# AUTOMATIC CLASSIFICATION OF CELLS DURING NEURODEVELOPMENT

SUDHA GOPALAKRISHNAN BRAIN CENTRE, IIT MADRAS

A Siddharth Reddy[1]

[1] Intern, SGBC, arvindsiddharthreddy@gmail.com

## Abstract

The research commenced with the implementation of a watershed-based cell detection algorithm designed to identify cells within specific tiles of JP2 images, focusing on the extraction of regions of interest (ROIs) and refining these selections based on a defined threshold area. This foundational step enabled the generation of data for more complex analyses. The feature vectors extracted from these ROIs were subsequently enhanced with additional attributes. To streamline the clustering process, Principal Component Analysis (PCA) was applied, reducing the feature space to two dimensions. Following this, KMeans clustering and Graph Cut Segmentation were employed to automatically cluster the vectors and identify distinct cell classes. Average feature values for each cluster or segment were then calculated and used to generate synthetic cell images.

The scope of the research was broadened by populating a comprehensive training dataset, involving feature extraction from multiple tiles across a significant portion of the fetal brain image. The relationship between the number of segments (cell classes) and changes in the number of training data vectors was systematically analyzed to identify a point of stability in segment count. Once this stability was established, average feature values for each cluster were recalculated, leading to the creation of a new set of synthetic cell images. An algorithm was developed to extract visible features from these synthetic images, forming the basis for training a semi-supervised model. This model is crucial for validating the accuracy of cell classification and ensuring the biological relevance of the identified cell classes. Additionally, visualization techniques were employed to observe changes across clusters, enhancing the understanding of cell differentiation patterns during neurodevelopment.

## 1 Introduction

The primary objective of this research is to develop an automated classifier that accurately classifies cells during neurodevelopment within the fetal brain. This process begins with the detection of cells and the extraction of regions of interest (ROIs) using a watershed algorithm, which lays the groundwork for further analysis. The feature vectors from these ROIs are then enhanced with additional attributes. Principal Component Analysis (PCA) is applied to reduce the feature space, facilitating easier clustering. KMeans clustering and Graph Cut Segmentation are employed to identify distinct cell classes based on these refined feature vectors.

This research also focuses on generating synthetic cell images derived from the average feature values of each cluster. A comprehensive training dataset is developed by extracting features from multiple tiles, enabling a robust analysis of segment count stability as the number of training data vectors changes. Finally, an algorithm is created to extract visible features from these synthetic images, which are then used to train a semi-supervised model. Visualization approaches are also integrated to monitor changes across clusters, providing insights into cell differentiation patterns during neurodevelopment. This model is essential for ensuring the biological relevance and accuracy of cell classification, ultimately contributing to a deeper understanding of neurodevelopment.

## 2 Motivation

The motivation behind this research is to enhance the classification of cells during fetal brain development, a crucial part of neurodevelopmental studies. Accurate classification can reveal important details about cell differentiation and growth patterns, aiding in the understanding of brain disorders and brain abnormalities. Traditional methods

are time-consuming and error-prone, highlighting the need for automated techniques. While neuron classification has been attempted [1] in adult brains, this is, to our knowledge, the first time such cell classification has been attempted on a developing human brain.

# 3 Cell Detection

## 3.1 Tile Extraction and Preprocessing

### 3.1.1 Tools Used:

- Python Libraries: **rasterio, cv2 (OpenCV), numpy, matplotlib.**

### 3.1.2 Steps:

- **Tile Extraction** - Extracted 512x512 pixel tile (Fig. 1) from the original JP2 image using rasterio. Then, converted extracted tile to a format suitable for OpenCV (BGR format).
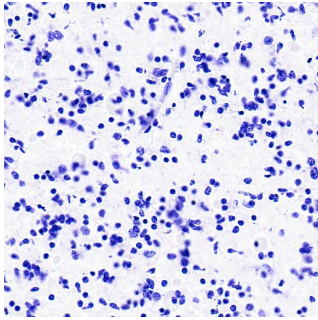


Fig. 1: Extracted 512x512 tile

- **Blue Channel Extraction** - Extracted blue channel from the image tile, as extracting blue channel enhances contrast and improves accuracy of edge detection (Fig. 2).
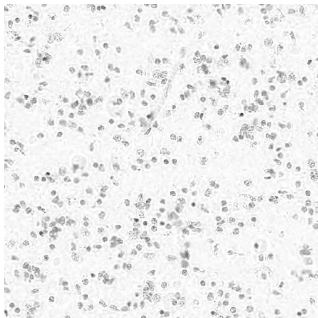


Fig. 2: Blue channel

- **Gaussian Blur** - Applied Gaussian blur to blue channel to reduce noise and smoothen the image. Experimentation revealed that a 9x9 mask was optimal for detecting the maximum number of cells (Fig. 3).
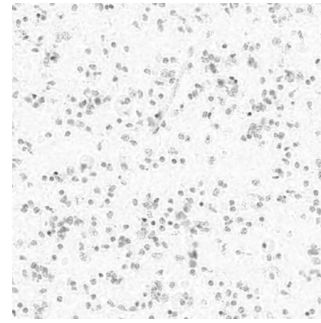


Fig. 3: With 9x9 gaussian blur

- **Edge Detection** - Applied Canny Edge Detection on blurred blue channel to detect edges more accurately (Fig. 4).
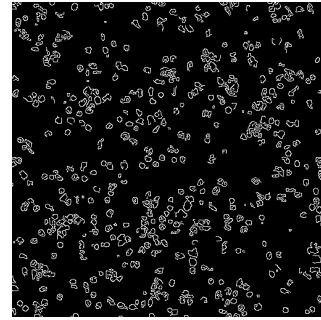


Fig. 4: Canny edge detection

## 3.2 Image Segmentation

### 3.2.1 Steps:

- Identified sure background and foreground regions for the Watershed algorithm (Fig. 5).

  1. **Sure Background**: Regions that are definitely back ground.
  2. **Distance Transform**: Calculated distances to the nearest zero pixel.
  3. **Sure Foreground**: Thresholded the distance tansform.
  4. **Unknown Area**: Subtracted sure foreground from sure background.

## 3.3 Watershed Algorithm

### 3.3.1 Steps:

- **Labeling** - Labeled the connected components in the sure foreground. Adjusted the label so that the background is not zero. Then, marked the unknown regions with zero.

- **Watershed Application** - Applied the Watershed Algorithm to segment the image into different regions (Fig. 6).
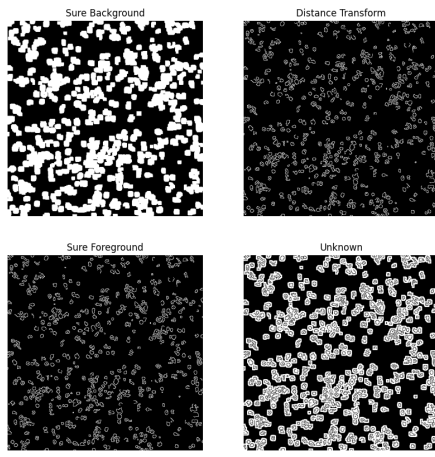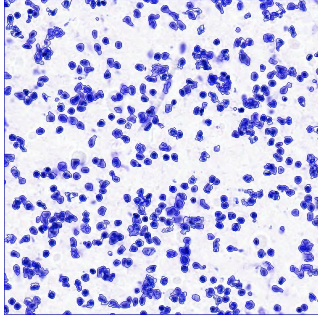
Fig. 5



Fig. 6: Marker labeling

# 4 Results

## 4.1 ROI Extraction and Display

- Extracted and saved all ROI's.

- Displayed all ROIs in a subplot(Fig. 7).

- Centroids are calculated and displayed for each cell detected in the tile(Fig. 8).
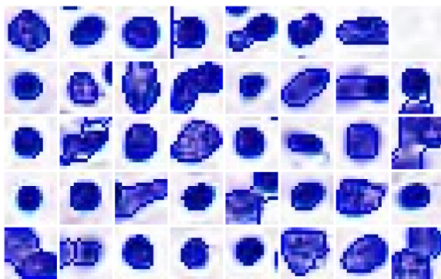
- Total 185 ROIs found (For 512x512 tile).



Fig. 7: Detected ROIs (A few...)

## 4.2 Choosing Threshold

A threshold area of 150 sq units was fixed, upon observing all the ROIs that were extracted so far. Therefore, all those ROIs with area less than 150
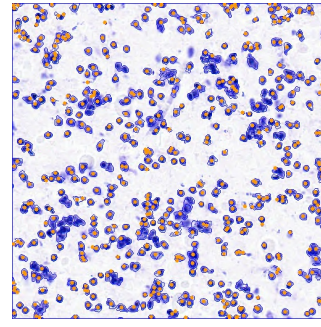


Fig. 8: Tile with centroid markings

units will be extracted. These ROIs will then be used for feature extraction which will also be used as training data for the cell auto-classifier. 180 ROIs below threshold (For the 512x512 tile used here).Example of an extracted ROI is shown in Fig. 9.



Fig. 9: Extracted ROI

# 5 Summary - Cell Detection

The cell detection algorithm performed so far performs the following tasks:

- Reads and preprocesses a JP2 image.

- Extracts the blue channel and applies Gaussian blur.

- Performs Canny Edge Detection.

- Identifies sure background and foreground areas.

- Applies the watershed algorithm to segment the image.

- Extracts and saves regions of interest (ROIs). Also displays them in subplots.

- Draws contours (for detected objects) and centroids on the original image.

- Extracts specific ROIs based on a fixed threshold.

# 6 Cell Feature Extraction

Now, for each of the extracted ROIs (the ones below the threshold area), a feature vector is prepared. This vector contains information about various features of a cell such as:

- Aspect Ratio
- Extent
- Solidity
- Contrast
- Local Homogeneity
- Global Homogeneity
- Perimeter
- Minor Axis
- Major Axis
- Equivalent Diameter
- Circularity
- Convexity

## 6.1 Definitions

### 6.1.1 Aspect Ratio

- Ratio of the width to the height of the bounding rectangle of the contour.

### 6.1.2 Extent

- Ratio of the contour area to the area of its bounding rectangle. It measures how much of the bounding rectangle is filled by the contour.

### 6.1.3 Solidity

- Ratio of the contour area to the area of it's convex hull. The convex hull is the smallest convex shape that can enclose the contour.

### 6.1.4 Contrast

- Calculated as the square of the variance of the greyscale pixel values within the contour. Indicates the level of detail or texture.

### 6.1.5 Local Homogeneity

- Measures the similarity of pixel values within a region (Inside the cell). Calculated using the gray-level co-occurrence matrix (GLCM) and derived from this matrix using the **greycoprops(glcm, 'homogeneity')** function.

### 6.1.6 Global Homogeneity

- Represents the overall uniformity of pixel values in an image. Lower entropy values indicate higher homogeneity.

### 6.1.7 Perimeter

- Total length of the contour of the cell. It is calculated using the **cv2.arcLength** function, which takes the contour and a flag indicating whether the contour is closed (True) or open (False).

### 6.1.8 Minor Axis

- Shorter side of the bounding rectangle around the contour. It is derived from the dimensions of the bounding rectangle (width w and height h).

### 6.1.9 Major Axis

- Longer side of the bounding rectangle around the contour. It is derived from the dimensions of the bounding rectangle (width w and height h).

### 6.1.10 Equivalent Diameter

- Diameter of a circle that has the same area as the contour.

### 6.1.11 Circularity

- Measure of how close the shape of the contour is to a perfect circle. A perfect circle has a circularity value of 1, while other shapes have values less than 1.

### 6.1.12 Convexity

- Binary measure that indicates whether the contour is convex or not. Determined using the **cv2.isContourConvex** function, which returns True if the contour is convex and False otherwise.
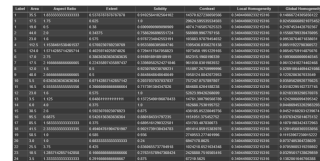
## 6.2 Results



Fig. 10: Extracted features

Fig. 11: Extracted features (Contd.)

# 7 KMeans Algorithm

The KMeans [2] algorithm is a clustering technique used to partition a dataset into K distinct, non-overlapping subsets (clusters). The algorithm iterates to find cluster centroids that minimize the within-cluster variance

## 7.1 Algorithm

### 7.1.1 Initialization

- Select K initial centroids randomly from the dataset.

### 7.1.2 Assignment Step

- Assign each data point to the nearest centroid, creating K clusters.

### 7.1.3 Update Step

- Calculate new centroids as the mean of all data points assigned to each cluster.

### 7.1.4 Repeat

- Repeat the assignment and update steps until the centroids no longer change or a maximum number of iterations is reached.

## 7.2 Implementation

### 7.2.1 PCA for Dimentionality Reduction

- Before applying KMeans, PCA was used to reduce the feature space to 2 dimensions for easier clustering.

### 7.2.2 KMeans Clustering

- We initialized the KMeans algorithm with 6 clusters.

- The the KMeans algorithm was applied to the PCA - reduced data and each data point was assigned to one of the 6 clusters. The cluster assignment was added as a new column 'Cluster' to the DataFrame.

### 7.2.3 Calculate Average Features

- The features used for clustering are then averaged within each cluster to obtain a representative feature vector for each cluster.

## 7.3 Results



Fig. 12: KMeans clustering on PCA components



Fig. 13: Average of each feature in each cluster

# 8 Graph Cut Segmentation

The graph cut segmentation [3] algorithm first constructs a graph where each pixel (or data point) is a node, and edges represent the similarity between neighboring pixels. It then applies label propagation to detect communities within this graph, effectively grouping similar pixels into segments. Finally, it assigns a unique label to each segment, ensuring the total number of segments does not exceed a specified limit. Overall,this function leverages graph-based methods and community detection to segment data effectively. It provides a flexible and efficient approach to identifying distinct regions or clusters in a dataset.

## 8.1 Algorithm

### 8.1.1 Graph Construction

- Construct a k-nearest neighbors graph from the input data, where nodes represent pixels

and edges represent similarity between neighbors.

### 8.1.2 Community Detection

- Apply the label propagation algorithm to detect communities (segments) within the graph.

### 8.1.3 Segment Limitation

- Ensure the number of detected segments does not exceed the specified maximum number.

### 8.1.4 Label Assignment

- Assign a unique label to each segment based on the detected communities.

### 8.1.5 Output Labels

- Return the array of segment labels for each pixel in the input data.

## 8.2 Results



Fig. 14: Graph Cut Segmentation on PCA components
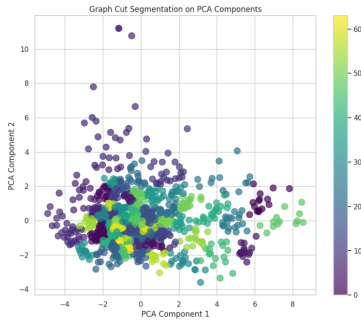


```
Average of each feature for each segment:
            Area  Aspect Ratio    Extent  Solidity      Contrast  \
Segment
0       104.400000      1.042442  0.383971  0.663806  1.151812e+06
1        31.787500      0.991716  0.628895  0.971243  4.391821e+05
2        19.947368      1.075313  0.490632  0.894187  7.295557e+05
3        67.875000      0.963991  0.616040  0.928288  1.050579e+06
4       135.821429      0.966663  0.536368  0.830429  9.195017e+05
...            ...           ...       ...       ...           ...
59       40.950000      1.231905  0.587467  0.952009  5.372128e+05
60       54.625000      0.927409  0.570585  0.910419  7.550679e+05
61       37.653846      1.013278  0.583599  0.944028  4.321689e+05
62       53.076923      0.991994  0.631333  0.962918  6.514644e+05
63       44.750000      0.946759  0.634404  0.964009  6.053291e+05

         Local Homogeneity  Global Homogeneity  Perimeter  Minor Axis  \
Segment
0                 0.343429            0.114030  61.126098   15.000000
1                 0.343429            0.077457  21.882159    6.775000
2                 0.343429            0.103922  19.838741    5.210526
3                 0.343429            0.067735  33.904699    9.187500
4                 0.343429            0.084620  55.020093   14.285714
...                    ...                 ...        ...         ...
59                0.343429            0.086578  25.606602    7.400000
60                0.343429            0.069558  31.038582    8.875000
61                0.343429            0.098369  24.180791    7.615385
62                0.343429            0.049918  28.236785    8.304615
63                0.343429            0.081071  25.859055    8.000000

         Major Axis  Equivalent Diameter  Circularity  Convexity
Segment
0         18.333333            11.523166     0.354122       0.00
1          7.500000             6.356053     0.834278       0.05
2          7.894737             5.026488     0.637784       0.00
3         12.125000             9.282264     0.742979       0.00
4         17.857143            13.146217     0.566626       0.00
...             ...                  ...          ...        ...
59         9.500000             7.218054     0.786683       0.00
60        10.875000             8.338684     0.713712       0.00
61         8.538462             6.922618     0.809234       0.00
62        10.076923             8.218552     0.836444       0.00
63         8.833333             7.548002     0.841448       0.00

[64 rows x 13 columns]
```

Fig. 15: Average of each feature in each segment

# 9 Synthetic Generation of Cell Images

## 9.1 Algorithm

The algorithm for generating synthetic cell images based on average features consists of several key steps. These steps are designed to create representative images that encapsulate the characteristic properties of clusters identified in the PCA - transformed feature space. Here's a detailed explanation:

### 9.1.1 Image Initialization

- The algorithm starts by creating a blank image (img) and a mask (mask) of a specified size (e.g., 20x20 pixels). These arrays are initialized to zeros, which represents the background.

### 9.1.2 Drawing the Cell Shape

- The center of the image is calculated to position the synthetic cell in the middle of the image.

- The major and minor axis lengths of the cell are derived from the average feature values of the cluster. These axes determine the shape and size of the ellipse that represents the cell.

- The **cv2.ellipse** function is used to draw the ellipse on the mask, which fills the cell shape with a white color (255). This mask is used to delineate the cell region from the background.

### 9.1.3 Assigning Gray Values

- The pixel values within the cell area (as defined by the mask) are set to a mid-gray value (128). This uniform color represents the baseline intensity of the cell.

### 9.1.4 Adding Texture (Homogeneity)

- To simulate local homogeneity within the cell, Gaussian noise is added to the cell area. The standard deviation of this noise is proportional to the Local Homogeneity feature, allowing the texture to reflect the variability observed in real cells.

- The **cv2.addWeighted** function blends the noise with the base cell intensity to create a realistic texture while ensuring that the background remains black.

### 9.1.5  Generating and Displaying Images

- For each cluster, a synthetic image is generated using the average features.

- The synthetic image is saved to disk and displayed in a grid layout using Matplotlib for visual inspection.
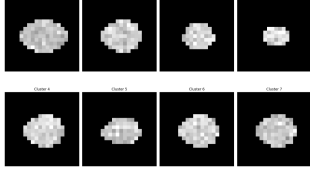
## 9.2  Results



Fig. 16: Based on KMeans cluster data



Fig. 17: Based on Graph Cut segment data



Fig. 18: Based on Graph Cut segment data
(Contd.)

## 10  Number of Cell Classes

The objective here was to identify the number of cell classes present in the fetal brain image, which depends on the number of segments detected. However, due to variations in the number of segments with changes in the training data size, accurately estimating the number of cell classes is challenging. By plotting the changes in training data size against the detected number of segments, a discernible trend emerged, allowing for an approximate determination of the number of cell classes.

| No. of Vectors | Segments |
|---|---|
| 1093 | 71 |
| 1289 | 77 |
| 2798 | 116 |
| 5166 | 108 |
| 7994 | 100 |
| 13106 | 105 |
| 23536 | 105 |
| 54297 | 108 |
| 83332 | 105 |
| 100451 | 105 |

Table 1: Number of Vectors vs. Segments



Fig. 19: Number of Vectors vs. Segments

Based on the graphical analysis and tabular representation, it can be concluded that the number of cell classes in the fetal brain image falls within the range of 100 to 108, with 105 being the most prominent estimate.

## 11  Distinctive Features of Cell Classes

We aim to explore the distinguishing characteristics of segments or cell classes by extracting specific visible features. The features chosen are as follows:

- **Area**: The number of pixels in the region

- **Perimeter**: The length of the boundary of the region.

- **Aspect Ratio**: The ratio of the major axis length to the minor axis length of the region. It describes the elongation of the region.

- **Extent**: The ratio of the area of the region to the area of the bounding box. It describes how much of the bounding box is filled by the region.

- **Solidity**: The ratio of the area of the region to the area of the convex hull. It describes how convex the region is.

- **Mean Intensity**: The average intensity of the pixels in the region.

- **Std Intensity**: The standard deviation of the intensity values of the pixels in the region.

- **Contrast**: A measure of the local variations in the gray-level co-occurrence matrix (GLCM).

- **Dissimilarity**: A measure of how different the elements of the GLCM are from each other

- **Homogeneity**: A measure of how similar the elements of the GLCM are to each other.

- **ASM (Angular Second Moment)**: Also known as energy, it measures the uniformity of the GLCM.

- **Energy**: The sum of squared elements in the GLCM, which is another measure of texture uniformity.



Fig. 20: Synthetically generated cell

Such a feature extraction was performed on a single synthetically generated image (Fig. 20) and the results were as follows:

- Area : 165

- Perimeter : 46.62741699796952

- Aspect Ratio : 1.3137962574389555

- Extent : 0.746606334841629

- Solidity : 0.9322033898305084

- Mean Intensity : 127.56363636363636

- Std Intensity : 11.018601852179772

- Contrast : 1710.3509615384617

- Dissimilarity : 20.418269230769234

- Homogeneity : 0.24878671916240513

- ASM : 0.03290264423076923

- Energy : 0.1813908603837835

We then applied this feature extraction algorithm to all the 105 synthetically generated images. Upon completion, a CSV file, as shown below, was compiled (Fig. 21).



Fig. 21: Extracted features from synthetic cells

Our next step involves refining these parameters to calculate the significant features from the un-biased clusters, followed by re-clustering using those significant parameters. This is a semi-supervised approach aimed at ensuring that the clusters obtained are valid and that the model is a real world applicable model.

## 12 What makes a cluster unique?

### 12.1 Visualization of Feature Variations Across Clusters

To analyze the variations in feature values across different clusters, we employed a bar plot visualization technique. This approach allowed us to compare how individual features change across various clusters, providing insights into the underlying patterns within the dataset.

- **Data Preparation**: The feature data was first reshaped from a wide to a long format, which made it easier to plot. Each feature was associated with its corresponding cluster, enabling a clear comparison across clusters.
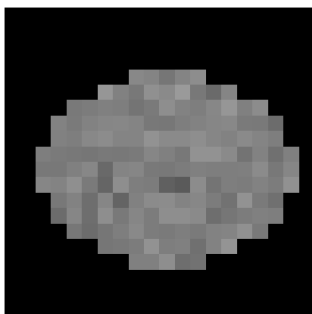
- **Normalization**: To ensure that all feature values were represented on a positive axis, we normalized the data by shifting the entire dataset upwards. Specifically, we added a constant value to each feature, which was calculated as the absolute value of the minimum feature value in the dataset, plus a small epsilon to avoid numerical issues. This normalization step ensured that all plotted values were positive, providing a clearer and more interpretable visualization.

- **Visualization**: We generated a bar plot with features on the x-axis and their corresponding normalized values on the y-axis. Different clusters were represented by distinct colors (hues), enabling a visual comparison of how each feature varied across clusters. The plot was enhanced by rotating the x-axis labels for better readability and setting an appropriate figure size for clarity (Fig. 22).
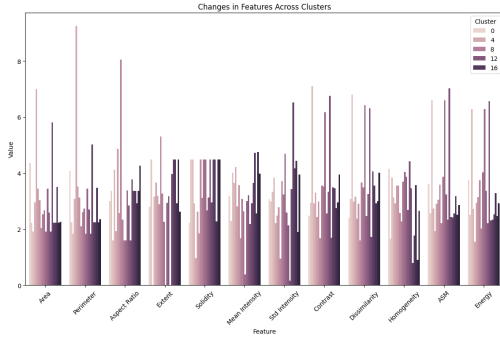


Fig. 22: Changes in features across clusters

This visualization helped to highlight key differences and similarities in feature values across clusters, serving as a valuable tool in the analysis of the data's structure and the effectiveness of the clustering algorithm.

## 12.2 Generation of Synthetic Image Equations

In our analysis, we sought to quantify the relationship between feature vectors and the corresponding synthetic images. To achieve this, we developed a method to generate mathematical equations that represent synthetic images as a linear combination of their feature vectors. This approach allowed us to express the contribution of each feature to the overall synthetic image in a clear and interpretable form.

- **Equation Generation**: For each cluster center, an equation was generated that defines the synthetic image associated with that cluster. The feature vector of the cluster center was used as the basis for this equation. Specifically, each feature in the vector was multiplied by its corresponding value, and these terms were summed to produce the final equation. The general form of the equation is as follows:

$$\text{Synthetic Image} = \sum_{i=1}^{n} \left( \text{Value}_i \times \text{Feature}_i \right)$$

Where $Value_i$ represents the numerical value of the i-th feature, and $Feature_i$ is the name of the feature. Each equation thus provides a linear combination of features that defines the synthetic image for a specific cluster.

- **Implementation**: The process was automated using a custom Python function, which constructs the equation string by iterating over the feature vector of each cluster center. For each feature, the function formats the corresponding value to two decimal places and appends it to the equation string. The final equation is then stored in a list, with each equation associated with its respective cluster.

## 12.3 Results

The generated equations (Fig. 23) were displayed as shown below, each corresponding to a different cluster. These equations provide a straightforward representation of how the features combine to form the synthetic images, offering insights into the structure and characteristics of each cluster.
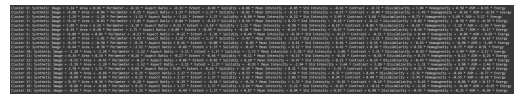


Fig. 23: Generated equations

This approach of representing synthetic images as equations is particularly useful for understanding the contributions of individual features and for comparing the relative importance of features across different clusters.

# 13 Comparison with Baseline

The baseline study [4] on neuron classification primarily focuses on adult human brain tissues, employing a combination of manual and semi-automated techniques [5] to distinguish between

neurons, glia, and endothelial cells. While effective, this approach is limited by its reliance on static histological markers and manual intervention, which can introduce variability and restrict scalability.

In contrast, our work presents a fully automated pipeline for cell classification, specifically designed for application to developing brain tissue—an area where, to our knowledge, no such classification has been attempted before. The dynamic nature of cellular development in the fetal brain presents unique challenges that our model addresses by integrating advanced image processing techniques and machine learning algorithms.

Our model automates the entire classification process, eliminating the need for manual intervention and allowing for the efficient processing of large datasets. This scalability is particularly crucial in the context of fetal brain research, where the cell populations are not only diverse but also rapidly evolving. Furthermore, our approach significantly advances classification accuracy by incorporating features such as Gaussian Blur, watershed algorithm, and PCA-based clustering. These enhancements enable a more nuanced distinction between cell types, even within the complex environment of the developing brain.

This work represents a novel contribution to the field, extending the scope of neuron classification from static adult tissues to the dynamic and formative stages of brain development. By applying our automated model to fetal brain tissue, we provide new insights into neurodevelopment, paving the way for future research into early brain formation and associated developmental disorders.

In summary, our model not only builds upon the foundation laid by previous studies but also expands the frontier of neurodevelopmental research, offering a robust and scalable solution for the automated classification of cells in the developing human brain.

## 14   Future Works

In the next phase of this project, I plan to enhance the robustness of the model by applying additional layers of synthetic image generation and feature extraction. This iterative approach will enable further refinement of the dataset, allowing the model to better capture the complexities of fetal brain cell classification. By generating more diverse synthetic images and extracting increasingly detailed features, the model's ability to generalize across various cell types and developmental stages will be significantly improved.

Furthermore, I aim to explore the incorporation of DALL-E, a cutting-edge generative model, into the workflow. DALL-E's ability to generate high-quality, detailed images from textual descriptions could be leveraged to create synthetic images that represent various cell types and conditions more accurately. This could serve as an additional source of training data, particularly for underrepresented cell types or rare developmental stages. By integrating DALL-E-generated images with the existing dataset, the model could benefit from a more diverse and enriched training set, potentially leading to improved classification accuracy and generalizability.

These enhancements will not only strengthen the current model but also open new avenues for research into the development of more sophisticated tools for neurodevelopmental studies. By continually refining the dataset and leveraging advanced generative models, the project will move closer to achieving highly accurate and reliable cell classification in developing human brains.

## 15   Codes

All the codes can be found in the following links:

- Part-1:Watershed Algorithm, ROI Extraction & Feature Extraction

- Part-2: Feature Extraction from Single Cell

- Part-3: Graph Cut Segmentation, KMeans Clustering & Synthetic Image Generation

- Part-4: Visible Feature Extraction & Cluster Visualization Techniques

## References

[1] Ofek Ophir, Orit Shefi and Ofir Lindenbaum, "Classifying Neuronal Cell Types Based on Shared Electrophysiological Information from Humans and Mice," 2024.

[2] Abiodun M. Ikotun, Absalom E. Ezugwu, Laith Abualigah, Belal Abuhaija, and Jia Heming, " K-means clustering algorithms: A comprehensive review, variants analysis, and advances in the era of big data," 2023.

[3] Chen Yu-ke, Wu Xiao-ming, Cai Ken and Ou Shan-xin, "CT Image Segmentation based on Clustering and Graph-Cuts," 2011.

[4] Miguel Á. García-Cabezas, Yohan J. John, Helen Barbas and Basilis Zikopoulos, "Distinction of Neurons, Glia and Endothelial

Cells in the Cerebral Cortex: An Algorithm Based on Cytological Features," 2016.

[5] Agata KOŁODZIEJCZYK, Magdalena ŁAD-NIAK and Adam PIORKOWSKI, "Constructing Software for analysis of Neuron, Glial and Endotheliaal Cell Numbers and Density in Histological NISSL-Stained Rodent Brain Tissue," 2014.