# Computational Methods for ROC Identification in Xenopus Tail Regeneration: Denoising and Batch Integration Analysis

David Benson

October 9, 2025

**Abstract**

This study evaluates computational methods for identifying Regeneration-Organizing Cells (ROCs) in *Xenopus* tail tissue using single-cell RNA sequencing data from 13,199 cells. I implemented two clustering algorithms (Walktrap and Leiden), two marker selection methods (logistic regression and Wilcoxon rank-sum test), two denoising techniques (PCA reconstruction and k-nearest neighbor smoothing), and two batch integration approaches (Harmony and BBKNN). Baseline clustering achieved silhouette coefficients of 0.046 (Walktrap) and 0.132 (Leiden) with 19-26 gene overlap with published ROC markers. Denoising methods showed contrasting effects: PCA reconstruction provided minimal marker improvement (+2 genes) while k-NN smoothing achieved substantial improvement (+14 genes, reaching 40 total overlapping markers). Batch integration demonstrated strong clustering improvements with Harmony achieving 10.3-fold silhouette increase (0.046 to 0.471) and BBKNN showing 8.0-fold improvement, though neither method improved marker overlap with published data. The analysis suggests that while batch correction dramatically improves clustering structure, marker validation requires domain-specific optimization that may be limited by experimental design differences rather than methodological deficiencies.

## 1 Introduction

Single-cell RNA sequencing enables identification of rare cell populations, but requires systematic evaluation of computational methods to optimize detection accuracy. This study focuses on Regeneration-Organizing Cells (ROCs) in *Xenopus* tail regeneration, using the published dataset from Aztekin et al.[1] to compare computational approaches across four methodological categories required for the analysis.

The primary objective is to evaluate how different preprocessing and analysis choices affect ROC identification accuracy, measured through clustering quality metrics and validation against published marker genes. I implemented multiple approaches within each required category: clustering algorithms, marker selection methods, denoising techniques, and batch integration approaches.

## 2 Methods

### 2.1 Data and Preprocessing

I analyzed the published *Xenopus* tail regeneration dataset (13,199 cells × 31,535 genes) across intact and regenerating conditions with four experimental batches. Following the original study's approach, I selected 2,308 highly variable genes using Fano factor filtering:

$$F_i = \frac{\sigma_i^2}{\mu_i} \qquad (1)$$

where $F_i$ is the Fano factor for gene $i$, $\sigma_i^2$ is variance, and $\mu_i$ is mean expression. Genes were retained if mean expression fell between 5th and 80th percentiles and Fano factor exceeded the 65th percentile. The expression matrix was then transformed as:

$$X_{\text{transformed}} = \log_2(1 + X_{\text{CP10K}}) \qquad (2)$$

### 2.2 Clustering Algorithms

**Walktrap Algorithm**: Graph-based clustering using random walks on the k-nearest neighbor graph constructed from UMAP's fuzzy simplicial set. The transition probability matrix is defined as:

$$P_{ij} = \frac{w_{ij}}{\sum_k w_{ik}} \qquad (3)$$

where $P_{ij}$ is the transition probability from node $i$ to $j$, and $w_{ij}$ are edge weights from UMAP's fuzzy set construction ($k = 10$ neighbors, 10 steps).

**Leiden Algorithm**: Modularity optimization clustering on PCA-reduced data using resolution=3.0. The modularity function optimized is:

$$Q = \frac{1}{2m} \sum_{ij} \left[ A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j) \qquad (4)$$

where $A_{ij}$ is the adjacency matrix, $k_i$ is the degree of node $i$, $m$ is total edge weight, and $\delta(c_i, c_j) = 1$ if nodes $i, j$ are in the same community.

1

## 2.3 Marker Selection Methods

**Logistic Regression**: L2-regularized one-vs-rest multiclass classification:

$$\min_w \frac{1}{2} w^T w + C \sum_{i=1}^{n} \log(1 + e^{-y_i(X_i w + b)}) \qquad (5)$$

where $w$ are feature weights, $C$ is the inverse regularization strength ($C = 1.0$), and $y_i \in \{-1, +1\}$ indicates ROC membership.

**Wilcoxon Rank-Sum Test**: Non-parametric test statistic for differential expression:

$$U = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} S(X_{1i}, X_{2j}) \qquad (6)$$

where $S(a, b) = 1$ if $a > b$, 0.5 if $a = b$, and 0 otherwise. Effect size calculated as log2 fold change between group means with FDR correction at 0.05.

## 2.4 Denoising Techniques

**PCA Reconstruction**: Noise reduction via truncated SVD:

$$X_{\text{denoised}} = U_r S_r V_r^T + \mu \qquad (7)$$

where $r = 20$ components, $U_r, S_r, V_r$ are truncated matrices, and $\mu$ is the mean vector. This preserves approximately 85% of variance while removing noise in lower components.

**k-Nearest Neighbor Smoothing**: Local averaging in PCA space:

$$\hat{x}_i = \frac{1}{k} \sum_{j \in N_k(i)} x_j \qquad (8)$$

where $N_k(i)$ are the $k = 10$ nearest neighbors of cell $i$ using cosine distance in 30-dimensional PCA space.

## 2.5 Batch Integration Methods

**Harmony**: Iterative clustering in PCA space with soft cluster assignments:

$$Y_{\text{corrected}} = Y - \sum_{k} R_k \theta_k \qquad (9)$$

where $Y$ is the original PCA embedding, $R_k$ are soft cluster assignments, and $\theta_k$ are cluster-specific correction vectors learned through iterative optimization with $\theta = 1.0$.

**BBKNN (Batch Balanced k-NN)**: Graph correction by modifying neighborhood connectivity:

$$k_{\text{batch}} = \frac{k \cdot n_{\text{batch}}}{N} \qquad (10)$$

where neighbors are selected proportionally from each batch to maintain $k$ total connections while balancing batch representation.

## 2.6 Validation Approach

I validated computational results against Supplementary Table 3 from the original publication, comparing overlap in top 100 and 150 marker genes after canonicalizing *Xenopus* gene names (removing .L/.S allele suffixes). Clustering quality was assessed using silhouette coefficient:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \qquad (11)$$

where $a(i)$ is mean intra-cluster distance and $b(i)$ is mean nearest-cluster distance using cosine metric. Adjusted Rand Index provided chance-corrected agreement between clustering methods.

# 3 Results

## 3.1 Baseline Clustering Performance

UMAP embedding revealed distinct cell populations across 13,199 cells. Clustering algorithm comparison:

**Table 1:** Clustering algorithm performance comparison

| Method | Silhouette | Clusters | ARI |
|---|---|---|---|
| Walktrap | 0.046 | 48 | 0.330 |
| Leiden | 0.132 | 23 | 0.330 |

Leiden clustering achieved 2.9-fold better silhouette score than Walktrap (0.132 vs 0.046), indicating superior cluster separation in the PCA-reduced feature space.

## 3.2 Marker Selection and Validation

Both marker selection methods identified ROC-enriched genes with substantial overlap with published data:

**Baseline Performance (No Preprocessing)**:

- Top 100 markers: 19 genes overlap (19% validation rate)

- Top 150 markers: 26 genes overlap (17% validation rate)

- Key overlapping genes: fgf7, fgf9, lef1, sp9, krt, bmp5, nid2, pltp

**Method Concordance**: Logistic regression and Wilcoxon approaches showed complementary identification patterns, with parametric methods capturing multivariate signatures and non-parametric methods identifying genes with largest effect sizes.

## 3.3 Denoising Method Evaluation

Denoising approaches showed contrasting effects on marker recovery:

**Table 2:** Denoising method impact on marker validation

| Method | Markers | Improvement |
|---|---|---|
| No Denoising | 26 genes | – |
| PCA Reconstruction | 28 genes | +2 (+8%) |
| k-NN Smoothing | 40 genes | +14 (+54%) |

**PCA Reconstruction** (20 components): Provided minimal improvement in marker overlap despite improving clustering silhouette scores. The truncation approach may have removed biological signal relevant for ROC identification.

**k-NN Smoothing** ($k = 10$): Achieved substantial improvement in marker validation, suggesting that local averaging in PCA space enhanced detection of ROC-specific expression patterns while preserving biological signals.

## 3.4 Batch Integration Impact

Batch integration methods demonstrated dramatic improvements in clustering structure:

**Table 3:** Batch integration method comparison

| Method | Walktrap Sil. | Leiden Sil. | Factor |
|---|---|---|---|
| No Correction | 0.046 | 0.230 | – |
| Harmony | 0.471 | 0.234 | 10.3× |
| BBKNN | 0.366 | 0.223 | 8.0× |

**Harmony Integration**: Achieved 10.3-fold improvement in Walktrap clustering silhouette (0.046 to 0.471) through iterative batch correction while maintaining biological structure.

**BBKNN Integration**: Provided 8.0-fold improvement with faster computation, using graph-based neighborhood modification for batch balance.

**Limitation Observed**: Despite dramatic clustering improvements, both batch integration methods showed no improvement in marker overlap with published data (maintained at 19 genes for top 100, 24 genes for top 150). This suggests that temporal batch confounding (severity = 0.75) may be inherent to the experimental design rather than a technical artifact correctable through computational methods.

## 3.5 ROC Signature Validation

Computational analysis identified ROC populations through enrichment testing:

- **Strict signature**: 16 genes including key regeneration regulators (fgf7, fgf9, lef1, sp8, sp9, tp63, wnt3a, wnt5a)

- **Expanded signature**: 96 genes capturing broader regeneration-associated expression patterns

- **Top ROC clusters**: Clusters 23, 33, 6, 41, 38 showed highest signature enrichment

# 4 Discussion

## 4.1 Technical Performance Assessment

This systematic evaluation demonstrates that computational method selection substantially impacts ROC identification accuracy, though improvements are concentrated in specific analytical stages.

**Clustering Algorithm Selection**: Leiden's modularity optimization in PCA space consistently outperformed Walktrap's random walk approach (2.9-fold silhouette improvement), confirming that dimensionality reduction before clustering benefits rare population detection.

**Denoising Strategy Impact**: The contrasting performance of PCA reconstruction versus k-NN smoothing highlights the importance of method selection. k-NN smoothing's 54% improvement in marker validation suggests that local averaging preserves relevant biological signals better than global variance-based truncation.

**Batch Integration Effectiveness**: Both Harmony and BBKNN achieved substantial clustering improvements (8-10 fold), but the absence of marker validation improvement indicates that batch effects in this dataset may reflect legitimate biological or temporal variation rather than technical artifacts.

## 4.2 Methodological Limitations

The batch integration results illustrate an important limitation: computational methods cannot distinguish between technical batch effects and legitimate biological variation when batches are confounded with experimental conditions. The temporal structure of this dataset (batch-time confounding severity = 0.75) may explain why dramatic clustering improvements did not translate to better marker validation.

This suggests that while batch correction methods perform their intended function of improving clustering structure, the validation against published markers may be limited by differences in experimental design, cell iso-

lation protocols, or analytical approaches between studies rather than computational deficiencies.

## 4.3 Assignment Compliance

The analysis successfully implements all required methodological categories:

- **Clustering algorithms**: 2 methods (Walktrap, Leiden)

- **Marker selection**: 2 methods (logistic regression, Wilcoxon)

- **Denoising techniques**: 2 methods (PCA reconstruction, k-NN smoothing)

- **Batch integration**: 2 methods (Harmony, BBKNN)

# 5 Conclusion

This computational evaluation identifies k-NN smoothing denoising as the most effective preprocessing step for ROC marker detection, improving validation overlap by 54% (26 to 40 genes). Batch integration methods provide substantial clustering improvements but show no marker validation benefit, likely reflecting experimental design constraints rather than methodological limitations.

Key findings:

1. **Leiden clustering consistently outperforms Walktrap** for rare cell detection ($2.9\times$ silhouette improvement)

2. **k-NN smoothing denoising provides substantial marker improvement** ($+54\%$ validation overlap)

3. **Batch integration improves clustering structure** ($8\text{-}10\times$ silhouette improvement) but not marker validation

4. **Method selection affects different analytical stages differently**, requiring stage-specific optimization

The results demonstrate that systematic method evaluation is essential for optimizing rare cell detection pipelines, though domain expertise remains crucial for interpreting validation results in the context of experimental design limitations.

## 5.1 Limitations and Future Directions

This analysis represents a computational methods comparison on a single dataset with inherent batch-time confounding. The validation approach relies on published markers that may not capture the complete ROC expression signature. Future work should evaluate these methods across multiple datasets with varying batch effect structures and validate results through orthogonal experimental approaches.

## 5.2 Code and Reproducibility

Complete computational pipeline with all methods and results available at: `https://github.com/db-d2/stat4243_proj1`. Analysis performed using Python 3.10+ with scanpy 1.9+, with fixed random seeds (42) for reproducibility. All figures, clustering visualizations, and method comparison plots are provided in the Supplementary Materials document.

## 5.3 AI Assistance Acknowledgment

Claude AI (Anthropic) was used for code debugging, formatting alignment with PEP8 standards, language simplification, LaTeX formatting, and condensing material into fewer pages. All computational analysis and result interpretation were performed independently.

# References

[1] Aztekin C, et al. (2019). Identification of a regeneration-organizing cell in the Xenopus tail. *Science* 364(6441):653-658.

[2] Luecken MD, et al. (2022). Benchmarking atlas-level data integration in single-cell genomics. *Nature Methods* 19(1):41-50.