g., 5-fold) may be employed for robust performance estimation.

## 4.4 Experimental Results:

### 4.4.1 Quantitative Performance:

The GAT-based model demonstrated strong performance across multiple evaluation metrics on the BraTS2021 validation dataset:

- **Dice Similarity Coefficient:**
  - Whole Tumor (WT): 0.91
  - Tumor Core (TC): 0.86
  - Enhancing Tumor (ET): 0.79

- **Hausdorff Distance (95th percentile):**
  - Whole Tumor: 5.91 mm
  - Tumor Core: 6.08 mm
  - Enhancing Tumor: 9.52 mm

These results indicate that the model performs best on whole tumor segmentation, which typically has clearer boundaries, while the enhancing tumor presents the greatest challenge due to its smaller size and irregular boundaries.

### 4.4.2 Comparison with Baseline Models:

The paper compared the GAT approach against traditional segmentation architectures:

- **CNN-based Methods:** Standard 3D CNNs without explicit relational modeling typically achieved Dice scores 3-5% lower than the GAT model, particularly struggling with tumor core and enhancing regions where spatial relationships are critical.

- **U-Net Architecture:** The ubiquitous U-Net baseline, while effective for general segmentation, showed limitations in capturing long-range dependencies and contextual relationships that the graph-based approach explicitly models. The GAT model outperformed U-Net variants by approximately 2-4% in Dice score across tumor regions.

- **Other Graph-Based Methods:** Compared to standard Graph Convolutional Networks (GCNs) without attention mechanisms, the GAT model's ability to dynamically weight neighbor importance led to improved performance, especially in heterogeneous tumor regions where not all neighboring regions are equally relevant.

### 4.4.3 Qualitative Analysis:

- **Segmentation Visualization:** Visual comparisons of segmentation outputs revealed that the GAT model produced:
  - Smoother and more coherent tumor boundaries compared to CNN baselines.
  - Better handling of small, disconnected tumor components.

- Improved detection of subtle enhancing regions that CNNs occasionally missed.

- More accurate delineation of tumor core boundaries where multiple tissue types intersect.

- **Attention Map Analysis:** Visualization of learned attention weights demonstrated that the model dynamically focuses on:

  - Regions with high contrast boundaries between tumor and healthy tissue.

  - Areas where multiple modalities provide complementary information.

  - Spatially distant but histologically similar regions, capturing non-local tumor patterns.

## 4.5 Key Findings and Insights:

### 4.5.1 How GAT Improved Segmentation Performance:

1. **Explicit Relational Modeling:** By representing MRI data as graphs, the GAT architecture explicitly captures spatial relationships and anatomical connectivity that are implicit or harder to learn in purely convolutional approaches.

2. **Attention-Based Feature Aggregation:** The attention mechanism allows the model to selectively focus on relevant neighboring regions while down-weighting less informative areas, leading to more discriminative feature representations.

3. **Multi-Scale Context Integration:** Graph structures naturally accommodate both local and global context, enabling the model to integrate information across different spatial scales without the architectural complexity of multi-scale CNN designs.

4. **Robustness to Irregular Structures:** Unlike grid-based convolutions, graph representations handle irregular tumor shapes and sizes more naturally, as the graph topology adapts to the underlying image structure.

### 4.5.2 Observed Limitations and Failure Cases:

1. **Small Tumor Detection:** The model occasionally struggled with very small enhancing tumor regions, particularly when they were isolated or had subtle contrast differences from surrounding tissue.

2. **Boundary Precision:** While overall segmentation quality was high, fine-grained boundary delineation remained challenging in regions with gradual intensity transitions or unclear anatomical boundaries.

3. **Computational Complexity:** Graph construction and GAT operations are more computationally expensive than standard convolutions, leading to longer training and inference times, particularly for high-resolution 3D volumes.

4. **Graph Construction Dependency:** Model performance is sensitive to the initial graph construction process, including the choice of super-pixel algorithm and parameters, which may require task-specific tuning.

5. **Class Imbalance:** Like many medical imaging tasks, the extreme imbalance between tumor and background classes (and between tumor sub-regions) required careful loss function design and training strategies.

### 4.5.3  Paper Conclusions and Future Directions:

- **Main Conclusions**: The authors concluded that Graph Attention Networks provide a promising alternative to purely convolutional approaches for brain tumor segmentation, effectively leveraging structural relationships in medical imaging data. The attention mechanism's ability to dynamically weight spatial relationships proved particularly valuable for handling heterogeneous tumor appearances.

- **Suggested Future Directions:**

  1. **Hybrid Architectures:** Combining the strengths of CNNs for local feature extraction with GATs for relational reasoning could yield further performance improvements.

  2. **Multi-Task Learning:** Extending the approach to simultaneously perform segmentation and other related tasks such as tumor grading or survival prediction.

  3. **Temporal Analysis:** Adapting the graph-based framework for longitudinal studies to track tumor progression over time by incorporating temporal edges in the graph structure.

  4. **Efficient Graph Construction:** Developing learnable or adaptive graph construction methods that optimize the graph topology during training rather than relying on fixed preprocessing steps.

  5. **Clinical Translation:** Investigating the model's performance on diverse clinical datasets beyond research challenges to assess real-world applicability and generalization.

  6. **Interpretability:** Further exploration of attention weight visualization to provide clinicians with interpretable insights into the model's decision-making process.

While the original paper focused on pixel-wise segmentation of tumor regions, our work adapts the GAT framework for tumor classification tasks. Instead of predicting class labels for each voxel, we modified the architecture to produce image-level predictions distinguishing between glioma, meningioma, pituitary tumors and healthy brain tissue. The core principles of graph representation and attention-based feature aggregation remain applicable, but with adjusted output layers and evaluation metrics appropriate for classification rather than segmentation.

# 5 Problem Definition:

The goal of this work is to perform multi-class brain tumor classification from magnetic resonance imaging (MRI) scans using a Graph Attention Network (GAT) framework. Given the complex structural dependencies in brain anatomy and the heterogeneous appearance of tumor regions, the task demands models capable of capturing both local visual patterns and non-local contextual relationships.

## 5.1 Formal Definition:

Let's denote the dataset by $\mathcal{D} = (X_i, y_i)_{i=1}^{N}$ where each ( $X_i$ ) is an MRI image and ( $y_i \in \{1, 2, 3, 4\}$ ) corresponds to one of four tumor classes:

- Glioma.

- Meningioma.

- Pituitary tumor.

- No tumor.

Each image ( $X_i$ ) is represented as a graph ( $G_i = (V_i, E_i)$ ), where:

- ( $V_i$ ) are nodes corresponding to superpixels or image regions.

- ( $E_i$ ) are edges capturing spatial adjacency or feature similarity.

- each node ( $v \in V_i$ ) is associated with a feature vector ( $\mathbf{x}_v \in \mathbb{R}^d$ ).

The objective of the GAT model ( $f_\theta$ ) is to learn a mapping $f_\theta : G_i \rightarrow y_i$.

that predicts the correct tumor class given the graph-structured representation of the MRI.

## 5.2 Input–Output Relationship:

- **Input:** Brain MRI image preprocessed into a graph structure with node-level features derived from pixel intensity, texture, or CNN-based embeddings.

- **Output:** A categorical prediction indicating the tumor type present in the image.

- **Intermediate Representation:** The GAT learns attention coefficients between connected nodes, reflecting the relative importance of spatially or semantically related regions in the classification process.

## 5.3 Challenges:

1. **Class Imbalance:** The dataset exhibits unequal class distribution, particularly between tumor and non-tumor categories, which can bias learning.

2. **Spatial Dependency:** Tumor appearance depends on complex spatial and anatomical contexts, simple pixel-level models struggle to capture these dependencies.

3. **Inter-patient Variability:** MRI scans vary significantly in orientation, intensity, and contrast due to acquisition differences.

4. **Small Sample Size:** Compared to natural image datasets, medical datasets are limited in scale, making deep learning prone to overfitting.

5. **Boundary Ambiguity:** Tumor borders often blend into healthy tissue, challenging precise region-based feature extraction.

These challenges motivate the adoption of graph-based approaches particularly attention-driven GNNs to model contextual relationships and improve the robustness of classification across heterogeneous MRI data.

# 6 Methodology:

## 6.1 Data Preprocessing and MRI Preparation:

Our preprocessing pipeline transforms standard MRI images into graph-structured data suitable for Graph Attention Network processing. The pipeline consists of several sequential stages:

- **Image Loading and Normalization:** Each brain MRI image is loaded and processed through the following steps:

  - Images are resized to a standardized dimension of 224×224 pixels to ensure uniform input size across the dataset.

  - Grayscale images are converted to RGB format for consistent processing.

  - Pixel intensity values are normalized to the $[0, 1]$ range by dividing by 255.

  - A corresponding grayscale version is maintained for texture feature extraction.

- **Superpixel Segmentation:** We employ the SLIC (Simple Linear Iterative Clustering) algorithm to partition each MRI image into 80 superpixels. This segmentation serves as the foundation for graph construction, where each superpixel becomes a node in our graph representation. The SLIC parameters are configured as follows:

  - **Number of segments:** 80 (balancing computational efficiency with spatial granularity).

  - **Compactness:** 10 (controlling the trade-off between color similarity and spatial proximity).

  - **Sigma:** 1 (Gaussian smoothing parameter for preprocessing).

This superpixel-based approach offers several advantages over pixel-level processing: reduced computational complexity, incorporation of local spatial context and more meaningful feature extraction from coherent image regions.

## 6.2 Graph Construction:

### 6.2.1 Node Feature Extraction:

For each superpixel (node) in the segmented image, we extract a comprehensive 25-dimensional feature vector encompassing intensity, color, texture, and shape characteristics:

- **Intensity Features (5 dimensions):** From the grayscale representation, we compute statistical descriptors of pixel intensities within each superpixel:

  - **Mean intensity:** average brightness of the region.

  - **Standard deviation:** intensity variability indicating homogeneity.

  - **Minimum and maximum intensity:** dynamic range of the region.

  - **Median intensity:** robust central tendency measure.

- **Color Features (6 dimensions):** From the RGB representation, we extract channel-wise statistics:

  - Mean values for R, G, and B channels, capturing average color composition.

– Standard deviations for R, G, and B channels, measuring color variability within the region.

- **Texture Features (5 dimensions):** We compute Gray Level Co-occurrence Matrix (GLCM) features to capture texture patterns, which are particularly informative for distinguishing tumor types:

  – **Contrast:** measures local intensity variations.

  – **Dissimilarity:** quantifies difference between adjacent pixels.

  – **Homogeneity:** assesses uniformity of texture.

  – **Energy:** indicates texture orderliness.

  – **Correlation:** describes linear dependency of gray levels

  The GLCM is computed using four orientations (0°, 45°, 90°, 135°) at distance 1, and the resulting texture descriptors are averaged across orientations for rotational invariance.

- **Shape Features (9 dimensions):** Geometric and morphological properties of each superpixel region:

  – **Area:** number of pixels in the region.

  – **Perimeter:** boundary length.

  – **Eccentricity:** deviation from circularity.

  – **Solidity:** proportion of convex hull area filled by the region.

  – **Extent:** ratio of region area to bounding box area.

  – **Compactness:** circularity measure computed as $4area/perimeter$.

  – **Major and minor axis lengths:** principal dimensions of the fitted ellipse.

  – **Aspect ratio:** ratio of major to minor axis lengths.

After extraction, all 25 features are standardized using z-score normalization (subtracting mean and dividing by standard deviation) to ensure equal contribution during model training.

### 6.2.2 Edge Construction and Graph Validation:

We construct a Region Adjacency Graph where edges connect spatially adjacent superpixels that share boundaries. This is accomplished using the rag_mean_color function, which creates edges between neighboring regions based on their spatial proximity in the original image. The resulting graph captures the spatial topology of the brain MRI, allowing the GAT model to propagate information along anatomically meaningful pathways. Edges are created bidirectionally to enable information flow in both directions during message passing. As a fallback mechanism, if no edges are detected (which rarely occurs), we create a chain connection between consecutive nodes to ensure graph connectivity.

Before feeding graphs to the model, we perform validation to ensure edge indices remain within valid bounds corresponding to the actual number of nodes, preventing potential indexing errors during training.

## 6.3 GAT Architecture:

Our Graph Attention Network architecture consists of multiple attention layers followed by graph-level pooling and classification heads:

- **Input Layer:**
  - GATConv layer: transforms 25-dimensional input features to 128 hidden dimensions.
  - Multi-head attention with 6 heads operating in parallel.
  - Head concatenation: outputs 768 dimensions (128×6).
  - Batch normalization applied for training stability.
  - ELU (Exponential Linear Unit) activation for smooth gradients.
  - Dropout (p=0.35) for regularization.

- **Hidden Layers:** Our architecture includes 1 intermediate GAT layer (resulting in 3 total layers):
  - Input: 768-dimensional features from previous layer.
  - GATConv with 128 dimensions per head and 6 attention heads.
  - Concatenated output: 768 dimensions.
  - Batch normalization and ELU activation.
  - Dropout (p=0.35).

- **Output GAT Layer:**
  - Input: 768-dimensional features.
  - GATConv with 128 dimensions and single attention head (no concatenation).
  - Output: 128-dimensional node embeddings.
  - Batch normalization and ELU activation.

- **Graph-Level Pooling:** To obtain a fixed-size representation for the entire graph (necessary for image-level classification), we apply two complementary global pooling operations:
  - Global mean pooling: averages node features across all nodes, capturing average characteristics.
  - Global max pooling: takes element-wise maximum across nodes, highlighting salient features.
  - Concatenation: combines both representations into a 256-dimensional graph embedding.

  This dual pooling strategy allows the model to capture both typical and distinctive features of the tumor regions.

- **Classification Head:**
  - Fully connected layer: 256 → 128 dimensions.
  - Batch normalization and ELU activation.
  - Dropout (p=0.35).

- Output layer: $128 \rightarrow 4$ dimensions (one per class).

- Log-softmax activation for classification probabilities.

The complete architecture contains 749,316 trainable parameters.

- **Attention Mechanism:** The GAT layers implement the attention mechanism as follows: for each node, attention coefficients are computed for all neighboring nodes, determining the importance of each neighbor's features. These coefficients are learned during training and allow the model to dynamically focus on relevant spatial relationships while suppressing less informative connections. Multi-head attention enables the model to attend to different aspects of node relationships simultaneously.

## 6.4 Training Strategy:

### 6.4.1 Loss Function:

We employ Negative Log-Likelihood (NLL) loss, which is the standard choice for multi-class classification with log-softmax outputs $L = -log(p(y|x))$ where $p(y|x)$ is the predicted probability of the correct class. This loss function penalizes incorrect predictions more heavily when the model is confident, encouraging well-calibrated probability estimates.

### 6.4.2 Optimizer:

AdamW (Adam with Weight Decay) optimizer with the following configuration:

- Initial learning rate: 0.001

- Weight decay: 1e-4 (L2 regularization to prevent overfitting)

- Default $\beta_1 = 0.9$, $\beta_2 = 0.999$ for adaptive moment estimation

AdamW was chosen over standard Adam for its decoupled weight decay, which has been shown to improve generalization in deep learning models.

### 6.4.3 Learning Rate Scheduling:

We implement ReduceLROnPlateau scheduling to adaptively adjust the learning rate:

- Monitoring metric: validation accuracy.

- Reduction factor: 0.5 (halves the learning rate).

- Patience: 7 epochs (waits for 7 epochs without improvement before reducing).

- Mode: maximize (since we're tracking accuracy).

This strategy prevents the model from getting stuck in local minima and enables fine-grained optimization as training progresses.

### 6.4.4 Gradient Clipping:

To prevent exploding gradients during backpropagation through the graph structure, we apply gradient norm clipping with a maximum norm of 1.0. This ensures stable training, particularly important for graph neural networks where message passing can amplify gradients.

### 6.4.5 Regularization Techniques:

Multiple regularization strategies are employed to combat overfitting:

1. Dropout (0.35) after each GAT layer and in the classification head.

2. Batch normalization after every layer to stabilize activations.

3. Weight decay (1e-4) in the optimizer.

4. Early stopping with patience of 15 epochs.

### 6.4.6 Early Stopping:

Training terminates if validation accuracy does not improve for 15 consecutive epochs. The model weights corresponding to the best validation accuracy are saved and used for final evaluation.

### 6.4.7 Hyperparameters Summary:

- **Image size:** 224×224 pixels.

- **Number of superpixels:** 80 per image.

- **Batch size:** 16 graphs.

- **Hidden dimensions:** 128.

- **Number of attention heads:** 6 (except output layer with 1 head).

- **Number of GAT layers:** 3.

- **Dropout rate:** 0.35.

- **Learning rate:** 0.001 (with adaptive scheduling).

- **Maximum epochs:** 300.

- **Early stopping patience:** 15 epochs.

# 7    Implementation Details:

## 7.1    Development Environment:

- **Framework and Libraries:** The implementation leverages the following Python ecosystem:
  - **PyTorch:** Version 1.x for deep learning framework.
  - **PyTorch Geometric (PyG):** Specialized library for graph neural networks, providing efficient implementations of GATConv, global_mean_pool, global_max_pool and graph batching utilities.
  - **scikit-image:** For SLIC superpixel segmentation, regionprops feature extraction, and GLCM texture analysis.
  - **OpenCV (cv2):** Image loading, resizing, and color space conversions.
  - **NumPy:** Numerical computations and array operations.
  - **scikit-learn:** Dataset splitting, evaluation metrics (accuracy, classification report).
  - **tqdm:** Progress bars for data processing visualization.
- **Computing Resources:**
  - Hardware: GPU-accelerated training (CUDA-enabled device).
  - The model automatically detects and utilizes available GPU resources, falling back to CPU if necessary.
  - Training time: Approximately 4-6 hours for 300 epochs on a modern GPU.

## 7.2    Dataset Handling:

1. **Dataset Source:** Brain Tumor MRI Dataset from Kaggle containing four classes:
   - **Glioma:** malignant brain tumor arising from glial cells.
   - **Meningioma:** typically benign tumor of the meninges.
   - **Pituitary:** tumor of the pituitary gland.
   - **No Tumor (Healthy):** normal brain MRI scans.

2. **Dataset Statistics:**
   (a) **Training Set:**
       - **Glioma:** 1,321 images.
       - **Meningioma:** 1,339 images.
       - **No Tumor:** 1,595 images.
       - **Pituitary:** 1,457 images.
       - **Total:** 5,712 images.
   (b) **Testing Set:**

- **Glioma:** 300 images.

- **Meningioma:** 306 images.

- **No Tumor:** 405 images.

- **Pituitary:** 300 images.

- **Total:** 1,311 images.

3. **Dataset Split:** The provided training set is further divided into training and validation subsets:

   - **Training:** 4,569 graphs (80

   - **Validation:** 1,143 graphs (20

   - **Testing:** 1,311 graphs (held-out test set).

   The split is stratified by class label to maintain class distribution across subsets, ensuring representative sampling of all tumor types in each partition.

4. **Data Loading:** Images are processed sequentially for each class, with progress tracking via tqdm. Each image undergoes:

   (a) Graph construction (superpixel segmentation + feature extraction).

   (b) Feature normalization.

   (c) Edge index creation.

   (d) Label assignment.

   (e) Conversion to PyTorch Geometric Data object.

   Error handling is implemented to skip corrupted or unreadable images without interrupting the entire process.

5. **Batching Strategy:** PyTorch Geometric's DataLoader automatically handles variable-size graphs by:

   - Creating mini-batches of multiple graphs.

   - Concatenating node features and adjusting edge indices.

   - Tracking graph membership via a batch vector.

   - Batch size: 16 graphs per mini-batch.

6. **Data Augmentation:** While not explicitly implemented in the current pipeline, the superpixel-based graph construction provides implicit regularization by creating slightly different graph structures due to SLIC's stochastic nature. Future improvements could incorporate:

   - Random rotation and flipping during image loading.

   - Elastic deformations.

   - Brightness and contrast adjustments.

   - Gaussian noise injection.

## 7.3   Training and Validation Procedures:

- **Training Loop:** For each epoch, the following steps are executed:

  1. Model set to training mode (enables dropout and batch norm updates).

  2. Iterate through training batches:

     - **Forward pass:** compute predictions.

     - **Loss calculation:** NLL loss between predictions and ground truth.

     - **Backward pass:** compute gradients.

     - **Gradient clipping:** prevent exploding gradients.

     - **Optimizer step:** update model parameters.

  3. Track cumulative loss and accuracy across all batches.

  4. Return epoch-level metrics.

- **Validation Loop:** After each training epoch:

  1. Model set to evaluation mode (disables dropout, uses running statistics for batch norm).

  2. Iterate through validation batches with gradient computation disabled.

  3. Collect predictions and ground truth labels.

  4. Compute validation accuracy.

  5. Update learning rate scheduler based on validation performance.

  6. Save model if validation accuracy improves (checkpoint best model).

  7. Increment early stopping counter if no improvement.

- **Evaluation Metrics:** During training and validation we used accuracy and loss but for final test evaluation we used accuracy, precision, recall, F1-score and support (number of true instances per class).

- **Model Checkpointing:** The best-performing model (based on validation accuracy) is saved to disk as best_gat_model.pt. This checkpoint is loaded for final test set evaluation, ensuring results reflect the model's optimal state rather than the final epoch's potentially overfit weights.

- **Monitoring and Logging:** Training progress is monitored through:

  - Epoch-level metrics printed to console (loss, train accuracy, validation accuracy).

  - Learning rate adjustments logged when triggered.

  - Best model notifications when validation accuracy improves.

  - Early stopping announcements.

- **Reproducibility:**
  - Random seed set to 42 for train-validation split.
  - Stratified splitting ensures consistent class distributions.
  - Deterministic operations where possible (though SLIC has inherent randomness).

# 8 Results and Discussion:

## 8.1 Quantitative Results:

Our Graph Attention Network achieved strong performance on the brain tumor classification task where the overall test accuracy was 96.43%. This result significantly outperforms many traditional machine learning and standard CNN approaches, demonstrating the effectiveness of graph-based representation for brain tumor classification.

| | | | | |
|---|---|---|---|---|
| Glioma | 0.99 | 0.88 | 0.93 | 300 |
| Meningioma | 0.95 | 0.94 | 0.94 | 306 |
| No Tumor | 0.96 | 0.97 | 0.97 | 405 |
| Pituitary | 0.93 | 0.96 | 0.95 | 300 |
| Macro Average | 0.96 | 0.84 | 0.95 | 1311 |
| Weighted Average | 0.96 | 0.86 | 0.96 | 1311 |

## 8.2 Qualitative Analysis:

### 8.2.1 Attention Mechanism Insights:

The Graph Attention Network's ability to learn meaningful attention weights is crucial to its performance. Through the multi-head attention mechanism:

1. **Spatial Relationship Modeling:** The model learns to weigh connections between adjacent superpixels based on their relevance to tumor classification. Superpixels within tumor regions likely receive higher attention from their neighbors compared to background regions.

2. **Multi-Scale Feature Capture:** With 6 attention heads in the initial layers, the model can simultaneously focus on different aspects:

   - Some heads may focus on intensity patterns characteristic of different tumor types.

   - Other heads may emphasize texture features distinguishing malignant from benign formations.

   - Additional heads could capture shape and boundary characteristics.

3. **Tumor-Background Discrimination:** The attention mechanism helps the model distinguish tumor-bearing regions from healthy tissue by learning to suppress attention weights for irrelevant background superpixels while amplifying weights for tumor-relevant regions.

### 8.2.2 Graph Representation Advantages:

The superpixel-based graph construction offers several qualitative benefits:

1. **Semantic Coherence:** Unlike pixel-wise processing, superpixels represent coherent anatomical regions, making the extracted features more interpretable and meaningful.

2. **Computational Efficiency:** Processing 80 nodes per image (vs. 224×224=50,176 pixels) dramatically reduces computational requirements while preserving critical information.

3. **Robustness to Noise:** Aggregating features over superpixel regions provides natural noise reduction compared to pixel-level processing.

### 8.2.3 Confusion Patterns:

Based on the classification results, we observe:

1. **Glioma-Meningioma Confusion:** Glioma's lower recall suggests occasional misclassification as meningioma or pituitary tumors, likely due to overlapping intensity and texture patterns in certain cases.

2. **Clean Separation of Healthy Cases:** The No Tumor class shows minimal confusion, indicating the model reliably distinguishes tumor presence from absence.

3. **Pituitary Reliability:** High recall for pituitary tumors reflects their distinctive location and morphological characteristics.

## 8.3 Comparison with Baseline Models:

### 8.3.1 Advantages Over Traditional CNNs:

While we don't have direct comparative experiments in this implementation, we can infer several theoretical and empirical advantages of our GAT approach:

1. **Explicit Spatial Modeling:** CNNs learn spatial relationships implicitly through convolutional filters, while GATs explicitly model spatial adjacencies through the graph structure. This explicit modeling can be more sample-efficient for medical imaging where spatial relationships carry diagnostic significance.

2. **Parameter Efficiency:** Our model contains 749,316 parameters, which is substantially fewer than typical ResNet or DenseNet architectures (commonly 10M-25M parameters) used for medical image classification, reducing the risk of overfitting.

3. **Interpretability:** Attention weights provide interpretable insights into which regions the model considers important for classification, whereas CNN feature maps are often less directly interpretable.

4. **Handling Irregular Structures:** Tumors exhibit irregular shapes and sizes; graph representations naturally accommodate this variability without the grid-based constraints of CNNs.

### 8.3.2 Comparison with U-Net (Segmentation Baseline):

While U-Net architectures are designed for segmentation rather than classification, adapted U-Nets are sometimes used for classification by:

- Adding global pooling layers after the encoder.

- Training end-to-end for image-level labels.

Compared to such approaches, our GAT model:

- Achieves comparable or superior accuracy with fewer parameters.

- Provides more explicit modeling of spatial relationships through graph edges.

- Requires less computational memory (processing 80 nodes vs. full resolution feature maps).

### 8.3.3 Comparison with the Source Paper:

The source segmentation paper achieved Dice scores of 0.91, 0.86, and 0.79 for different tumor regions. While these metrics are not directly comparable to classification accuracy, our 96.43% classification accuracy demonstrates that the graph-based approach successfully transfers to the classification domain, validating the architectural principles.

## 8.4 Limitations and Challenges:

### 8.4.1 Current Limitations:

1. **Computational Overhead of Graph Construction:** The SLIC segmentation and feature extraction process is computationally expensive (2-3 seconds per image), which limits real-time application potential. Graph construction requires more processing time than simple image resizing used in CNN pipelines.

2. **Fixed Graph Structure:** Once constructed, the graph topology remains fixed during training. An ideal system might learn to adaptively refine the graph structure, emphasizing relevant superpixels while coarsening irrelevant regions.

3. **Superpixel Parameter Sensitivity:** The choice of 80 superpixels is somewhat arbitrary. Too few superpixels might lose important details, while too many increase computational cost. The optimal number may vary by tumor type or image characteristics.

4. **Limited Data Augmentation:** The current pipeline lacks robust data augmentation strategies. Traditional image augmentations (rotation, flipping, intensity jittering) could improve generalization, but require careful integration with graph construction to maintain structural consistency.

5. **Class Imbalance Handling:** While our dataset is reasonably balanced, real-world clinical data often exhibits severe class imbalance. The current approach doesn't explicitly address this through techniques like class weighting or focal loss.

6. **Interpretability Gaps:** Although attention weights provide some interpretability, we don't currently visualize them or analyze which superpixels receive high attention. Adding attention visualization could enhance clinical trust and model understanding.

7. **Single-Modality Input:** Clinical MRI protocols typically acquire multiple sequences (T1, T2, FLAIR, T1-contrast). Our current implementation uses only single-channel or RGB images, not leveraging multi-modal information that could improve accuracy.

### 8.4.2 Failure Cases and Challenges:

1. **Glioma Recall:** The model's 88% recall for glioma indicates 12% of glioma cases are misclassified. These false negatives likely correspond to:

   - Low-grade gliomas with subtle appearance.
   - Gliomas with atypical presentations mimicking other tumor types.
   - Boundary cases where tumor characteristics overlap with meningioma.

2. **Small Tumor Detection:** The superpixel approach may struggle with very small tumors that occupy only a few superpixels, as limited spatial context could hinder classification.

3. **Artifact Sensitivity:** MRI artifacts (motion artifacts, intensity inhomogeneity) could affect superpixel segmentation quality, propagating errors into the graph representation.

## 8.5 Potential Improvements and Future Directions:

### 8.5.1 Short-Term Enhancements:

1. **Data Augmentation:** Implement comprehensive augmentation strategies including:

   - Geometric transformations (rotation, scaling, flipping).

   - Intensity manipulations (brightness, contrast, gamma correction).

   - Elastic deformations to simulate biological variability.

   - MixUp or CutMix for graph data.

2. **Attention Visualization:** Develop visualization tools to display attention weights overlaid on original MRI images, helping clinicians understand model decisions.

3. **Ensemble Methods:** Combine multiple GAT models trained with different random seeds or hyperparameters to improve robustness and accuracy.

4. **Class Balancing:** Implement focal loss or class-weighted loss functions to handle potential class imbalances in deployment scenarios.

5. **Hyperparameter Optimization:** Conduct systematic grid search or Bayesian optimization to find optimal values for:

   - Number of superpixels

   - Number of GAT layers

   - Hidden dimensions

   - Number of attention heads

   - Dropout rates

### 8.5.2 Medium-Term Directions:

1. **Multi-Scale Graphs:** Construct hierarchical graphs with multiple resolution levels, enabling the model to capture both fine-grained local details and coarse global structure.

2. **Learnable Graph Construction:** Replace fixed SLIC-based graphs with learnable graph construction mechanisms that can adapt node placement and edge connections during training.

3. **Multi-Modal Integration:** Extend the approach to handle multiple MRI sequences by:

   - Creating separate graphs per modality

   - Fusing features from different modalities at node level

   - Learning cross-modal attention mechanisms

4. **Edge Feature Learning:** Currently, edges have no features. Augmenting edges with learned or handcrafted features (e.g., boundary strength, intensity gradients) could improve performance.

5. **Graph Pooling Strategies:** Explore learnable pooling methods (e.g., DiffPool, TopKPooling) to hierarchically coarsen graphs while preserving important structures.

### 8.5.3 Long-Term Research Directions:

1. **Weakly Supervised Segmentation:** Extend the classification model to provide rough segmentation masks using Class Activation Mapping (CAM) or attention-based localization, enabling tumor localization without pixel-level annotations.

2. **Multi-Task Learning:** Jointly train the model for classification, segmentation, and tumor grading or survival prediction, leveraging shared representations across tasks.

3. **3D Graph Construction:** Extend to true 3D volumetric analysis by constructing graphs over 3D MRI volumes rather than 2D slices, capturing inter-slice relationships.

4. **Temporal Analysis:** For longitudinal studies, incorporate temporal edges connecting graphs from different time points to track tumor progression or treatment response.

5. **Uncertainty Quantification:** Implement Bayesian GATs or ensemble methods to provide confidence estimates alongside predictions, crucial for clinical decision support.

6. **Clinical Deployment:** Develop a streamlined inference pipeline optimized for speed, potentially using:

   - GPU-accelerated graph construction

   - Model quantization or pruning

   - Cached superpixel computations for repeated scans

7. **Explainable AI Integration:** Combine attention mechanisms with counterfactual explanations or causal analysis to provide clinicians with actionable insights: "The model predicts glioma because regions X and Y exhibit patterns similar to known glioma cases."

8. **Cross-Dataset Generalization:** Evaluate and improve model performance on external datasets (e.g., BRATS, TCGA) to assess generalization beyond the Kaggle dataset.

# 9   Conclusion:

Our Graph Attention Network approach successfully adapted a segmentation-focused architecture to the classification domain, achieving 96.43% accuracy on the brain tumor MRI classification task. The graph-based representation, combined with attention mechanisms, provides an effective and parameter-efficient alternative to traditional CNN architectures. The model demonstrates balanced performance across all four tumor classes, with particularly strong results for healthy brain detection (97% F1-score) and reliable identification of meningioma and pituitary tumors (94-95% F1-scores).

The key innovation lies in transforming MRI images into graph structures where nodes represent coherent anatomical regions and edges encode spatial relationships. This explicit structural modeling, combined with learnable attention weights, enables the model to focus on diagnostically relevant features while suppressing irrelevant background information.

Github repo where you can find the code:
`https://github.com/amaliahm/BrainGAT.git`

# 10    References:

- The original paper by Dhrumil Patel, Dhruv Patel, Rudra Saxena and Thangarajah Akilan, *Multi-class Brain Tumor Segmentation using Graph Attention Network* on `https://arxiv.org/pdf/2302.05598`

- Brain Tumor MRI Dataset, available at: `https://www.kaggle.com/datasets/masoudnickparvar/brain-tumor-mri-dataset`

- Medium articles on `https://medium.com/`