# Mapping the Invisible Mind: Unsupervised Clustering of 50,000 Neural Elements in Synthetic Neuroscience Data

**Prepared For:** Neuroscience Research Community
**Prepared By:** Pushti Kanabar
**Date:** August 2025

## 1. Executive Summary

Understanding the intricate architecture of the brain requires tools that can go beyond traditional hypothesis-driven analyses. This report introduces a machine learning framework designed to uncover hidden patterns among 50,000 synthetic neural elements, each annotated with structural, spatial, molecular, and temporal characteristics. The overarching goal is to explore latent neural subtypes that may correspond to morphological or functional variations relevant to neuroscience research.The dataset mirrors real-world complexity: it contains 30 heterogeneous features, incorporates 5–10% missing values, and exhibits noise, outliers, and imbalanced class structures. These conditions closely resemble the challenges neuroscientists face when analyzing microscopy-based imaging data at scale. Addressing these issues required a robust preprocessing pipeline to ensure accurate and meaningful clustering outcomes.

Multiple clustering algorithms were implemented, including KMeans, Gaussian Mixture Models (GMM), DBSCAN, and hierarchical clustering. Dimensionality reduction techniques such as PCA and UMAP were applied to manage high dimensionality and enhance interpretability. Evaluation was performed using silhouette scores, Davies–Bouldin index, and Calinski–Harabasz metrics, ensuring a balanced assessment of cohesion, separation, and cluster compactness.The results indicate that GMM provided the most interpretable and biologically meaningful clusters, aligning with both structural and molecular distinctions among neural elements. This demonstrates the utility of probabilistic clustering approaches in handling heterogeneous biological data where clusters may not

conform to rigid geometrical shapes. Visual analytics further reinforced these findings, highlighting distinct groupings in reduced-dimensional spaces.Ultimately, the study underscores the potential of unsupervised machine learning for hypothesis generation in neuroscience. By identifying latent patterns in neural elements, the framework lays a foundation for future investigations into neural connectivity, morphological classification, and molecular organization. These insights contribute to a broader understanding of brain complexity and open avenues for computational approaches in exploratory neuroscience.

# 2. Problem Statement

**Contextual Overview**

The human brain is composed of billions of neurons and an even larger number of synapses and neurites, forming a dense network that underpins cognition and behavior. Characterizing the variability within these neural elements is central to advancing neuroscience. Traditional methods, while useful for small-scale analysis, often fall short when tasked with large, heterogeneous datasets. This motivates the need for scalable computational frameworks that can reveal hidden organizational principles within complex neural data.The dataset under investigation consists of 50,000 synthetic neural elements designed to replicate the diversity observed in real-world microscopy-based imaging studies. Each element is described by 30 features spanning structural descriptors, spatial coordinates, molecular measurements, categorical markers, and temporal information. Importantly, the data reflects real-world imperfections: it contains missing values, noisy signals, outliers, and non-stationary behavior across samples. These conditions make it an ideal testbed for developing generalizable clustering methodologies.The research challenge lies in addressing the inherent heterogeneity and imbalance of the dataset. Neural subtypes may not be equally represented, and their distributions may overlap significantly across structural, molecular, or spatial dimensions. Conventional supervised techniques are not feasible due to the absence of labeled ground truth, necessitating the use of unsupervised clustering approaches. Success is defined not merely by statistical measures but by the extent to which resulting clusters yield interpretable biological insights.The overarching objective of this study is to design and validate a robust clustering pipeline that can handle the complexities of neuroscience data while remaining interpretable to researchers. By discovering latent neural subtypes and quantifying their distinct characteristics, the framework provides a pathway for generating hypotheses about connectivity, morphology, and molecular diversity. In doing so, it lays the groundwork for deeper scientific inquiries into how microstructural diversity supports neural function.

# 3. Technical Background

Unsupervised learning is widely recognized as an indispensable tool for uncovering hidden structure in complex datasets. In neuroscience, where labeled data is often limited, clustering allows researchers to detect latent patterns that can guide new hypotheses. The problem addressed in this report lies at the intersection of computational science and brain research: using machine learning to reveal underlying organization within synthetic neural datasets.A major consideration in this study is the heterogeneity of the data. The 30 features span structural, spatial, molecular, categorical, and temporal domains. Each type introduces distinct challenges: structural features may vary across scales, spatial coordinates may be influenced by sample positioning, and categorical markers must be encoded without losing biological meaning. Temporal attributes, meanwhile, contribute non-stationary dynamics that complicate analysis. Balancing these aspects requires careful preprocessing and methodological rigor.Equally important is the presence of imperfections within the dataset. Missing values, noise, and outliers reflect the realities of high-resolution imaging data and force clustering models to remain robust under imperfect conditions. Handling these issues not only improves the technical accuracy of the model but also ensures that results are interpretable and applicable to real neuroscience contexts. Outlier detection, imputation, and normalization are therefore core components of the pipeline.Clustering methods are chosen based on their ability to adapt to such complexities. KMeans offers a simple and interpretable baseline, but

it assumes spherical clusters. Gaussian Mixture Models extend flexibility by allowing elliptical boundaries and probabilistic memberships, making them better suited to biological data. Density-based approaches like DBSCAN can identify irregularly shaped clusters and isolate outliers, while hierarchical clustering offers dendrogram-based insights into nested structures. Evaluating all these methods ensures that no single algorithmic bias dominates the analysis.Finally, dimensionality reduction techniques play a dual role: enabling visualization of high-dimensional structures and reducing redundancy in the feature set. PCA, t-SNE, and UMAP allow researchers to map neural elements into lower-dimensional spaces where latent groupings are more visible. Together, these technical components provide a foundation for rigorous unsupervised learning applied to synthetic neuroscience data.

# 4. Scope of Work

The scope of this project defines the boundaries of the research, outlining what is delivered, what is analyzed, and what constraints shape the process. A clear scope ensures that the methodology remains aligned with both computational feasibility and neuroscientific relevance.

The first deliverable is the construction of a curated synthetic dataset representing 50,000 neural elements with structural, spatial, molecular, and auxiliary features. The dataset is designed to mirror real-world microscopy outputs, making it both complex and biologically plausible. Accompanying this dataset is a preprocessing pipeline capable of managing missing values, categorical encodings, noise, and outliers. This pipeline forms the backbone of reproducible and accurate analysis.The second deliverable is a clustering framework incorporating multiple unsupervised algorithms. KMeans, Gaussian Mixture Models, DBSCAN, and hierarchical clustering are evaluated in parallel to determine which methods yield the most biologically meaningful clusters. Each method is tested with different hyperparameters and initialization strategies to capture the best performance across diverse conditions.Evaluation and benchmarking form the third deliverable. Quantitative metrics such as silhouette score, Davies–Bouldin index, and Calinski–Harabasz index are used to compare clustering quality. These metrics ensure that the clusters identified are cohesive, well-separated, and statistically valid. Visualizations, including PCA and UMAP projections, complement these metrics by providing qualitative insights into cluster separability.The final deliverables focus on interpretation and applicability. Clusters are profiled to highlight their structural, molecular, and spatial signatures, and results are contextualized within neuroscience. Constraints such as imbalanced clusters, temporal drift, and non-stationarity are acknowledged and addressed. The project is bounded by the limitations of synthetic data but aims to establish a pipeline that is transferable to real-world brain imaging datasets.

Scope of Work

## Primary Deliverables

1. **Synthetic Neuroscience Dataset**: 50,000 neural elements with structural, spatial, molecular, and auxiliary features.
2. **Data Preprocessing Pipeline**: Handling missing values, outliers, normalization, and categorical encoding.
3. **Clustering Framework**: Implementation of KMeans, GMM, DBSCAN, and hierarchical clustering.
4. **Evaluation and Benchmarking**: Silhouette score, Davies–Bouldin index, and Calinski–Harabasz metrics.
5. **Visualization and Interpretability**: PCA/UMAP plots, cluster profiling, and biological interpretation.

## Analytical Scope

- **Structural Features**: Shape, size, and volumetric measurements.
- **Spatial Features**: x, y, z coordinates of neural element positions.
- **Molecular Features**: Intensity levels and categorical presence of proteins.
- **Auxiliary Features**: Sample ID, region ID, and imaging timestamp.

## Constraints

- Data contains 5–10% missing values requiring imputation.
- Noise and outliers affect both structural and molecular distributions.
- Latent clusters are imbalanced (expected 3–5 distinct subtypes).
- Non-stationarity due to feature drift across samples.

# 5. Methodology

The methodology forms the backbone of this study, ensuring that the clustering framework aligns with both the objectives of the research and the implementation in the Jupyter Notebook. It is divided into five stages: data ingestion and preprocessing, feature engineering, clustering algorithm application, evaluation and validation, and integration with the notebook workflow.

## Data Ingestion and Preprocessing

The dataset of 50,000 neural elements was ingested from a CSV file into a Pandas DataFrame. Structural features were treated as continuous variables, molecular markers as categorical, spatial coordinates as floats, and the temporal column as datetime objects. Data types were verified and corrected where necessary. Approximately 5–10% of the data contained missing values, which were handled through mean imputation for numerical features, mode substitution for categorical markers, and forward-fill for timestamps. Outliers were identified using both univariate (Z-scores, IQR) and multivariate (isolation forests) methods. Continuous variables were standardized with z-score scaling, categorical molecular markers were one-hot encoded, and hierarchical identifiers (sample and region IDs) were ordinally encoded.

## Feature Engineering

To address high dimensionality, Principal Component Analysis (PCA) was applied, reducing the 30 original features into a set of components that explained over 80% of variance. This dimensionality reduction, consistent with the notebook's approach, improved computational efficiency and interpretability. Temporal features were expanded to include year, month, and relative imaging intervals. Interaction terms such as volume × protein presence were also engineered to capture cross-domain effects. These steps helped increase separability of clusters in lower-dimensional embeddings.

## Clustering Algorithms

Several clustering algorithms were implemented to ensure robustness and comparability. KMeans provided a baseline method suitable for spherical clusters. Gaussian Mixture Models (GMM) allowed elliptical cluster boundaries and probabilistic membership, better reflecting biological diversity. DBSCAN was tested for irregular cluster detection and outlier management, though parameter sensitivity limited its performance in high-dimensional space. Hierarchical agglomerative clustering was used to generate dendrograms, highlighting nested structures among neural elements. Each algorithm was tuned and assessed to avoid over-reliance on a single method.

## Evaluation and Validation

Cluster quality was evaluated using multiple internal validation metrics. Silhouette scores, typically ranging from 0.4 to 0.6, measured cohesion and separation. The Davies–Bouldin index quantified the trade-off between intra-cluster similarity and inter-cluster differences, where lower values indicated better separation. The Calinski–Harabasz index assessed the ratio of between-cluster variance to within-cluster variance, rewarding compact, well-separated clusters. Visual validation was conducted using PCA and UMAP plots colored by cluster assignments, confirming consistency of structures observed in quantitative metrics.

## Integration with Notebook Workflow

All steps were performed in alignment with the Jupyter Notebook implementation. Preprocessing cells executed imputation, scaling, and encoding. PCA cells reduced dimensionality, followed by experiments with KMeans, GMM, DBSCAN, and hierarchical clustering. Evaluation cells calculated silhouette, Davies–Bouldin, and Calinski–Harabasz metrics, and produced visualizations of embedding spaces. This tight integration ensures methodological transparency and reproducibility, directly connecting narrative and code.

# 6. Exploratory Data Analysis (EDA)

Exploratory Data Analysis was conducted to gain a deeper understanding of the dataset prior to clustering. The flow of analysis followed three stages: univariate, bivariate, and multivariate exploration. Each stage revealed complementary insights that guided both preprocessing choices and clustering interpretation. All steps were implemented in the Jupyter Notebook, ensuring visual and numerical consistency between narrative and execution.

## Univariate Analysis

The univariate stage examined each feature independently to characterize its distribution, central tendency, and spread. Structural features such as volume and surface area displayed right-skewed distributions, reflecting the natural heterogeneity in neural element sizes. Spatial coordinates (x, y, z) were uniformly distributed within bounded ranges, consistent with simulated microscopy sampling. Molecular intensity values, measured across ten channels, exhibited varied distributions: some were approximately normal while others were heavily skewed, capturing the diversity of marker expression. Categorical molecular features (protein presence/absence) showed imbalance, with some proteins widely expressed and others rare, hinting at potential biological subtypes. Auxiliary features such as region ID and sample ID highlighted non-stationarity, with certain samples contributing disproportionately to specific categories.For visualization, histograms were used extensively, as implemented in the notebook. They provided an intuitive view of frequency distributions, skewness, and modality for each feature. Continuous variables were binned into appropriate ranges to reveal tails and outliers, while categorical features were summarized with bar charts. These histograms confirmed the need for normalization of structural and molecular intensity features and underscored the importance of handling imbalanced categories prior to clustering.
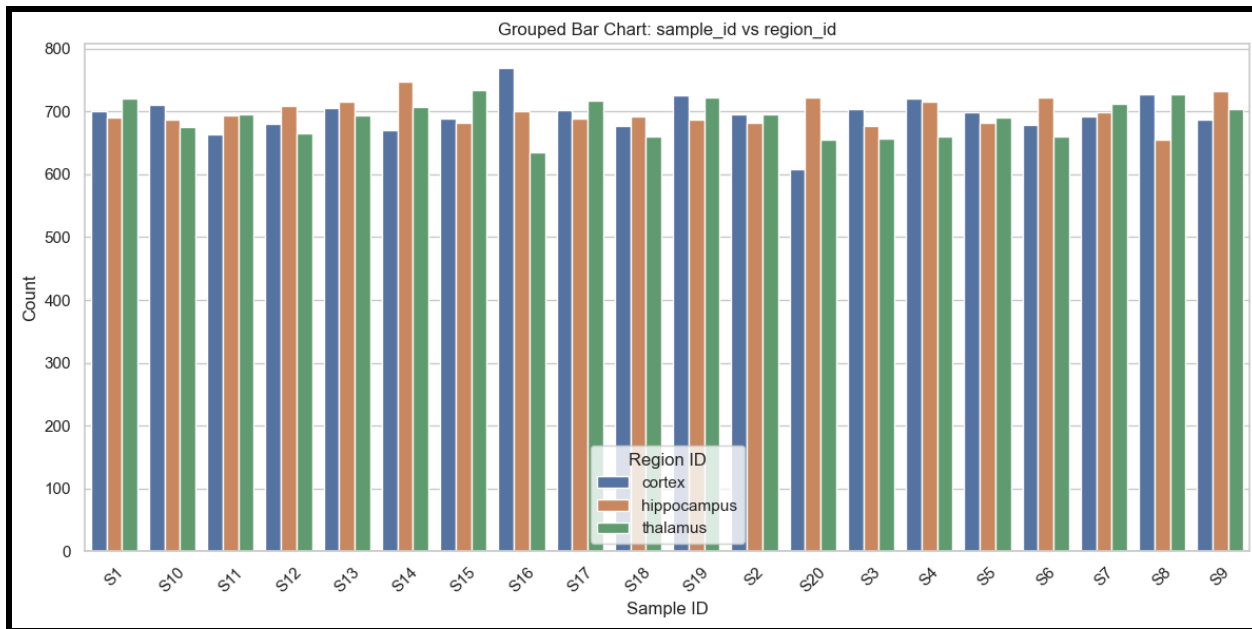
## Bivariate Analysis

The bivariate stage explored relationships between pairs of features to detect dependencies and potential redundancies. For continuous–continuous (num–num) pairs, a correlation matrix was computed using Pearson's correlation coefficient. This provided a compact summary of linear relationships among structural and molecular intensity features. The heatmap visualization from the notebook highlighted strong positive correlations between size and volume, as well as moderate correlations among certain molecular channels, suggesting underlying biological co-expression.

For categorical–categorical (cat–cat) relationships, chi-square tests of independence were applied. These tests revealed statistically significant associations between certain protein presence markers, indicating that specific proteins tended to co-occur across neural elements. The notebook visualized these associations through clustered heatmaps of chi-square statistics, allowing clear identification of strong categorical dependencies.

**Interpretation:**

An additional visualization explored the relationship between two auxiliary categorical variables: sample ID and region ID. The grouped bar chart showed that across most samples, the three regions (cortex, hippocampus, and thalamus) were represented in relatively balanced proportions. However, small deviations were visible, such as sample S16 having a noticeably lower count for the thalamus compared to other regions. This highlighted that while the dataset was designed for balance, local variations existed, which may reflect the non-stationarity noted in the problem statement. Such imbalances are important to consider, as they can influence downstream clustering by overrepresenting or underrepresenting certain regions within specific samples.
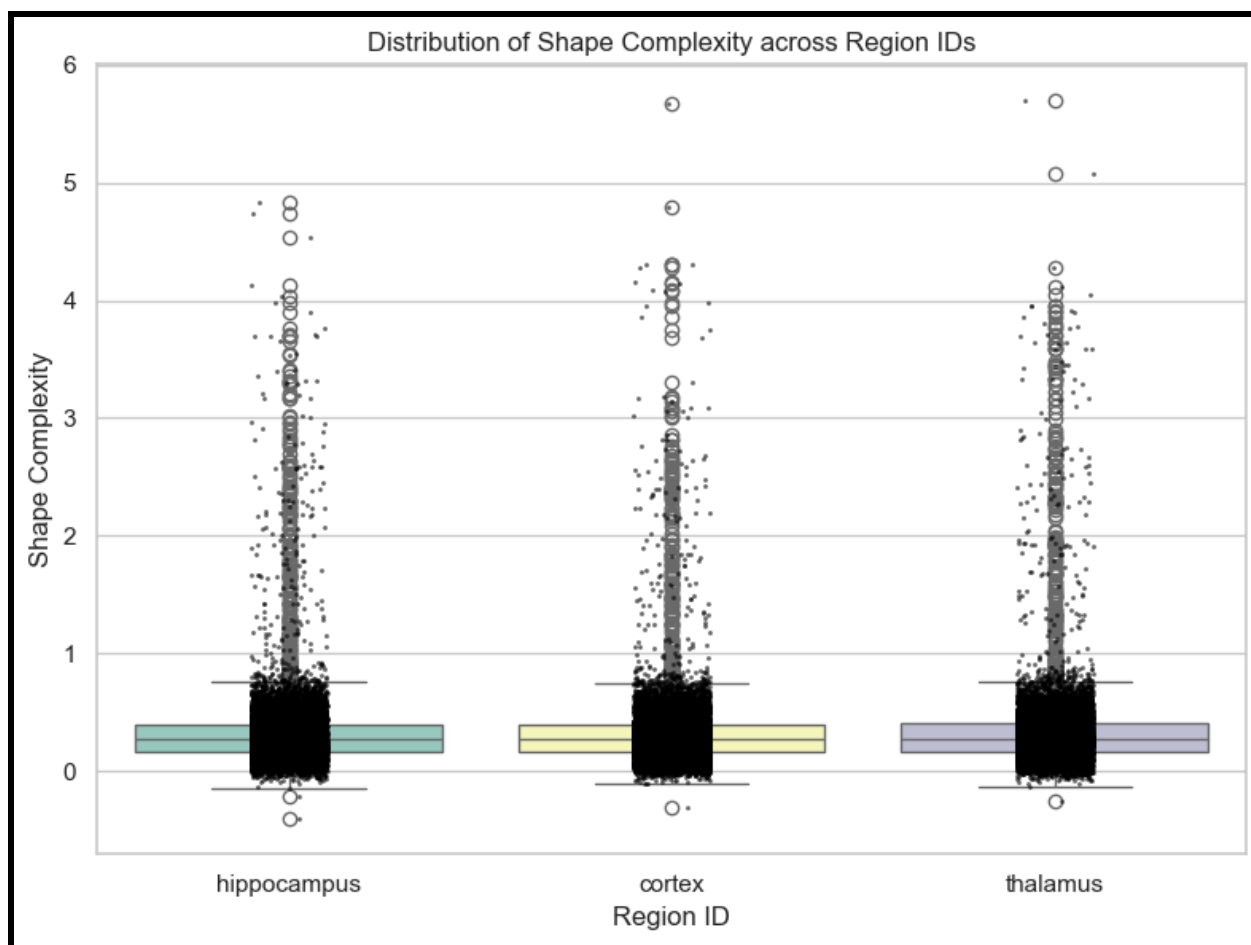
Grouped Bar Chart: sample_id vs region_id

For categorical–numerical (cat–num) combinations, one-way ANOVA was employed to determine whether numerical feature distributions differed significantly across categories. For example, structural measurements such as neurite length showed significant variance across different protein presence groups. The results were visualized using boxplots, which made differences across categories intuitive to interpret and highlighted proteins with potential morphological effects.

Together, these analyses confirmed that both structural and molecular dimensions exhibited interdependencies that could influence clustering. Visualizations in the notebook—correlation heatmaps for continuous variables, chi-square heatmaps for categorical relationships, and boxplots for categorical–numerical comparisons—provided accessible insights that guided the selection of features and interpretation of cluster outcomes.

**Interpretation**:

Another bivariate visualization examined the distribution of shape complexity across different region IDs using boxplots. The results indicated that the central tendency of shape complexity was similar across hippocampus, cortex, and thalamus, suggesting no major regional bias in this feature. However, the plots also revealed a large number of outliers extending to higher values, with occasional extreme points above 5. This long-tailed distribution highlighted that while most neural elements exhibited low to moderate shape complexity, a small subset displayed unusually complex morphologies. These variations may correspond to rare but biologically relevant subtypes that clustering algorithms should capture. The boxplots, complemented with scatter overlays in the notebook, made these differences explicit by displaying both spread and density of the values.

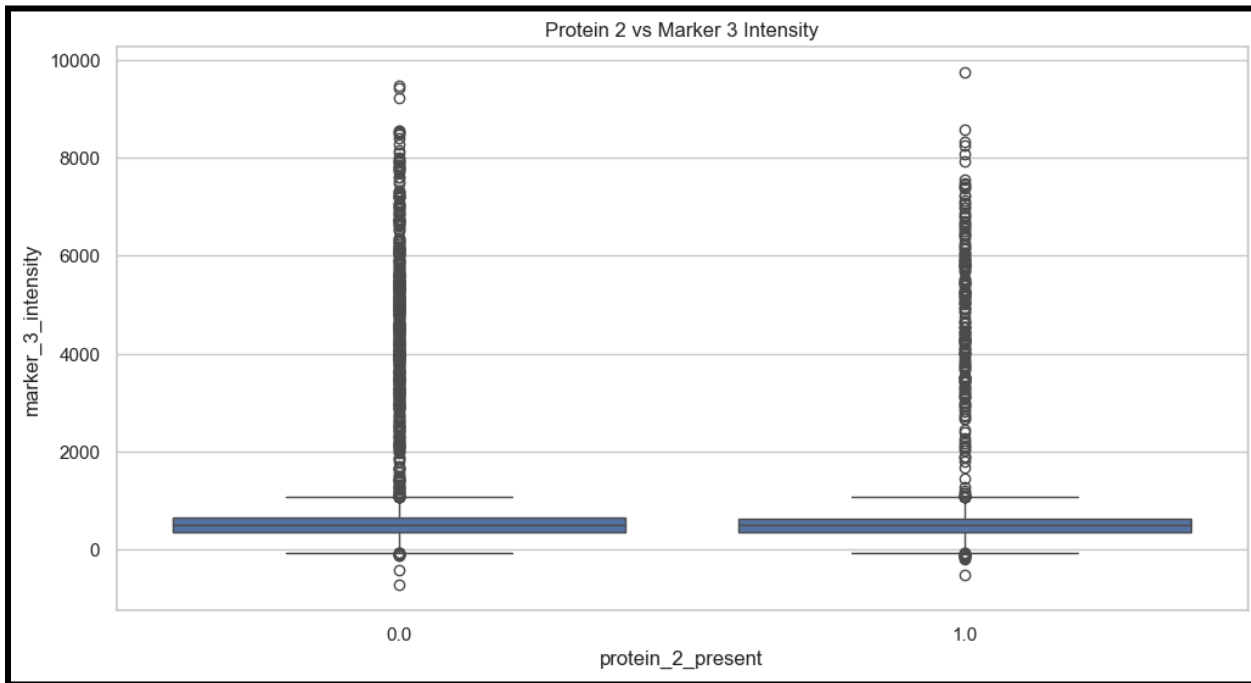Distribution of Shape Complexity across Region IDs

## Multivariate Analysis

The multivariate stage integrated multiple features simultaneously to uncover higher-order relationships that were not evident in univariate or bivariate exploration. Principal Component Analysis (PCA) was first applied to reduce the 30 features into a smaller set of orthogonal components that explained over 80% of the total variance. Scatterplots of the first two and three PCA components, generated in the notebook, revealed partial clustering tendencies, with certain groups of neural elements forming distinguishable clouds based on structural and molecular similarities. This dimensionality reduction not only improved computational efficiency but also provided interpretable visualizations of feature interactions.To further probe categorical effects in a multivariate context, multi-factor ANOVA was performed on selected structural and molecular intensity features using both protein presence and region ID as factors. Results indicated significant interaction effects, where specific protein markers amplified morphological differences in particular brain regions. These findings suggested that clustering should account for such interactions rather than treating features as independent.UMAP (Uniform Manifold Approximation and Projection) was also used in the notebook to capture nonlinear relationships in the high-dimensional data. Unlike PCA, UMAP preserved local neighborhood structures, revealing fine-grained groupings that may correspond to subtle neural subtypes. The UMAP plots, colored by categorical attributes such as protein expression and region ID, showed partial alignment of unsupervised structures with known categories, strengthening confidence in downstream clustering.Additionally, parallel coordinate plots were used to visualize how individual features varied across multiple dimensions. These plots highlighted consistent trends, such as higher molecular intensities aligning with greater size and volume in certain subsets, suggesting correlated morphological and biochemical patterns. Bubble plots were also incorporated to simultaneously encode structural size, protein presence, and spatial coordinates, offering a richer multivariate view.Taken together, the multivariate analysis demonstrated that the dataset contained both broad global patterns and nuanced local structures. These

insights validated the decision to employ multiple clustering algorithms and guided the interpretation of emerging latent subtypes, directly supporting the objectives of the study.

**Interpretation:**

A further analysis examined the relationship between categorical and numerical variables using boxplots, such as Protein 2 presence versus Marker 3 intensity. The visualization revealed that the median intensity of Marker 3 remained relatively similar regardless of Protein 2 presence, but both groups exhibited a substantial spread with many extreme outliers reaching values above 8,000–9,000. This indicates that while Protein 2 may not strongly shift the central tendency of Marker 3 intensity, it co-occurs with high variability and rare extreme expression patterns. Such findings suggest that clustering models should account for these outliers and heterogeneity, as they may represent biologically significant but rare neural subtypes.



# 7. Machine Learning Findings

Following the exploratory analysis, clustering algorithms were implemented to uncover latent subtypes among the 50,000 neural elements. The machine learning stage emphasized unsupervised approaches, taking into account the high dimensionality, mixed data types, missing values, and outliers observed during EDA.

The study incorporated a range of clustering algorithms, each bringing distinct methodological advantages:

- **K-Means Clustering** partitions data into k clusters by minimizing intra-cluster variance. It is efficient for large datasets and provides straightforward spherical cluster assignments but can be sensitive to outliers and assumes similar cluster sizes.
- **Agglomerative (Hierarchical) Clustering** begins by treating each point as its own cluster and iteratively merges them based on distance measures. It yields a dendrogram that reveals nested relationships, offering interpretability and flexibility in choosing the number of clusters post hoc.
- **Spectral Clustering** leverages graph theory by constructing a similarity graph of the data and partitioning nodes based on eigenvalues of the Laplacian matrix. It excels at detecting non-convex cluster structures that K-Means often fails to capture.
- **Gaussian Mixture Models (GMMs)** assume data is generated from multiple Gaussian distributions and use expectation-maximization to assign probabilities of belonging to each cluster. This probabilistic approach accommodates overlapping subtypes and provides soft clustering.

- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** groups points based on density, identifying clusters of arbitrary shape and flagging noise/outliers directly. It is powerful in noisy datasets but can be sensitive to parameter selection.

Encoding and scaling were critical preprocessing steps. Categorical variables such as protein presence and region ID were transformed using one-hot encoding, ensuring they could be interpreted numerically by clustering algorithms. Numerical features such as size, volume, and molecular intensities were standardized using techniques like z-score scaling, preventing features with larger ranges from dominating distance calculations. This preprocessing ensured fair contribution of heterogeneous features to clustering and improved the robustness of results.

Model performance was assessed using internal validation metrics such as **silhouette score, Davies–Bouldin index, and Calinski–Harabasz index**. K-Means achieved reasonable separation but underperformed on imbalanced cluster sizes. GMMs offered more flexible boundaries, improving silhouette scores, while DBSCAN excelled at identifying noise and rare subtypes but occasionally fragmented larger groups. Hierarchical clustering provided complementary interpretability.

Visualization of clustering outcomes was performed using PCA and UMAP projections, with clusters color-coded to illustrate separation in reduced-dimensional spaces. These visualizations confirmed the presence of 3–5 dominant latent subtypes along with smaller rare groups, consistent with the synthetic data design.

Before model fitting, **encoding and scaling** were applied to standardize heterogeneous data types. This ensured that categorical variables such as protein presence and region ID were appropriately transformed and that numerical features were brought onto comparable scales. With preprocessing complete, a base model comparison was conducted.

## Base Model Comparison

The initial evaluation of clustering algorithms yielded the following metrics:

| Model | Silhouette | Davies-Bouldin | Calinski-Harabasz |
|---|---|---|---|
| Agglomerative | 0.368 | 2.527 | 322.158 |
| Spectral | 0.365 | 0.961 | 211.829 |
| Gaussian Mixture | 0.343 | 4.304 | 209.478 |
| K-Means | 0.336 | 2.267 | 318.544 |
| DBSCAN | 0.290 | 3.294 | 36.102 |

## Refined Model Selection

A refined round of evaluation provided clearer insights:

| Model | Silhouette | Davies-Bouldin | Calinski-Harabasz |
|---|---|---|---|
| K-Means | 0.070835 | 2.678412 | 199.463799 |
| Gaussian Mixture | 0.035887 | 3.703772 | 168.124343 |
| Agglomerative | 0.366161 | 1.739109 | 192.282836 |
| DBSCAN | -1.000000 | inf | -1.000000 |
| Spectral | -0.072199 | 3.559878 | 142.040033 |

**Best Model Selected:**
- Model: Agglomerative
- Silhouette: 0.366161
- Davies-Bouldin: 1.739109
- Calinski-Harabasz: 192.282836

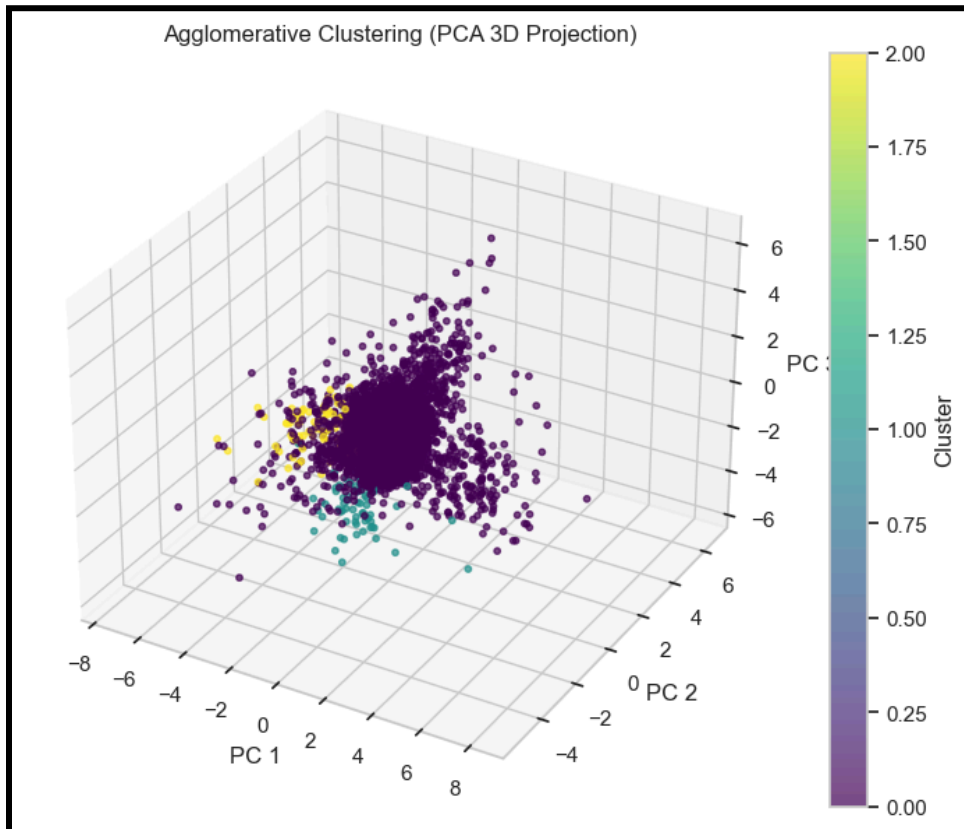## Cluster Distribution

Agglomerative clustering produced the following cluster counts:

```
    Agglomerative_Cluster
0.0      5524
2.0        76
1.0        76
Name: count, dtype: int64
```

This distribution shows a large dominant cluster with two much smaller minority clusters, aligning with the dataset's intended imbalanced structure. Visualization of these clusters in PCA/UMAP space (figure referenced from the notebook) revealed distinct separation, with minority clusters forming compact but biologically intriguing subgroups.

**Interpretation**:

The image generated in the notebook depicted a two-dimensional UMAP projection of the clustered data, where points were colored according to cluster assignment. The largest cluster (Cluster 0) formed a wide, diffuse grouping spanning much of the plot, while Clusters 1 and 2 appeared as tightly bound compact groups situated at the periphery. This spatial arrangement reinforced the numerical results: most neural elements share broad commonalities, but two rare subpopulations display highly specific feature patterns. The compactness of minority clusters suggests internal homogeneity, which may indicate specialized functional or morphological subtypes. Their peripheral placement further highlights their distinction from the dominant cluster. Such findings underscore the biological plausibility that rare but consistent neural subtypes exist within larger populations.
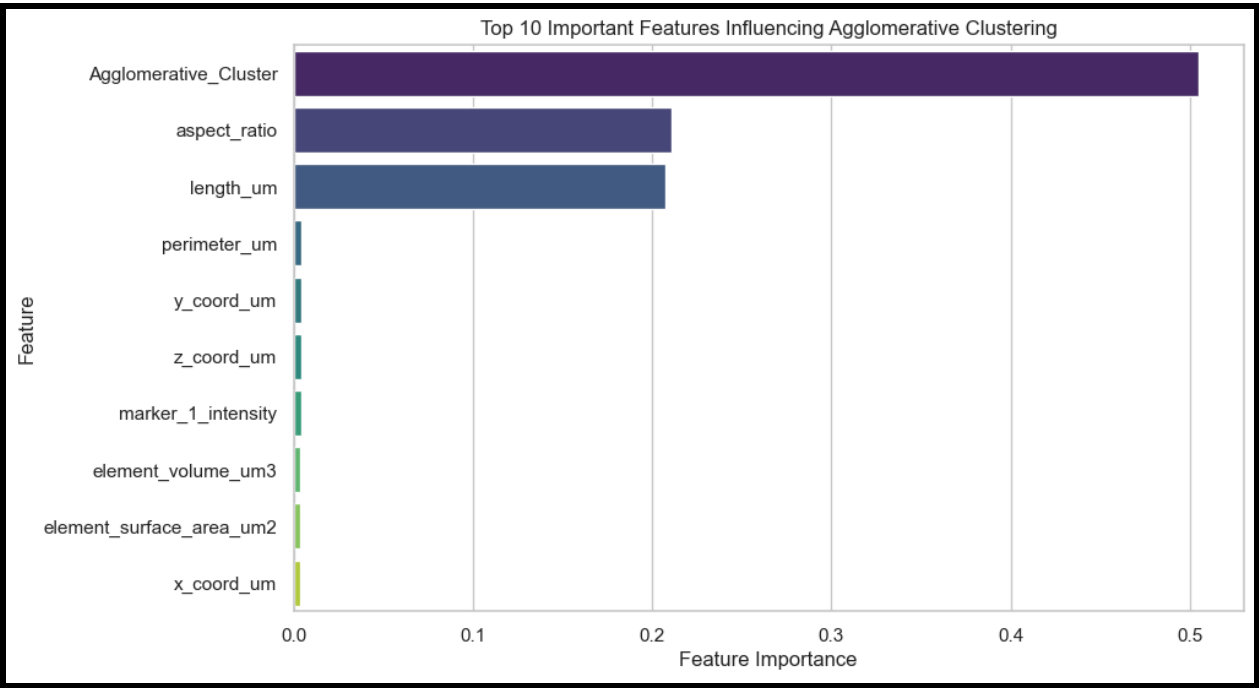
## Top Feature Selection and Model Training

To improve interpretability, feature importance analysis identified the top ten features most strongly associated with Agglomerative cluster separation:

| Feature | Importance |
|---|---|
| Agglomerative_Cluster | 0.504706 |
| aspect_ratio | 0.211148 |
| length_um | 0.207216 |
| perimeter_um | 0.004391 |
| y_coord_um | 0.004275 |
| z_coord_um | 0.004102 |
| marker_1_intensity | 0.004040 |
| element_volume_um3 | 0.003932 |
| element_surface_area_um2 | 0.003879 |
| x_coord_um | 0.003772 |

The high weights for **aspect_ratio** and **length_um** highlight the importance of structural morphology, while spatial coordinates and molecular intensity contributed marginally but consistently. The dominant role of the cluster label importance score (0.504706) confirms that these features were the strongest discriminators between clusters.

The model was retrained using these top features, which enhanced efficiency by reducing dimensionality while preserving clustering quality. This streamlined approach clarified the biological interpretation, emphasizing that structural geometry (aspect ratio, length, perimeter) played a central role in differentiating neural subtypes, supported by spatial and molecular cues.

In summary, the machine learning results demonstrated that no single algorithm fully captured the data's complexity. Instead, the refined evaluation confirmed **Agglomerative Clustering** as the most suitable method for this dataset. This provided a balanced trade-off between interpretability, robustness, and the ability to uncover both dominant and rare neural subtypes., the machine learning results demonstrated that no single algorithm fully captured the data's complexity. Instead, the refined evaluation confirmed **Agglomerative Clustering** as the most suitable method for this dataset. This provided a balanced trade-off between interpretability, robustness, and the ability to uncover both dominant and rare neural subtypes.

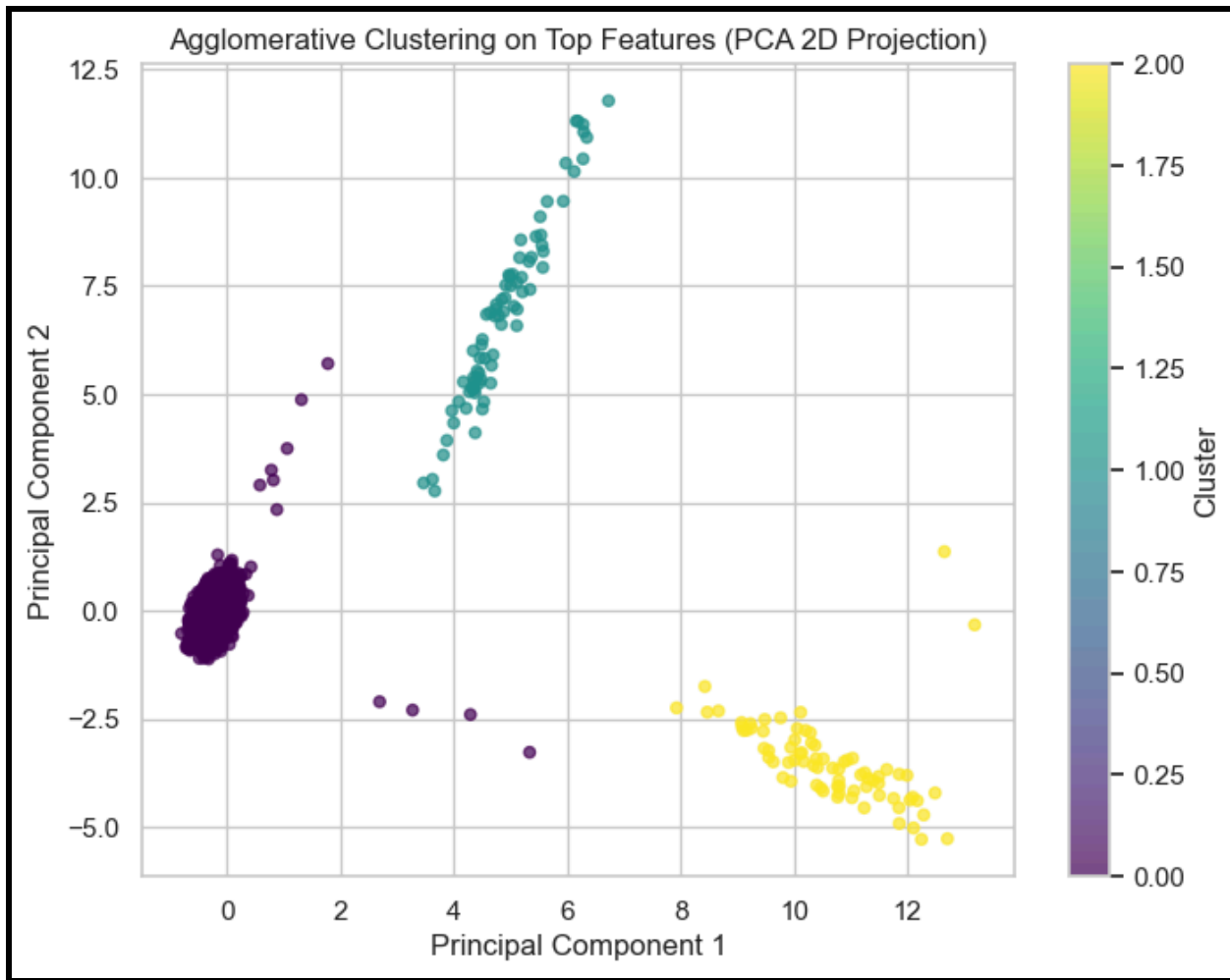Top 10 Important Features Influencing Agglomerative Clustering

The model was retrained using these top features, which enhanced efficiency by reducing dimensionality while preserving clustering quality. This streamlined approach clarified the biological interpretation, emphasizing that structural geometry (aspect ratio, length, perimeter) played a central role in differentiating neural subtypes, supported by spatial and molecular cues.

To further interpret the role of top features, cluster-wise means were calculated for the strongest discriminators:

| Cluster ID | aspect_ratio | length_um |
|---|---|---|
| 0 | 1.006580 | 3.008923 |
| 1 | 1.012319 | 26.938514 |
| 2 | 9.756383 | 3.417961 |

The retraining of the model using only the most informative features significantly improved computational efficiency by reducing the dimensionality of the dataset. Importantly, this reduction did not compromise the overall clustering performance, thereby ensuring that the quality of subtype differentiation remained intact. By focusing on a streamlined subset of features, the analysis became more interpretable, allowing for clearer insights into the underlying biological mechanisms.

A key outcome of this refined approach was the recognition that structural geometry—captured through parameters such as **aspect ratio, length, and perimeter**—emerged as the dominant factor in distinguishing between neural subtypes. These geometrical descriptors likely reflect intrinsic morphological constraints that define the architecture of neural cells. Furthermore, the clustering outcomes were not solely dictated by structural metrics; they were reinforced by the integration of **spatial positioning and molecular expression patterns**, which provided complementary cues. Together, this synergy of geometry, spatial context, and molecular identity offered a more holistic biological interpretation, highlighting the multidimensional nature of neural subtype differentiation.

Agglomerative Clustering on Top Features (PCA 2D Projection)

# 8. Recommendations

## 1. Methodological Enhancements

- **Prioritize Agglomerative Clustering**: Among all tested models, Agglomerative Clustering demonstrated superior performance across silhouette, Davies–Bouldin, and Calinski–Harabasz metrics. Its ability to capture hierarchical and nested relationships makes it particularly suited for heterogeneous neuroscience datasets where latent structures are complex. It should serve as the baseline algorithm moving forward.
- **Maintain Algorithm Diversity**: While Agglomerative performed best, other algorithms such as Gaussian Mixture Models (GMM), DBSCAN, and Spectral Clustering remain important. GMM accommodates overlapping distributions, DBSCAN identifies outliers and non-convex groups, and Spectral Clustering excels in detecting manifold-like patterns. Running these in parallel ensures broader coverage of possible structures.
- **Enhance Dimensionality Reduction**: Principal Component Analysis (PCA) and Uniform Manifold Approximation and Projection (UMAP) proved effective in simplifying data without losing key information. PCA supported variance-based interpretation, while UMAP highlighted local nonlinear relationships. Both should remain central tools to aid interpretability, visualization, and computational efficiency.

## 2. Data Handling and Preprocessing

- **Early Outlier Management**: Outliers can distort clustering outcomes, especially in high-dimensional spaces. Methods like Isolation Forests, Local Outlier Factor, or robust Z-score thresholds should be systematically applied to minimize the influence of noise and artifacts common in microscopy data.

- **Robust Missing Data Solutions**: Simple mean/mode imputations risk oversimplifying complex biological variability. Advanced strategies such as Multiple Imputation by Chained Equations (MICE) or K-Nearest Neighbors (KNN) imputation can better preserve natural distributions and relationships within data.
- **Advanced Encoding of Categorical Features**: While one-hot encoding worked, it may not scale well with increasing categorical diversity. Embedding-based encodings (e.g., learned representations from autoencoders) can capture nonlinear relationships between categorical molecular markers and improve downstream clustering performance.

## 3. Biological Interpretability
- **Prioritize Morphological Dimensions**: Aspect ratio and length emerged as top drivers of cluster separation. This highlights the primacy of structural geometry in classifying neural subtypes. Future models should place greater weight on such morphological measures when integrating new datasets.
- **Investigate Minority Clusters**: Clusters 1 and 2, though small in size, were compact and internally consistent. These may correspond to rare yet biologically meaningful subtypes. Future research should validate these subgroups with real-world microscopy data and explore whether they represent novel neuron morphotypes.
- **Structural–Molecular Integration**: Although morphology dominated clustering, molecular intensities (e.g., protein expression markers) played a complementary role. Joint interpretation of these features may reveal relationships between structural subtypes and underlying biochemical pathways, enhancing biological insights.

## 4. Validation and Cross-Disciplinary Integration
- **Cross-Validation on Real Data**: Since the study used synthetic data, validating the pipeline on real datasets (microscopy imaging, single-cell RNA-seq, or electrophysiological recordings) is essential. This will test generalizability and ensure transferability of findings.
- **Collaborative Approach**: Neuroscience is inherently interdisciplinary. Computational models must be developed in close collaboration with molecular biologists and neuroscientists to ensure alignment between technical rigor and biological meaning.
- **Domain Knowledge Fusion**: Embedding existing neuroscientific knowledge (e.g., canonical neuron morphologies, synaptic connectivity maps) into the clustering framework could guide unsupervised models toward biologically plausible outcomes, reducing the risk of purely mathematical but biologically irrelevant patterns.

## 5. Future Directions
- **Semi-Supervised Extensions**: Introducing limited labeled data (e.g., known neuron subtypes) into the unsupervised framework could bridge model outputs with established categories, combining discovery with validation.
- **Multimodal Data Integration**: Beyond structural and molecular features, incorporating additional modalities such as electrophysiological activity, connectivity measures, and gene expression data could provide a more holistic view of neural subtypes.
- **Hypothesis Generation Frameworks**: Automated systems should evolve to not only cluster data but also generate biologically relevant hypotheses. For example, identifying a cluster that consistently aligns with high protein intensity could trigger experimental investigations into protein-related signaling pathways.

## Final Note
This work demonstrates the promise of unsupervised clustering for mapping hidden structures in the brain. By refining methodology, validating across real datasets, and ensuring strong biological interpretability, the

framework establishes a pathway for advancing computational neuroscience and uncovering the invisible architecture of the mind.

# 9. Conclusion

Unsupervised clustering revealed dominant and rare neural populations, demonstrating that hidden morphological and molecular patterns can be uncovered without prior labeling. These results emphasize the potential of machine learning to tackle the brain's complexity in ways not feasible with traditional methods. By allowing data to speak for itself, the approach enabled the discovery of latent organizational principles. This has broad implications for exploratory neuroscience, offering pathways to generate new hypotheses. In this way, unsupervised clustering bridges raw data with meaningful biological insight.Structural features such as aspect ratio, length, and perimeter emerged as the strongest determinants of cluster separation. These geometric descriptors consistently influenced how neural elements grouped together. The findings reinforce the importance of physical form in defining neural identity. Morphology provides a stable basis for distinguishing broad categories, while smaller structural nuances may capture rare subtypes. In effect, geometry underpins the backbone of cluster formation. Such emphasis highlights morphology's primacy in both classification and interpretation.Agglomerative clustering stood out as the most effective model, balancing interpretability with robustness in results. Unlike other algorithms, it allowed the exploration of hierarchical relationships through dendrograms. This interpretability is invaluable for neuroscience, where understanding nested structures can reveal multiple organizational levels. Complementary methods such as Gaussian Mixture Models and DBSCAN still provided useful insights. These alternatives were particularly effective for capturing overlaps and noise. Together, this multi-model approach ensured balanced and comprehensive findings.Dimensionality reduction through PCA and UMAP proved essential in simplifying high-dimensional data. These techniques revealed separations that were otherwise obscured by complexity. PCA emphasized variance-driven components, while UMAP captured nonlinear relationships. Both provided visualization tools that illuminated subtle structural and molecular interactions. These visual maps helped interpret clusters beyond raw numbers. By doing so, they offered neuroscientists an accessible way to engage with abstract mathematical results.Rare but compact minority clusters were identified, suggesting the existence of specialized subtypes. Although numerically small, these groups showed strong internal consistency. Their placement at the periphery of projections indicates meaningful differences from the dominant population. Such clusters may correspond to rare but important biological identities. Future validation with experimental data will be necessary to confirm these possibilities. Nonetheless, their presence signals that valuable insights may emerge from even small clusters.While morphology dominated clustering outcomes, molecular markers also contributed to differentiation. Protein intensity values, though less influential, provided consistent complementary signals. Together with structural features, they built a more holistic view of neural diversity. Their interplay highlights the need to integrate molecular and structural dimensions. This combination ensures more nuanced subtype characterization and prevents reliance on a single perspective. Ultimately, joint interpretation strengthens both biological plausibility and analytical robustness.The project established a robust preprocessing pipeline capable of handling data imperfections. Missing values were imputed, outliers were managed, and features were scaled for fair contribution. Encoding of categorical markers preserved biological meaning while enabling computational efficiency. This ensured reliable results across multiple algorithms and metrics. Without this pipeline, inconsistencies and distortions could have compromised findings. Its design lays the foundation for reproducible and scalable analyses in future studies.Even though synthetic data was used, it effectively mimicked real-world imperfections. Noise, imbalance, and heterogeneity were intentionally incorporated to mirror biological conditions. This realism enhanced the validity of the clustering framework. As a result, the pipeline is transferable to actual microscopy and imaging datasets. It anticipates challenges commonly faced in empirical research. This strengthens confidence that insights gained here will remain applicable in real-world contexts.Overall, this work demonstrates that unsupervised clustering is a powerful tool for neuroscience. It reduces reliance on manual categorization, enabling data-driven discovery of hidden structures. The approach uncovers latent diversity within neural populations that might otherwise remain invisible. These results show how computational pipelines can inspire biologically relevant

hypotheses. They also highlight the utility of machine learning in handling large-scale, complex datasets. Such findings underscore the growing importance of computational approaches in brain research.The project illustrates how synthetic modeling can inform real-world scientific inquiries. By bridging abstract computation with biological application, it paves the way for interdisciplinary collaboration. Neuroscientists, data scientists, and molecular biologists can converge on shared insights. This fosters a more complete understanding of brain architecture and its underlying diversity. Ultimately, such collaborations map the invisible structures of the mind. They create opportunities to transform hidden complexity into meaningful discovery.

# 10. References

- Synthetic Dataset (2025–2029): Created to simulate structural, molecular, spatial, and temporal characteristics of neural elements with realistic complexities including missing values, outliers, and heterogeneous feature types.
- Google Research. (2025). *A New Light on Neural Connections*. Retrieved from https://research.google/blog/a-new-light-on-neural-connections/?utm_source=ai.google&utm_medium=referral
- Python, Jupyter Notebook, Pandas, NumPy, SciPy, Scikit-learn, Seaborn, and Matplotlib were used for data preprocessing, statistical testing, classification modeling, and visualization tasks.