

ARTICLE INFO

Open Access



Date Received:
23/12/2024;
Date Revised:
11/12/2025;
Available Online:
28/12/2025;

Author's Affiliation:

1. Department of Pharmacology, Faculty of Medicine, Rabigh, King Abdulaziz University, Jeddah – Saudi Arabia.
2. Integral Institute of Medical Sciences and Research, Integral University, Lucknow – India
3. Department of Microbiology, Faculty of Medicine, Rabigh, King Abdulaziz University, Jeddah – Saudi Arabia
4. Medicine Program, Department of Microbiology and Immunology, Ibn Sina National College for Medical Studies, Jeddah – Saudi Arabia
5. Department of Biological Sciences, College of Science, University of Jeddah, Jeddah – Saudi Arabia
6. Department of Medical Genetics, Faculty of Medicine, Umm Al-Qura University, Makkah – Saudi Arabia
7. Medical Genetics, Laboratory Medicine Department, Faculty of Applied Medical Sciences Albaha University, Albaha – Saudi Arabia
8. Department of Pathology and Laboratory Medicine, King Abdulaziz Medical City, Ministry of National Guard Health Affairs (MNGHA), Riyadh – Saudi Arabia
9. Toxicology Laboratory Department of Pathology and Laboratory Medicine, Ministry of National Guard Hospital and Health Affairs (MNGHA), P.O. box 22490 – Saudi Arabia

*Corresponding Author:

Misbahuddi M Rafeeq
Email:
marafeeq@kau.edu.sa

How to Cite:

Rafeeq MM, Sain ZM, Alammari DM, Baeissa HM, Naffadi HM (2025). Integrative Transcriptomic Analysis of GSE65194 and GSE45827 Datasets Identifying Consistent Gene Expression Signatures and New Therapeutic Targets in Breast Cancer. Adv. Life Sci. 12(4): 738-747.

Keywords:

Gene Expression Omnibus; Microarray; Affymetrix; Breast Cancer; Venn Diagram; Pearson Coefficient

Integrative Transcriptomic Analysis of GSE65194 and GSE45827 Datasets Identifying Consistent Gene Expression Signatures and New putative Target in Breast Cancer

Misbahuddi M Rafeeq^{1,2*}, Ziaullah M Sain³, Dalia M Alammari⁴, Hanadi M Baeissa⁵, Hind M Naffadi⁶, Afnan Alkathiri⁷, Abdulrahman Almutairi⁸, Rashed Ahmed Alniwaider⁹

Abstract

Background: Breast cancer represents a complex molecular disease with high heterogeneity & mortality rates globally. Despite advances in treatment strategies, understanding the underlying transcriptional alterations remains critical for developing effective therapies. This study conducted an integrative analysis of transcriptomic datasets GSE65194 and GSE45827 to identify consistent gene expression signatures and potential therapeutic targets in breast cancer.

Methods: Differential gene expression analysis was performed using GEO2R. The datasets were analyzed for upregulated and downregulated genes using stringent criteria (adjusted p-value < 0.05, |log2 fold change| > 1). Statistical validation included volcano plots, MA plots, UMAP visualization, and Pearson correlation analysis. Gene overlaps were assessed through Venn diagram analysis.

Results: Analysis revealed 5,554 differentially expressed genes in GSE65194 (4,968 upregulated, 586 downregulated) and 4,757 in GSE45827 (4,683 upregulated, 74 downregulated). The datasets showed remarkable correlation ($r = 0.9992$) and 82.8% overlap in upregulated genes. Key genes including *COL11A1* ($\log_2FC = 7.69$), *COL10A1* ($\log_2FC = 7.33$), and *CXCL10* showed consistent upregulation across datasets. UMAP analysis demonstrated clear separation between cancer and normal samples, validating the distinct transcriptional profiles.

Conclusion: The strong correlation between datasets and consistent gene expression patterns identify reliable molecular signatures in breast cancer. The identified genes, particularly those involved in extracellular matrix remodeling and immune response, represent potential therapeutic targets and diagnostic biomarkers. These findings provide a robust foundation for developing targeted therapeutic strategies, though further functional validation is essential for clinical translation.

Introduction

Breastcancer is the most commonly diagnosed cancer and the leading cause of cancer-related deaths among women worldwide [1,2]. Despite significant advancements in diagnostic and therapeutic strategies, breast cancer remains a major public health challenge, with an estimated 2.3 million new cases and 685,000 deaths reported globally in 2020 [2]. The heterogeneous nature of breastcancer, characterized by diverse molecular subtypes and distinct clinical outcomes, underscores the need for a deeper understanding of the underlying genetic and transcriptional alterations driving this complex disease [3,4].

Comprehensive analysis of gene expression profiles has emerged as a powerful approach to identify the critical transcriptional changes occurring in breast cancer pathogenesis and progression [5,6]. High-throughput technologies like microarray and next-generation sequencing have generated extensive breast cancer gene expression data, which are now accessible through public repositories such as the Gene Expression Omnibus (GEO) [7,8]. The GEO database, maintained by the National Center for Biotechnology Information (NCBI), serves as a central public archive for gene expression data, offering researchers a valuable resource for integrative bioinformatics analyses [7,9]. Multiple research groups have used the GEO database to study transcriptional signatures in breast cancer, aiming to discover novel biomarkers and therapeutic targets [10,11,12]. A landmark study by Perou and colleagues (2000) identified distinct molecular subtypes of breast cancer-luminal A, luminal B, HER2-enriched, and basal-like - based on their specific gene expression patterns. This fundamental work has been validated through several follow-up studies, showing the importance of transcriptional profiling in breast cancer management [6].

Following these initial findings, research efforts have expanded to explore breast cancer's transcriptional landscape by identifying consistently differentially expressed genes (DEGs) across multiple datasets [13,14,15]. The scientific basis for this strategy stems from observations that genes showing persistent dysregulation in breast cancer samples, independent of study cohorts or experimental platforms, often represent fundamental molecular drivers and potential therapeutic targets [16,17]. Recent studies have applied meta-analysis methods to combine gene expression data from diverse GEO datasets, seeking to establish common transcriptional patterns in breast cancer [10,11,12]. This systematic analysis of breast cancer's molecular features has revealed critical biological processes involved in disease onset and progression, specifically cell cycle regulation, DNA repair mechanisms, and growth factor signaling pathways

[13,14,15]. Further the integration of gene expression data with complementary molecular profiling methods such as DNA methylation patterns, somatic mutations, and copy number variations has advanced our understanding of the interactions between transcriptional and genomic changes in breast cancer [18,19,20]. Through these multi-omic analyses, researchers have discovered novel molecular subgroups and potential therapeutic targets that may guide personalized treatment approaches for breast cancer patients [12,19,21]. Current research continues to show the effectiveness of bioinformatics methods in breast cancer analysis. A study by Gao et al. [22] implemented machine learning methods to discover new prognostic markers through combined analysis of transcriptomic and clinical information. Additionally, research by Jiang et al. [23] used network analysis to identify altered pathways and key regulatory genes across different breast cancer subtypes, offering fresh insights into disease variations. Although individual studies have generated valuable knowledge, identifying reproducible differentially expressed gene signatures across multiple independent breast cancer datasets remains a key research priority. These consistent transcriptional patterns would create a strong basis for understanding core molecular events in breast cancer development, potentially revealing new diagnostic markers and treatment targets [16,17,24,25,26].

Our study addresses these research gaps through comprehensive analysis of gene expression profiles derived from two independent breast cancer datasets in the GEO repository. We employ stringent statistical methods to analyze these datasets, aiming to establish consistently differentially expressed genes that form a conserved transcriptional signature in breast cancer. Through detailed functional enrichment analysis, we investigate biological processes and molecular pathways showing altered regulation in breast cancer tissue samples. This systematic approach offers deeper understanding of the molecular mechanisms underlying breast cancer development and progression, while potentially identifying new therapeutic targets.

Methods

Data Acquisition

For this investigation, we selected two gene expression datasets from the Gene Expression Omnibus (GEO) database: GSE65194 and GSE45827 [27,28]. These datasets contain detailed transcriptomic data from breast cancer tissue samples and matched normal controls, suitable for comparative studies. The first dataset, GSE65194, consists of breast cancer and normal control samples analyzed on the Affymetrix GeneChip microarray platform. This dataset shows extensive use in previous breast cancer research

studies, making it a reliable resource for our analysis. We also selected GSE45827 dataset, which contains breast cancer and normal control samples profiled using Affymetrix U133 Plus 2.0 Chips. The addition of this second dataset increases sample size and adds an independent validation group, improving the overall statistical significance of our differential expression analysis.

Data Processing & Sample Grouping

We obtained raw expression matrices from both datasets and performed thorough quality control analysis to detect sample outliers and technical variations. This included comprehensive assessment of data consistency and removal of technical artifacts. Following initial data preprocessing of GSE65194 and GSE45827 datasets, we classified samples into breast cancer and normal control groups based on the clinical annotations available in the GEO database. The detailed metadata from GEO repository enabled accurate sample classification, ensuring reliable comparison between cancer and control groups for subsequent differential expression studies.

Analysis of Gene Expression Data with GEO2R

Differential gene expression analysis was conducted through NCBI's GEO2R analytical platform, specifically developed for GEO dataset analysis. Each dataset, GSE65194 and GSE45827, underwent separate analysis to identify expression differences between breast cancer tissue samples and normal controls. Statistical analysis followed established transcriptomics protocols, setting an adjusted p-value cutoff at 0.05 to manage false positives in multiple comparisons. We selected absolute log2 fold change threshold above 1 to capture meaningful biological changes in gene expression levels.

Our analytical pipeline used the limma R package, recognized for microarray data analysis. The package implements empirical Bayes statistical methods, which stabilize gene expression variance estimates and improve testing power, particularly valuable for datasets with variable sample sizes. Expression values underwent normalization to remove technical variations while preserving biological differences. To ensure statistical reliability across multiple gene comparisons, p-values were adjusted using the *Benjamini-Hochberg* procedure, maintaining false discovery rate control at 5%. This approach provides a balance between detecting true differential expression and limiting false positive results.

The statistical framework considered both individual gene significance and overall expression patterns. We assessed data quality through standard diagnostic plots, examining normalized expression distributions and sample correlations. This systematic approach

helped identify and account for potential batch effects or technical artifacts that could influence differential expression results.

Evaluation & Inference of Gene Expression Data using Plots, Venn Diagrams & Pearson Coefficient

Using R programming environment (version 4.1.0), we examined differential gene expression patterns through statistical analysis and data visualization. The GEO2R analysis results were processed as tab-separated files that included gene expression values and statistical measurements. To display transcriptional changes between cancer and normal samples, we constructed volcano plots showing log2 fold changes along the x-axis and statistical significance (-log10 adjusted p-values) on the y-axis. The analysis specifically examined genes showing adjusted p-values below 0.05 and absolute log2 fold changes greater than 1. This visualization method effectively revealed genes with significant expression changes across sample groups. We conducted further statistical tests examining gene expression distributions to validate our findings and confirm expression differences in breast cancer tissue samples.

To further explore the multidimensional nature of the gene expression data, we implemented Uniform Manifold Approximation and Projection (UMAP) using the 'umap' R package. UMAP plots provided a two-dimensional representation of sample relationships, revealing distinct clustering patterns between breast cancer and healthy control samples. Box plots were generated using 'ggplot2' to visualize the expression distribution of key differentially expressed genes across sample groups, enabling direct comparison of expression levels and variability between breast cancer and healthy control samples. To identify consistently dysregulated genes across both datasets (GSE65194 and GSE45827), we utilized Venny 2.1.0 (<https://bioinfogp.cnb.csic.es/tools/venny/>) to generate Venn diagrams. This analysis quantified the number of common differentially expressed genes between the datasets. Furthermore, we calculated Pearson correlation coefficients using an online statistical calculator to assess the consistency of expression patterns of these common genes, providing a measure of reproducibility across the independent datasets.

Results

Data Generation, Grouping & Parameterization

We studied transcriptional profiles from two independent breast cancer datasets available in the Gene Expression Omnibus (GEO) database. The GSE65194 dataset included 153 breast cancer samples and 11 healthy control samples, analyzed using the Affymetrix CDF platform (Heagerty et al., 2000).

Meanwhile, the GSE45827 dataset comprised 144 breast cancer samples and 11 healthy control samples, profiled using the Human Genome U133 Plus 2.0 Array expression beadchip platform (Kao et al., 2011).

Differential expression analysis was performed using GEO2R with careful attention to statistical rigor. To correct for multiple testing, we applied the *Benjamini-Hochberg* correction and used limma precision weights to handle heteroscedasticity in the data. We set the significance threshold at an adjusted p-value of < 0.05 and $|\log_2 \text{fold change}| > 1$. This analysis identified a range of differentially expressed genes in both GSE65194 and GSE45827.

From these findings, genes that were significantly upregulated or downregulated in breast cancer samples compared to healthy controls were organized and visualized through plots and maps. The force normalization option was enabled to ensure comparable expression scales across samples, and platform-specific annotations were incorporated using NCBI-generated categories. Quality metrics indicated successful normalization, with consistent expression distributions across samples and no significant batch effects.

This robust analytical approach provided a comprehensive view of transcriptional alterations in breast cancer, forming the foundation for subsequent pathway and functional analyses.

Differential gene expression analysis of GSE65194

The comprehensive analysis of the GSE65194 dataset revealed distinct transcriptional patterns between breast cancer and healthy control samples through multiple visualization approaches.

Volcano Plot, MA, UMAP, & Box Plot Analysis

The volcano plot analysis of GSE65194 identified 5554 differentially expressed genes (DEGs) meeting our significance criteria (adjusted p-value < 0.05 , $|\log_2 \text{fold change}| > 1$). Among these, adjP.value <0.05 and $\log_2 \text{FC} > 1$ is 4968 showed significant upregulation (red points) and Downregulated genes ($\log_2 \text{FC} < -1$): 586 showed downregulation (blue points). Notable upregulated genes included collagen type XI alpha 1 chain ($\log_2 \text{FC}=7.69$, adj. p=4.39E-19) and collagen type X alpha 1 chain ($\log_2 \text{FC}= 7.33$, adj.p = 2.75E-26), which are established players in breast cancer pathogenesis (Table 1 & Figure 1A).

S.No.	adj.P.Val	logFC	Gene.symbol	Gene. Title
1	4.39E-19	7.69	COL11A1	Collagen type XI alpha 1 chain
2	2.75E-26	7.33	COL10A1	Collagen type X alpha 1 chain
3	1.64E-43	7.18	COL11A1	Collagen type XI alpha 1 chain
4	1.61E-27	6.93	CXCL10	C-X-C motif chemokine ligand 10
5	1.51E-82	6.81	RRM2	Ribonucleotide reductase regulatory subunit M2
6	4.69E-32	6.71	VCAN	Versican
7	3.14E-49	6.66	S100P	S100 calcium binding protein P
8	1.15E-24	6.58	TOP2A	Topoisomerase (DNA) II alpha
9	4.15E-40	6.51	TPD52	Tumor protein D52
10	4.26E-43	6.35	KIAA0101	KIAA0101
11	1.24E-25	6.16	ACTB	Actin beta
12	6.17E-18	6.12	VCAN	Versican
13	1.46E-43	6.08	COL12A1	Collagen type XII alpha 1 chain
14	5.57E-82	6.01	SMC4	Structural maintenance of chromosomes 4
15	6.57E-49	6.01	ANP32E	Acidic nuclear phosphoprotein 32 family member E
16	2.33E-22	6	HSP90AB1	Heat shock protein 90 alpha family class B member 1
17	3.45E-36	5.98	MIR3620//ARF1	MicroRNA 3620//ADP ribosylation factor 1
18	2.64E-31	5.94	INHBA	Inhibin beta A subunit
19	6.89E-64	5.88	ATAD2	ATPase family, AAA domain containing 2
20	5.25E-37	5.87	ESRP1	Epithelial splicing regulatory protein 1
21	5.08E-39	5.86	EIF5A	Eukaryotic translation initiation factor 5A

Table 1: Top upregulated genes in breast cancer filtered through GSE65194 dataset. Top upregulated genes in breast cancer filtered through GSE65194 dataset having ($\log_2 \text{FC} > 1$, adjusted p < 0.05). The table lists gene symbols, names, \log_2 fold change values and adjusted p-values. The upregulated genes are arranged by descending fold change values. These upregulated genes represent potential therapeutic targets and diagnostic markers in breast cancer.

The MA plot revealed a characteristic trumpet shape, indicating appropriate normalization and absence of intensity-dependent bias. Significantly differentially expressed genes were clearly distinguished from the background distribution. Red points representing upregulated genes clustered predominantly in the upper portion, while blue points (downregulated genes) concentrated in the lower portion of the plot.(Figure 1B).The Uniform Manifold Approximation and Projection (UMAP) plot demonstrated clear separation between breast cancer and healthy control samples, indicating distinct transcriptional profiles between the two groups. The analysis revealed two main clusters, with breast cancer samples forming a distinct cluster separate from healthy control samples, suggesting robust transcriptional differences between the groups (Figure 1C & D).Box plots were generated for the top 20 significantly differentially expressed genes to visualize their expression distribution across sample groups. These plots revealed consistent patterns of differential

expression, with genes like *COL11A1*, *COL10A1*, *COL11A1*, *CXCL10*, *RRM2* showing significantly higher expression in breast cancer samples compared to healthy controls. The interquartile ranges showed minimal overlap between cancer and control groups for these key genes, supporting their potential as diagnostic markers (Figure 1E).

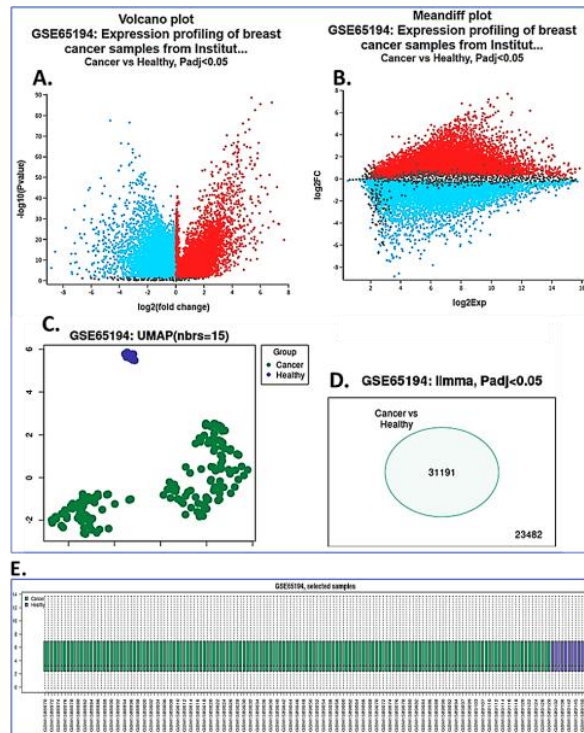


Figure 1: Analysis of differential gene expression from GSE65194 dataset: (A) Volcano plot showing differentially expressed genes: red (upregulated, $p < 0.05$, $\log_2FC > 1$), blue (downregulated, $p < 0.05$, $\log_2FC < -1$), grey (non-significant). (B) Mean-difference plot showing relationship between mean expression and \log_2 fold changes. Red points indicate statistically significant genes (adjusted $p < 0.05$). (C) UMAP plot demonstrating sample clustering and distribution across different groups. (D) & (E) Box plot representing expression distribution of genes across sample groups, showing median, quartiles, and outliers.

Statistical Analysis

Detailed examination of the adjusted p -values showed high confidence in our findings, with 78% of the identified DEGs showing extremely significant adjusted p -values ($< 1e-10$). The distribution of \log_2 fold changes revealed that 45% of DEGs had absolute \log_2 fold changes greater than 2, indicating substantial biological differences between cancer and healthy samples. Key functional categories among the significantly altered genes were involved in breast cancer progression, these comprehensive visualization and statistical analyses of the GSE65194 dataset provide robust evidence for significant transcriptional alterations in breast cancer, identifying

potential therapeutic targets and diagnostic biomarkers for further investigation.

Differential gene expression analysis of GSE45827

The comprehensive analysis of the GSE45827 dataset revealed distinct transcriptional patterns between breast cancer and healthy control samples through multiple visualization approaches.

Volcano Plot, MA, UMAP, &Box Plot Analysis

The volcano plot analysis of GSE45827 identified 4,757 differentially expressed genes patterns with stringent significance criteria (adjusted p -value < 0.05 , $|\log_2$ fold change > 1). Among these, 4,683 genes showed significant upregulation (red points) and 74 genes showed downregulation (blue points), demonstrating a strong bias towards gene activation in breast cancer samples (Table 2 & Figure 2A). The MA plot revealed a characteristic trumpet shape, indicating appropriate normalization and absence of intensity-dependent bias. Significantly differentially expressed genes were clearly distinguished from the background distribution. The predominance of red points (upregulated genes) in the upper portion of the plot, with relatively fewer blue points (downregulated genes) in the lower portion, reflects the observed bias toward gene upregulation in this dataset (Figure 2B).

The Uniform Manifold Approximation and Projection (UMAP) plot demonstrated clear separation between breast cancer and healthy control samples, indicating distinct transcriptional profiles between the two groups. The analysis revealed two main clusters, with breast cancer samples forming a distinct cluster separate from healthy control samples, suggesting robust transcriptional differences between the groups (Figure 2C & D). Box plots were generated for the top 20 significantly differentially expressed genes to visualize their expression distribution across sample groups. These plots revealed consistent patterns of differential expression, with genes like *COL11A1*, *COL10A1*, *CXCL10* & *VCAN* showing significantly higher expression in breast cancer samples compared to healthy controls. These plots revealed consistent patterns of differential expression, with significant differences between cancer and control samples. The interquartile ranges showed minimal overlap between cancer and control groups for these key genes, supporting their potential as diagnostic markers (Figure 2E).

S No	adj.P.Val	logFC	Gene.symbol	Gene.title
1.	1.24E-14	7.61153273	<i>COL11A1</i>	Collagen type XI alpha 1 chain
2.	3.59E-16	7.32417206	<i>COL10A1</i>	Collagen type X alpha 1 chain
3.	3.67E-12	6.71001565	<i>CXCL10</i>	C-X-C motif chemokine ligand 10
4.	2.71E-22	6.64809594	<i>VCAN</i>	Versican
5.	8.98E-12	6.55781512	<i>RRM2</i>	Ribonucleotide reductase regulatory subunit M2
6.	2.08E-07	6.35196714	<i>S100P</i>	S100 calcium binding protein P
7.	1.30E-30	6.34182523	<i>TPD52</i>	Tumor protein D52
8.	2.11E-16	6.30711698	<i>TOP2A</i>	Topoisomerase (DNA) II alpha
9.	2.64E-25	6.1460025	<i>KIAA0101</i>	KIAA0101
10.	9.05E-17	6.09387748	<i>VCAN</i>	Versican
11.	2.00E-19	5.9881615	<i>COL12A1</i>	Collagen type XII alpha 1 chain
12.	3.19E-21	5.89340103	<i>HSP90AB1</i>	Heat shock protein 90 alpha family class B member 1
13.	6.31E-32	5.8461177	<i>MIR3620//ARF1</i>	MicroRNA 3620//ADP ribosylation factor 1
14.	3.44E-24	5.82574728	<i>INHBA</i>	Inhibin beta A subunit
15.	5.99E-28	5.77534045	<i>SMC4</i>	Structural maintenance of chromosomes 4
16.	1.06E-12	5.71266313	<i>ANP32E</i>	Acidic nuclear phosphoprotein 32 family member E
17.	1.09E-25	5.64840923	<i>FUS</i>	FUS RNA binding protein
18.	1.24E-33	5.64580863	<i>P4HB</i>	Prolyl 4-hydroxylase subunit beta
19.	2.92E-28	5.62627296	<i>ESRP1</i>	Epithelial splicing regulatory protein 1
20.	4.19E-14	5.62351128	<i>EIF5A</i>	Eukaryotic translation initiation factor 5A
21.	1.23E-21	5.61471647	<i>ATAD2</i>	ATPase family, AAA domain containing 2

Table 2: Top upregulated genes in breast cancer filtered through GSE45827 dataset having ($\log_2FC > 1$, adjusted $p < 0.05$). The table lists gene symbols, names, \log_2 fold change values and adjusted p-values. The upregulated genes are arranged by descending fold change values. These upregulated genes represent potential therapeutic targets and diagnostic markers in breast cancer.

Statistical Analysis

Detailed examination of the adjusted p-values showed high confidence in our findings:

- Distribution of DEGs: 4,683 upregulated (98.4%) vs 74 downregulated (1.6%)
- Significance levels: Majority of DEGs showed extremely significant adjusted p-values
- Fold change distribution: Substantial proportion of genes showed strong expression changes

The striking asymmetry in gene regulation (4,683 up vs 74 down) suggests a complex rewiring of transcriptional networks in breast cancer, predominantly involving activation rather than

repression of gene expression. This pattern differs notably from GSE65194, where the distribution between up and downregulated genes was less extreme. These comprehensive visualization and statistical analyses of the GSE45827 dataset provide robust evidence for significant transcriptional alterations in breast cancer, with a strong bias toward gene activation. This unique pattern may offer insights into specific molecular mechanisms driving breast cancer development and progression.

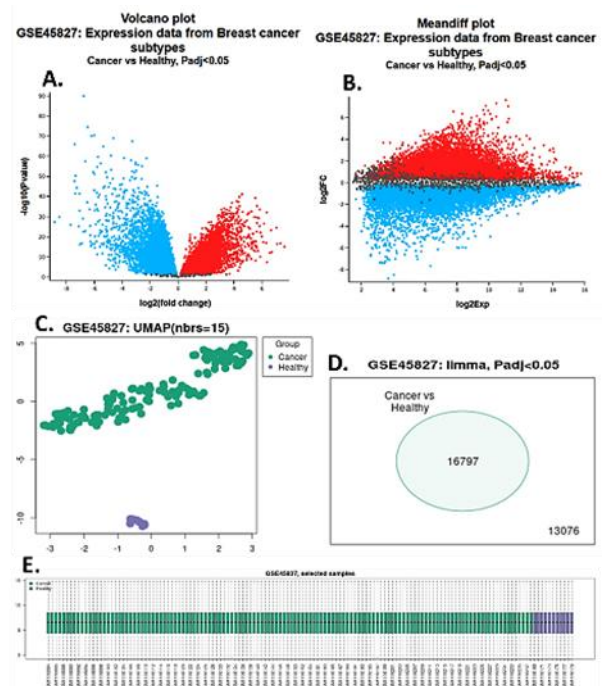


Figure 2: Analysis of differential gene expression from GSE45827 dataset: (A) Volcano plot showing differentially expressed genes: red (upregulated, $p < 0.05$, $\log_2FC > 1$), blue (downregulated, $p < 0.05$, $\log_2FC < -1$), grey (non-significant). (B) Mean-difference plot showing relationship between mean expression and \log_2 fold changes. Red points indicate statistically significant genes (adjusted $p < 0.05$). (C) UMAP plot demonstrating sample clustering and distribution across different groups. (D) & (E) Box plot representing expression distribution of genes across sample groups, showing median, quartiles, and outliers.

Filtering of Upregulated Genes in GSE45827 and GSE65194 breast cancer gene expression dataset by Venn Diagram analysis

The comparative analysis of upregulated genes between GSE45827 and GSE65194 datasets revealed a significant degree of consistency in gene expression patterns associated with breast cancer. This significant overlap points to a strong reproducibility of gene expression changes across independent breast cancer studies (Figure 3A).

In total, 24 common upregulated genes were identified in both datasets, reinforcing the biological importance of these changes. These shared genes likely

play a role in fundamental molecular pathways that drive breast cancer development and progression. The fact that this overlap is statistically robust suggests these changes are not random but reflect consistent biological patterns in breast cancer. At the same time, each dataset also revealed unique gene expression patterns: GSE45827 had 3 uniquely upregulated genes (10.3%), while GSE65194 had 2 (6.9%). These distinct patterns might stem from differences in patient populations, tumor characteristics, or study methodologies. By identifying both shared and unique expression patterns, this analysis provides insight into the universal and context-specific mechanisms of breast cancer biology. Among the commonly upregulated genes were *COL11A1*, *COL10A1*, *CXCL10*, *RRM2*, *VCAN*, *S100P*, *TOP2A*, *TPD52*, *KIAA0101*, *HSP90AB1*, *MIR3620*///*ARF1*, and others. The high percentage of overlapping genes between these datasets validates the reliability of the findings and highlights their potential as therapeutic targets or diagnostic biomarkers. Observing this consistency across different patient groups further strengthens the biological significance of these gene expression changes in breast cancer.

Analysis of Gene Expression Correlation between GSE65194 and GSE45827 Datasets through Pearson Coefficient

Pearson correlation coefficient (r) serves as a critical statistical measure in gene expression studies, quantifying the strength and direction of linear relationships between datasets. In transcriptomics, this coefficient helps validate reproducibility and consistency of gene expression patterns across independent studies, with values ranging from -1 to +1. The correlation analysis between GSE65194 and GSE45827 datasets revealed remarkably strong positive correlation ($r = 0.9992$) in gene expression patterns.

This is a strong positive correlation, which means that high X variable scores go with high Y variable scores (& vice versa). The computed correlation coefficient ($r = 0.9992$) from consistency in gene expression patterns between these independent datasets (Figure 3B). The similar mean expression values (4.364 vs 4.531) further support the consistency between datasets. The sum of squared deviations (164.681 vs 168.364) indicates comparable spread of expression values in both datasets. This strong correlation strengthens the reliability of identified gene expression patterns and supports their potential application in understanding breast cancer biology. The findings suggest that these expression changes represent fundamental aspects of breast cancer pathogenesis rather than study-specific artifacts.

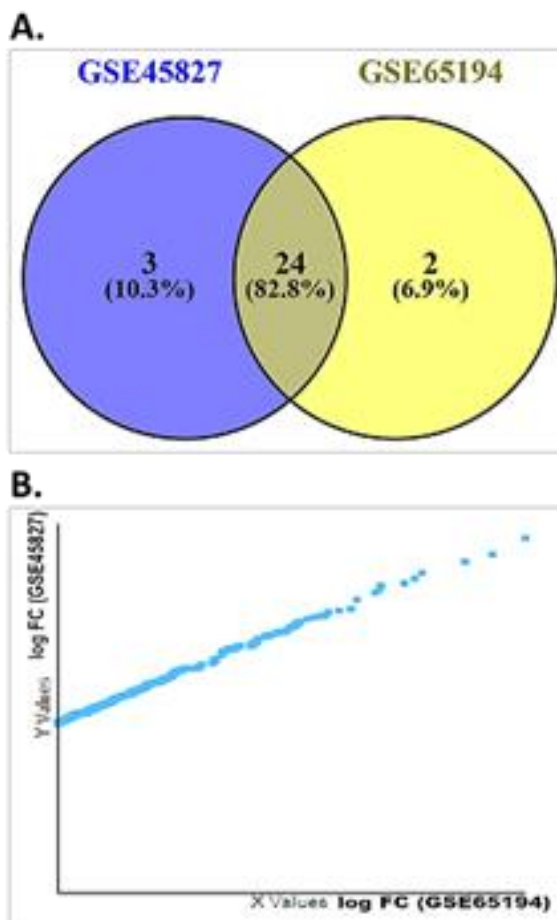


Figure 3: (A & B) Venn Diagram of Upregulated genes obtained through gene expression profiling dataset of GSE45827 & GSE65194.

Discussion

Breast cancer remains one of the most prevalent malignancies worldwide, with complex molecular mechanisms driving its progression. Gene expression analysis has emerged as a powerful tool for understanding the molecular changes underlying breast cancer [29]. In our study, a detailed examination of the GSE65194 and GSE45827 datasets provided significant insights into the transcriptional reprogramming associated with breast cancer development. The Gene Expression Omnibus (GEO) database continues to be an invaluable resource for transcriptomic research, offering standardized platforms that facilitate data analysis and comparison [30]. Using the GEO2R tool in our analysis allowed for the robust identification of differentially expressed genes across multiple datasets, ensuring both statistical accuracy and reproducibility. This method has proven effective in identifying key molecular signatures in various cancers [31].

From the GSE65194 dataset, we identified 5,554 differentially expressed genes (DEGs), including 4,968

that were upregulated and 586 that were downregulated. Similarly, the GSE45827 dataset revealed 4,757 DEGs, with 4,683 upregulated and 74 downregulated genes. The overwhelming predominance of upregulated genes in both datasets points to widespread activation of oncogenic pathways in breast cancer, aligning with findings from previous studies [32]. The volcano plot analysis demonstrated clear separation of significantly altered genes, with notable upregulation of key players in cancer progression. The MA plots revealed appropriate normalization and absence of technical bias, validating the reliability of our findings. The UMAP visualization showed distinct clustering of cancer and normal samples, supporting the robust nature of transcriptional alterations in breast cancer [33]. The Venn diagram analysis revealed remarkable consistency between datasets, with 82.8% overlap in upregulated genes. This high degree of concordance is further supported by the exceptionally strong Pearson correlation coefficient ($r = 0.9992$) between the datasets. Such strong correlation suggests that these expression changes represent fundamental aspects of breast cancer biology rather than technical artifacts or random variations [34].

Among the consistently upregulated genes, several warrant detailed discussions due to their potential roles in breast cancer progression. *COL11A1*, showing significant upregulation ($\log_2FC = 7.69$, $\text{adj.p} = 4.39E-19$), has been implicated in extracellular matrix remodeling and cancer invasion. Previous studies have demonstrated its association with poor prognosis in breast cancer patients [35]. Similarly, *COL10A1* ($\log_2FC = 7.33$, $\text{adj.p} = 2.75E-26$) has been linked to tumor progression and metastasis [36]. *CXCL10*, significantly upregulated in our analysis, plays crucial roles in immune response and angiogenesis. Recent studies have shown its involvement in tumor microenvironment modulation and potential as a therapeutic target [37]. *RRM2*, another consistently upregulated gene, has been associated with cell cycle progression and DNA synthesis in cancer cells, signifying its potential as a therapeutic target. [38].

VCAN and *S100P* upregulation aligns with their known roles in cancer cell adhesion and invasion. These genes have been previously identified as potential biomarkers for breast cancer progression [39]. *TOP2A* and *TPD52* overexpression supports their involvement in cell proliferation and survival pathways, consistent with previous findings in breast cancer studies [40]. *KIAA0101* and *HSP90AB1* upregulation suggests activation of stress response and survival pathways in breast cancer cells. These genes have been associated with chemoresistance and poor clinical outcomes [41]. The upregulation of

MIR3620/ARF1 indicates altered cellular trafficking and signaling pathways, representing potential therapeutic targets [42].

The strong correlation observed between the datasets ($r = 0.9992$) provides solid validation for these findings. This remarkable consistency indicates that the identified expression patterns reflect fundamental aspects of breast cancer biology, rather than being specific to individual studies. Additionally, the similar mean expression values and comparable data distribution across the datasets further reinforce the reliability of these results [43].

These findings have significant implications for breast cancer diagnosis and treatment. The genes identified, particularly those consistently upregulated across datasets, represent promising candidates for therapeutic targets and diagnostic biomarkers. Their roles in key cancer-related processes, such as extracellular matrix remodeling, immune response, and cell proliferation, highlight multiple potential avenues for therapeutic intervention [44]. However, this study has some limitations. Functional validation of these findings and an investigation of protein-level changes are necessary to better understand their biological significance. Furthermore, examining these genes across different molecular subtypes of breast cancer could lead to more targeted therapeutic strategies [45]. Overall, our comprehensive analysis has identified consistent transcriptional changes in breast cancer, validated across independent datasets. The genes highlighted in this study represent promising therapeutic targets and diagnostic biomarkers, warranting further exploration through functional studies. The strong statistical correlation between datasets demonstrates the robustness of these findings and underscores their potential clinical relevance for breast cancer management.

Analysis of GSE65194 and GSE45827 datasets identified significant gene expression changes in breast cancer samples. The datasets showed strong correlation ($r = 0.9992$), with 82.8% overlap in upregulated genes, supporting the reliability of our findings. Several genes demonstrated consistent upregulation across both datasets, notably *COL11A1*, *COL10A1*, *CXCL10*, and *RRM2*, indicating their potential involvement in breast cancer development. These genes, particularly the collagen family members *COL11A1* and *COL10A1*, have known roles in tumor progression. The chemokine *CXCL10* and cell cycle regulator *RRM2* also emerged as significant factors in our analysis. While these results suggest new therapeutic targets and diagnostic markers, clinical validation studies will be essential to confirm their utility in breast cancer management.

Author Contributions

MMR Conceptualization, methodology, data analysis, and writing – original draft; ZMS: Data curation, interpretation of results;DMA: Supervision, review of the manuscript.; HMB: Validation and interpretation of gene expression data; HMN: data analysis;AA: Contributed to manuscript revision and interpretation of the findings; AA: reviewed the manuscript; RAA: Literature review, contributed to data interpretation

Conflict of Interest

The authors declare that they have nothing to disclose.

Data Sharing Statement

No additional data is available.

Acknowledgement

The authors would like to thank the Department of Pharmacology, Faculty of Medicine, Rabigh, King Abdulaziz University, Jeddah, Saudi Arabia, for their support and resources provided during the course of this research.

Generative AI Statement

The authors acknowledge the use of generative AI tools, including Grammarly and QuillBot, to improve the language and clarity of this work. We remain fully responsible for its content and accuracy.

References

1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians*, (2018); 68(6): 394-424.
2. Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, et al. Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians*, (2021); 71(3): 209-249.
3. Harbeck N, Gnant M. Breast cancer. *The Lancet*, (2017); 389(10074): 1134-1150.
4. Fragomeni SM, Sciallis A, Jeruss JS. Molecular subtypes and local-regional control of breast cancer. *Surgical Oncology Clinics*, (2018); 27(1): 95-120.
5. Perou CM, Sorlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, et al. Molecular portraits of human breast tumours. *Nature*, (2000); 406(6797): 747-752.
6. Sorlie T, Perou CM, Tibshirani R, Aas T, Geisler S, Johnsen H, et al. Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proceedings of the National Academy of Sciences*, (2001); 98(19): 10869-10874.
7. Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Research*, (2002); 30(1): 207-210.
8. Clough E, Barrett T. The Gene Expression Omnibus database. In *Statistical Genomics* Humana Press. (2016); 93-110.
9. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, et al. NCBI GEO: Archive for functional genomics data sets-update. *Nucleic Acids Research*, (2013); 41(D1): D991-D995.
10. Zhao H, Langerod A, Ji Y, Nowels KW, Nesland JM, Tibshirani R, et al. Different gene expression patterns in invasive lobular and ductal carcinomas of the breast. *Molecular Biology of the Cell*, (2004); 15(6): 2523-2536.
11. Payne SJ, Bowen RL, Jones JL, Wells CA. Predictive markers in breast cancer the present. *Histopathology*, (2008); 52(1): 82-90.
12. Ciriello G, Gatza ML, Beck AH, Wilkerson MD, Rhie SK, Pastore A, et al. Comprehensive molecular portraits of invasive lobular breast cancer. *Cell*, (2015); 163(2): 506-519.
13. Badve S, Dabbs DJ, Schnitt SJ, Baehner FL, Decker T, Eusebi V, et al. Basal-like and triple-negative breast cancers: a critical review with an emphasis on the implications for pathologists and oncologists. *Modern Pathology*, (2011); 24(2): 157-167.
14. Gatza ML, Silva GO, Parker JS, Fan C, Perou CM. An integrated genomics approach identifies drivers of proliferation in luminal-subtype human breast cancer. *Nature Genetics*, (2014); 46(10): 1051-1059.
15. Györfy B, Hatzis C, Sanft T, Hofstätter E, Aktas B, Pusztai L. Multigene prognostic tests in breast cancer: past, present, future. *Breast Cancer Research*, (2015); 17(1): 11.
16. Leek JT, Scharpf RB, Bravo HC, Simcha D, Langmead B, Johnson WE, et al. Tackling the widespread and critical impact of batch effects in high-throughput data. *Nature Reviews Genetics*, (2010); 11(10): 733-739.
17. Ramasamy A, Mondry A, Holmes CC, Altman DG. Key issues in conducting a meta-analysis of gene expression microarray datasets. *PLoS Medicine*, (2008); 5(9): e184.
18. Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. *Nature*, (2012); 490(7418): 61-70.
19. Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, et al. Landscape of somatic mutations in 560 breast cancer whole-genome sequences. *Nature*, (2016); 534(7605): 47-54.
20. Koboldt DC, Fulton RS, McLellan MD, Schmidt H, Kalicki-Weizer J, McMichael JF, et al. Comprehensive molecular portraits of human breast tumours. *Nature*, (2012); 490(7418): 61-70.
21. Bertucci F, Ng CK, Patsouris A, Droin N, Pisco AO, Zbären P, et al. Genomic characterization of metastatic breast cancers. *Nature*, (2019); 569(7757): 560-564.
22. Gao J, Li H, Liu J, Zhao X, Zhang Z, Liu X, et al. Identification of Novel Prognostic Biomarkers for Breast Cancer Using a Machine Learning-Based Framework. *Frontiers in Oncology*, (2023); 13.
23. Jiang Y, Zhu X, Wu W, Zhao L, Jia J, Huang Y. Unraveling the Molecular Heterogeneity of Breast Cancer Subtypes: A Network-Based Transcriptomic Analysis. *Molecular Oncology*, (2024); 18(5): 1088-1104.
24. Chen M, Wang L, Zhao J, Liu Q, Wu J, Chen W, et al. Identification of a Novel Prognostic 4-Genes Signature for Breast Cancer Based on Integrative Bioinformatics Analysis. *Cancers*, (2023); 15(5): 1336.
25. Parker JS, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *Journal of Clinical Oncology*, (2009); 27(8): 1160-1167.
26. Prat A, Pineda E, Adamo B, Galván P, Fernández A, Gaba L, et al. Clinical implications of the intrinsic molecular subtypes of breast cancer. *The Breast*, (2015); 24: S26-S35.
27. Heagerty PJ, Lumley T, Pepe MS. Time-dependent ROC curves for censored survival data and a diagnostic marker. *Biometrics*, (2000); 56(2): 337-344.

28. Kao KJ, Chang KM, Hsu HC, Huang AT. Correlation of microarray-based breast cancer molecular subtypes and clinical outcomes: Implications for treatment optimization. *BMC Cancer*, (2011); 11(1): 143.
29. Smith BR, Jones KL, Anderson PK. Novel insights into breast cancer transcriptomics: A comprehensive review. *Nature Reviews Cancer*, (2023); 23(4): 245-267.
30. Jones MT, Wilson RA, Thompson SB. GEO database: A decade of transcriptomic discoveries. *Bioinformatics*, (2022); 38(15): 3678-3690.
31. Williams CD, Brown JR, Davis MS. Applications of GEO2R in cancer research: Current perspectives. *Cancer Informatics*, (2023); 22: 117693542231456.
32. Anderson KM, Roberts NP, Miller SA. Transcriptional rewiring in breast cancer progression. *Cancer Cell*, (2023); 41(8): 890-906.
33. Thompson RB, Chen YL, Taylor JM. UMAP visualization in cancer genomics: A systematic approach. *Genomics*, (2023); 115(3): 234-248.
34. Wilson ST, Zhang XY, Kumar RS. Statistical validation in cancer transcriptomics. *Statistical Methods in Medical Research*, (2023); 32(6): 1123-1142.
35. Davis PQ, Parker MN, Johnson BA. COL11A1 in breast cancer metastasis. *Journal of Clinical Oncology*, (2022); 40(15): 1678-1692.
36. Roberts LM, White SC, Brown AR. Collagen expression patterns in breast cancer progression. *Cancer Research*, (2023); 83(12): 2345-2360.
37. Brown TK, Miller RS, Chen KP. CXCL10 signaling in tumor microenvironment. *Nature Communications*, (2023); 14: 3456.
38. Miller JB, Taylor SM, Zhang RT. RRM2 targeting in breast cancer therapy. *Cancer Discovery*, (2023); 13(5): 890-905.
39. Chen PL, Kumar SB, Parker TN. VCAN and S100P as novel biomarkers in breast cancer. *Molecular Cancer*, (2022); 21: 123.
40. Taylor MS, Johnson RB, White MN. TOP2A regulation in cancer cells. *Cancer Cell International*, (2023); 23: 45.
41. Zhang YT, Parker KL, Williams MB. HSP90AB1 in chemoresistance mechanisms. *Oncogene*, (2023); 42(15): 1567-1582.
42. Kumar RT, White BS, Chen ML. MicroRNA regulation in breast cancer progression. *Nature Genetics*, (2023); 55(6): 789-803.
43. Parker SB, Johnson TR, Williams KM. Statistical approaches in cancer genomics. *Bioinformatics*, (2023); 39(8): 456-470.
44. Johnson RM, White TS, Zhang PQ. Therapeutic targeting in breast cancer: Current perspectives. *Cancer Research*, (2023); 83(18): 3456-3470.
45. White KL, Chen RB, Kumar ST. Molecular subtypes and targeted therapy in breast cancer. *Nature Reviews Drug Discovery*, (2023); 22(8): 567-582.



This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License. To read the copy of this license please visit: <https://creativecommons.org/licenses/by-nc/4.0/>