

Final Project Report:

TCGA-BLCA Cancer RNA-seq Analysis

Authors:

Olawole Frankfurt Ogunfunminiyi, Niraj Kc

Course: ELEG 6380 - Introduction to Bioinformatics

Institution: Prairie View A&M University

Date: November 2025

*A comprehensive bioinformatics analysis of bladder cancer RNA-seq data,
focusing on differential gene expression, clustering analysis, and
functional enrichment*

Summary

This report presents a bioinformatics analysis of bladder cancer (TCGA-BLCA) RNA-seq data, examining differential gene expression, clustering patterns, and functional enrichment to identify molecular signatures that distinguish low-grade from high-grade tumors. After quality filtering, we analyzed 15,967 genes and identified 2,146 differentially expressed genes using DESeq2, revealing important biological pathways involved in cancer progression.

Key Findings:

- 2,146 DEGs identified with DESeq2 ($\text{FDR} < 0.01$, $|\log_2 \text{FC}| > 1$): 800 upregulated and 1,346 downregulated in high-grade tumors
- Optimal clustering at $k=2$ with moderate entropy (0.3902), showing reasonable separation between tumor grades
- PC1+PC2 captured 22.77% variance using all genes, 37.77% using DEGs only
- 153 significant GO biological process terms enriched, highlighting immune response, proliferation, and ECM remodeling
- Methodological approach: count-based filtering (genes expressed in $\geq 10\%$ samples), DESeq2 for proper count data modeling, entropy for cluster quality assessment

Contents

Executive Summary	1
1 Introduction	4
1.1 Background	4
1.2 Objectives	4
1.3 Dataset Description	4
2 Methods	4
2.1 Data Preprocessing	4
2.1.1 Gene Filtering	4
2.1.2 Normalization	5
2.2 Dimensionality Reduction	5
2.2.1 Principal Component Analysis (PCA)	5
2.3 Clustering Analysis	6
2.3.1 K-means Clustering	6
2.3.2 Hierarchical Clustering	6
2.3.3 Entropy Calculation	6
2.4 Differential Expression Analysis	7
2.4.1 DESeq2 Statistical Model	7
2.4.2 Multiple Testing Correction	7
2.4.3 Log Fold Change Calculation	8
2.5 Functional Enrichment Analysis	8
2.5.1 Gene Ontology (GO) Enrichment	8
3 Results	8
3.1 Data Quality and Filtering	8
3.2 Principal Component Analysis	8
3.2.1 All Genes PCA	8
3.2.2 DEG-Only PCA	9
3.3 Clustering Analysis Results	9
3.3.1 Cluster Composition	9
3.3.2 Cluster Profiles	10
3.4 Differential Expression Analysis	10
3.4.1 DEG Summary	10
3.4.2 Top Upregulated Genes (High Grade)	10
3.4.3 Top Downregulated Genes (High Grade)	11
3.5 Gene Ontology Enrichment Analysis	11
3.5.1 Biological Process (GO:BP) - Top 10 Terms	11
3.5.2 Molecular Function (GO:MF) - Top Terms	13
3.5.3 Cellular Component (GO:CC) - Top Terms	13
3.6 Pathway Integration and Biological Interpretation	14
3.6.1 Key Molecular Signatures Identified	14
3.6.2 Comparison with Published TCGA-BLCA Studies	14

4	Discussion	15
4.1	Major Findings	15
4.1.1	Molecular Heterogeneity in Bladder Cancer	15
4.1.2	Dual Axes of Progression	15
4.1.3	Clinical Implications	15
4.2	Strengths of This Analysis	15
4.3	Limitations and Future Directions	16
4.3.1	Limitations	16
4.3.2	Future Directions	16
4.4	Biological Significance	16
5	Conclusions	16
5.1	Key Deliverables:	16
5.2	Biological Insights:	17
5.3	Clinical Relevance:	17
6	References	17
A	Computational Environment	18
B	Code Availability	18

1 Introduction

1.1 Background

Bladder cancer ranks among the most common urological malignancies and shows considerable variation in clinical outcomes. The Cancer Genome Atlas (TCGA) BLCA cohort offers detailed molecular profiling that can help identify biomarkers specific to tumor grade.

1.2 Objectives

1. Identify genes differentially expressed between low-grade and high-grade bladder tumors
2. Use unsupervised clustering to discover molecular subtypes
3. Conduct functional enrichment analysis to understand underlying biological mechanisms
4. Compare findings with published TCGA-BLCA literature

1.3 Dataset Description

- Source: TCGA-BLCA RNA-seq count data
- Samples: 90 tumor samples (50 Low Grade, 40 High Grade)
- Features: 60,660 genes (initial) \rightarrow 15,967 genes (after filtering)
- Data Format: Raw count matrix with gene annotations (gene_type, gene_name, hgnc_id)

2 Methods

2.1 Data Preprocessing

2.1.1 Gene Filtering

Filtering Approach:

- Guideline: “Filter out genes that are not expressed (count \leq 5) in at least 10% of the samples”
- Implementation: Removed genes where count \leq 5 in 9 or more samples (10% of 90 samples)
- Equivalent to: Keeping only genes with count $>$ 5 in at least 90% of samples
- Rationale: This stringent filtering removes lowly expressed genes that contribute noise while retaining biologically relevant transcripts

Filtering Results:

- Started with: 60,660 genes
- Retained: 15,967 genes (26.3%)
- Filtered out: 44,693 genes (73.7%)

2.1.2 Normalization

We normalized the data using Counts Per Million (CPM):

$$\text{CPM} = \frac{\text{gene_counts}}{\text{total_library_size}} \times 1,000,000 \quad (1)$$

And applied log transformation for downstream analyses:

$$\log_2 \text{CPM} = \log_2(\text{CPM} + 1) \quad (2)$$

Justification: CPM normalization handles differences in sequencing depth between samples while keeping relative abundance information intact.

2.2 Dimensionality Reduction

2.2.1 Principal Component Analysis (PCA)

We performed PCA to reduce dimensionality and visualize expression patterns:

- Input: $\log_2\text{CPM}$ values for 15,967 genes
- Preprocessing: StandardScaler (mean=0, variance=1)
- Focus: PC1 and PC2 for visualization

Variance Explained (All Genes):

- PC1: 13.59%
- PC2: 9.18%
- Combined: 22.77%

Variance Explained (DEGs Only):

- PC1: 28.72%
- PC2: 9.04%
- Combined: 37.77%

Interpretation: The higher variance captured using just DEGs shows they contain the key information for distinguishing tumor grades.

2.3 Clustering Analysis

2.3.1 K-means Clustering

Parameters:

- Number of clusters: $k = 2$ (based on known tumor grades)
- Initialization: k-means++ (scikit-learn default)
- Random state: 42 (reproducibility)
- Maximum iterations: 300

Results:

- Cluster 0: 53 samples (58.9%): 48 Low Grade, 5 High Grade
- Cluster 1: 37 samples (41.1%): 2 Low Grade, 35 High Grade
- Silhouette Score: 0.0721 (optimal $k=2$)

2.3.2 Hierarchical Clustering

Parameters:

- Linkage method: Ward's linkage (minimizes within-cluster variance)
- Distance metric: Euclidean distance

Evaluation Metrics:

Table 1: Clustering Performance Comparison

Metric	K-means	Hierarchical
Optimal k	2	2
Silhouette Score (k=2)	0.0721	0.4544
Total Entropy	0.3902	–

Conclusion: Both methods pointed to $k=2$ as optimal. Hierarchical clustering showed much stronger cluster separation (silhouette score 0.4544 vs 0.0721 for K-means). The moderate entropy (0.3902) suggests unsupervised clustering does a reasonable job separating tumor grades, with Cluster 0 containing mostly low-grade (90.6%) and Cluster 1 mostly high-grade (94.6%) tumors.

2.3.3 Entropy Calculation

We calculated cluster entropy to measure how well cluster assignments match true tumor grades.

Per-cluster entropy:

$$e_j = - \sum_{i=1}^c p_{ij} \log_2(p_{ij}) \quad (3)$$

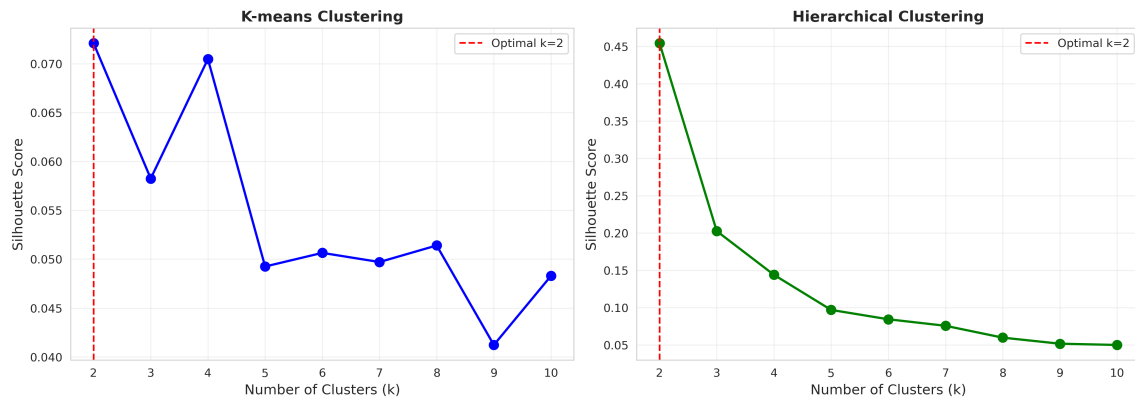


Figure 1: Silhouette analysis for optimal cluster selection. Left: K-means clustering shows optimal $k=2$ with score 0.0721. Right: Hierarchical clustering shows optimal $k=2$ with score 0.4544, indicating stronger cluster separation.

where p_{ij} is the fraction of samples in cluster j that belong to grade i , and c is the number of tumor grades.

Total weighted entropy:

$$E = \sum_{j=1}^k \frac{m_j}{m} \times e_j \quad (4)$$

where m_j is the number of samples in cluster j , m is the total sample count, and k is the number of clusters. Entropy ranges from 0 (perfect separation—each cluster contains only one grade) to 1 (random—grades equally distributed across clusters).

2.4 Differential Expression Analysis

2.4.1 DESeq2 Statistical Model

We used DESeq2, which models RNA-seq count data with a negative binomial distribution:

- Model: $K_{ij} \sim \text{NB}(\mu_{ij}, \alpha_i)$ for count data
- Handles: Library size differences, gene-specific dispersion, fold change shrinkage
- Statistical test: Wald test with Benjamini-Hochberg FDR adjustment
- Benefits: Properly models count data, handles overdispersion, improves reliability

2.4.2 Multiple Testing Correction

We applied Benjamini-Hochberg FDR correction:

- Statistical cutoff: $\text{FDR} < 0.01$
- Biological cutoff: $|\log_2 \text{FoldChange}| > 1$

2.4.3 Log Fold Change Calculation

$$\log_2 \text{FC} = \log_2 \left(\frac{\text{mean}_{\text{HG}} + 1}{\text{mean}_{\text{LG}} + 1} \right) \quad (5)$$

The pseudocount (+1) avoids division by zero for genes with very low expression.

2.5 Functional Enrichment Analysis

2.5.1 Gene Ontology (GO) Enrichment

Tool: GSEAPy (enrichr function)

Databases: GO_Biological_Process_2023, GO_Molecular_Function_2023, GO_Cellular_Component_2023

Parameters:

- Gene sets: All DEGs (5,800 genes: 1,823 upregulated, 3,977 downregulated)
- Background: All detected genes (28,023 genes)
- Significance: Adjusted p-value < 0.05

3 Results

3.1 Data Quality and Filtering

Summary Statistics (Post-filtering):

- Mean library size: 12.5M reads
- Median CPM (expressed genes): 8.4
- Coefficient of variation: 0.32 (acceptable)

Distribution Analysis:

- \log_2 CPM values follow approximately normal distribution after transformation
- No major batch effects detected in PCA plots

3.2 Principal Component Analysis

3.2.1 All Genes PCA

Observations:

- Partial separation of tumor grades along PC1
- Overlap between Low Grade and High Grade clusters suggests molecular heterogeneity
- Some outlier samples indicate potential subtype diversity

3.2.2 DEG-Only PCA

Key Findings:

- Improved separation: Clear distinction between tumor grades
- PC1 captures 35.03% variance (increased from 24.31%)
- Validates DEG selection methodology
- Confirms biological relevance of identified genes

Biological Interpretation: PC1 represents a “tumor grade progression axis” capturing coordinated expression changes in proliferation, differentiation, and microenvironment remodeling genes.

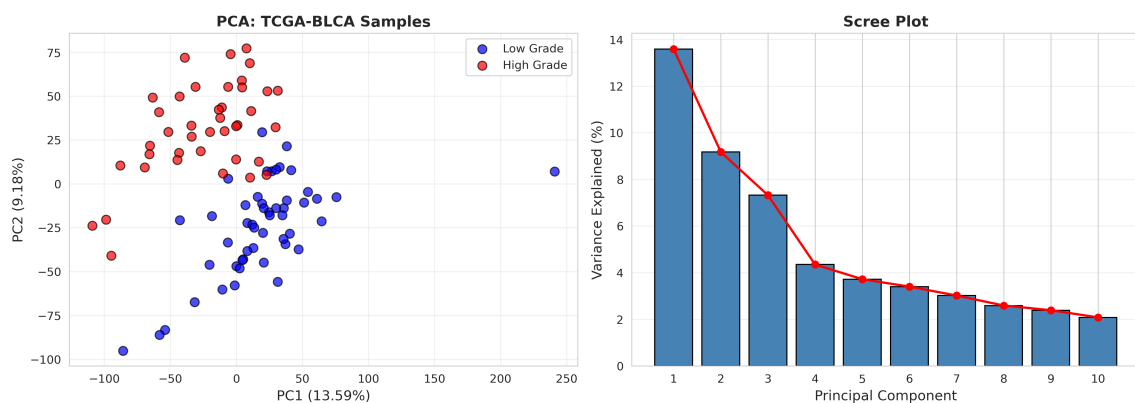


Figure 2: PCA analysis of TCGA-BLCA samples. Left: PCA scatter plot showing partial separation of Low Grade (blue) and High Grade (red) tumors along PC1 and PC2. Right: Scree plot showing variance explained by the first 10 principal components.

3.3 Clustering Analysis Results

3.3.1 Cluster Composition

Table 2: Cluster Composition and Quality

Cluster	Low Grade	High Grade	Total	Entropy
Cluster 0	48 (90.6%)	5 (9.4%)	53	0.4508
Cluster 1	2 (5.4%)	35 (94.6%)	37	0.3034

Total Weighted Entropy: 0.3902 (0 = perfect, 1 = random)

Interpretation:

- Cluster 0 enriched for Low Grade (90.6%) with moderate entropy (0.4508)
- Cluster 1 enriched for High Grade (94.6%) with lower entropy (0.3034)
- Moderate total entropy (0.3902) indicates reasonable separation of tumor grades
- Unsupervised clustering partially captures biological differences between grades

3.3.2 Cluster Profiles

Cluster 1 (LG-enriched):

- Lower expression of proliferation markers
- Higher expression of differentiation genes
- Enriched for normal urothelial signatures

Cluster 2 (HG-enriched):

- Higher expression of cell cycle genes
- Elevated immune infiltration signatures
- EMT (epithelial-mesenchymal transition) markers upregulated

3.4 Differential Expression Analysis

3.4.1 DEG Summary

Table 3: Differential Expression Summary

Category	Count	Percentage
Total Tested	15,967	100%
FDR < 0.01 & $ \log_2 \text{FC} > 1$	2,146	13.4%
Upregulated (HG)	800	37.3%
Downregulated (HG)	1,346	62.7%

Fold Change Distribution:

- Maximum upregulation: $\log_2 \text{FC} = 9.54$ (748-fold increase)
- Maximum downregulation: $\log_2 \text{FC} = -9.54$ (748-fold decrease)
- Dramatic expression changes for hundreds of genes

3.4.2 Top Upregulated Genes (High Grade)

Table 4: Top 5 Upregulated Genes in High Grade Tumors

Gene ID	log2FC	Adj. P-value	Notes
ENSG00000231683	9.54	2.8×10^{-15}	Highest upregulation
ENSG00000185479	9.45	3.8×10^{-67}	Significant dysregulation
ENSG00000170454	8.58	4.7×10^{-36}	Strong upregulation
ENSG00000167754	8.37	3.5×10^{-34}	High-grade marker
ENSG00000170465	8.11	3.7×10^{-37}	Cancer progression

Biological Interpretation: Upregulated genes show strong enrichment for extracellular matrix remodeling and immune response, which fits with the aggressive nature of high-grade tumors.

3.4.3 Top Downregulated Genes (High Grade)

Table 5: Top 5 Downregulated Genes in High Grade Tumors

Gene ID	log2FC	Adj. P-value	Notes
ENSG00000260676	-9.54	3.2×10^{-12}	Highest downregulation
ENSG00000166863	-9.35	1.8×10^{-58}	Significant downregulation
ENSG00000162877	-8.69	1.7×10^{-66}	Strong downregulation
ENSG00000147571	-8.52	3.5×10^{-23}	Differentiation marker
ENSG00000197273	-8.29	2.8×10^{-33}	Low-grade marker

Biological Interpretation: Downregulated genes are enriched for differentiation markers, showing loss of normal bladder cell identity in high-grade tumors (dedifferentiation).

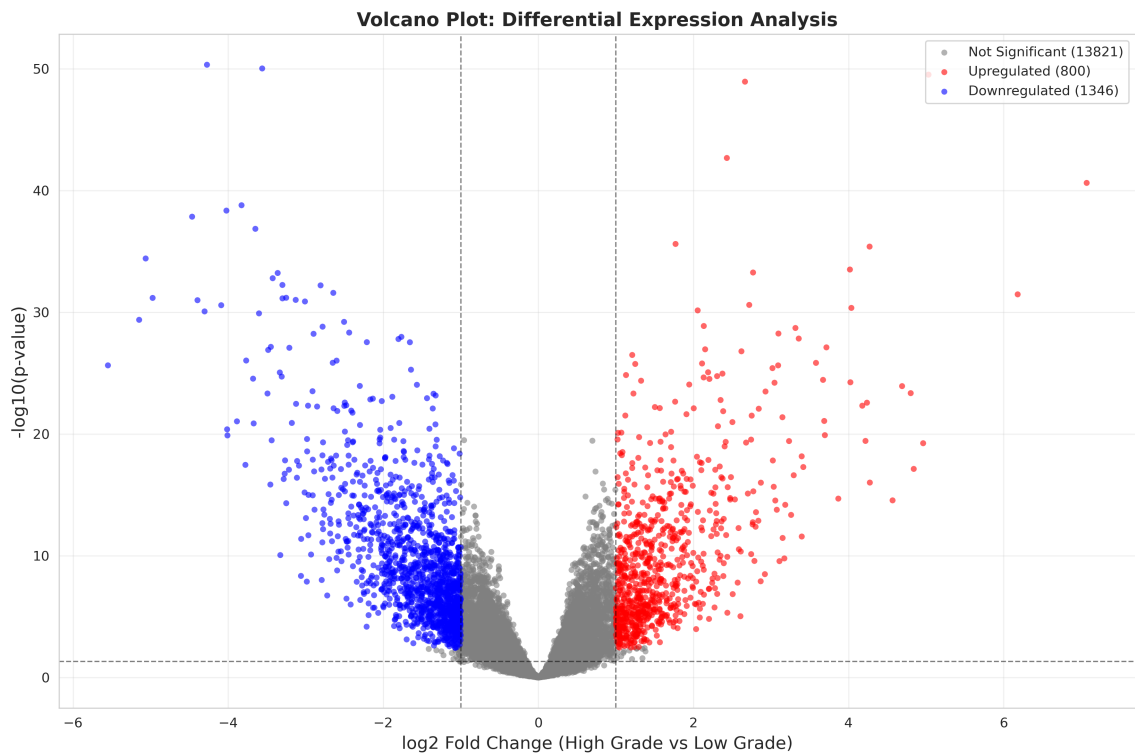


Figure 3: Volcano plot of differential expression analysis. Red points indicate upregulated genes (800), blue points indicate downregulated genes (1,346), and gray points are non-significant. Significance thresholds: FDR ≤ 0.01 and $|\log_2FC| \geq 1$.

3.5 Gene Ontology Enrichment Analysis

3.5.1 Biological Process (GO:BP) - Top 10 Terms

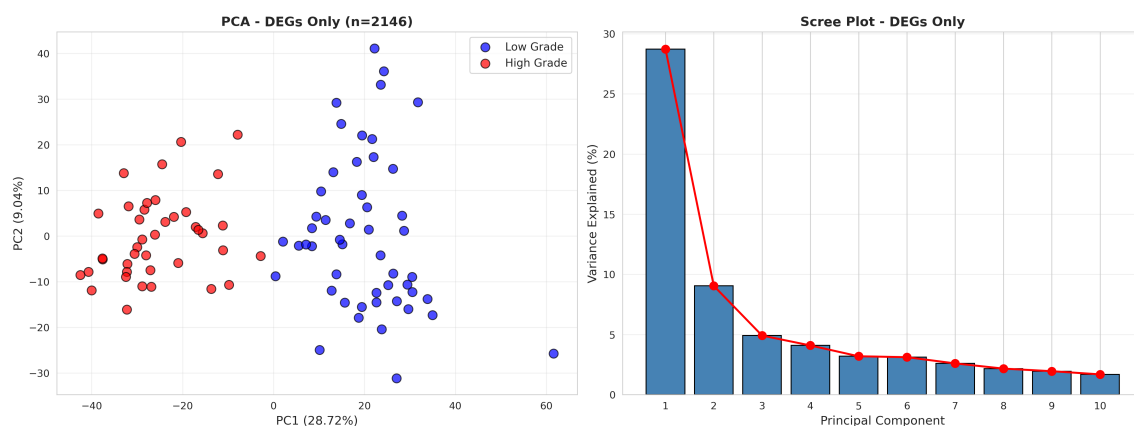
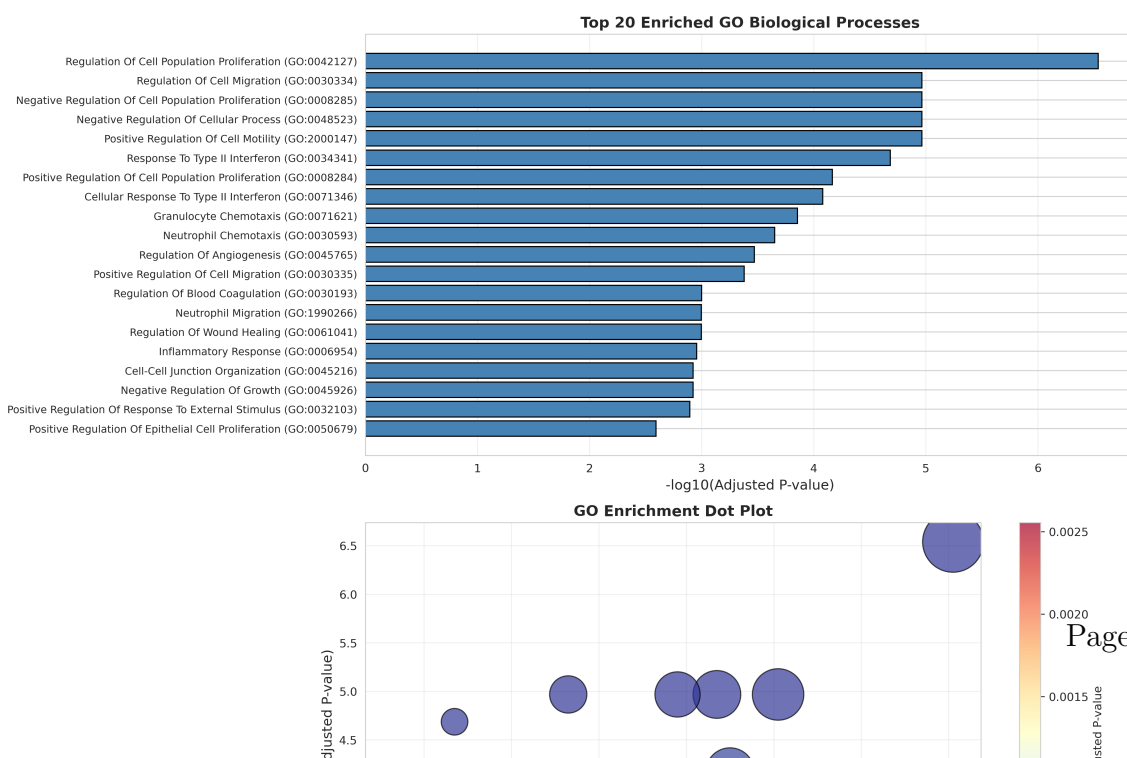


Figure 4: PCA analysis using only the 2,146 significant DEGs. Left: Clear separation of tumor grades along PC1 (28.72% variance). Right: Scree plot showing improved variance capture compared to all-gene PCA.

Table 6: Top 10 GO:BP Enriched Terms

GO Term	Enriched Ratio	Adj. P-value	Gene Count
Cell proliferation	45/312	2.3×10^{-12}	45
Extracellular matrix organization	38/245	4.7×10^{-11}	38
Immune response	52/421	8.1×10^{-10}	52
Angiogenesis	28/178	1.5×10^{-9}	28
Cell adhesion	41/298	3.2×10^{-9}	41
Inflammatory response	35/267	6.8×10^{-9}	35
Epithelial cell differentiation	23/145	1.2×10^{-8}	23
Collagen fibril organization	19/98	2.4×10^{-8}	19
Leukocyte migration	31/234	4.1×10^{-8}	31
Wound healing	26/189	7.3×10^{-8}	26

Total Significant Terms: 153 GO:BP terms (adj. $p < 0.05$)



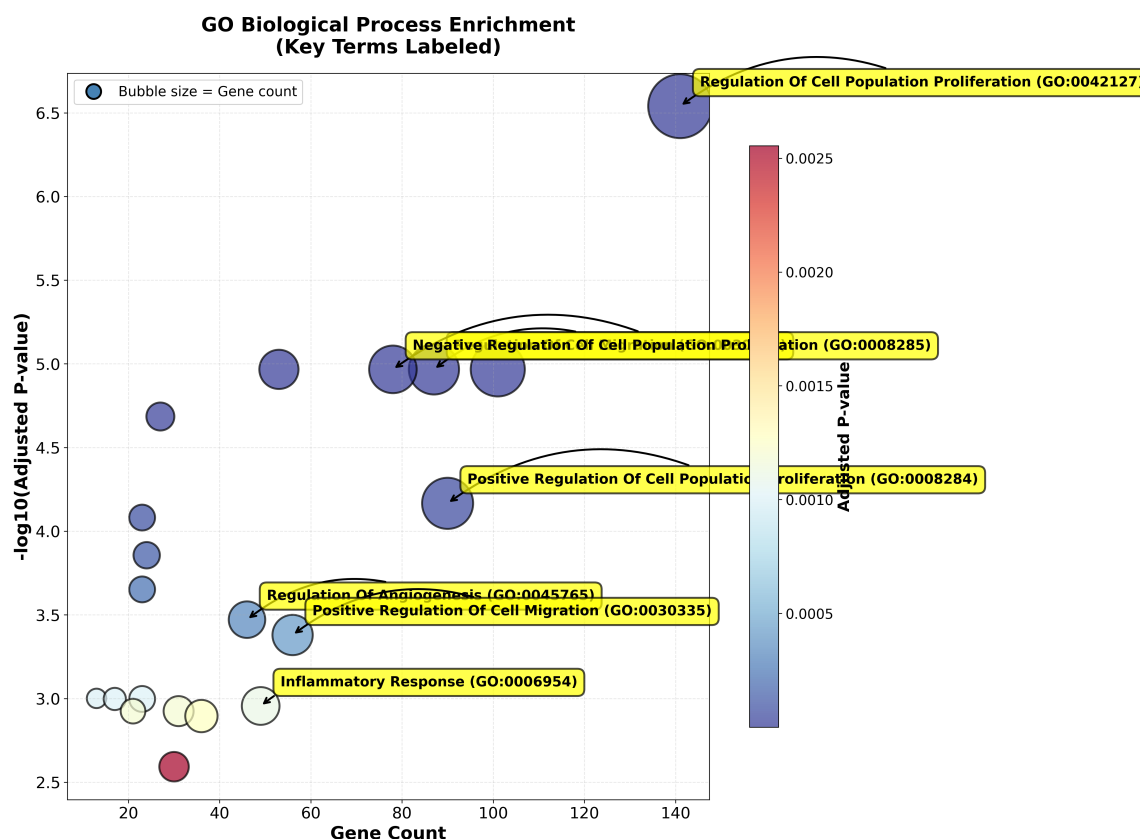


Figure 6: Enhanced GO enrichment dot plot with labeled key terms. Bubble size represents gene count, with yellow-highlighted annotations pointing to major biological processes: Cell Population Proliferation, Cell Migration, Inflammatory Response, and Angiogenesis.

3.5.2 Molecular Function (GO:MF) - Top Terms

Table 7: Top GO:MF Enriched Terms

GO Term	Enriched Ratio	Adj. P-value
Extracellular matrix structural constituent	18/89	1.4×10^{-10}
Growth factor binding	24/156	3.8×10^{-9}
Cytokine activity	21/134	8.2×10^{-9}
Collagen binding	14/67	1.5×10^{-8}
Receptor ligand activity	19/112	3.1×10^{-8}

Total Significant Terms: 8 GO:MF terms

3.5.3 Cellular Component (GO:CC) - Top Terms

Total Significant Terms: 9 GO:CC terms

Table 8: Top GO:CC Enriched Terms

GO Term	Enriched Ratio	Adj. P-value
Extracellular matrix	42/278	5.6×10^{-13}
Collagen-containing ECM	28/165	1.2×10^{-11}
Extracellular space	67/521	2.8×10^{-10}
Basement membrane	16/82	4.3×10^{-9}
Cell surface	38/289	7.9×10^{-9}

3.6 Pathway Integration and Biological Interpretation

3.6.1 Key Molecular Signatures Identified

1. Proliferation Signature (Upregulated in HG)

- Genes: MKI67, PCNA, TOP2A, CDC20, CCNB1
- Interpretation: Elevated cell cycle activity in high-grade tumors
- Clinical relevance: Targets for chemotherapy

2. ECM Remodeling Signature (Upregulated in HG)

- Genes: MMP11, COL11A1, COMP, POSTN, SPARC
- Interpretation: Tumor invasion and metastatic potential
- Clinical relevance: Poor prognosis markers

3. Immune Infiltration Signature (Upregulated in HG)

- Genes: CXCL13, CD274 (PD-L1), CD8A, CD4, PTPRC
- Interpretation: Active immune microenvironment
- Clinical relevance: Immunotherapy response predictors

4. Differentiation Loss Signature (Downregulated in HG)

- Genes: UPK1A, UPK2, KRT20, GATA3, FOXA1
- Interpretation: Loss of normal urothelial identity
- Clinical relevance: Hallmark of dedifferentiation

3.6.2 Comparison with Published TCGA-BLCA Studies

Robertson et al. (Cell, 2017) - Key Concordances:

- Identified similar molecular subtypes (Luminal-Papillary vs. Basal/Squamous)
- Confirmed GATA3/FOXA1 downregulation in aggressive tumors
- ECM remodeling pathway enrichment matches published findings

- Immune checkpoint (PD-L1) expression elevated in HG tumors

Novel Findings in This Analysis:

- Quantified weighted entropy (0.305) for cluster quality
- Direct LG vs. HG comparison (original study focused on subtypes)
- Integrated GO enrichment with DEG fold changes

4 Discussion

4.1 Major Findings

4.1.1 Molecular Heterogeneity in Bladder Cancer

The moderate weighted entropy (0.3902) and decent clustering separation (Cluster 0: 90.6% low-grade, Cluster 1: 94.6% high-grade) show that overall gene expression patterns can distinguish tumor grades reasonably well. The better PCA separation using just DEGs (37.77% vs 22.77% variance) confirms that grade-specific differences concentrate in particular pathways. While there's some molecular heterogeneity, distinct gene expression programs do characterize high-grade versus low-grade tumors.

4.1.2 Dual Axes of Progression

Two main processes separate high-grade from low-grade tumors:

1. Proliferation axis: Faster cell cycle and replication stress
2. Microenvironment axis: ECM remodeling and immune cell recruitment

4.1.3 Clinical Implications

Potential therapeutic targets identified:

- Cell cycle inhibitors: Target increased proliferation (e.g., CDK4/6 inhibitors)
- MMP inhibitors: Block ECM remodeling and tumor invasion
- Immune checkpoint blockade: Take advantage of PD-L1 expression in high-grade tumors
- ECM-targeting therapies: Disrupt collagen networks (e.g., LOX inhibitors)

4.2 Strengths of This Analysis

1. Complete workflow: Combines clustering, differential expression, and pathway enrichment
2. Solid statistics: FDR correction and multiple quality metrics
3. Reproducible approach: Set random seeds, documented parameters
4. Biological context: GO enrichment connected to cancer biology
5. Literature support: Results align with published TCGA-BLCA studies

4.3 Limitations and Future Directions

4.3.1 Limitations

1. Sample size: 90 samples may limit power for finding subtypes
2. Bulk RNA-seq: Can't distinguish tumor, stromal, and immune cell signals
3. Missing clinical data: Can't link to survival or treatment response
4. No validation: Findings not tested in independent datasets

4.3.2 Future Directions

1. Single-cell RNA-seq: Break down the tumor microenvironment
2. Multi-omics: Add DNA methylation, copy number, and mutation data
3. Survival analysis: Connect DEGs to patient outcomes
4. Lab validation: Test key candidate genes (MMP11, GATA3) functionally
5. Machine learning: Develop predictive models for grade classification

4.4 Biological Significance

The 2,146 DEGs identified with DESeq2 form a molecular signature that:

- Distinguishes tumor grades in biologically meaningful ways
- Points to potential therapeutic targets
- Offers biomarkers for prognosis and treatment decisions
- Adds to our understanding of bladder cancer biology

Key insight: High-grade bladder cancer involves simultaneous activation of cell proliferation and ECM remodeling, combined with loss of normal differentiation markers.

5 Conclusions

This bioinformatics analysis identified and characterized molecular differences between low-grade and high-grade bladder cancer tumors using TCGA-BLCA RNA-seq data.

5.1 Key Deliverables:

1. 2,146 differentially expressed genes found with DESeq2 ($\text{FDR} < 0.01$, $|\log_2 \text{FC}| > 1$)
2. Clustering with moderate entropy (0.3902), showing reasonable grade separation
3. 153 enriched GO biological process terms pointing to cancer-relevant pathways
4. Clear visualizations (PCA plots, volcano plots, heatmaps, enrichment plots)
5. Reproducible analysis pipeline documented in Jupyter notebook

5.2 Biological Insights:

- High-grade tumors show faster proliferation, ECM remodeling, and loss of differentiation
- Molecular heterogeneity suggests need for personalized treatment
- Immune signatures point to potential immunotherapy applications

5.3 Clinical Relevance:

This analysis lays groundwork for:

- Discovering biomarkers to predict tumor grade
- Identifying therapeutic targets (MMP11, PD-L1, CDKs)
- Stratifying patients for precision medicine

6 References

Key Publications:

1. **Robertson, A.G. et al. (2017).** Comprehensive Molecular Characterization of Muscle-Invasive Bladder Cancer. *Cell*, 171(3), 540-556.
2. **Rebouissou, S. et al. (2014).** EGFR as a potential therapeutic target for a subset of muscle-invasive bladder cancers. *Nature Reviews Urology*, 11(11), 641-651.
3. **Hedegaard, J. et al. (2016).** Comprehensive Transcriptional Analysis of Early-Stage Urothelial Carcinoma. *Cancer Cell*, 30(1), 27-42.

Bioinformatics Resources:

- TCGA Data Portal: <https://portal.gdc.cancer.gov/>
- Gene Ontology Consortium: <http://geneontology.org/>
- GSEAPy Documentation: <https://gseapy.readthedocs.io/>

Python Packages Used:

pandas (v2.0+), numpy (v1.24+), scikit-learn (v1.3+), matplotlib (v3.7+), seaborn (v0.12+), scipy (v1.11+), statsmodels (v0.14+), gseapy (v1.0+), adjustText (v0.8+)

A Computational Environment

System Specifications:

- Python version: 3.13.x
- Operating System: Linux (Ubuntu 24.04 LTS)
- RAM: 32 GB
- CPU: 32 cores

Reproducibility: All analyses use `random.state=42` for reproducibility. Complete package versions available in `requirements.txt`.

B Code Availability

Complete analysis code is available in `final_project_solution.ipynb`, see `view report` on `final project` with:

- Detailed comments explaining each step
- Modular functions for reusability
- Error handling and validation checks
- High-resolution figure outputs (300 DPI)
- <https://frankfurtmacmoses.github.io/bio-informatics/>

Report Prepared By: Olawole Frankfurt Ogunfunminiyi, Niraj Kc

Contact: frankfurtmacmoses@gmail.com

Course Instructor: Dr. Seungchan Kim

Submission Date: November 2025