



## OPEN ACCESS

## EDITED BY

Kishan Nyati,  
Osaka University, Japan

## REVIEWED BY

Ravi Ramalingam,  
The University of Texas Health Science  
Center at San Antonio, United States  
Sarvagya Sharma,  
Saveetha Dental College And Hospitals, India

## \*CORRESPONDENCE

Junying Yang

✉ yangjuny@mail.sysu.edu.cn

Shan Chen

✉ chensh53@mail.sysu.edu.cn

†These authors share first authorship

RECEIVED 23 March 2025

ACCEPTED 19 May 2025

PUBLISHED 05 June 2025

## CITATION

Sun S, Ren J, Zeng X, Chen Y, Zhou Q, Yang J and Chen S (2025) Integrated machine learning and single-cell RNA sequencing reveal COL4A2 and CXCL6 as oxidative stress-associated biomarkers in periodontitis. *Front. Immunol.* 16:1598642. doi: 10.3389/fimmu.2025.1598642

## COPYRIGHT

© 2025 Sun, Ren, Zeng, Chen, Zhou, Yang and Chen. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

# Integrated machine learning and single-cell RNA sequencing reveal COL4A2 and CXCL6 as oxidative stress-associated biomarkers in periodontitis

Siyu Sun<sup>1,2†</sup>, Jing Ren<sup>1,2†</sup>, Xiujuan Zeng<sup>1</sup>, Yanbin Chen<sup>1</sup>, Qianbing Zhou<sup>3</sup>, Junying Yang<sup>1\*</sup> and Shan Chen<sup>1\*</sup>

<sup>1</sup>Department of Stomatology, The First Affiliated Hospital of Sun Yat-sen University, Guangzhou, China, <sup>2</sup>National-Guangdong Joint Engineering Laboratory for Diagnosis and Treatment of Vascular Diseases, Guangzhou, Guangdong, China, <sup>3</sup>Department of Stomatology, The First Affiliated Hospital of Guangdong Pharmaceutical University, Guangzhou, China

**Background:** Periodontitis, recognized as the second most prevalent oral disease globally, is strongly linked to systemic disorders like diabetes and cardiovascular diseases, highlighting the critical need for effective prevention and treatment strategies. Oxidative stress plays an important role in periodontitis pathogenesis and progression, yet their specific association remains unclear. This study aims to explore the association between oxidative stress and periodontitis pathogenesis while identifying potential diagnostic biomarkers and therapeutic targets for this condition.

**Methods:** Transcriptomic data from gingival tissues of periodontitis patients and controls were obtained from the Gene Expression Omnibus (GEO) database. Key genes linked to oxidative stress in periodontitis were identified through a comprehensive analytical approach, including differential expression analysis, weighted gene co-expression network analysis (WGCNA), gene set enrichment analysis (GSEA), and functional enrichment analyses (GO and KEGG). Machine learning algorithms were subsequently employed to refine the selection of key genes. The relationship between oxidative stress and the expression of these key genes was validated using external datasets and a periodontitis rat model. Additionally, single-cell RNA sequencing (scRNA-seq) data were interrogated to delineate the cellular subpopulations expressing the key genes, leveraging clustering and annotation approaches.

**Results:** Comprehensive bioinformatics analysis identified COL4A2, CYR61, and CXCL6 as key genes associated with oxidative stress in periodontitis. Among these genes, COL4A2 and CXCL6 showed elevated expression levels in the gingival tissues of periodontitis rats. Single-cell RNA-seq analysis further demonstrated that COL4A2 exhibited predominant expression within endothelial and stromal cell clusters, whereas CXCL6 was predominantly localized to epithelial cell clusters.

**Conclusions:** This study demonstrates a correlation between oxidative stress and the progression of periodontitis. COL4A2 and CXCL6 were identified as potential therapeutic targets for the treatment of periodontitis.

#### KEYWORDS

periodontitis, oxidative stress, single-cell RNA sequencing, COL4A2, CXCL6

## 1 Introduction

Periodontitis is a chronic inflammatory disease caused by bacterial infection. It is characterized by chronic gingival inflammation and progressive destruction of periodontal structures, including alveolar bone, periodontal ligament, and cementum. This condition is also one of the primary causes of tooth loss (1). In recent years, oxidative stress has been identified as a key mechanism in the development of periodontitis. Oxidative stress exacerbates periodontal tissue damage by promoting inflammatory responses, disrupting bone metabolism, and regulating gene expression (2). Oxidative stress (OS) occurs when there is an imbalance between the production of reactive oxygen species (ROS) and the body's antioxidant defenses (3). During periodontitis, oxidative stress causes direct damage to periodontal tissues. It induces oxidative damage to mitochondrial DNA and degrades the extracellular matrix (4, 5). This process exacerbates inflammatory responses and triggers apoptosis in periodontal ligament cells. Additionally, oxidative stress enhances the activity of matrix metalloproteinases (MMPs) (6). This indirectly regulates signaling pathways related to inflammation and apoptosis, further increasing periodontal tissue damage and alveolar bone resorption. Over the past decades, therapeutic strategies targeting oxidative stress have been explored for periodontitis treatment. For example, some approaches aim to reduce oxidative damage by scavenging reactive oxygen species or activating antioxidant signaling pathways (7, 8). However, their clinical efficacy remains limited, partly due to the absence of biomarkers capable of precisely indicating individual oxidative stress status, leading to a lack of targeted therapeutic strategies. The cellular mechanisms underlying the imbalance of oxidative stress in periodontitis remain incompletely understood, particularly as the cell type-specific responses to oxidative microenvironments have not been systematically characterized, necessitating further research to elucidate these processes.

Sequencing technology has become a widely used tool in the biomedical field. RNA sequencing (RNA-seq) datasets, which integrate transcriptomic profiling and microarray data analysis, represent a high-throughput technique for genome-wide gene expression quantification and analysis (9). It enables the simultaneous detection of expression levels for a large number of genes, offering significant advantages such as cost-effectiveness and well-established technological reliability. However, RNA-seq is limited in its ability to reveal intercellular heterogeneity (10). In contrast, single-cell RNA sequencing (scRNA-seq) allows for gene expression analysis at the

single-cell level. This method provides significant advantages, including high resolution and the ability to identify new cell types. As a result, scRNA-seq has emerged as a powerful tool for studying cellular heterogeneity and gene expression patterns in complex tissues (11).

In this study, RNA sequencing data, single-cell RNA sequencing (scRNA-seq) data, and associated clinical metadata were acquired from the Gene Expression Omnibus (GEO) database. Gene Set Enrichment Analysis (GSEA) is a bioinformatics method that evaluates the correlation between gene sets and clinical variables across different samples (12). Weighted Gene Co-expression Network Analysis (WGCNA) is a systems biology method used to construct gene co-expression networks, identify modules of highly correlated genes, and explore their relationships with phenotypic traits or external conditions (13). We used the WGCNA and GSEA to investigate the relationship between oxidative stress (OS) gene expression and periodontitis. However, traditional bioinformatics methods are prone to limitations from multicollinearity, gene redundancy, and linear assumptions when processing high-dimensional genomic data (14), which struggle to analyze the complex nonlinear relationships between genes and phenotypes (15, 16). By employing three machine learning algorithms, we systematically screened a set of oxidative stress (OS)-related genes associated with periodontitis progression. The combination strategy of multiple algorithms ensured robustness in biomarker screening and improved the accuracy and reliability of our analysis. Additionally, the key genes were validated through external dataset analysis and experimental validation in a rat periodontitis model, thereby confirming their reliability and biological significance. Finally, we validated the expression of key genes at the single-cell level. By integrating RNA-seq and scRNA-seq data, this study provides a novel approach to understanding cell type-specific and functional regulatory networks in periodontal tissues. This integration also helps clarify the role of oxidative stress in the mechanisms underlying periodontitis.

## 2 Materials and methods

### 2.1 Dataset acquisition and preprocessing

RNA-seq datasets (GSE10334, GSE16134) and scRNA-seq dataset (GSE171213) were obtained from the Gene Expression Omnibus (GEO) database (<https://www.ncbi.nlm.nih.gov/geo/>).

The training dataset, GSE10334, included transcriptome data from 183 gingival samples of periodontitis lesion sites and 64 samples from healthy sites. For validation, we used the GSE16134 dataset, which contains gingival tissue samples from 120 patients who underwent periodontal surgery. Additionally, 404 genes related to oxidative stress were retrieved from the Gene Set Enrichment Analysis (GSEA) database (<https://www.gsea-msigdb.org/gsea/index.jsp>). Data analysis was performed using R language (version 4.3.2) (<https://www.r-project.org/>).

## 2.2 Differentially expressed genes analysis

The “limma” package (version 3.58.1) (17) was utilized to analyze differentially expressed genes, aiming to compare gene expression profiles between periodontitis patients and controls. The “normalizeBetweenArrays” function was applied to correct potential technical errors and minimize batch effects across multiple samples. Differentially expressed genes (DEGs) were identified based on the following criteria:  $|\log_2(\text{fold-change})| \geq 1$  and adjusted P-value  $\leq 0.05$ . The “ggplot2” package (version 3.5.1) was utilized to generate heat and volcano maps. The pcomp function was used to perform PCA analyses of the samples, and the “factoextra” package (version 1.0.7) was employed to generate scatter plots.

## 2.3 GSEA and WGCNA

Gene Set Enrichment Analysis (GSEA) was primarily used to identify gene sets significantly enriched under specific biological conditions, such as disease states or drug treatments (18). The “GSEABase” package (version 1.64.0) and “GSVA” package (version 1.50.1) were employed in this study to assess the enrichment of oxidative stress-related gene sets in transcriptome data expression profiles through GSEA, and barcode enrichment maps were generated. The “WGCNA” package (version 1.72) was utilized to construct gene co-expression networks (13). Initially, hierarchical clustering analysis was applied to the transcriptome dataset to detect and eliminate outlier samples. Subsequently, the optimal soft threshold power was determined using the “pickSoftThreshold” function in the “WGCNA” package. Dynamic hybrid cutting was employed to identify co-expression modules, and a hierarchical clustering dendrogram was constructed to visualize the module structure (minModuleSize = 50, mergeCutHeight = 0.25, the colors representing different modules). Finally, the oxidative stress-related gene sets identified through GSEA were integrated with results obtained from WGCNA. Pearson correlation analysis was performed to evaluate the relationships between these gene sets and the modules, with correlation heatmaps generated to visualize the results. In the heatmaps, rows represented modules, columns represented traits, and the corresponding boxes displayed correlation coefficients and P-values. This approach identified gene co-expression modules closely associated with oxidative stress.

## 2.4 Identification of oxidative stress-related DEGs

The list of DEG was cross-referenced with the list of genes from the oxidative stress-related co-expression module identified through WGCNA. Genes that appeared in both lists were designated as oxidative stress-related DEGs. The “VennDiagram” package (version 1.12) (19) was used to create a Venn diagram, visually representing the intersection results.

## 2.5 Functional enrichment analysis

Functional enrichment analysis of oxidative stress-associated DEGs was conducted using the Gene Ontology (GO) and the Kyoto Encyclopedia of Genes and Genomes (KEGG). GO enrichment analysis included three categories: molecular function (MF), biological process (BP), and cellular component (CC). KEGG pathway analysis was also performed to identify relevant signaling pathways. The “clusterProfiler” package (version 4.10.1) and the “org.Hs.eg.db” package (version 3.18.0) (20) were used to perform these analyses. The “GOplot” package (version 1.0.2) and the “ggplot2” package (version 3.5.1) were utilized to visualize the results, generating bar charts and pie charts to illustrate the enrichment results.

## 2.6 Machine learning screens for key genes

We further analyzed the previously screened genes to identify the optimal key genes associated with OS. Three machine learning algorithms were employed for feature selection and key gene identification: Least Absolute Shrinkage and Selection Operator (LASSO) (21), Gradient Boosted Tree (GBM) (22), and Extreme Gradient Boosting (XGBoost) (23). Genes consistently identified by all three algorithms were considered optimal key genes. LASSO regression overcomes the limitations of gene redundancy in WGCNA modules by introducing a regularization term into the loss function (15), thereby enhancing predictive performance. The R package “glmnet” (version 4.1) (24) was used to identify the optimal tuning parameters through 10-fold cross-validation, ensuring robust parameter tuning and model evaluation. GBM can capture nonlinear interactions between genes through iterative residual optimization. The “gbm” package (version 2.1.9) was utilized to implement the GBM algorithm. XGBoost employs an iterative approach to construct decision trees (16), rectifying the errors from preceding iterations, and enhances accuracy via an early stopping mechanism to avert overfitting. The “xgboost” package (version 1.7.1) was used to perform the XGBoost algorithm. Finally, the UpSetR package (version 1.4.0) was used to generate UpSet plots, which visually represent the key DEGs identified by the intersection of the three algorithms.

## 2.7 Validation of key genes and diagnostic performance assessment

The differential expression of the key genes was validated using the external dataset GSE16134. Differential box line plots were generated to visualize the expression patterns using the “ggpubr” package (version 1.7.7). The “pROC” package (version 1.18.5) was used to perform receiver operating characteristic (ROC) curve analysis for evaluating the diagnostic performance of the key genes. The area under the curve (AUC) was calculated to quantify the diagnostic efficacy of the key genes.

## 2.8 Construction of a rat model of ligature-induced periodontitis

The animal experimental protocol was approved by the Ethics Committee of the First Affiliated Hospital of Guangdong Pharmaceutical University (No. G2R2024012), and all experimental animals were housed in the institution’s Animal Experiment Center. Six-week-old male Sprague-Dawley (SD) rats were randomly assigned to a periodontitis group (n=4, ligation-induced) and a control group (n=4, untreated). Following established methods from previous studies (8), periodontitis was induced by placing orthodontic steel ligatures (diameter: 0.2 mm) around the necks of the bilateral maxillary second molars in the periodontitis group. Prior to ligation, rats were anesthetized via intraperitoneal injection of 2% pentobarbital sodium (0.2 ml/100 g). Ligatures were checked daily to ensure retention.

After 28 days, the rats were fasted for 2 hours and euthanized by cervical dislocation. Alveolar bone and gingival tissue samples were then collected for further analysis. Micro-CT technology was employed to perform three-dimensional reconstruction image analysis, allowing for the observation of structural changes in periodontal tissues. Additionally, histological sectioning and pathological analysis were conducted using hematoxylin and eosin (H&E) staining to assess the degree of inflammation in the periodontal tissues.

## 2.9 Measurement of oxidative stress markers in periodontal tissues

After homogenizing and grinding the periodontal tissues, the protein concentration of the samples was quantified using a BCA assay kit (Biosharp, BL521A). Subsequently, malondialdehyde (MDA) levels were quantified using the MDA kit (Beyotime,

S0131S). The procedure involved adding the TBA working solution and incubating the mixture at 100°C for 45 minutes. After cooling, the supernatant was collected via centrifugation, and its absorbance was measured at 532 nm. The MDA concentration in the samples was then calculated based on the standard curve. A total SOD assay kit (Beyotime, S0101S) was used to measure the concentration of superoxide dismutase (SOD). After homogenizing the periodontal tissues and determining the protein concentration of the samples, the samples were mixed with the prepared WST-8 working solution and incubated at 37°C for 30 minutes. The absorbance was measured at 560 nm, and the SOD concentration was calculated accordingly.

## 2.10 RNA extraction and quantitative real-time PCR

Gingival tissue samples were homogenized in Trizol (Sigma-Aldrich, T9424) to extract total RNA. The RNA was reverse-transcribed into complementary DNA (cDNA) using the HiScript III RT SuperMix for qPCR kit (Vazyme, R323-01). Subsequently, quantitative real-time PCR (qPCR) was performed to quantify the expression of key genes using the ChamQ Blue Universal SYBR qPCR Master Mix (Vazyme, Q312-02). The relative mRNA expression levels of the target genes were calculated using the  $2^{-\Delta\Delta Ct}$  method, with  $\beta$ -actin serving as the internal reference gene. The primer sequences used in this study are provided in Table 1.

## 2.11 ScRNA-seq data analysis

The scRNA-seq data (GSE171213) (25) were processed using the “Seurat” software package (version 5.1.0). The filter conditions were as follows: nFeature\_RNA > 300, nFeature\_RNA < 10,000, percent.mt < 10, and nCount\_RNA > 600. Gene expression data were normalized and scaled using the “LogNormalize” method. Principal component analysis (PCA) was performed to identify principal components (PCs), and the batch correction was applied using the “harmony” software package (version 1.2.0) to mitigate batch effects. Cell clustering and sub-clustering analyses were performed using the FindClusters function in the Seurat package with a resolution parameter of 1 (resolution = 1). This analysis classified cells into 25 distinct clusters. Cell types were manually annotated using cellMarker, and the expression patterns of key genes were subsequently identified and visualized through Uniform Manifold Approximation and Projection (UMAP) and violin plots (VlnPlot).

TABLE 1 Primers used for qPCR.

Gene name	Forward primer (5'-3')	Reverse primer (5'-3')
$\beta$ -actin	AACACAGTGCTGTCTGGTG	GTAACAGTCCGCTAGAAAGC
Col4a2	GGGACCTGCCATTACTTCGCTAAC	GGATGGTGTGCTCTGGAAGTTCTG
CYR61	CGGTGCGAAGATGGCGAGATG	GGGATGCGGGCAGTTGTAGTTAC
CXCL6	GTCTTGACCCAGAAGCTCCGTTG	GGCTGATCTGACCAGTGCAAGTG

## 2.12 Statistical analysis

All statistical analyses were conducted using R software (version 4.3.2). All tests were two-sided, and a  $p$ -value  $< 0.05$  was considered statistically significant.

## 3 Results

### 3.1 Identification of DEGs in periodontitis

The raw RNA-seq data were normalized to minimize expression differences arising from technical variability and batch effects (Figure 1A). Dimensionality reduction was performed on the normalized data using principal component analysis (PCA). Dimensionality reduction was performed on the normalized data using PCA, and the results revealed a clear separation between the periodontitis and the controls in principal component space (Figure 1B), indicating significant differences in gene expression profiles between the two groups. To further investigate these differences, heatmaps were constructed to visualize gene expression patterns (Figure 1C). Additionally, the statistical significance of the

DEGs was assessed using volcano plots (Figure 1D). The analysis revealed a total of 139 DEGs between the periodontitis and control groups, with 111 genes up-regulated and 28 genes down-regulated.

### 3.2 Identifying co-expression modules of genes associated with oxidative stress

GSEA demonstrated the enrichment of oxidative stress-related genes within the list of genes associated with periodontitis. The peaks and significance of the running enrichment score curves indicated a statistically significant correlation ( $P < 0.05$ ) between the oxidative stress-related genes and the periodontitis transcriptome data (Figure 2A). To further analyze the transcriptome data, we employed WGCNA to identify gene co-expression modules linked to oxidative stress. A soft threshold (power) of 10 was selected based on the optimal scale-free fit and average connectivity (Figure 2B), and this approach identified 15 distinct gene co-expression modules (Figure 2C). A correlation heatmap was generated to visualize module-trait relationships, with color intensity indicating correlation strength—red for positive and blue for negative. Each row represents a gene module, and each column represents a

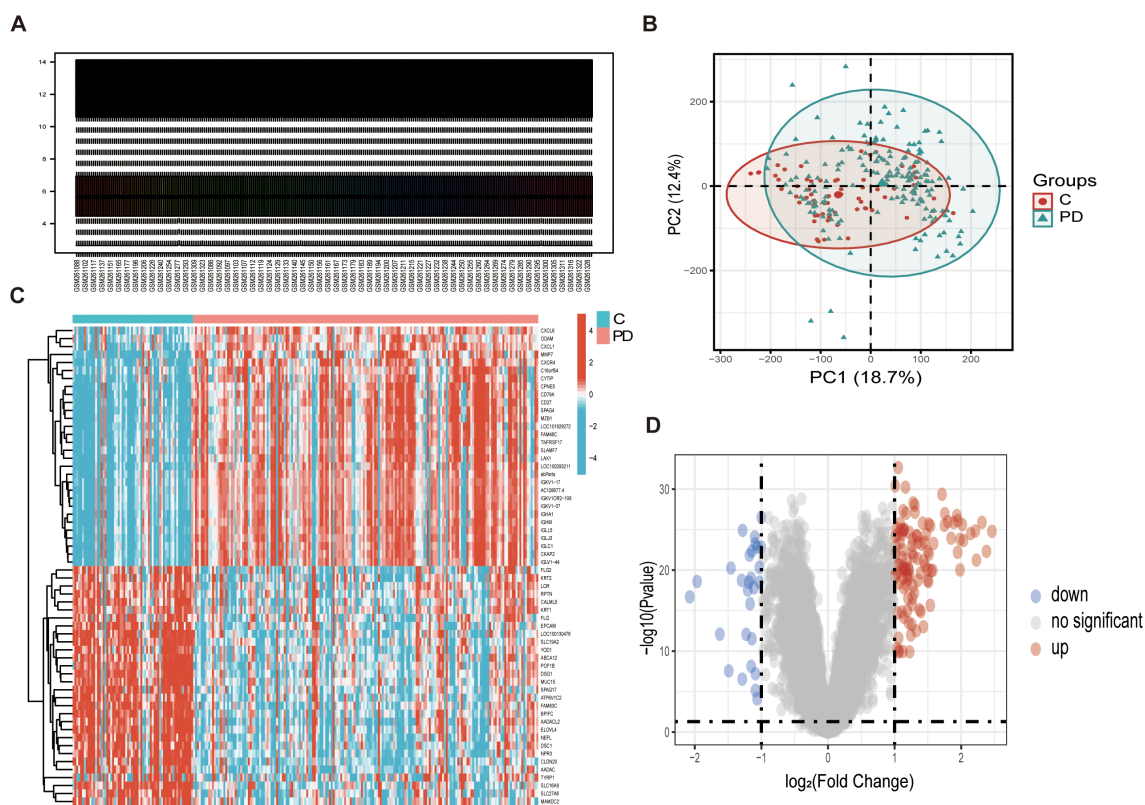


FIGURE 1

Differential expression analysis of the Periodontitis dataset. (A) A standardized box plot was generated to visualize the distribution of gene expression data from GEO datasets, with each box representing the expression profile of an individual sample. (B) The PCA plot illustrates the distribution of the controls (red) and the periodontitis (blue) in the principal component space. (C) Heatmap depicting the expression patterns of DEG. The intensity of the color corresponds to the expression level, with red representing up-regulated genes and blue representing down-regulated genes. (D) Volcano plot displaying the results of the statistical significance analysis of DEGs. The horizontal axis represents the logarithmic fold change in gene expression ( $\log_2(\text{FC})$ ), while the vertical axis represents the statistical significance ( $-\log_{10}(P\text{-value})$ ). Significantly changed genes are highlighted: up-regulated in red, down-regulated in blue, and non-significant in gray.

clinical trait. Among these, the white module, containing 154 genes, exhibited the most significant correlation with the oxidative stress gene set (Cor = 0.7, P = 4e-38) (Figure 2D).

### 3.3 Identification of DEGs associated with oxidative stress and their functions

The 139 DEGs were identified through a comparison of gene expression profiles between periodontitis samples and controls, as presented in Figure 1B. By intersecting these DEGs with oxidative stress-related modules derived from WGCNA, we identified 13 genes

highly associated with oxidative stress (Figure 3A). Subsequent Gene Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) enrichment analyses were performed on these 13 genes. GO analysis revealed significant enrichment in biological processes and molecular functions related to the extracellular matrix, chemokine activity, and cytokine activity (Figure 3B). Meanwhile, KEGG pathway analysis highlighted significant enrichment in key pathways such as IL-17 signaling, AGE-RAGE signaling, and cytokine-cytokine receptor interaction (Figure 3C). In summary, this analysis identified 13 genes with significant expression changes in periodontitis, closely linked to oxidative stress, offering insights into the molecular mechanisms of periodontitis pathogenesis.

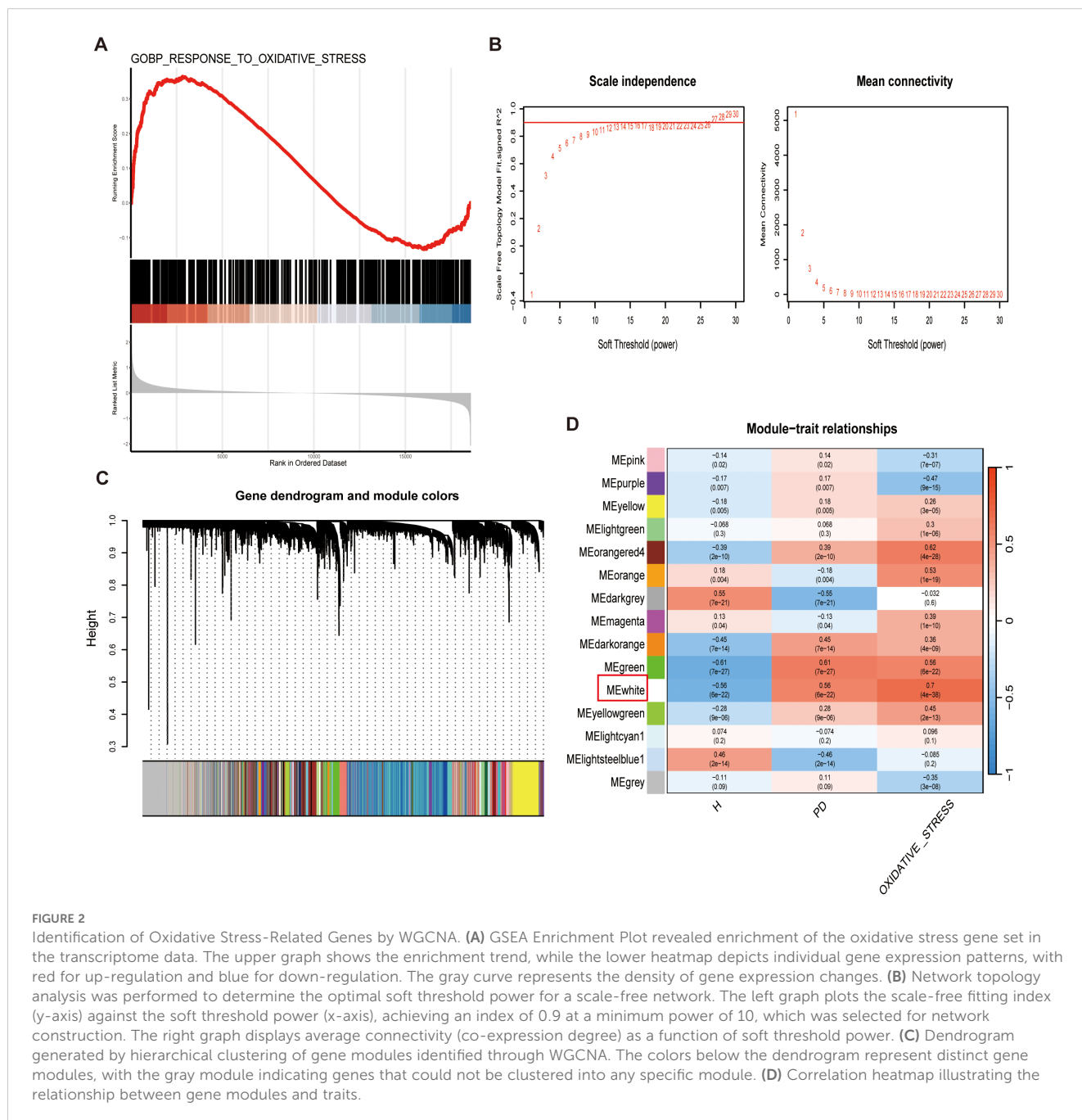


FIGURE 2

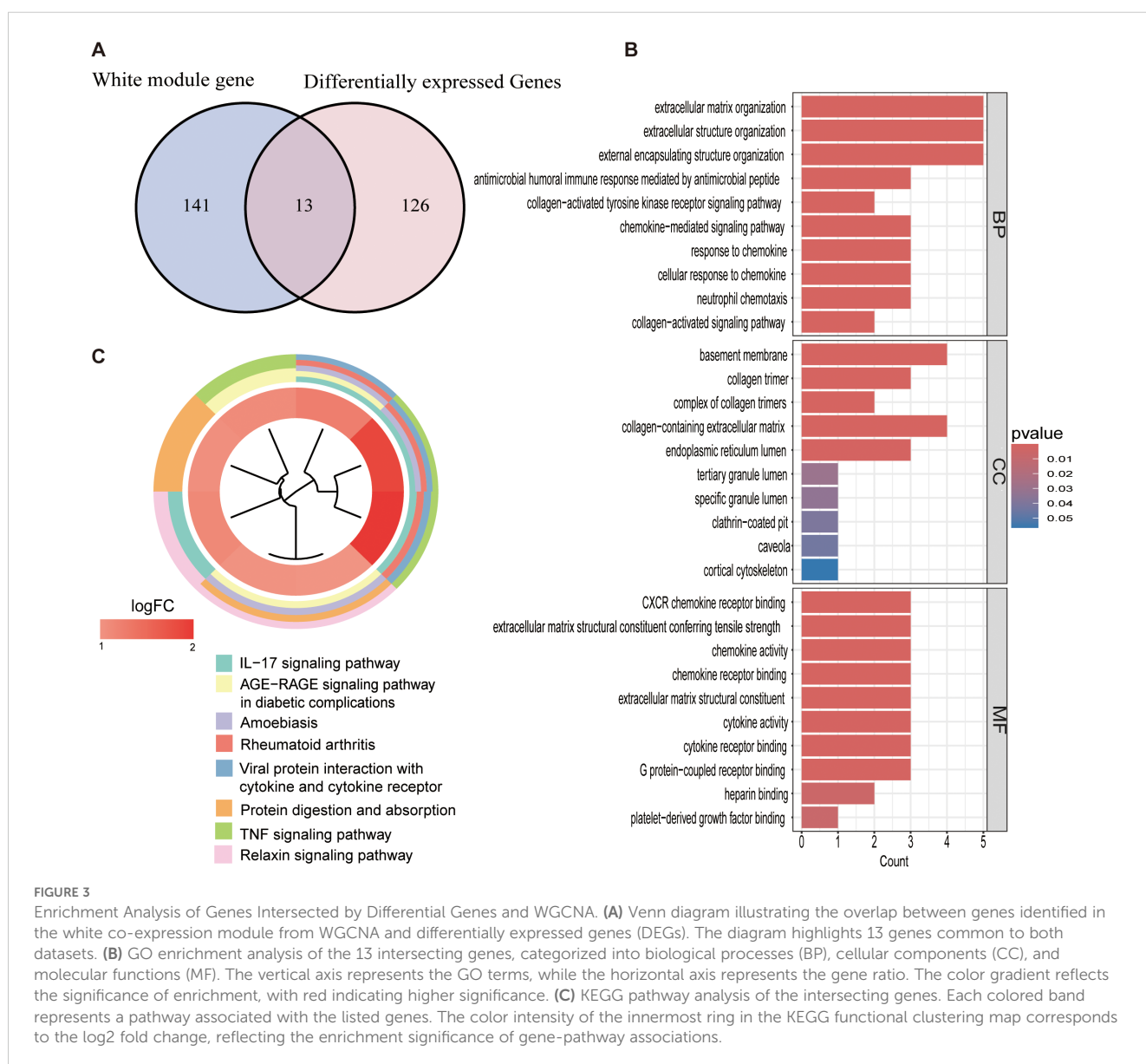
Identification of Oxidative Stress-Related Genes by WGCNA. (A) GSEA Enrichment Plot revealed enrichment of the oxidative stress gene set in the transcriptome data. The upper graph shows the enrichment trend, while the lower heatmap depicts individual gene expression patterns, with red for up-regulation and blue for down-regulation. The gray curve represents the density of gene expression changes. (B) Network topology analysis was performed to determine the optimal soft threshold power for a scale-free network. The left graph plots the scale-free fitting index (y-axis) against the soft threshold power (x-axis), achieving an index of 0.9 at a minimum power of 10, which was selected for network construction. The right graph displays average connectivity (co-expression degree) as a function of soft threshold power. (C) Dendrogram generated by hierarchical clustering of gene modules identified through WGCNA. The colors below the dendrogram represent distinct gene modules, with the gray module indicating genes that could not be clustered into any specific module. (D) Correlation heatmap illustrating the relationship between gene modules and traits.

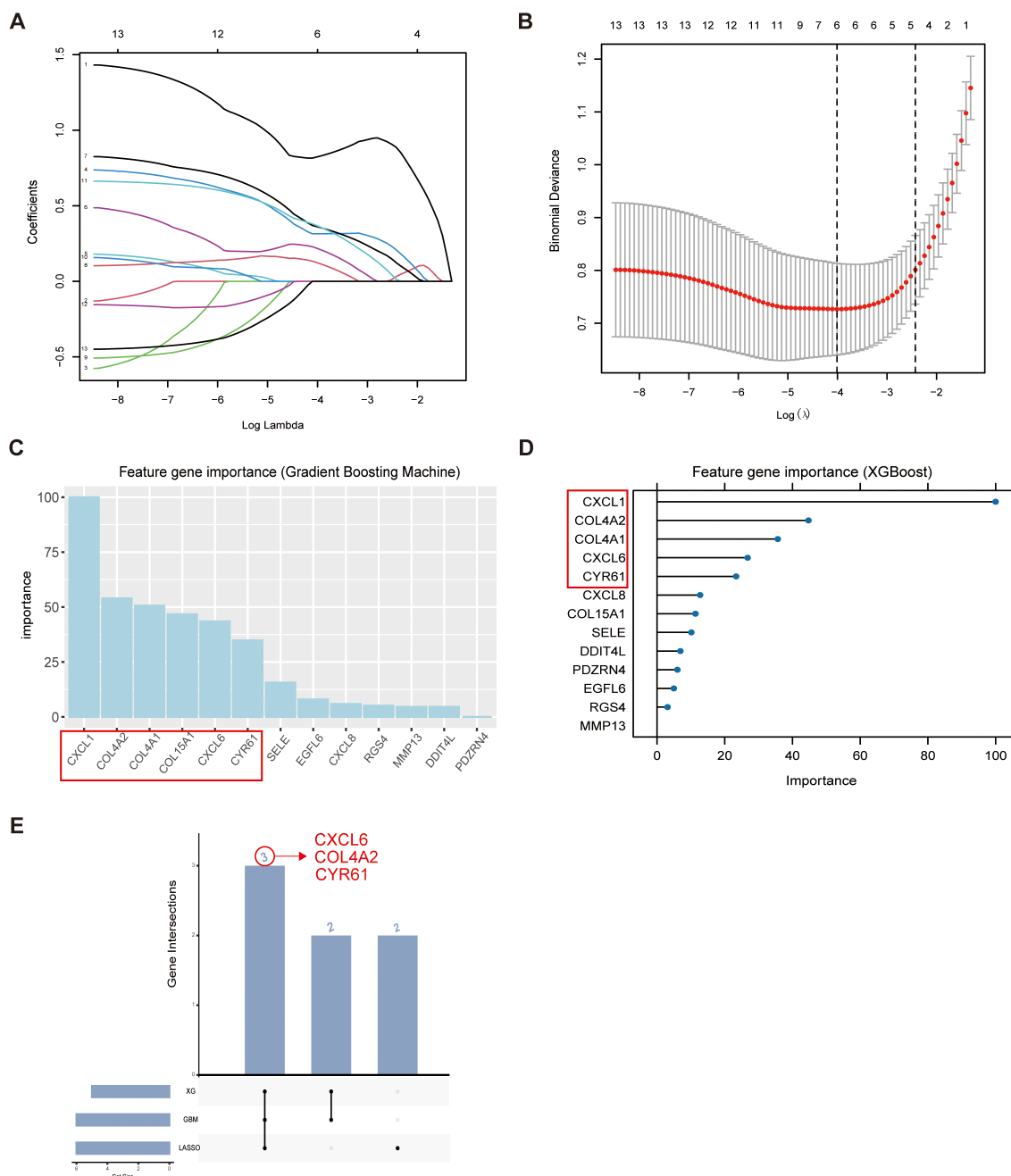
### 3.4 Identification of the key genes based on machine learning

Machine learning algorithms are widely recognized as powerful tools for identifying biomarkers associated with complex diseases. To further analyze the previously screened genes, we utilized three different machine learning algorithms to identify oxidative stress-related genes. Specifically, we utilized the LASSO regression (Figures 4A, B), Gradient Boosted Tree (GBM) (Figure 4C), and XGBoost (Figure 4D) algorithms for this analysis. The LASSO regression model was applied to refine the gene set by penalizing less relevant features, while the GBM and XGBoost algorithms were used to evaluate gene importance based on their contribution to the predictive models. Through cross-validation and comparative analysis of these algorithms, we identified three key genes: CXCL6, COL4A2, and CYR61 (Figure 4E). These genes were consistently highlighted across all three methods, suggesting their potential significance in the pathogenesis of periodontitis.

### 3.5 Validation of key genes and assessment of their diagnostic performance

The expression of the candidate key genes was verified using an external dataset (GSE16134). The results confirmed that the expression levels of COL4A2, CYR61, and CXCL6 were significantly different between periodontitis patients and the control population ( $P < 0.05$ ) (Figures 5A-C). The diagnostic performance of these three key genes was evaluated using receiver operating characteristic (ROC) curves. The areas under the ROC curves (AUC) were calculated as follows: COL4A2 (AUC = 0.874), CYR61 (AUC = 0.793), and CXCL6 (AUC = 0.838) (Figures 5D-F). These AUC values indicate that the key genes exhibit high predictive accuracy and are effective in distinguishing between periodontitis and controls. Overall, the findings demonstrate that COL4A2, CYR61, and CXCL6 have strong diagnostic potential, highlighting their utility as biomarkers for periodontitis.





**FIGURE 4** Machine learning screening of key genes in intersecting genes. **(A, B)** The lasso algorithm determined 6 feature genes based on lambda. min values. **(C)** The GBM algorithm analyzed intersecting genes, identifying 6 significant genes with importance scores >25. **(D)** The XGBoost algorithm was employed to further refine the analysis, identifying 5 genes with importance scores>20. **(E)** The upset plot was generated to visualize the overlap of key genes identified by the three algorithms, revealing three consistently highlighted intersecting genes - CXCL6, COL4A2, and CYR61.

### 3.6 Periodontitis rat model for validation of key gene expression

In this study, We established a ligature-induced rat model of periodontitis. Micro-CT analysis revealed significant bone resorption in the periodontitis group compared to controls (Figure 6A). Additionally, both the periodontal pocket depth

(Figure 6B) and the distance between the cemento-enamel junction and the alveolar bone crest (Figure 6C) were markedly increased in the periodontitis group. These differences were statistically significant (p < 0.05). Histopathological examination using H&E staining further corroborated the successful induction of periodontitis. The periodontitis group exhibited pronounced bone resorption and inflammatory cell infiltration (Figure 6D).



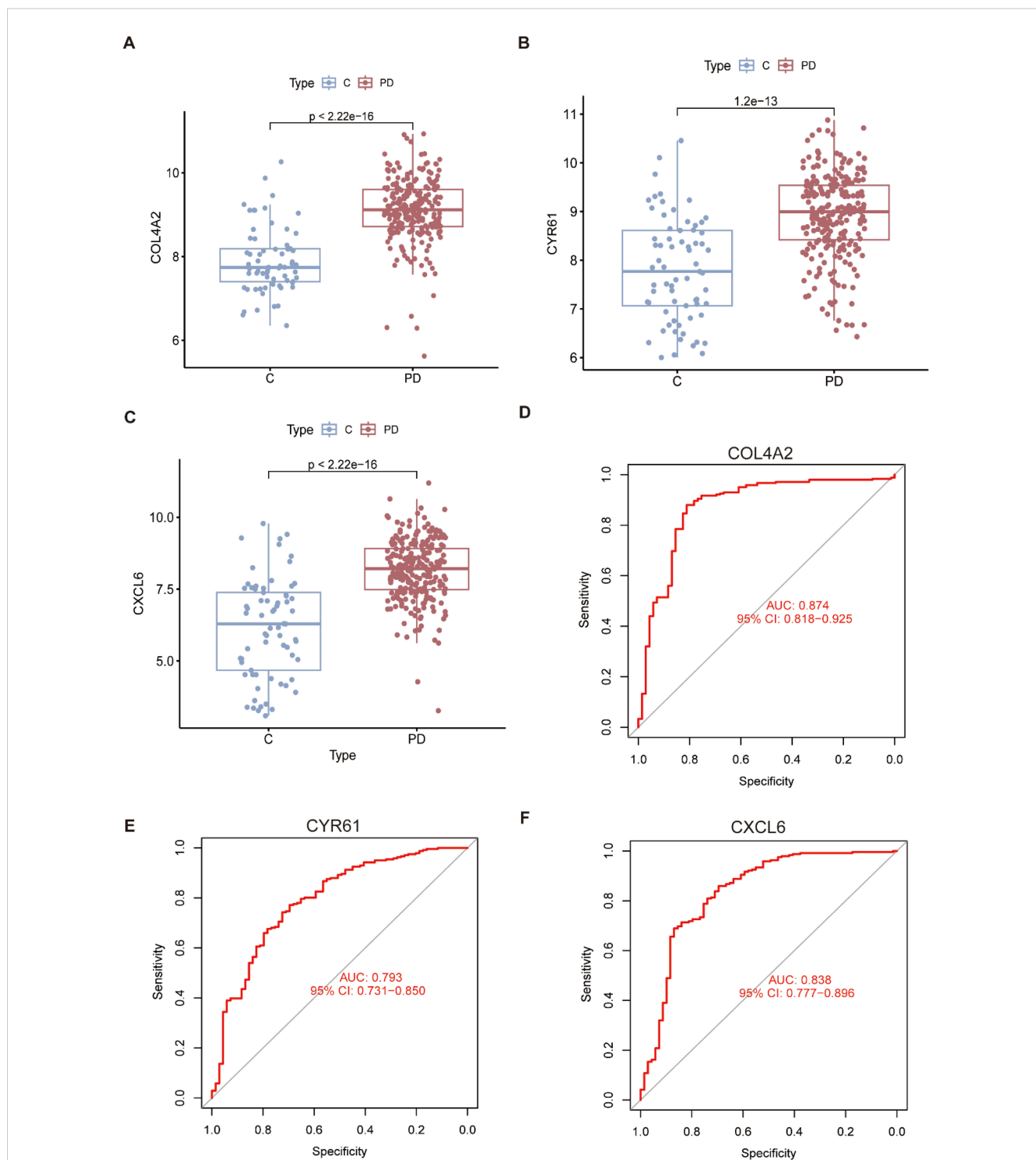


FIGURE 5

Expression and ROC profiles of key genes. (A–C) The box plot illustrates the expression differences of three key genes based on an external dataset. (D–F) ROC curve analysis of the three key genes, evaluating their diagnostic and prognostic significance.

Regarding oxidative stress levels, the periodontitis group demonstrated significantly higher levels of MDA and SOD compared to the controls (Figures 6E, F). These differences were also statistically significant ( $p < 0.05$ ), indicating elevated oxidative stress in the periodontal tissues of rats with periodontitis. Analysis of mRNA levels in gingival tissues revealed significant upregulation

of COL4A2 and CXCL6 in the periodontitis group compared to the controls (Figures 6G, I) ( $p < 0.05$ ), while no significant difference in CYR61 expression was observed between the two groups (Figure 6H). Consequently, COL4A2 and CXCL6 were identified as marker genes for further investigation into the specific cell types associated with their effects.

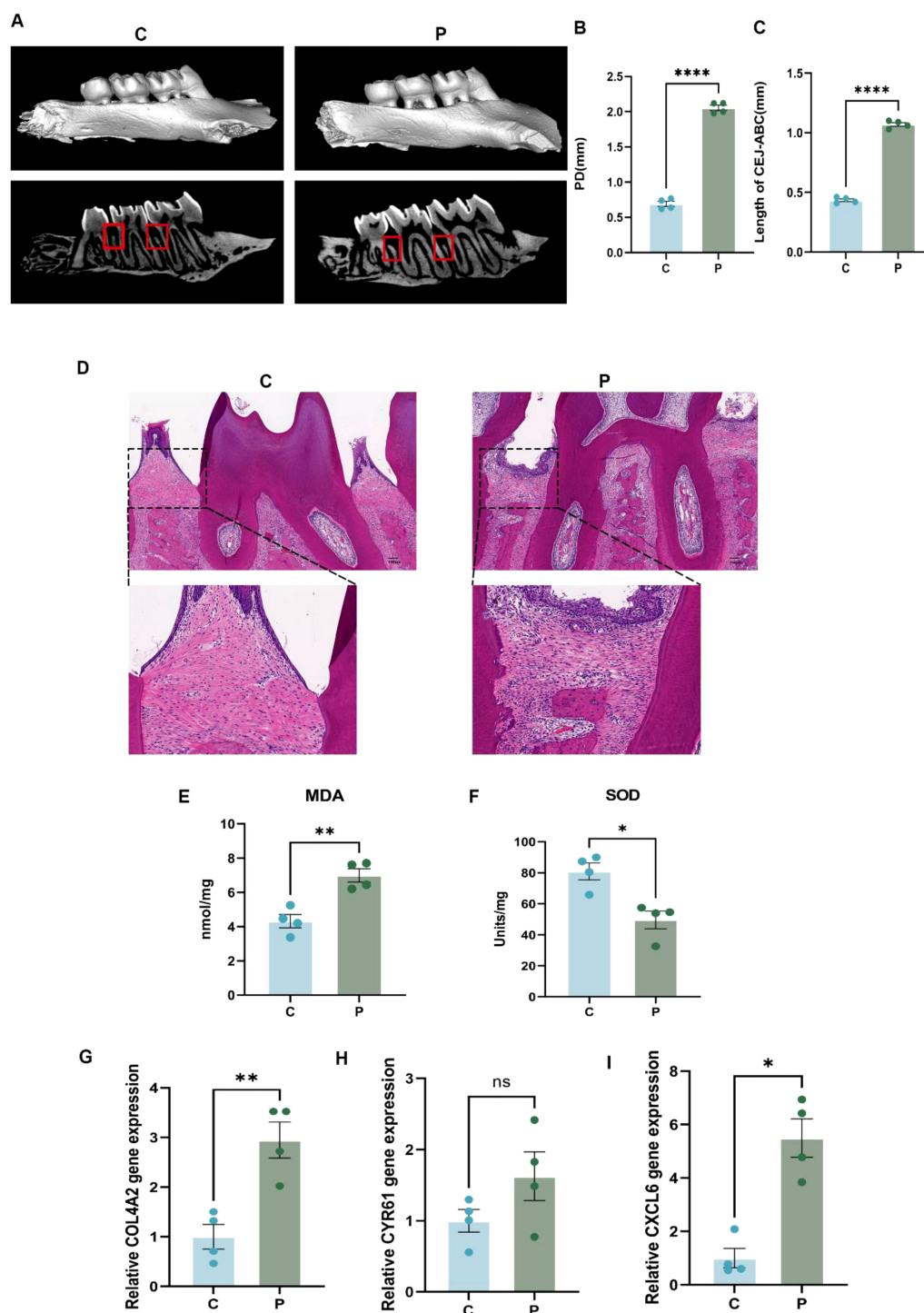


FIGURE 6

Rat experimental periodontitis model to verify the expression of key genes. (A) Three-dimensional reconstructed images of the periodontitis group (P) and the controls (C) obtained through Micro-CT scanning. (B) Bar graph illustrating periodontal pocket depth (PD), demonstrating a significant increase in the periodontitis group (P) compared to the control group (C). (C) Bar graph depicting the CEJ-ABC, indicating a significant increase in bone resorption in the periodontitis group (P) compared to the control group (C). (D) H&E stained sections of rat alveolar bone, showing the periodontal membrane, alveolar bone, and inflammatory cell infiltration. The scale bar represents 100 $\mu$ m. (E, F) Bar graphs showing the levels of MDA and SOD in the periodontitis group (P) compared to the controls (C), respectively, to evaluate differences in oxidative stress levels. (G-I) Relative mRNA expression levels of the genes COL4A2, CYR61, and CXCL6 were determined by qPCR in both the control and periodontitis groups. Data are expressed as mean  $\pm$  SD, n=4. \*p < 0.05; \*\*p < 0.01; \*\*\*\*p < 0.0001; ns, not significant.

### 3.7 Single-cell analysis reveals key gene expression patterns within cell clusters

To identify the primary sources and target cell types of two candidate genes, we analyzed scRNA-seq data. After filtering out low-quality cells, unsupervised cell clustering analysis was performed using UMAP, which annotated the cell clusters into 25 distinct cell clusters (Figure 7A). Following manual single-cell annotation, these clusters were further refined into seven major cell types: T cells, B cells, endothelial cells, epithelial cells, stromal cells, NK cells, fibroblasts, and mast cells (Figure 7B). The reliability of the cell type annotation was confirmed using bubble plots, which illustrated the average expression levels of marker genes within each cell subpopulation (Figure 7C). Subsequently, UMAP plots and violin plots were employed to visualize the spatial distribution and expression intensities of COL4A2 (Figures 7D, E) and CXCL6 (Figures 7F, G) across the cell populations. The results revealed that the expression of the COL4A2 gene was significantly upregulated in endothelial and stromal cells. In contrast, the CXCL6 gene exhibited pronounced upregulation primarily in epithelial cells.

## 4 Discussion

In this study, we developed a comprehensive analytical framework to identify biomarkers with potential diagnostic value by integrating bioinformatics analysis, machine learning algorithms, and Single cell analysis. Specifically, we identified two genes, COL4A2 and CXCL6, which are closely associated with oxidative stress and may play a critical role in the pathogenesis of periodontitis. Additionally, qPCR analysis confirmed the upregulated expression of COL4A2 and CXCL6 in gingival tissues of periodontitis-induced rats. Through scRNA-seq analysis, we further elucidated the distinct expression patterns of COL4A2 and CXCL6 across different cell types. COL4A2 was predominantly expressed in endothelial cells and stromal cells, while CXCL6 showed significant expression in epithelial cells. These findings suggest that these genes may contribute to the disease process through cell type-specific mechanisms, highlighting their potential as therapeutic targets or diagnostic markers in periodontitis. This integrated approach not only advances our understanding of the molecular mechanisms underlying periodontitis but also provides a robust methodology for identifying and validating biomarkers in complex diseases.

COL4A2 is the gene that encodes the  $\alpha 2$  chain of type IV collagen. Type IV collagen is an essential component of the basement membrane, providing structural support for cells such as endothelial cells and contributing to the stability of the extracellular matrix (26). Previous studies on diseases such as cerebral hemorrhage and ischemic brain injury have demonstrated that mutations in COL4A2 may increase the vulnerability of cerebral blood vessels by disrupting the structure and function of collagen IV (27, 28). Although no

direct evidence links COL4A2 to periodontitis, its significant role in endothelial cells suggests a potential mechanism by which it may mitigate oxidative stress-induced damage in periodontal tissues. We hypothesize that COL4A2 may protect periodontal tissues from reactive oxygen species (ROS) by maintaining the integrity of the endothelial basement membrane in microvessels (26, 29). Oxidative stress and the subsequent inflammatory response are key features of periodontitis pathogenesis. ScRNA-seq analysis demonstrated predominant COL4A2 localization in endothelial cells (Figures 7D, E). As mentioned in previous studies (27, 28), endothelial COL4A2 expression alters vascular permeability, thereby regulating the release of pro-inflammatory mediators including TNF- $\alpha$  and IL-1 $\beta$ , and consequently influencing inflammatory response severity. Additionally, COL4A2 expression in stromal cells contributes to extracellular matrix stability (30), which supports tissue repair and regeneration. In summary, the elevated expression of COL4A2 in endothelial and stromal cells may reflect a defense mechanism against oxidative stress. As a gene encoding extracellular matrix proteins, COL4A2 likely protects endothelial cells from oxidative damage by modulating extracellular matrix stability and intercellular signaling. Furthermore, studies have shown that COL4A2 promotes osteogenic differentiation of periodontal ligament stem cells (PDLSCs) by negatively regulating the Wnt/ $\beta$ -catenin signaling pathway, offering a potential therapeutic strategy for bone defect repair (31). However, the specific mechanisms linking COL4A2 to oxidative stress in periodontitis remain unclear. Further research is needed to determine whether COL4A2 can serve as a potential therapeutic target or a focus for mechanistic studies in periodontitis.

CXCL6, also known as GCP-2, is an ELR+ CXC chemokine that primarily mediates neutrophil chemotaxis by binding to CXCR1 and CXCR2 receptors (32). The results of scRNA-seq analysis demonstrated predominant CXCL6 expression in epithelial cells (Figures 7F, G), consistent with prior studies (33), CXCL6 can be induced in multiple cell types under inflammatory conditions, including epithelial cells. CXCL6 exhibits pro-inflammatory, pro-angiogenic, and antimicrobial properties, playing a critical role in modulating immune responses (34). Dysregulation of CXCL6 function and expression has been strongly linked to a range of disorders, particularly cancers, fibrosis, and inflammatory diseases (35–37). In the context of periodontitis, CXCL6 expression is markedly elevated in the gingival tissues of patients, where it is closely associated with inflammatory cell infiltration and tissue damage. Studies have demonstrated (38) that CXCL6 acts synergistically with IL-8 to enhance the inflammatory response by promoting neutrophil chemotaxis, thereby driving the pathological progression of periodontitis. Oxidative stress, a key factor in the inflammatory response, can induce CXCL6 expression through the activation of multiple signaling pathways, such as NF- $\kappa$ B (39). For instance, in models of ischemia-reperfusion injury, oxidative stress has been shown to regulate cell permeability, proliferation, and apoptosis by activating the AKT/FOXO3a signaling pathway. This pathway modulates the expression of Sirt3, which subsequently

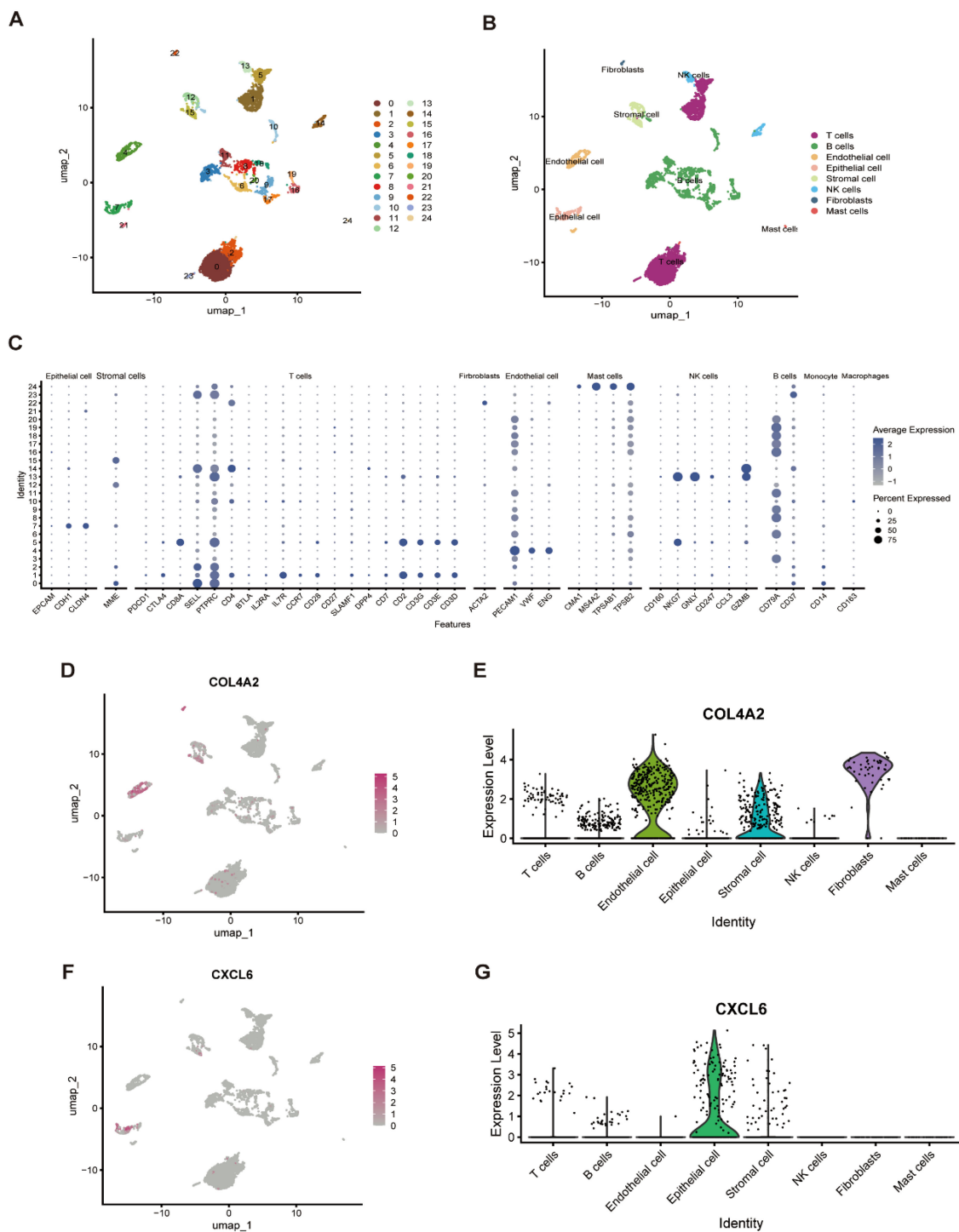


FIGURE 7

The expression patterns of key genes resolved at the single-cell level. **(A)** UMAP plot of single-cell data following dimensionality reduction analysis, where cells were categorized into 25 distinct clusters. **(B)** UMAP plot visualization displays 8 manually annotated cell clusters, providing a refined classification of the cell types. **(C)** Bubble plot illustrating the expression profiles of selected marker genes across the identified cell types. Each small bubble represents the distribution of gene expression within a specific cell type. **(D, E)** Expression pattern of COL4A2 at the single-cell level, visualized using a UMAP plot and a violin plot. **(F, G)** Expression pattern of CXCL6 at the single-cell level, depicted through a UMAP plot and a violin plot.

influences CXCL6 secretion (40). These findings suggest that oxidative stress may indirectly regulate CXCL6 expression through inflammatory signaling pathways, thereby influencing the trajectory of the inflammatory response.

Our study introduces an integrated strategy for the comprehensive characterization of oxidative stress-related gene expression and cellular heterogeneity in periodontitis. By leveraging this multi-faceted methodology, our study addresses the limitations inherent in relying on a single technique, thereby offering a more holistic understanding of the transcriptomic landscape and the molecular mechanisms underlying specific cell types.

Although this study offers novel insights and employs animal model experiments to validate the relevance of specific genes in periodontitis, several limitations warrant further investigation. First, the relatively small sample size ( $n = 4$ ) in the current rat model experimental may compromise the statistical power and reliability of the results. Additionally, while we identified a correlation between oxidative stress activity and the expression of COL4A2 and CXCL6, the precise mechanisms by which these genes influence oxidative stress and contribute to the progression of periodontitis remain unclear. More critically, the diagnostic value of these genes as biomarkers and their potential as therapeutic targets require validation in clinical cohort studies. Future investigations should expand experimental sample sizes, systematically elucidate the specific signaling pathways and cellular processes regulated by COL4A2 and CXCL6 under oxidative stress contexts, and evaluate the feasibility of their translational application to human diseases via multi-center clinical studies.

## 5 Conclusion

In this study, we employed an innovative approach to screen and identify two key genes, COL4A2 and CXCL6, by integrating machine learning algorithms with scRNA-seq analysis. Notably, our results demonstrated that the upregulation of these genes was closely associated with oxidative stress activity, with significant expression observed primarily in endothelial, stromal, and epithelial cells. These findings underscore the potential of COL4A2 and CXCL6 as biomarkers for periodontitis. Furthermore, they provide a promising foundation for the development of personalized and effective therapeutic strategies aimed at improving patient prognosis.

## Data availability statement

The original contributions presented in the study are included in the article/supplementary material. Further inquiries can be directed to the corresponding authors.

## Ethics statement

The animal study was approved by Experimental Animal Ethics Committee of the First Affiliated Hospital of Guangdong Pharmaceutical University. The study was conducted in accordance with the local legislation and institutional requirements.

## Author contributions

SS: Writing – original draft, Software, Visualization, Formal analysis, Methodology, Conceptualization, Data curation. JR: Investigation, Writing – review & editing. XZ: Writing – original draft. SC: Data curation, Formal analysis, Writing – review & editing. YC: Data curation, Investigation, Writing – review & editing. QZ: Writing – review & editing, Methodology, Investigation. JY: Writing – review & editing, Funding acquisition.

## Funding

The author(s) declare that financial support was received for the research and/or publication of this article. This study was partly supported by the Clinical Specialty Capacity Building Support Program of SYSU First Affiliated Hospital (Grant No. R70039), Guangzhou Municipal Science and Technology Project (2025A04J4077), and the Traditional Chinese Medicine Bureau of Guangdong Province (20251056).

## Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

## Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

## References

- Pihlstrom BL, Michalowicz BS, Johnson NW. Periodontal diseases. *Lancet*. (2005) 366:1809–20. doi: 10.1016/s0140-6736(05)67728-8
- Sczepanik FSC, Grossi ML, Casati M, Goldberg M, Glogauer M, Fine N, et al. Periodontitis is an inflammatory disease of oxidative stress: We should treat it that way. *Periodontol* 2000. (2020) 84:45–68. doi: 10.1111/prd.12342
- van der Pol A, van Gilst WH, Voors AA, van der Meer P. Treating oxidative stress in heart failure: past, present and future. *Eur J Heart Fail*. (2019) 21:425–35. doi: 10.1002/ehf.1320
- Wu Y, Hu H, Wang T, Guo W, Zhao S, Wei R. Characterizing mitochondrial features in osteoarthritis through integrative multi-omics and machine learning analysis. *Front Immunol*. (2024) 15:1414301. doi: 10.3389/fimmu.2024.1414301
- Wen P, Sun Z, Gou F, Wang J, Fan Q, Zhao D, et al. Oxidative stress and mitochondrial impairment: Key drivers in neurodegenerative disorders. *Ageing Res Rev*. (2025) 104:102667. doi: 10.1016/j.arr.2025.102667
- Radzki D, Negri A, Kusiak A, Obuchowski M. Matrix metalloproteinases in the periodontium-vital in tissue turnover and unfortunate in periodontitis. *Int J Mol Sci*. (2024) 25:2763. doi: 10.3390/ijms25052763
- Xu Y, Luo Y, Weng Z, Xu H, Zhang W, Li Q, et al. Microenvironment-responsive metal-phenolic nanozyme release platform with antibacterial, ROS scavenging, and osteogenic for periodontitis. *ACS Nano*. (2023) 17:18732–46. doi: 10.1021/acsnano.3c01940
- Yuan Z, Li J, Xiao F, Wu Y, Zhang Z, Shi J, et al. Sinensetin protects against periodontitis through binding to Bach1 enhancing its ubiquitination degradation and improving oxidative stress. *Int J Sci*. (2024) 16:38. doi: 10.1038/s41368-024-00305-z
- Bhandari N, Walambe R, Kotecha K, Khare SP. A comprehensive survey on computational learning methods for analysis of gene expression data. *Front Mol Biosci*. (2022) 9:907150. doi: 10.3389/fmolb.2022.907150
- Song B, Liu D, Dai W, McMyn NF, Wang Q, Yang D, et al. Decoding heterogeneous single-cell perturbation responses. *Nat Cell Biol*. (2025) 27:493–504. doi: 10.1038/s41556-025-01626-9
- Potter SS. Single-cell RNA sequencing for the study of development, physiology and disease. *Nat Rev Nephrol*. (2018) 14:479–92. doi: 10.1038/s41581-018-0021-7
- Gao X, Guo Z, Wang P, Liu Z, Wang Z. Transcriptomic analysis reveals the potential crosstalk genes and immune relationship between IgA nephropathy and periodontitis. *Front Immunol*. (2023) 14:1062590. doi: 10.3389/fimmu.2023.1062590
- Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinf*. (2008) 9:559. doi: 10.1186/1471-2105-9-559
- Tarca AL, Bhatti G, Romero R. A comparison of gene set analysis methods in terms of sensitivity, prioritization and specificity. *PLoS One*. (2013) 8:e79217. doi: 10.1371/journal.pone.0079217
- Frost HR, Amos CI. Gene set selection via LASSO penalized regression (SLPR). *Nucleic Acids Res*. (2017) 45:e114. doi: 10.1093/nar/gkx291
- Yang C, Chen M, Yuan Q. The application of XGBoost and SHAP to examining the factors in freight truck-related crashes: An exploratory analysis. *Accid Anal Prev*. (2021) 158:106153. doi: 10.1016/j.aap.2021.106153
- Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res*. (2015) 43:e47. doi: 10.1093/nar/gkv007
- Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A*. (2005) 102:15545–50. doi: 10.1073/pnas.0506580102
- Chen H, Boutros PC. VennDiagram: a package for the generation of highly-customizable Venn and Euler diagrams in R. *BMC Bioinf*. (2011) 12:35. doi: 10.1186/1471-2105-12-35
- Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing biological themes among gene clusters. *Omic*. (2012) 16:284–7. doi: 10.1089/omi.2011.0118
- Hsu CL, Wu PC, Wu FZ, Yu HC. LASSO-derived model for the prediction of lean-non-alcoholic fatty liver disease in examinees attending a routine health check-up. *Ann Med*. (2024) 56:2317348. doi: 10.1080/07853890.2024.2317348
- Mall R, Cerulo L, Garofano L, Frattini V, Kunji K, Bensmail H, et al. RGBM: regularized gradient boosting machines for identification of the transcriptional regulators of discrete glioma subtypes. *Nucleic Acids Res*. (2018) 46:e39. doi: 10.1093/nar/gky015
- Hou N, Li M, He L, Xie B, Wang L, Zhang R, et al. Predicting 30-days mortality for MIMIC-III patients with sepsis-3: a machine learning approach using XGboost. *J Transl Med*. (2020) 18:462. doi: 10.1186/s12967-020-02620-5
- Engelbrechts S, Bohlin J. Statistical predictions with glmnet. *Clin Epigenetics*. (2019) 11:123. doi: 10.1186/s13148-019-0730-1
- Cao Y, Fu L, Wu J, Peng Q, Nie Q, Zhang J, et al. Integrated analysis of multimodal single-cell data with structural similarity. *Nucleic Acids Res*. (2022) 50:e121. doi: 10.1093/nar/gkac781
- Steffensen LB, Rasmussen LM. A role for collagen type IV in cardiovascular disease? *Am J Physiol Heart Circ Physiol*. (2018) 315(3):H610–h25. doi: 10.1152/ajpheart.00070.2018
- McNeilly S, Thomson CR, Gonzalez-Trueba L, Sin YY, Granata A, Hamilton G, et al. Collagen IV deficiency causes hypertrophic remodeling and endothelium-dependent hyperpolarization in small vessel disease with intracerebral hemorrhage. *EBioMedicine*. (2024) 107:105315. doi: 10.1016/j.ebiom.2024.105315
- Maurice P, Guilbaud L, Garel J, Mine M, Dugas A, Friszer S, et al. Prevalence of COL4A1 and COL4A2 mutations in severe fetal multifocal hemorrhagic and/or ischemic cerebral lesions. *Ultrasound Obstet Gynecol*. (2021) 57:783–9. doi: 10.1002/uog.22106
- Kalluri R. The biology and function of fibroblasts in cancer. *Nat Rev Cancer*. (2016) 16:582–98. doi: 10.1038/nrc.2016.73
- Rannikmäe K, Sivakumaran V, Millar H, Malik R, Anderson CD, Chong M, et al. COL4A2 is associated with lacunar ischemic stroke and deep ICH: Meta-analyses among 21,500 cases and 40,600 controls. *Neurology*. (2017) 89:1829–39. doi: 10.1212/wnl.0000000000004560
- Wen Y, Yang H, Wu J, Wang A, Chen X, Hu S, et al. COL4A2 in the tissue-specific extracellular matrix plays important role on osteogenic differentiation of periodontal ligament stem cells. *Theranostics*. (2019) 9:4265–86. doi: 10.7150/thno.35914
- Cai X, Li Z, Zhang Q, Qu Y, Xu M, Wan X, et al. CXCL6-EGFR-induced Kupffer cells secrete TGF- $\beta$ 1 promoting hepatic stellate cell activation via the SMAD2/BRD4/C-MYC/EZH2 pathway in liver fibrosis. *J Cell Mol Med*. (2018) 22:5050–61. doi: 10.1111/jcmm.13787
- Wuyts A, Struyf S, Gijssbers K, Schutysse E, Put W, Conings R, et al. The CXC chemokine GCP-2/CXCL6 is predominantly induced in mesenchymal cells by interleukin-1 $\beta$  and is down-regulated by interferon- $\gamma$ : comparison with interleukin-8/CXCL8. *Lab Invest*. (2003) 83:23–34. doi: 10.1097/01.lab.0000048719.53282.00
- Dai CL, Yang HX, Liu QP, Rahman K, Zhang H. CXCL6: A potential therapeutic target for inflammation and cancer. *Clin Exp Med*. (2023) 23:4413–27. doi: 10.1007/s10238-023-01152-8
- Song M, He J, Pan QZ, Yang J, Zhao J, Zhang YJ, et al. Cancer-associated fibroblast-mediated cellular crosstalk supports hepatocellular carcinoma progression. *Hepatology*. (2021) 73:1717–35. doi: 10.1002/hep.31792
- Bahudhanapati H, Tan J, Apel RM, Seeliger B, Schupp J, Li X, et al. Increased expression of CXCL6 in secretory cells drives fibroblast collagen synthesis and is associated with increased mortality in idiopathic pulmonary fibrosis. *Eur Respir J*. (2024) 63:2300088. doi: 10.1183/13993003.00088-2023
- Caxaria S, Kouvatso N, Eldridge SE, Alvarez-Fallas M, Thorup AS, Cici D, et al. Disease modification and symptom relief in osteoarthritis using a mutated GCP-2/CXCL6 chemokine. *EMBO Mol Med*. (2023) 15:e16218. doi: 10.15252/emmm.20216218
- Plemmenos G, Evangelidou E, Polizogopoulos N, Chalazias A, Deligianni M, Piperi C. Central regulatory role of cytokines in periodontitis and targeting options. *Curr Med Chem*. (2021) 28:3032–58. doi: 10.2174/0929867327666200824112732
- Korbecki J, Kojder K, Kapczuk P, Kupnicka P, Gawrońska-Szklarz B, Gutowska I, et al. The effect of hypoxia on the expression of CXC chemokines and CXC chemokine receptors-A review of literature. *Int J Mol Sci*. (2021) 22:843. doi: 10.3390/ijms22020843
- Wang X, Dai Y, Zhang X, Pan K, Deng Y, Wang J, et al. CXCL6 regulates cell permeability, proliferation, and apoptosis after ischemia-reperfusion injury by modulating Sirt3 expression via AKT/FOXO3a activation. *Cancer Biol Ther*. (2021) 22:30–9. doi: 10.1080/15384047.2020.1842705