

omicML: An Integrative Bioinformatics and Machine Learning Framework for Transcriptomic Biomarker Identification

Joy Prokash Debnath^{1†}, Kabir Hossen^{1†}, Md. Sayeam Khandaker^{1, 3‡}, Shawon Majid², Md Mehrajul Islam², Siam Arefin², Preonath Chondrow Dev^{1*}, Saifuddin Sarker^{1*}, Tanvir Hossain^{1*}

1 Department of Biochemistry and Molecular Biology, Shahjalal University of Science and
Technology, Sylhet 3114, Bangladesh

2 Department of Software Engineering, Shahjalal University of Science and Technology, Sylhet
3114, Bangladesh

3 Bangladesh Medical University, Dhaka, Bangladesh

† Equal contribution

* Correspondence:

Tanvir Hossain

Email: tanvir-bmb@sust.edu

Saifuddin Sarker

Email: saifuddinker@gmail.com

Preonath Chondrow Dev

Email: preonath2838@gmail.com

21 Abstract

22 Introduction

23 Transcriptomic biomarker discovery has been a challenge due to variation in datasets
24 and platforms, complexity in statistical and computational methods, integration of multiple
25 programming languages, and intricacy of ML workflow to evaluate biomarkers. Standard
26 workflows necessitate several stages (quality control, normalization, differential expression),
27 typically executed in R or Python, resulting in bottlenecks for non-experts. Existing
28 platforms have alleviated certain challenges by offering graphical interfaces for data
29 loading, normalization, differential gene expression analysis, and functional analysis;
30 nevertheless, they typically do not incorporate integrated machine learning procedures for
31 biomarker selection.

32 Method

33 In this regard, we present omicML, an intuitive graphical user interface (GUI) that combines
34 transcriptomic data analysis with machine learning (ML)-based classification via integrating R
35 and Python packages/libraries. It supports both RNA-Seq and microarray data,
36 automating preprocessing (data import, quality control, and normalization) and differential
37 expression analysis. The tool annotates differentially expressed genes (DEGs) with
38 descriptions, gene ontology, and pathway information and incorporates comparative analysis.
39 Our extensive ML pipeline enables both supervised and unsupervised learning, integrates
40 various datasets based on candidate gene signatures, standardizes and eliminates less
41 significant features, benchmarks multiple ML classifiers with robust performance metrics (e.g.,
42 AUROC, AUPRC), assesses feature importance, develops single-gene and multi-gene predictive
43 models, and systematically finalizes the biomarker algorithm. All functionalities are available in
44 omicML, hence reducing the barrier for biologists without computational proficiency.

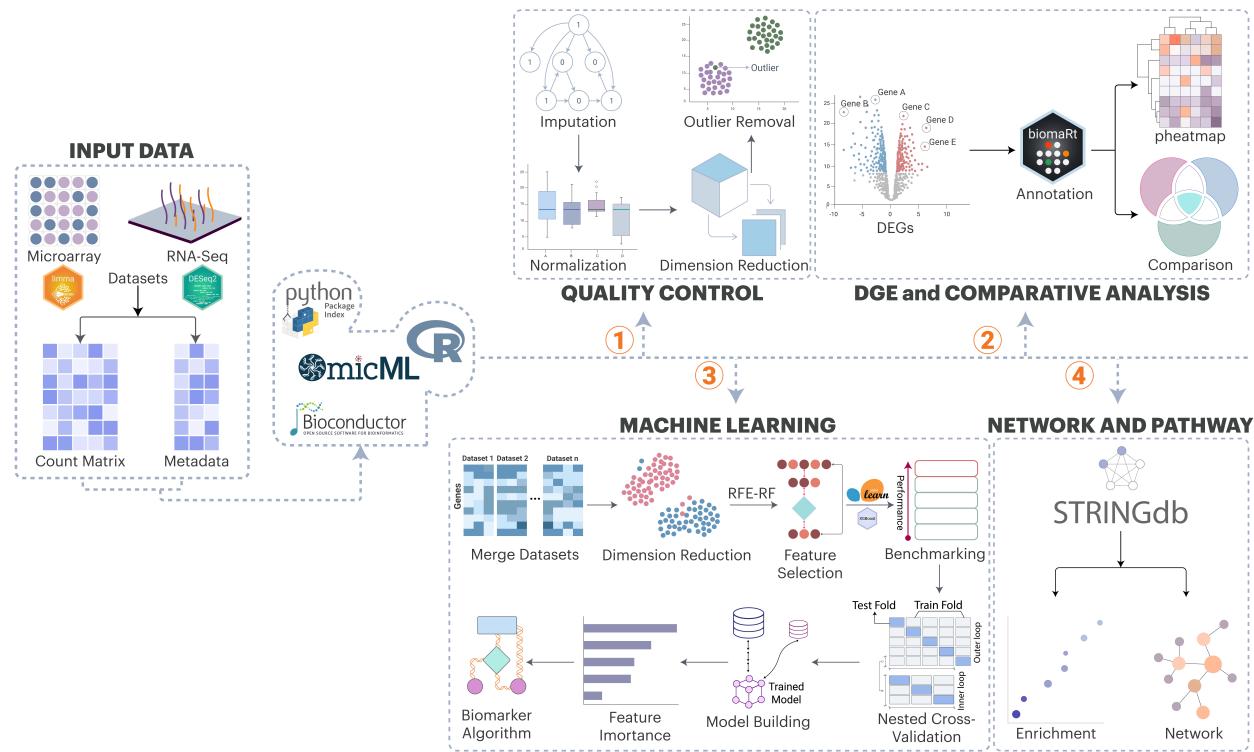
45 Result

46 In a case study, omicML identified a six-gene diagnostic model that distinguishes Mpox
47 (monkeypox virus) infections from those caused by other viruses, including SARS-CoV-2, HIV,
48 Ebola, and varicella-zoster. These results illustrate omicML's capacity to discern clinically
49 relevant biomarkers from complex transcriptome data.

50 Conclusion

51 Through the unified system, omicML (<https://omicml.org>), integrating data
52 preprocessing, differential gene expression analysis, annotation, heatmap analysis, dataset
53 integration, batch effect correction, machine learning approach, and functional analysis can
54 diminish technical barriers and accelerates the conversion of expression data into diagnostic
55 insights for clinicians and bench scientists.

56 **Keywords:** Machine Learning (ML), Biomarker, Mpox, Differentially Expressed Genes
57 (DEGs), Graphical User Interface (GUI), Transcriptomics, Bioinformatics, Web-based platform.



58 **Figure 1** Overview of omicML and its modules.

59 **Introduction**

60 Microarrays and RNA sequencing have been extensively utilized to investigate gene behavior in
61 diseases, environmental stresses, or infections. These techniques generate mountains of data and
62 translating high-dimensional expression data into robust biomarkers poses significant challenges,
63 particularly for researchers without computational expertise. A biomarker discovery workflow
64 demands proficiency in diverse tools for preprocessing, differential expression analysis, functional
65 annotation, and advanced machine learning (ML), each requires specialized programming skills in
66 R, Python, or Bash. Existing platforms streamline differential gene expression (DGE) analysis and
67 pathway enrichment, but they typically lack ML pipelines for predictive biomarker discovery. To
68 address these gaps, we present omicML, an interactive web application that integrates
69 bioinformatics and machine learning workflows for the development of transcriptome biomarkers
70 with clicks.

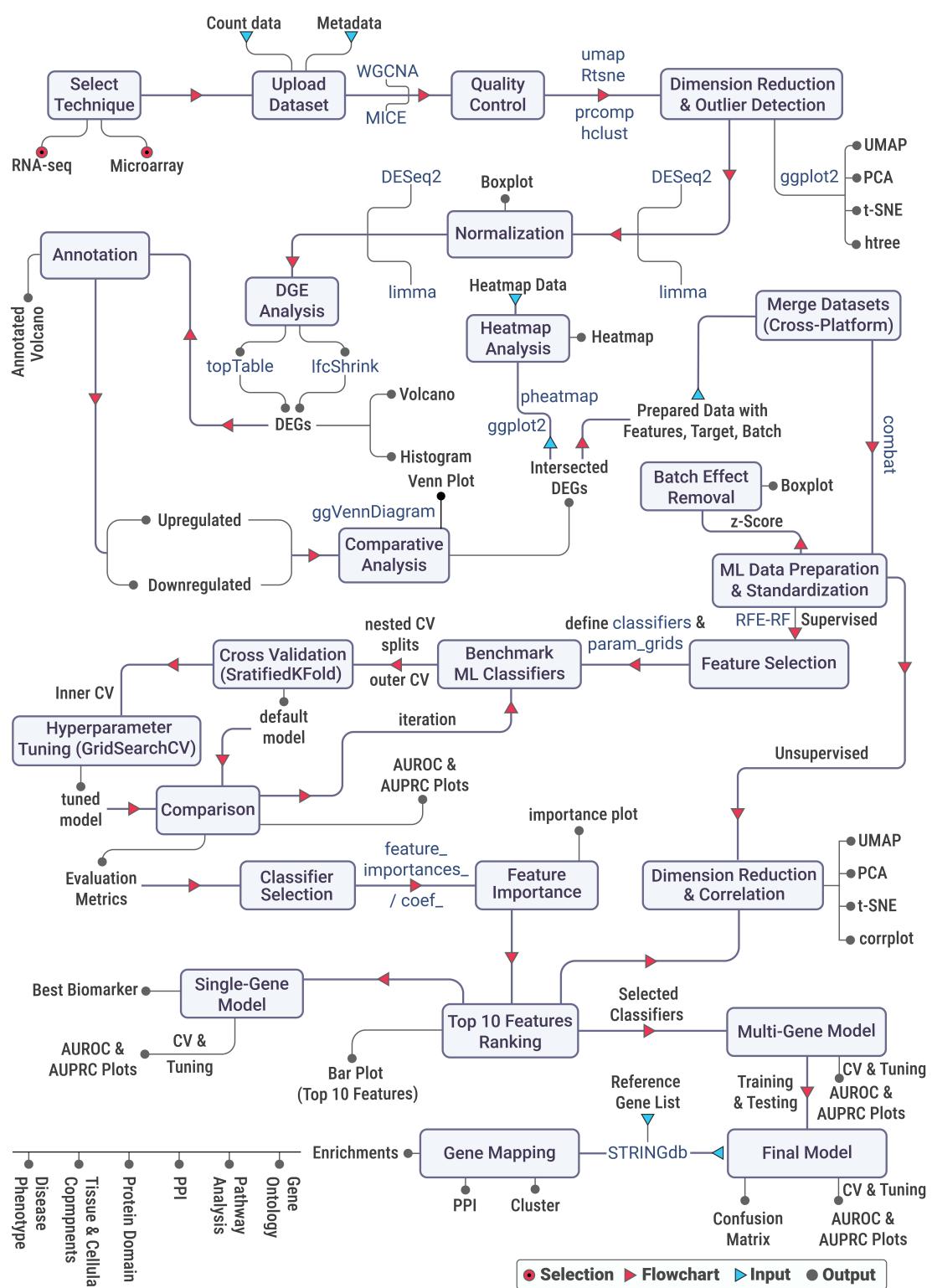
71 Over the years, numerous web-based platforms have been developed to optimize genomic and
72 transcriptomic investigations, providing intuitive interfaces that enable complex data interrogation
73 with minimal user burden. For example, iDEP (integrated Differential Expression and Pathway
74 analysis) automates extensive gene ID conversion, provides comprehensive gene annotation, and
75 integrates statistical and visualization techniques, including principal component analysis (PCA),
76 DGE analysis, and pathway enrichment [1]. Similar platforms, such as START App (Shiny
77 Transcriptome Analysis Resource Tool), Degust, and ShinyNGS, offer intuitive interfaces for
78 clustering, PCA, expression visualization, and DGE analysis [2–4]. Furthermore, an open-source
79 platform, Galaxy, provides a broad suite of web-based tools for genomic and differential analysis

80 [5]. These platforms facilitate preliminary data exploration and DGE identification but do not
81 inherently identify biomarkers using advanced ML. Additional specialized instruments also tackle
82 many facets of gene expression analysis. IRIS-EDA facilitates DGE analysis through discovery-
83 driven methodologies, including correlation analysis, heatmap creation, clustering, and PCA [6].
84 In contrast, Phantauus provides cross-platform datasets integration facilities from the Gene
85 Expression Omnibus (GEO), allowing for data normalization and filtering before differential
86 analysis [7]. The GEO2R tool in the NCBI, GEO portal, enables pairwise comparisons of
87 expression data; however, it depends on the pre-normalized data provided by submitters and does
88 not execute batch correction or standardize microarray and RNA-Seq processing. GEO2R
89 excludes any gene that displays at least one NA value in any of the samples in the comparison [8].
90 Despite the advancements in transcriptome analysis provided by the aforementioned technologies,
91 none are specifically engineered for predictive biomarker modeling utilizing comprehensive ML
92 methodologies. To address this significant deficiency, we offer a universal algorithmic pipeline
93 capable of predicting biomarkers linked to various diseases, environmental pressures, and other
94 situations across several species, based on the differentially expressed genes revealed in
95 preliminary research. omicML extends existing capabilities with the following features:(1)
96 automated preprocessing of input data (imputation of missing values, batch effect correction), (2)
97 compatibility with both RNA-Seq and microarray datasets, enabling cross-platform analysis, (3)
98 annotation support for 367 species, providing broad taxonomic coverage, (4) integrated ML
99 pipelines for data standardization, feature selection, benchmarking, nested cross-validation,
100 hyperparameter tuning, feature importance, single-gene model building, multi-gene model
101 (biomarker algorithm) building, model finalization, and (5) embedded network analysis and
102 functional enrichment to contextualize candidate biomarkers within biological pathways. By
103 integrating these processes, omicML distinctly connects transcriptomic analysis with predictive
104 biomarker modeling, enabling users to convert DEG lists into actionable diagnostic or therapeutic
105 targets without requiring coding expertise. omicML has been designed using six general modules
106 (**Figure 01**).

107 This paper also illustrates the effectiveness of omicML through a case study on monkeypox virus
108 (MPXV) infections. In the light of the identification of a novel clade of 2022 Mpox and their
109 biomarkers in our previous study [9], we have now used heterogeneous transcriptomic datasets to
110 identify the conserved biomarkers across different clades and cell models. omicML identified 34
111 shared DEGs through comparative analysis and employed a machine learning pipeline to prioritize
112 six high-confidence biomarkers. Among these, RRAD emerged as the most robust single-gene
113 predictor of MPXV infection. This example highlights omicML's capacity to democratize ML-
114 driven biomarker discovery, providing a reproducible, scalable, cross-platform workflow for users
115 navigating complex transcriptomic data.

116 Materials and Methods

117 Key steps are summarized in **Figure 2** which outlines the omicML workflow within a
118 comprehensive bioinformatics pipeline for transcriptomic data analysis, focusing on biomarker
119 discovery and functional interpretation.



120 **Figure 2** Detailed workflow of omicML

121 **Data Acquisition and Input**

122 The analytical pipeline within the omicML framework commences with the acquisition of primary
123 input datasets of raw sequencing derived from either microarray or RNA-Seq experiments. Initial
124 data can either be generated *de novo* via web-based analytical suites (e.g., Galaxy [5]) or retrieved
125 from curated repositories such as Gene Expression Omnibus (GEO) [10], ArrayExpress [11], The
126 Cancer Genome Atlas (TCGA) [12], the Sequence Read Archive [13]. The structured input
127 system—an expression matrix and a metadata—is designed to minimize errors during processing
128 and provide a robust foundation for ML-driven biomarkers discovery.

129 **Preprocessing, Dimensionality Reduction, and Outlier Detection**

130 To ensure data integrity, initial preprocessing is employed by the filtration through WGCNA [14]
131 and the imputation technique to mitigate low-variance genes and missing values respectively.
132 Subsequently, dimension reduction techniques are used to visualize the high dimensional and
133 complex gene expression data into lower dimension spaces. Additionally, hierarchical clustering
134 (dendrogram-based view) is conducted to explore the grouping of samples facilitating the detection
135 of potential outliers (user exclusion if necessary).

136 **Normalization and Batch Correction**

137 Normalization corrects the systematic biases and "uninteresting" factors, ensuring that observed
138 differences between experimental conditions accurately reflect true biological variation. For
139 microarray data, quantile normalization is applied for harmonizing intensity distributions across
140 arrays. RNA-Seq data is normalized in median-of-ratios approach.

141 **Differential Gene Expression (DGE) Analysis**

142 In DGE analysis, differentially expressed genes (DEGs) are pointed out through pairwise
143 comparison between samples with two different types of conditions. DGE analysis is employed
144 through platform-specific statistical frameworks. DEGs are defined by $|LFC| > 1$ and $padj < 0.05$,
145 visualized via volcano plots annotated with gene IDs.

146 **Gene ID Conversion and Annotation**

147 Gene (DEGs) IDs (Ensembl and Entrez) resulted in DGE analysis are converted to gene symbol
148 and annotated with description using the R package the *biomaRt* [15]—incorporates gene
149 annotation data for 367 organisms (214 from Ensembl and 153 from Ensembl Plants). Gene ID
150 conversion ensures the visualization of volcano plots with respective gene symbols (if
151 available) instead of gene IDs (Ensembl or Entrez).

152 **Comparative Transcriptomics and Integrative Analysis**

153 Comparative analysis of the outcomes in DGE analysis enhances our understanding of shared and
154 unique gene expression patterns, revealing condition-specific molecular signatures. The result is
155 visualized in a venn diagram using the *ggVennDiagram* [16] package, which can efficiently handle
156 up to 7 gene sets. For analyses involving more than 7 gene lists, the results are visualized using an
157 upset plot, offering a clear and scalable representation of complex gene set intersections. Users
158 can analyze one or multiple datasets and compare the DEGs to identify intersecting genes shared
159 across different contrasts as well as unique gene sets specific to individual contrasts. Log-fold
160 metrics of both microarray and RNA-seq are integrated to generate the data for heatmap analysis.

161 **Machine Learning Analysis**

162 A suite of advanced ML frameworks was leveraged to systematically evaluate the classification
163 efficacy of the DEGs as candidate biomarkers derived from transcriptomic profiles. The predictive

164 performance of the DEGs was quantified through iterative feature optimization and cross-
165 validation, identifying genes with robust discriminatory power. This integrative pipeline bridges
166 transcriptomic discovery and computational validation, thereby ensuring the robustness of putative
167 biomarkers in complex biological matrices.

168 **Data Preprocessing, Merging, and Batch-Effect Correction**

169 Raw multi-platform datasets are subjected to a systematic preprocessing workflow encompassing
170 variance stabilization, and global normalization to attenuate platform-derived technical variability
171 and facilitate cross-study comparability. Heterogeneous datasets (e.g., RNA-Seq and microarray)
172 were integrated into a unified expression matrix, annotated with metadata columns for conditions
173 and batch identifiers. Platform-induced technical artifacts are corrected via the *ComBat* [17],
174 followed by Z-score to standardize feature distributions prior to downstream analyses.

175 **Dimensionality Reduction and Unsupervised Analysis**

176 Unsupervised learning techniques, including Principal Component Analysis (PCA), t-distributed
177 Stochastic Neighbor Embedding (t-SNE), and Uniform Manifold Approximation and Projection
178 (UMAP), are applied to the normalized dataset to reduce dimensionality and reveal intrinsic data
179 structures. Sample-to-sample relationships are further assessed by computing a Pearson correlation
180 matrix, visualized as a correlation heatmap.

181 **Feature Selection**

182 Feature selection is performed using the Recursive Feature Elimination (RFE) algorithm [18] with
183 a Random Forest classifier as the base estimator. In this analysis, the top 50% of are retained by
184 default. Finally, the dataset is reduced to retain only the selected features for subsequent modeling
185 and analysis.

186 **Model Selection, Nested Cross-Validation, Hyperparameter Tuning and Model 187 Benchmarking**

188 Six ML classifiers, linear and tree-based methods, including Logistic Regression (LR), Extra Trees
189 (ET), Random Forest (RF), XGBoost (XGB), Gradient Boosting (GB), and AdaBoost (AB), are
190 selected for benchmarking experiment using Python libraries *scikit-learn* [19] and *XGBoost* [20].
191 To evaluate model performance and extract optimal hyperparameters, two-tiered nested cross-
192 validation (CV) [21] is implemented. The 5-fold outer CV is used for model evaluation to maintain
193 original class balance. In each of the five iterations, 4-folds (80% of the data) are utilized for
194 training and hyperparameter selection; the remaining 1-fold (20%) is employed for testing. This
195 process is repeated in all the 5-folds. To further optimize the model's performance, hyperparameter
196 tuning is conducted. Within each outer training set, an inner 3-fold Stratified CV (2-fold for
197 training and 1-fold for testing) is performed for hyperparameter tuning, evaluating all
198 combinations of hyperparameter using *GridSearchCV* [20]. A baseline evaluation is also
199 performed using each classifier's default hyperparameters under the same outer CV framework.
200 For every held-out test fold, several metrics—including Accuracy (ACC), balanced accuracy
201 (BACC), precision (PREC), recall (REC), F1 score (F1), AUROC, area under the precision-recall
202 curve (AUPRC), Matthews correlation coefficient (MCC), Cohen's kappa (KAPPA), and log loss
203 (LOGLOSS)—are computed for both the default and tuned pipelines.

204 **Performance evaluation**

205 The classification performance of each model is assessed through the calculation of the AUROC
206 and AUPRC. By default, the model with the highest AUPRC and AUROC scores is selected for
207 further analysis. But users can select any of the models which aligns well with their study.

208 **Feature Importance**

209 For the selected model, feature importance is evaluated by calculating the mean decrease in
210 accuracy, which measures the decrement of a model's accuracy for the removal of individual
211 feature. Based on this metric, the top 10 features are ranked, reflecting the features that contributed
212 the most to the model's accuracy in classifying samples.

213 **Single-Gene Model Building to Determine the Best Biomarker**

214 For each of the top 10 features, a separate single-gene model is trained, and cross-validated using
215 nested CV. The evaluation metrics (e.g., AUROC, AUPRC, ACC) are calculated accordingly.
216 AUROC value indicates the ability of a feature (biomarker) to distinguish between, while the
217 AUPRC value detects the ability of a feature to detect true positive case of the target condition.
218 The single-gene model with the highest AUPRC and AUROC is considered as the best biomarker
219 for the contrast.

220 **Feature Ablation Study to Finalize Multi-Gene Model (Biomarker Algorithm)**

221 To finalize the multi gene-model, the top ten features are incrementally reduced by removing the
222 least important feature one at a time. The multi-gene model yielding the best AUPRC along with
223 other matrices is selected as the prediction model. Nested CV, hyperparameter tuning, and model
224 evaluation (e.g., AUROC, AUPRC, ACC) are performed to finalize the multi-gene model
225 (biomarker algorithm).

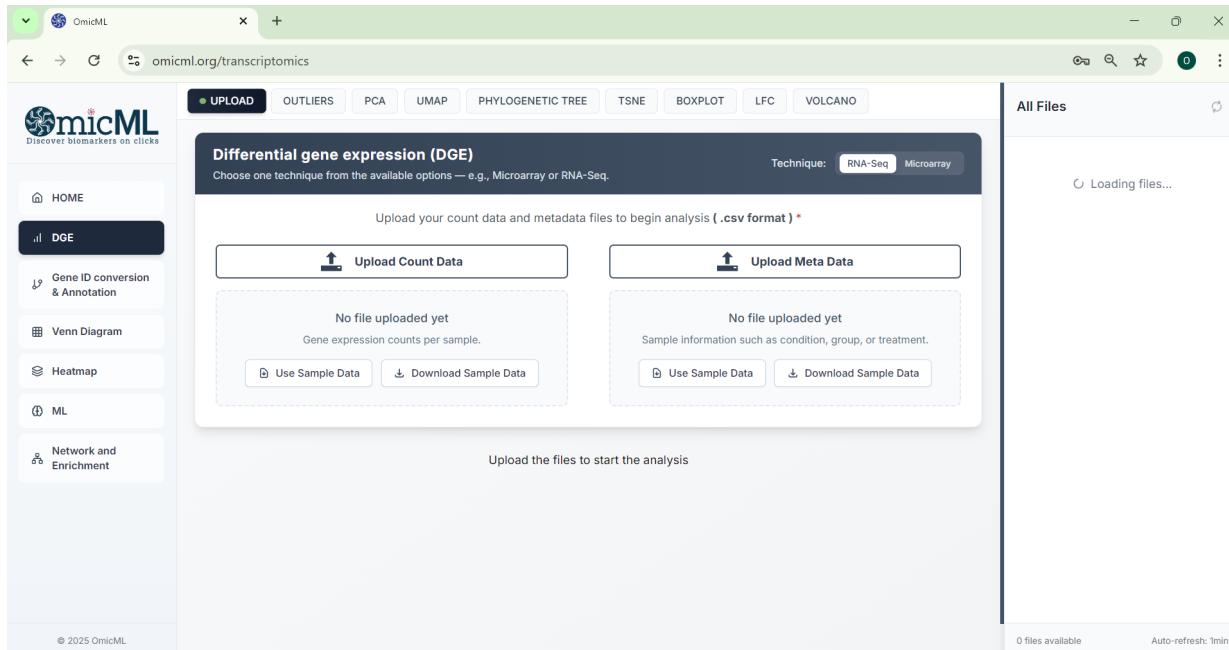
226 **Network and Functional Enrichment Framework**

227 A set of genes' symbols are uploaded and selected their corresponding organism to identify the
228 PPI network and other enrichments. Initially, PPI networks are constructed using **STRING-db** [22]
229 for the input gene symbols. Consequently, the common genes between the neighbor genes and
230 query genes (need to be uploaded by user) are selected and their interaction network is visualized.
231 The common genes are then analyzed to obtain their enrichments including Component, Process,
232 Function, WikiPathways, KEGG, Reactome, HPO, Diseases, Pfam, SMART, InterPro, TISSUES,
233 Compartments, NetworkNeighbors, PMID, Keyword with FDR<0.05. The top enriched terms are
234 visualized using **ggplot2**. Moreover, users can identify neighbor genes and upload multiple genes
235 directly to obtain their interaction networks and enrichments.

236 **Backend and Frontend Development**

237 The analytical pipelines are developed as APIs using *FastAPI* [23] to handle HTTP requests and
238 enable communication with Python. R-based computations are executed using the *rpy2*
239 (<https://rpy2.github.io/>) interface. This architecture enables seamless interoperability between
240 Python and R. The entire workflow is containerized using *docker* [24] to configure and manage
241 the server. The frontend is developed using *Next.js* (<https://nextjs.org/>) and *TypeScript*
242 (<https://www.typescriptlang.org/>) to create an interactive Graphical User Interface (GUI). *Next.js*
243 is chosen for its server-side rendering (SSR) capabilities, improved SEO, and efficient static site
244 generation (SSG). *TypeScript* is incorporated to enhance code maintainability and reliability. To
245 Deploy the full server, EC2 (<https://aws.amazon.com/ec2/>) service of AWS is utilized.

246 Results



247

248 **Figure 3** The snapshot represents the interface of the omicML web tool

249 We developed omicML, an intuitive web application integrating R and Python libraries to
250 streamline transcriptomic analysis and biomarker prediction through a six-phase workflow: (i)
251 differential gene expression analysis, (ii) gene annotation, (iii) comparative analysis, (iv) heatmap
252 analysis, (v) learning analysis and validation, and (vi) functional analysis. The platform automates
253 preprocessing and employs dimensionality reduction alongside incorporates exploratory tools
254 (histograms, volcano plots) to identify DEGs. Machine learning workflows rigorously prioritize
255 and validate biomarkers via data standardization, feature selection, benchmarking, nested cross-
256 validation, hyperparameter tuning, feature importance, single-gene model building, multi-gene
257 model (biomarker algorithm) building, model finalization, culminating in predictive models
258 selected by evaluation metrics (e.g., AUROC, AUPRC, accuracy). By uniting statistical, ML, and
259 functional analyses within an intuitive graphical user interface, omicML enables researchers
260 without rigorous programming expertise to rapidly translate expression data into validated
261 biomarkers and novel biological hypotheses with some clicks only.

262 Case Study

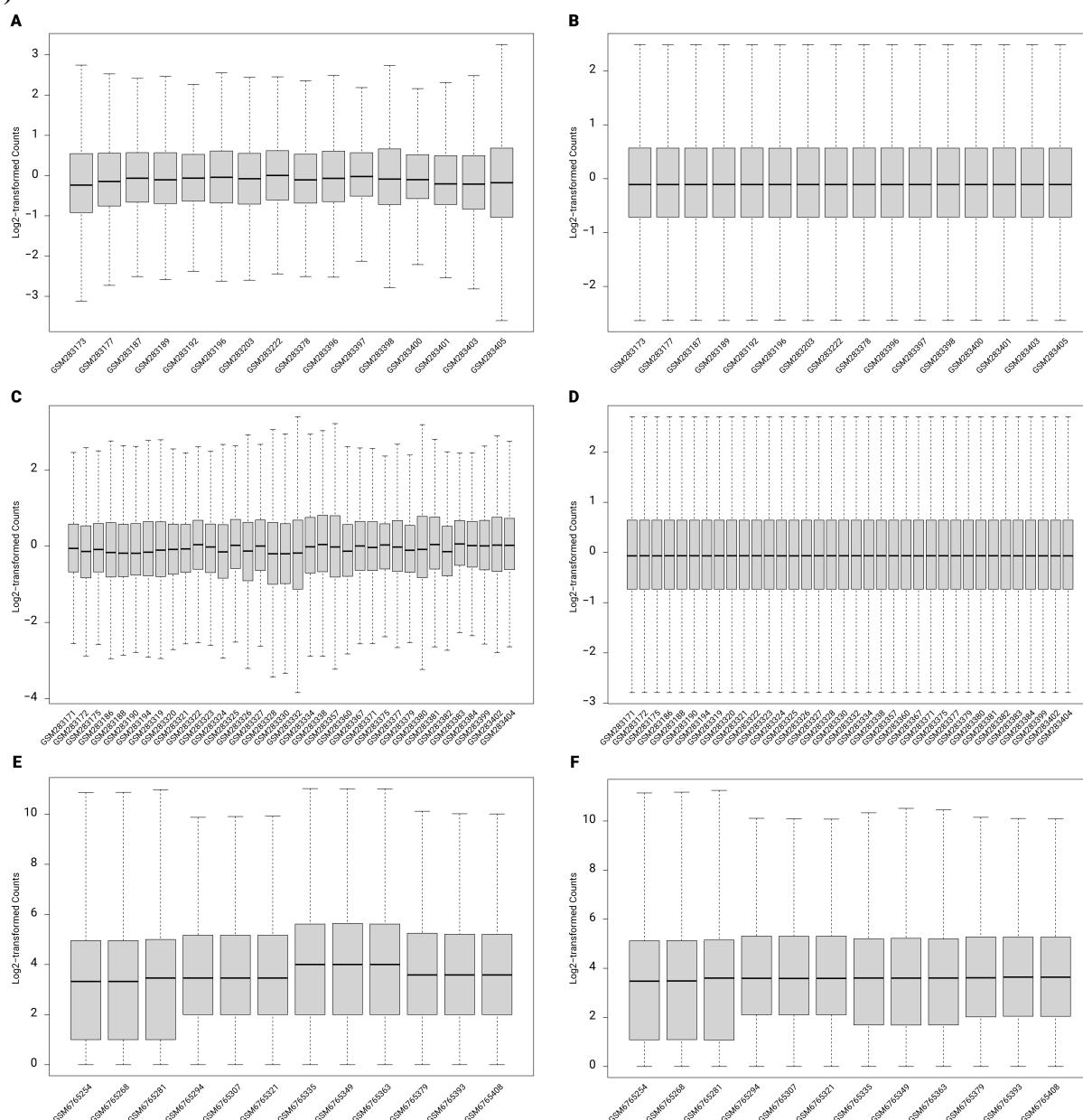
263 **markerMPXV: Upregulation of RRAD and Building of a Biomarker Algorithm** 264 **for Mpox Virus Infection**

265 To demonstrate the potential of omicML, we applied it to analyze real biological data, selected
266 from Gene Expression Omnibus (GEO), NCBI, of human cell models infected with the
267 monkeypox virus (MPXV) strains. The dataset, GSE11234 [25], comprises gene expression data
268 generated via expression profiling by array experiment of MPXV (Zaire strain)-infected dermal
269 fibroblasts and monocytes. The other dataset, GSE219036 [26], employs high-throughput
270 sequencing of transcriptomes of human keratinocytes infected with three distinct MPXV strains:

271 Clade I (historically endemic), Clade IIa (prior endemic), and Clade IIb (exclusively classified
272 during the 2022 global outbreak). Clade IIb shows distinguished expression pattern than the other
273 clades, we identified in our previous article [9].

274 **Pre-processing, and Normalization of Datasets**

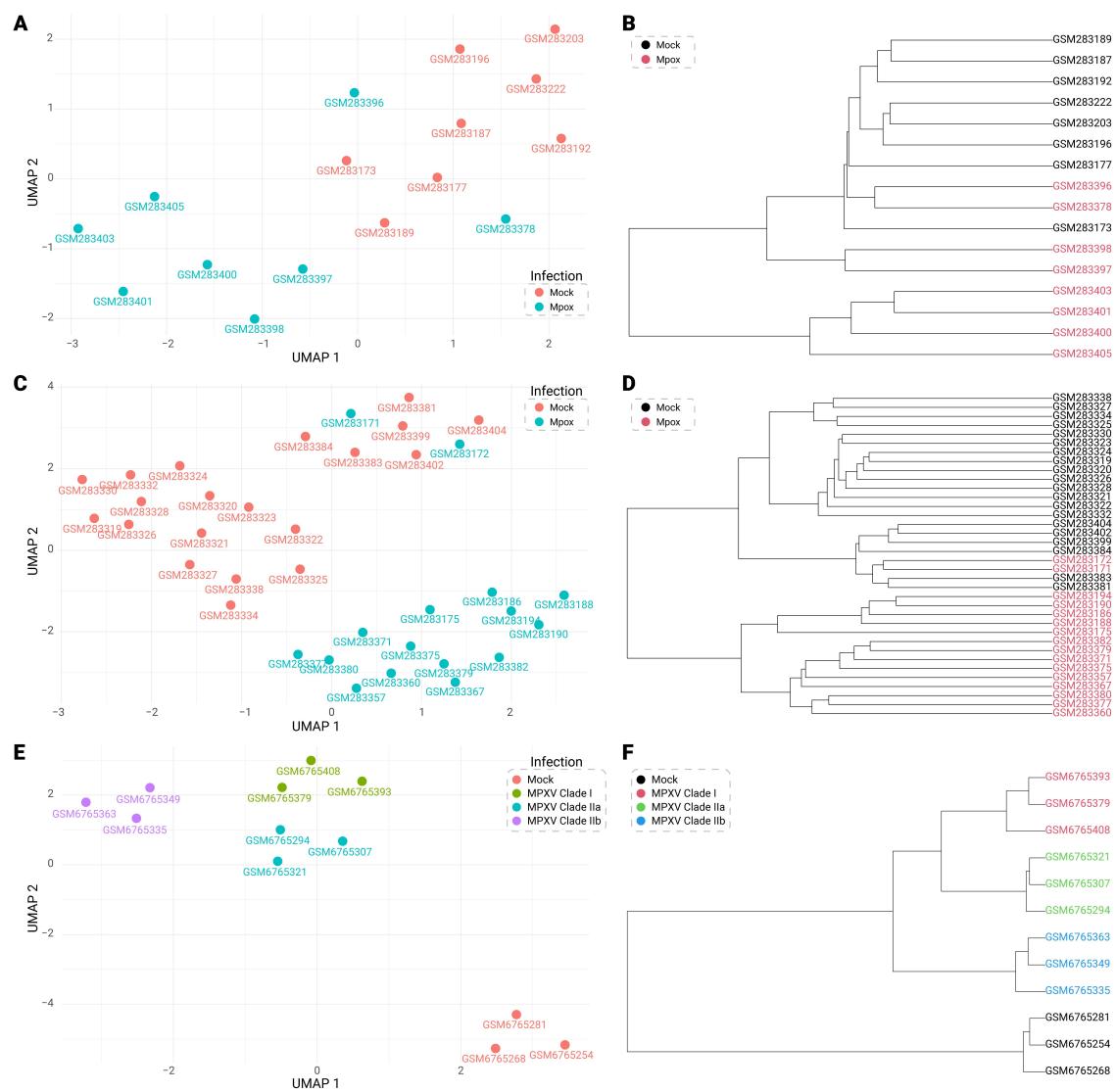
275 In the microarray dataset, data of fibroblast (16 samples) and monocyte (36 samples) cell lines
276 were normalized to mitigate technical variability and ensure cross-sample comparability, with
277 boxplots (**Figure 4A-D**) visualizing the reduction in technical biases and improved post-
278 normalization distribution consistency. For the RNA-Seq dataset, expression distributions before
279 and after normalization of keratinocyte cell-line samples were visualized in boxplots (**Figure 4E-F**).
280



281 **Figure 4 Expression pattern of datasets before and after normalization. The distribution of**
282 **expression data in the fibroblast (A-B), monocyte (C-D), and keratinocyte (E-F) cell lines are**

283 shown. Each of the boxes represents 50 percent (median) of the data. The lower boundary (Q1)
 284 marks the 25th percentile, and the upper boundary (Q3) denotes the 75th percentile, with the
 285 central line indicating the median. 50% of data points lie above the median, and 50% fall below
 286 it, offering a clear statistical summary of gene expression variability within each cell type and
 287 normalization state.

288 UMAP (**Figure 5A, C, E**) and hierarchical clustering (**Figure 5B, D, F**) analyses revealed distinct
 289 expression patterns between MPXV-infected and mock cell types. Hierarchical clustering
 290 identified four and two MPXV-infected samples in fibroblast (**Figure 5B**) and monocytes (**Figure**
 291 **5D**) respectively as outliers, underscoring technical or biological variability. Likewise,
 292 keratinocyte cell-lines (**Figure 5E-F**) exhibited high homogeneity, tight clustering with no
 293 misclassification between groups.

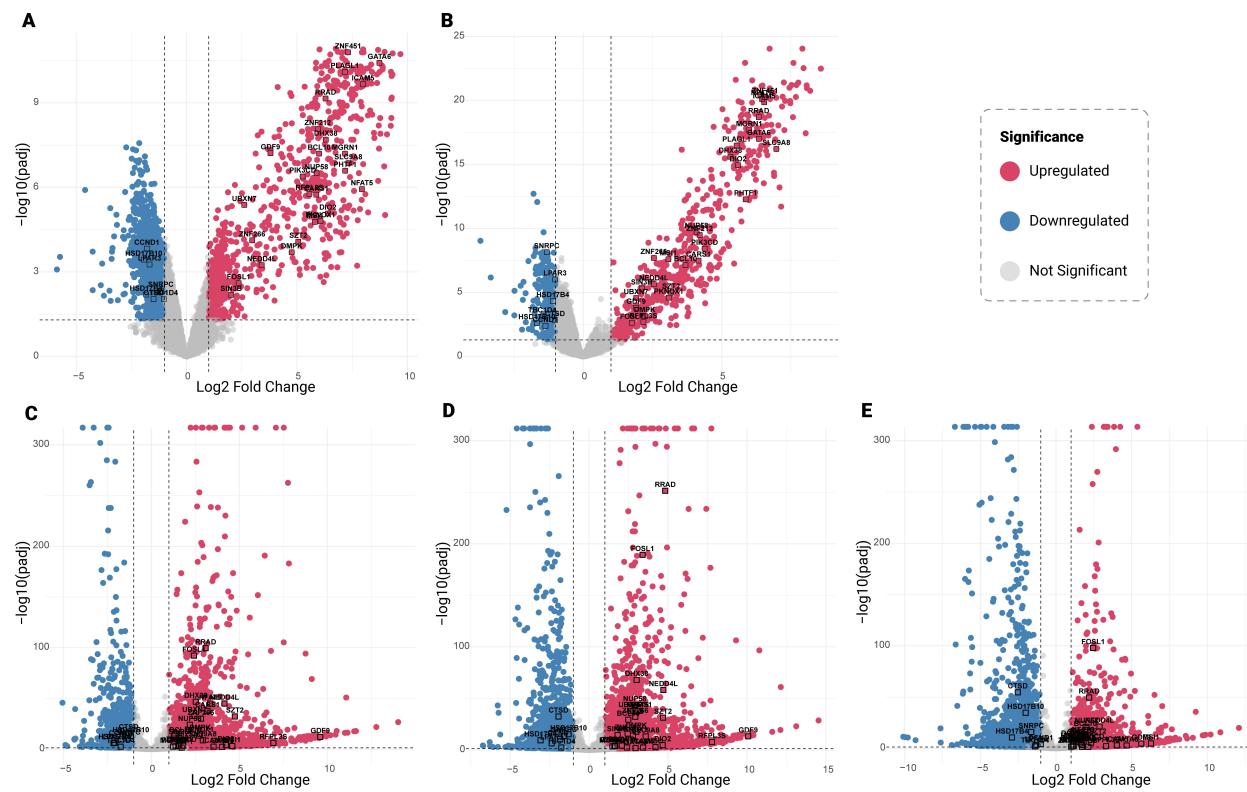


294 **Figure 5** Sample distribution patterns and outlier identification through dimensionality reduction
 295 (*UMAP, t-SNE*) and phylogenetic analysis. (A, C, and E) Clusters among samples identified from
 296 UMAPs. X-axis represents the first component (C1), which captures the highest variation in gene
 297 expression while Y-axis represents the second component (C2) which delineates the second most

298 variation in gene expression across the datasets. Each circle indicates a sample while varying
299 colors indicating different treatments on the samples. (B, D, F) Phylogenetic trees visualized
300 between the classes of samples to pinpoint the outlier samples.

301 Differential Gene Expression Analysis, Annotation, and Identification of DEGs

302 DGE analysis compared cell-line specific MPXV-infected samples to mock controls. In microarray
303 data, 5,520 significant ($\text{padj} < 0.05$) annotated genes were evident in fibroblasts (4MPXV vs
304 8Mock), while 3,548 in monocytes (14 MPXV vs. 20 mock). Consequently, 922 up- and 1,849
305 down-regulated genes in fibroblasts (Figure 6A), and 590 up- and 277 down-regulated genes in
306 monocytes (Figure 6B) were identified. Analyzing RNA-Seq data, substantial significant genes
307 and DEGs were found in keratinocytes. After annotating significant Ensembl IDs, the following
308 keratinocytes data were shown: (i) Clade I vs. mock: 2,631 upregulated and 2,212 downregulated
309 (Figure 6C), (ii) Clade IIa vs. mock: 3,108 upregulated, 2,735 downregulated (Figure 6D), and
310 (iii) Clade IIb vs. mock: 2,167 upregulated, 2,156 downregulated (Figure 6E).

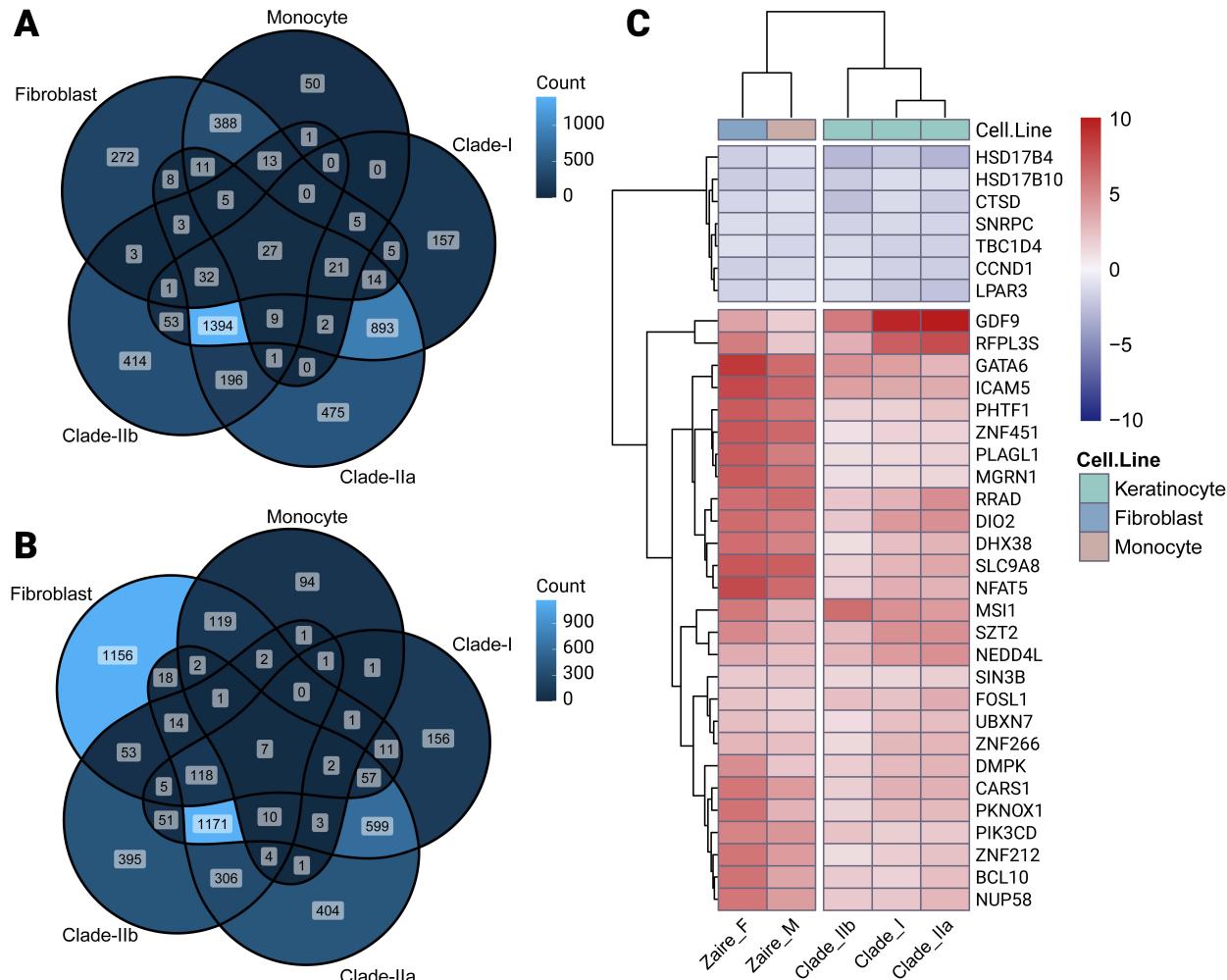


311
312 **Figure 6 Upregulated and downregulated genes.** X-axis denotes the log₂ fold change (LFC) in
313 gene expression. Positive values indicate upregulation whereas negative values indicate
314 downregulation. Y-axis shows the negative log-transformed adjusted P-value. The red circles and
315 the blue circles point out upregulated genes and downregulated genes, respectively. The horizontal
316 line represents the threshold value of FDR smaller than 0.05, ascertaining the significance of the
317 genes. However, the vertical lines represent the range of LFC less than -1 and greater than +1,
318 nominating the significant genes as differentially expressed genes.

319 Shared DEGs and Conserved Expression Patterns Across Cell Types and Clades

320 Comparative analysis of DEGs revealed a conserved transcriptional response across MPXV clades
321 (I, IIa, IIb) and cell types (fibroblasts, monocytes, and keratinocytes). Venn diagrams (Figure 7A-

322 **B)** identified 27 up- and 7 downregulated genes common to all clades' infection. Heatmap
323 visualization (**Figure 7C**) of these 34 shared DEGs further delineated clade- and cell type-specific
324 expression variability. This conserved signature underscores key genes are pivotal to MPXV
325 pathogenesis, irrespective of viral lineage or host cell type.



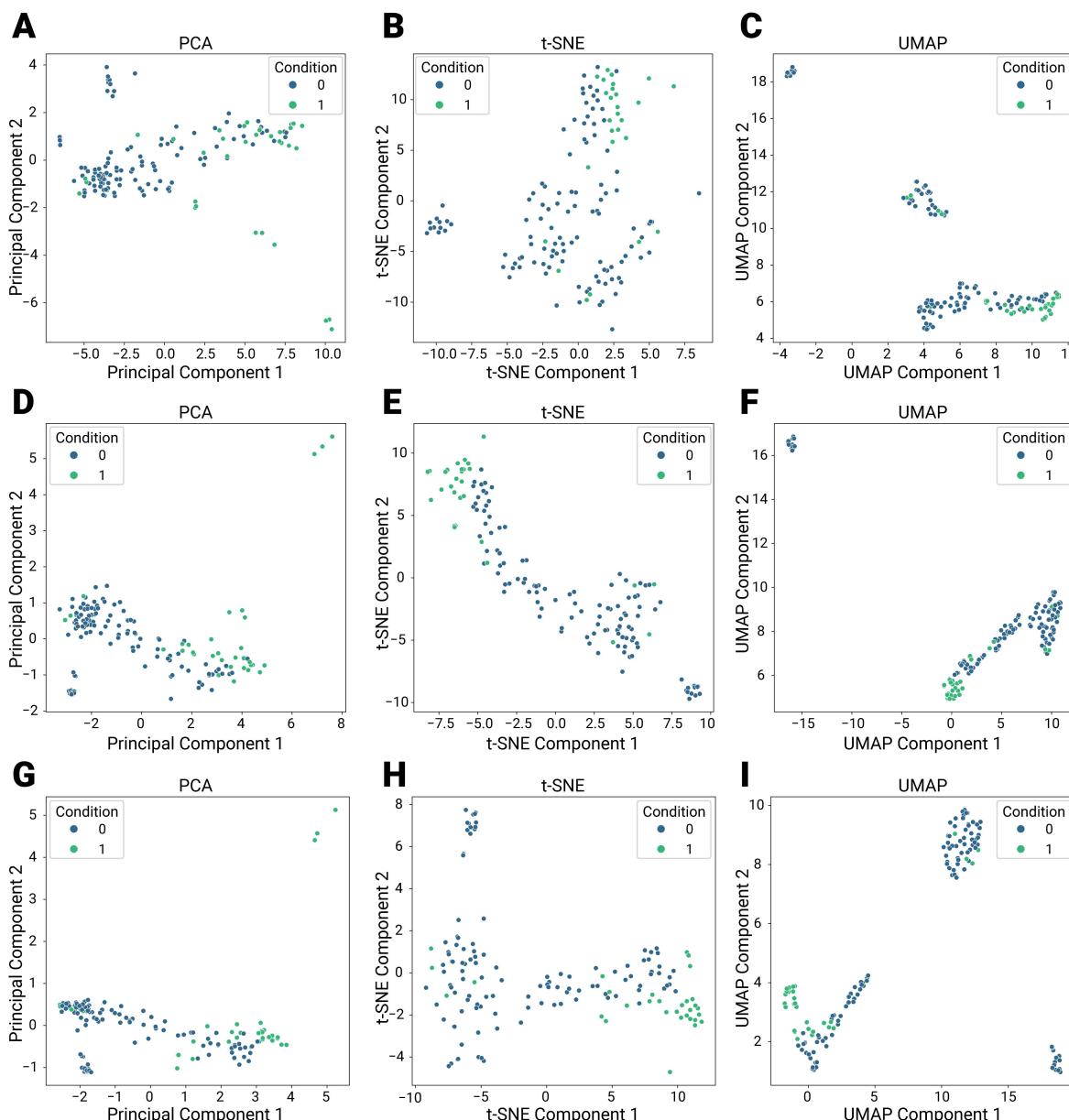
326

327 **Figure 7** Comparative analysis and Heatmap. Venn diagrams depict the commonness and
328 intersection in upregulated (A) and downregulated (B) genes. The heatmap (C) displays the
329 expression patterns of the 34 shared DEGs based on log fold-change (LFC) values across different
330 clades.

331 Integration of Machine Learning Framework

332 To enhance the reliability of biomarker discovery, automated machine learning (ML) algorithms
333 were employed to evaluate the discriminatory power of identified DEGs in distinguishing MPXV-
334 infected samples from controls. From comparative transcriptomic analyses, 34 DEGs conserved
335 across all MPXV clades (I, IIa, IIb) were selected as initial features.

336 Primary data was constructed by merging normalized RNA-Seq (24 samples) and microarray
337 datasets (124 samples), yielding 148 samples. Dimensionality reduction (PCA, t-SNE, UMAP),
338 post batch effect correction and Z-score standardization, visualized the distribution of 148 samples
339 based on the 34-feature expression profile (**Figure 8A-C**), while correlation analysis (**Figure 9**)
340 assessed interdependencies among features. Feature selection refined the 34 DEGs to 17 non-
341 redundant biomarkers using recursive feature elimination with random forests (RF-RFE). Using

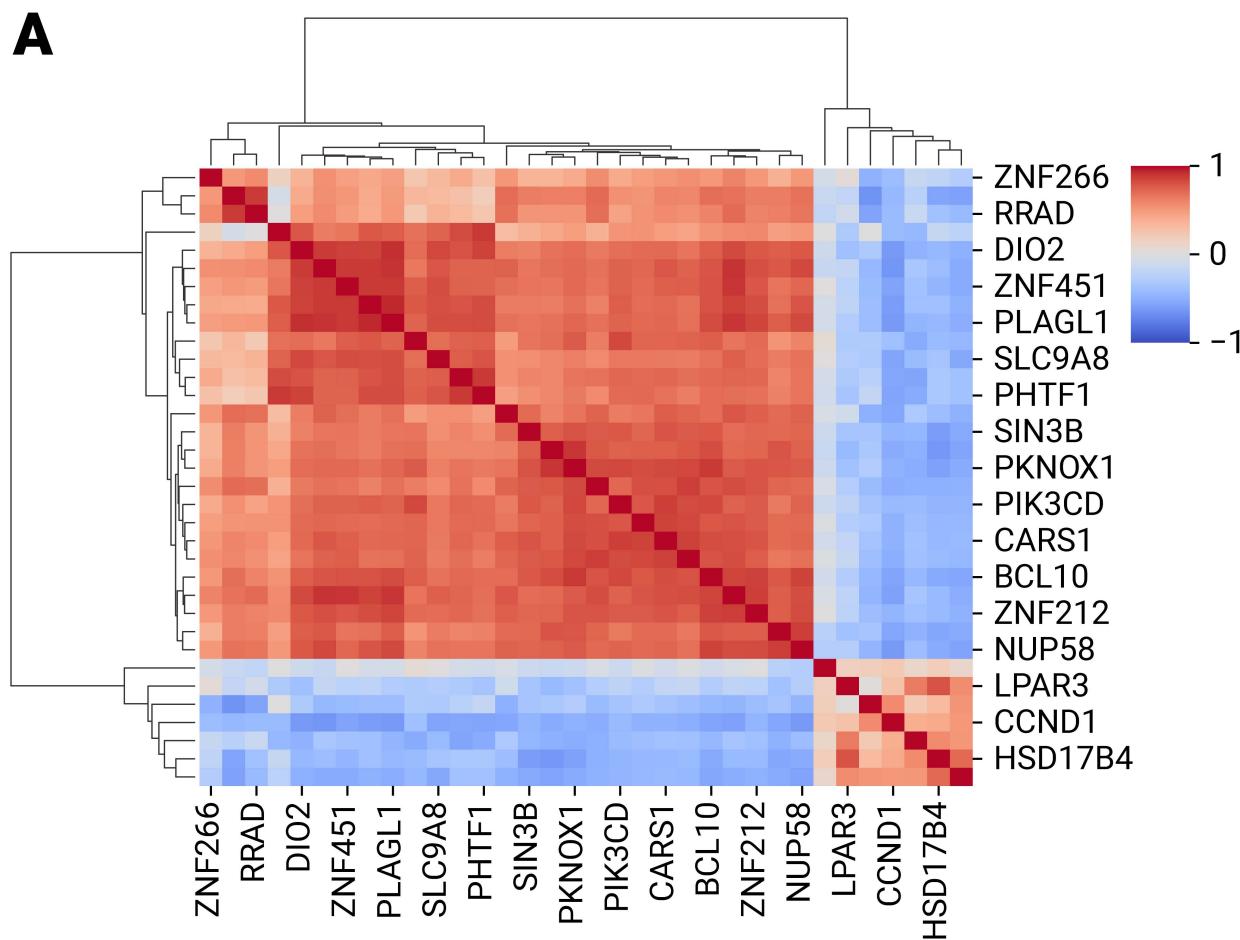


342 these 17 features, a reduced dataset was structured with the target classification (dependent
343 variable).

344 **Figure 8** *Assessment of the classification capabilities of features to cluster among the samples*
345 *according to dimension reduction. The sample distributions are depicted using PCA (A), t-SNE*
346 *(B), and UMAP (C) based on the 34 intersecting genes across infection of all clades of the Mpox*
347 *virus. Similarly, distributions of samples are illustrated utilizing PCA (D), t-SNE (E), and UMAP*
348 *(F) according to the top 10 genes ranked by importance score. Finally, discriminative power of*
349 *the six genes, selected for the final model, is highlighted through PCA (G), t-SNE (H), and*
350 *UMAP(I), plotted based on their expression levels. In the plots, each circle represents a sample*
351 *while "1" represents Mpox samples, and "0" represents samples infected by other pathogens.*

352 **Model Development and Benchmarking**

353 By combining the most important features associated with MPXV infection into a Composite
354 metric (disease indicator), a machine learning model, **markerMPXV**, was built. Twelve
355 classification algorithms were rigorously tested iteratively. Using nested cross-validation and
356 hyperparameter tuning, the Extra Trees (ET) classifier emerged as the top performer, achieving an
357 accuracy of 0.95, AUROC of 0.97, AUPRC of 0.94, and F1 score of 0.90 (**Figure 10A-B**).
358 Benchmarking results (mean ± std across outer folds) for all models are determined.



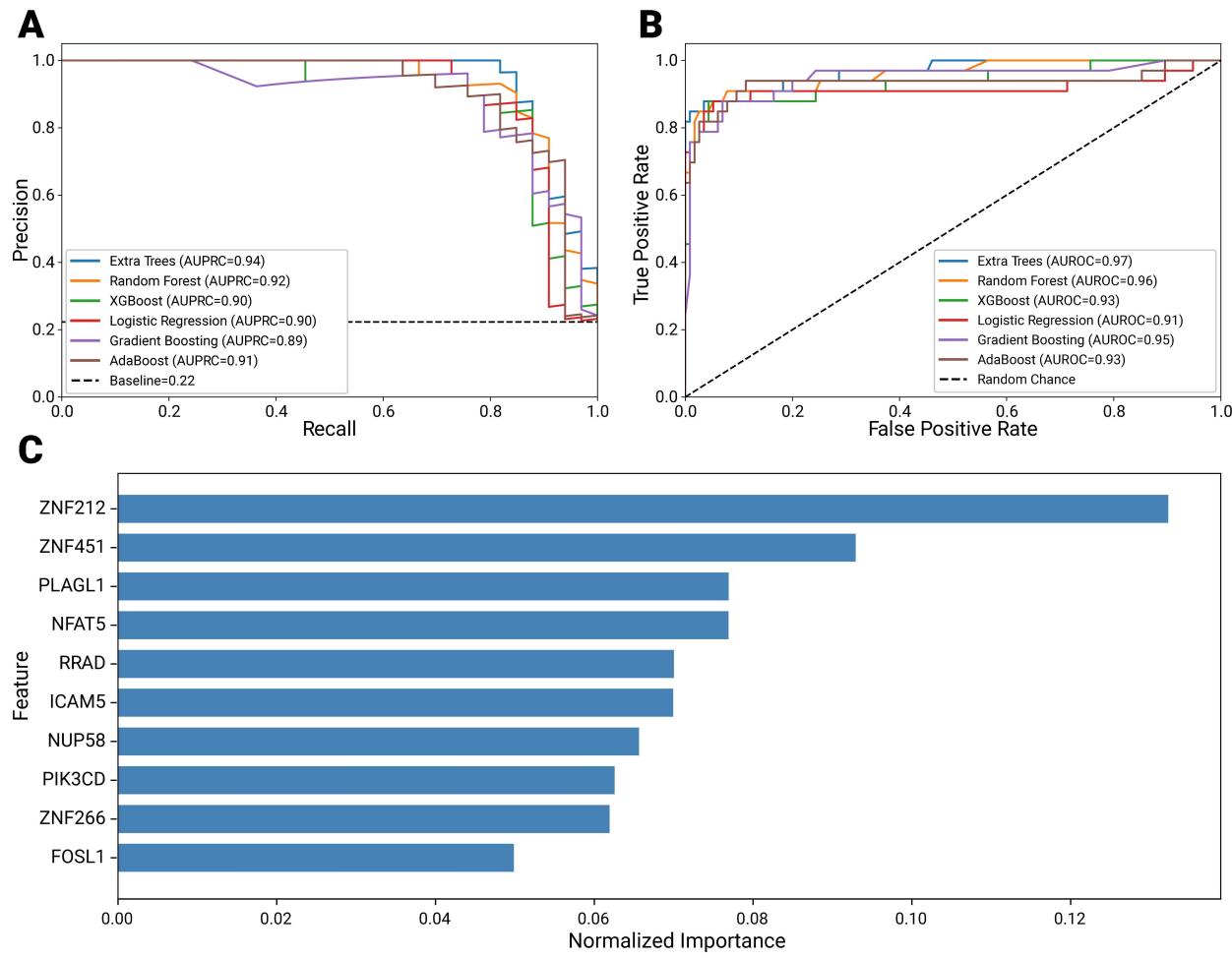
359

360

361 **Figure 9 Relationships among the features shared to infection of all MPXV clades.** The correlation
362 plot depicts the interactions among the 34 features common to all clade infections. A correlation
363 value of +1 indicates the highest positive correlation, while -1 represents the strongest negative
364 correlation between two features. Positive correlations are represented by red, whereas blue
365 signifies negative correlations.

366 **Feature Importance analysis and Single-Gene Models Facilitated Biomarker Prioritization
367 by Ranking**

368 The ET model ranked top 10 features by importance (**Figure 10C**). Single-gene Models ranked
369 *RRAD* as the top among 10 biomarkers selected by feature importance. *RRAD* achieved robust
370 performance (AUROC: 0.90; AUPRC: 0.85; F1: 0.76; accuracy: 0.91), demonstrating strong
371 discriminative power between Mpox-infected and control samples (**Figure 11A-B**).

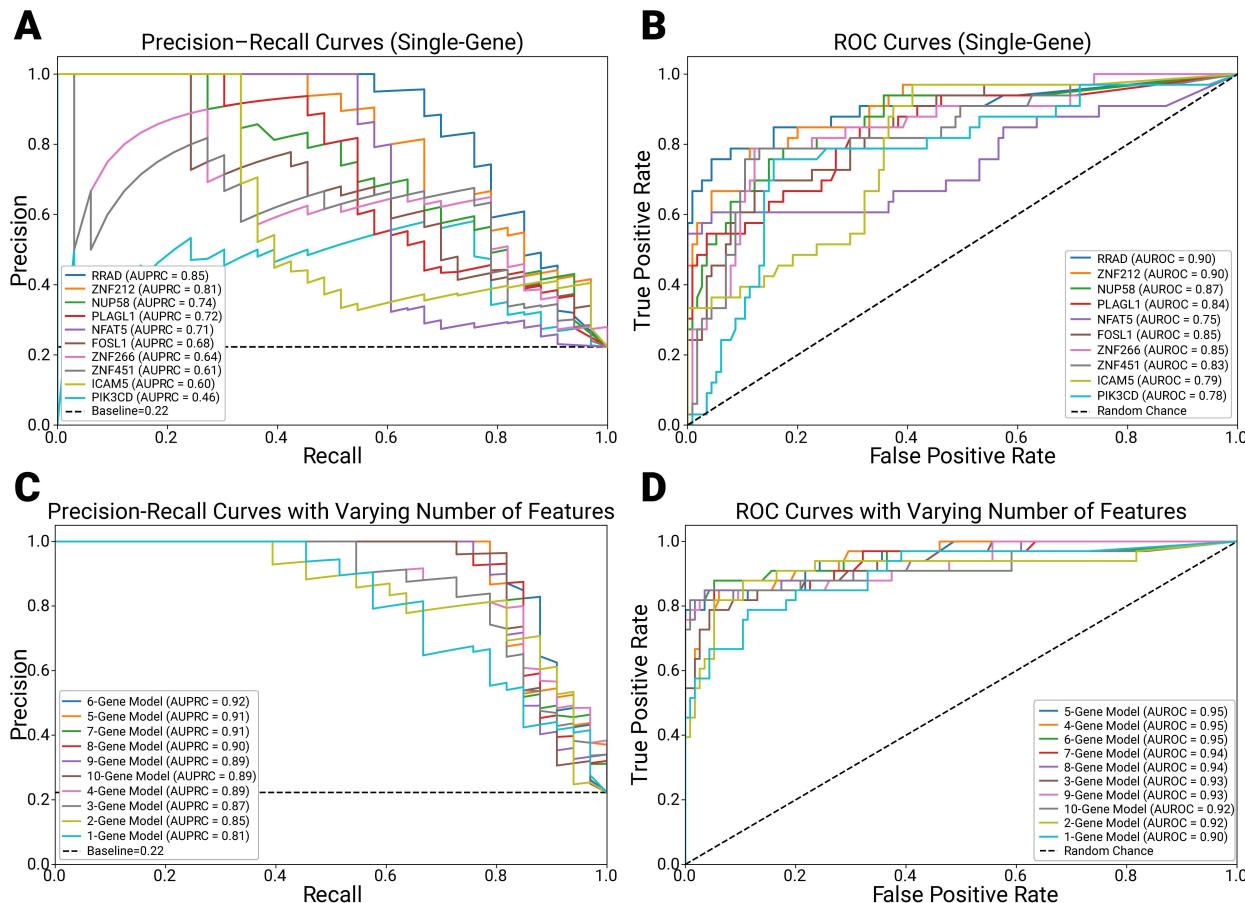


372

373 **Figure 10** The benchmarking experiment identifies the best-fitting classifier for the dataset. (A)
374 The AUPRC plot, also known as the precision-recall plot, illustrates the performance of six
375 classifiers based on nested cross-validation. (B) The AUROC plot represents the accuracy of the
376 classifiers to build models. (C) Highlights the ranking of the top 10 features based on their
377 importance scores.

378 **Feature Ablation and Optimal Multi-Gene Panel**

379 A feature ablation study revealed that combining six genes (*ZNF212*, *ZNF451*, *PLAGL1*, *NFAT5*,
380 *ICAM5*, *RRAD*) surpassed single-gene models, achieving superior performance (AUROC: 0.95;
381 AUPRC: 0.92; F1: 0.84; accuracy: 0.94) (**Figure 11 C-D**). This six-gene panel was selected for
382 final model refinement.



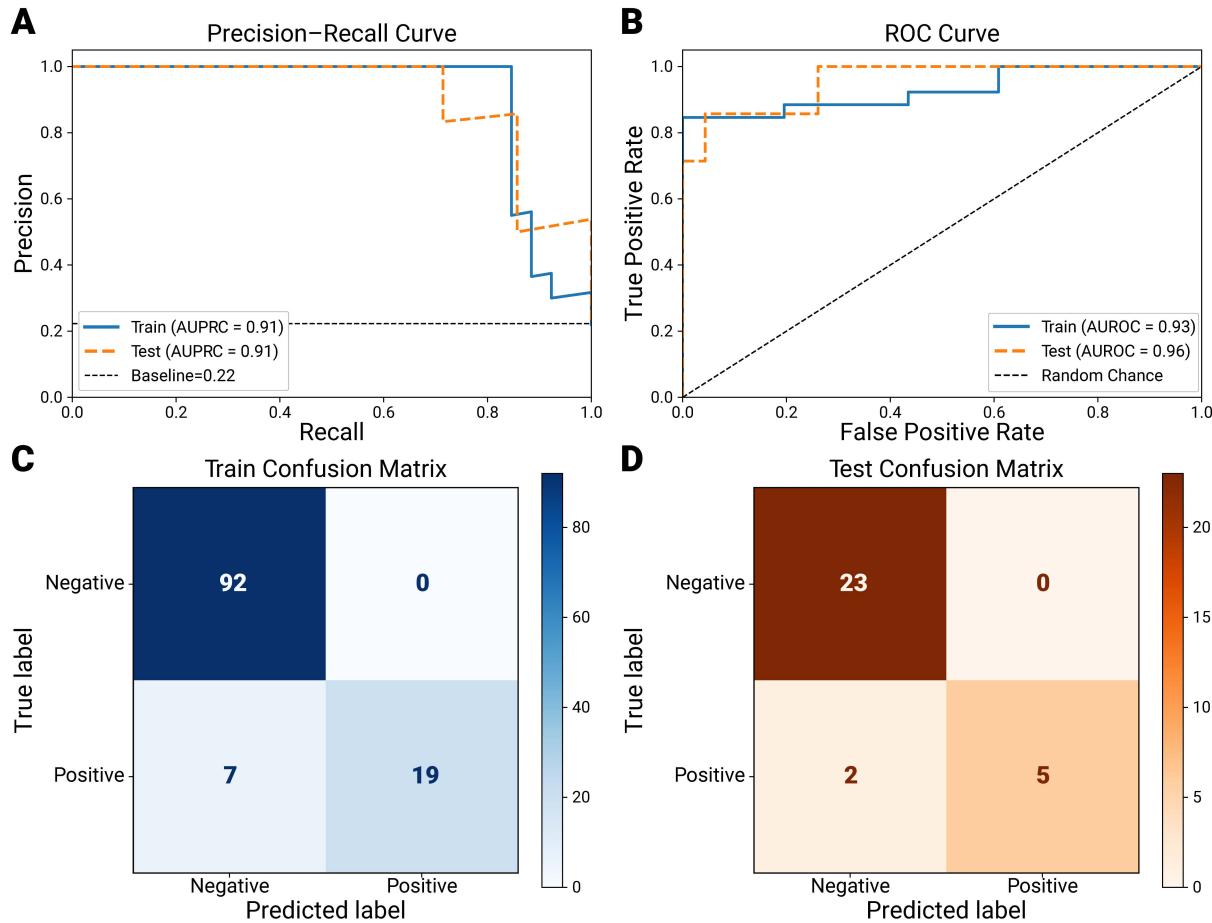
383

384 **Figure 11 Evaluation of single- and multi-gene models.** (A, B) The AUPRC and AUROC plots
385 depict the performance and class-distinguishing ability of individual features. (C, D) illustrate the
386 optimal combination of features for model building, determined by the ranking of performance
387 (PR) and AUC scores. The evaluation ranking is based on AUPRC and AUROC values.

388

389 **Final Model Performance and Validation**

390 Consequently, the predictive model was built and finalized using the expression of the six genes
391 (*ZNF212*, *ZNF451*, *PLAGL1*, *NFAT5*, *ICAM5*, *RRAD*), achieving an accuracy of 0.94 on train data
392 (AUPRC: 0.91, AUROC: 0.93, F1: 0.84) and 0.93 on test data (AUPRC 0.91, AUROC 0.96, F1
393 0.83) (**Figure 12A-B**). The confusion matrices of train and test data for the finalized models are
394 shown in **Figure 12C-D**. These results underscore the model's reliability in diagnosing Mpox
395 infection and its potential for clinical translation.



396

397 **Figure 12 Performances of the Final Model X.** (A, B) Model evaluation scores of AUPRC and
398 AUROC, representing the performance of both training and testing of the final model. (C)
399 Prediction outcomes on the test (C) and train (D) data are presented using the confusion matrix.

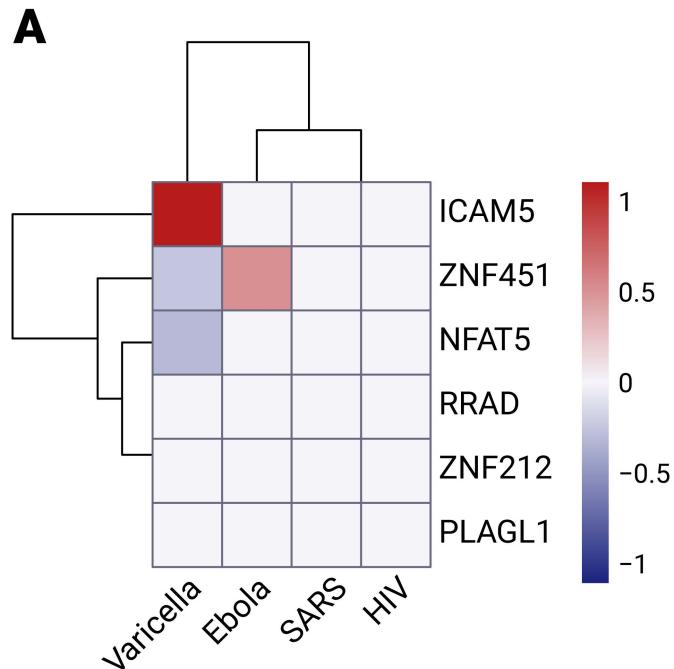
400 Network, Functional, and Enrichment Analysis

401 To elucidate the functional roles of the six predicted biomarkers in MPXV infection, omicML
402 mapped their interaction networks using the STRING database. 3,880 neighbor genes interacting
403 with the six-gene panel (*ZNF212*, *ZNF451*, *PLAGL1*, *NFAT5*, *ICAM5*, *RRAD*) were identified.
404 Each of the six genes were also analyzed individually and top 20 interacting neighbor genes'
405 network was plotted. Genes overlapping between the neighbor network and DEG lists were
406 identified with enrichment analysis revealing their involvement in key biological processes and
407 molecular functions.

408 Uniqueness of the Six-Gene Model as MPXV Biomarkers

409 To validate the hallmark signatures of Mpx infection of identified six-gene model including
410 *RRAD*, transcriptomic datasets (GSE141932 [27], GSE157103 [28], GSE184320 [29], and
411 GSE11234 [25]) of other potential viruses including varicella, ebola, HIV, and SARS-CoV-2 have
412 been analyzed via omicML. Only *ICAM5* (varicella) and *ZNF451* (Ebola) exhibited marginal
413 upregulation near significance thresholds, while other genes showed no differential expression.
414 Interestingly, none of the six genes met statistical significance ($\text{padj} < 0.05$, $|\text{LFC}| > 1$) in smallpox-
415 related varicella-zoster, SARS-CoV-2, HIV, or Ebola infections, confirming their specificity to

416 MPXV. This lack of cross-viral relevance underscores the panel's uniqueness as a robust MPXV
417 signature.



418 **Figure 13 Expression level of selected biomarkers in other virus infection in human. The heatmap**
419 **highlighting the scaled expression levels of six selected genes (ICAM5, ZNF451, NFAT5, RRAD,**
420 **ZNF212, and PLAGL1) across four viral infection conditions: Varicella, Ebola, SARS, and HIV.**
421 **The intensity of the colors reflects LFC values, where red indicates upregulation, blue indicates**
422 **downregulation, and gray indicates no significance. Cladogram was applied to both genes and**
423 **infection conditions to comprehend relationships pattern**

424 Discussions

425 The advent of high-throughput omics data and biomarker discovery techniques has resulted in a
426 fragmented and specialized ecology of isolated platforms, rendering end-to-end analysis laborious.
427 In practice, researchers frequently need to integrate disparate software for each task (e.g., DGE
428 analysis in R/Python, annotation in external databases, GO pathway analysis in another tool) via
429 manual file transfers and custom scripting, resulting in inefficiencies, errors, and reproducibility
430 issues for non-programmers [30]. Most conventional pipelines (e.g., DESeq2, edgeR, limma) and
431 annotation services (biomaRt, DAVID) are available solely as code libraries or standalone web
432 applications, whereas point-and-click platforms (DEBrowser, GenePattern, GEO2R) generally
433 cater to only a limited range of tasks, neglecting advanced procedures such as cross-study meta-
434 analyses, batch correction, dataset integration, or machine learning (ML) investigations.
435 Consequently, even standard procedures like quality control, normalization, and batch correction
436 necessitate bioinformatics assistance, thereby constraining scalability and impeding workflow
437 efficiency.

438 Moreover, workflows that use machine learning are typically absent. Traditional biomarker
439 investigations often culminate with the compilation of a list of differentially expressed candidates,
440 often evaluating them individually. Few tools offer built-in ML pipelines (feature selection, model
441 training, benchmarking) to rigorously evaluate candidate biomarkers. However, Leclercq et al.

442 indicated that current auto-ML systems are either inadequately designed for biological datasets or
443 excessively complex for individuals lacking expertise in machine learning to utilize effectively
444 [31]. In addition, the majority of machine learning technologies are inaccessible to biologists due
445 to their support for a restricted array of algorithms, the necessity for manual hyperparameter
446 optimization, or the assumption of programming proficiency [31]. Finally, inflexible, hard-coded
447 pipelines lack the modularity and graphical interfaces necessary for seamless adaptation across
448 various omics modalities or research strategies. The discipline presently needs a unified, adaptable
449 platform that seamlessly integrates (1) DGE analysis, (2) annotation, (3) cross-study comparison,
450 (4) ML-based validation, and (5) functional enrichment in an automated, user-friendly manner.
451 Bench scientists encounter obstacles at every stage of biomarker discovery [30,32]. We introduce
452 omicML, specifically engineered to address these deficiencies by offering a comprehensive,
453 graphical platform for transcriptome biomarker identification.
454 All fundamental preprocessing and differential expression procedures are consolidated into a
455 single platform, eliminating the necessity for users to transfer data across programs. Investigators
456 can define experimental groups with minimal clicks and promptly obtain differentially expressed
457 genes (DEGs) along with corresponding graphs. Integrated annotation (ID conversion, pathway
458 mapping) and comparative modules (e.g., Venn overlaps across conditions) obviate the need for
459 manual scripting or file transfers. The new approach eliminates the "silos" of disparate
460 technologies, allowing for seamless transitions of outputs from data extraction to annotation to
461 comparative modules. This immediately tackles the fragmentation problem identified in the
462 literature.
463 omicML incorporates an extensive, machine learning-driven validation suite to overcome the
464 shortcomings of traditional biomarker procedures. The platform transcends conventional methods
465 that only identify statistical connections by automating feature selection, model benchmarking
466 using nested cross-validation, and conducting ablation studies to thoroughly evaluate biomarker
467 stability and significance. omicML employs tools such as BioDiscML for comprehensive searches
468 and cross-validated classifiers, ensuring that biomarker candidates are evaluated using advanced
469 methodologies without necessitating user coding [31].
470 provided as a graphical, no-code interface designed for anyone without a background in
471 bioinformatics. Like BIOMEX, which illustrated the use of an interactive multi-omics platform
472 for laboratory researchers, the new program offers menus and wizards in lieu of command lines.
473 Non-experts may upload their data, configure parameters, and examine outcomes presented as
474 publication-quality graphs and tables. The technology automates laborious activities such as batch
475 correction and file merging, guaranteeing reproducibility without manual involvement. Bench
476 scientists can go from raw data to validated biomarker panels solely within the GUI, eliminating
477 the necessity for Python or R coding. This modular approach guarantees flexibility and
478 reproducibility, allowing pipelines to be re-executed or modified for new datasets.
479 In our case study, we conducted a comparative transcriptomic analysis to identify DEGs and
480 predictive biomarkers across multiple MPXV clades' infection because of limited therapeutic
481 options and the growing threat of a broader pandemic of monkey pox viruses. Using keratinocytes,
482 dermal fibroblast, and monocyte cell-types infected with various MPXV clades, we found a higher
483 number of DEGs in skin-derived cells compared to monocytes, a finding consistent with prior
484 observations of increased viral load in keratinocytes **Figure 5 (A-E)**. Notably, recent clades
485 appeared to elicit broader gene dysregulation compared to the older Zaire strain, with 34 DEGs
486 (27 upregulated, 7 downregulated) consistently expressed across all three cell types irrespective of
487 clade, suggesting their potential relevance in distinguishing MPXV pathogenesis **Figure 6 (A &**

488 **B).** To evaluate the diagnostic efficacy of these DEGs, we implemented machine learning models
489 that classified MPXV-infected versus control samples using integrated RNA-Seq and microarray
490 data. After batch effect correction, feature selection and benchmarking experiment, the Extra Trees
491 Classifier uncovered *RRAD* as the most potent single-gene biomarker (AUROC: 0.90; AUPRC:
492 0.85; F1: 0.76; accuracy: 0.91). Furthermore, a six-gene panel (*ZNF212*, *ZNF451*, *PLAGL1*,
493 *NFAT5*, *ICAM5*, *RRAD*) exhibited superior classification performance (AUROC: 0.95; AUPRC:
494 0.92; F1: 0.84; accuracy: 0.94), underscoring its utility for robust biomarker-based mpox detection.
495 The identified biomarkers are pivotal in orchestrating host responses to MPXV infection,
496 interacting with both upregulated and downregulated genes. Among the six key genes, five regulate
497 critical cellular processes: *ZNF212*, *ZNF451*, *NFAT5*, and *PLAGL1* govern gene expression and
498 biological pathways, while *RRAD* modulates molecular functions. The sixth gene, *ICAM5*, is
499 central to cellular adhesion. Collectively, these genes form an interconnected network influencing
500 signal transduction and immune responses, highlighting their systemic role in host-pathogen
501 interactions.

502 *RRAD* and *ICAM5* are central to immune evasion [33,34]. *RRAD* suppresses NF-κB signaling by
503 binding to its p50/p65 heterodimer, blocking inflammatory protein synthesis and cytokine
504 production [34,35]. This inhibition dampens immune activation, potentially aiding MPXV
505 survival. Notably, *RRAD* overexpression is linked to oncogenesis in skin cells and glucose
506 metabolism dysregulation, contributing to type II diabetes [36,37]. Similarly, *ICAM5*, a neuronal
507 immune modulator, is upregulated in MPXV-infected cells, impairing phagocytosis and T-cell
508 responses [33,38]. Its overexpression may suppress innate and adaptive immunity, enhancing viral
509 persistence and disease severity.

510 In contrast, *NFAT5* and *ZNF451* activate immune defenses. *NFAT5* promotes immune cell survival,
511 proliferation, and differentiation (e.g., macrophages, T-cells) while regulating NF-κB and Treg/Th
512 cell pathways. However, its overexpression risks rheumatoid arthritis and tumor progression, and
513 may stimulate viral replication [39,40]. *ZNF451* enhances immunity by inhibiting TGF-β
514 signaling, which otherwise suppresses NK cells, T-cells, and antigen-presenting cells [39]. By
515 countering TGF-β, *ZNF451* amplifies immune activation, though its role in MPXV-specific
516 responses warrants further study.

517 *PLAGL1* governs apoptosis, cell cycle control, and TP53-mediated transcription. As a tumor
518 suppressor, its overexpression regulates aberrant proliferation yet is paradoxically associated with
519 oncogenesis [41,42]. In MPXV infection, *PLAGL1*-induced apoptosis may restrict viral
520 dissemination, while its multiple functions in cancer underscore context-dependent effects on host-
521 pathogen interactions.

522 **Limitations**

523 While omicML currently provides a comprehensive GUI-driven pipeline for transcriptomics-based
524 biomarker discovery, more sophisticated functionalities are yet to be integrated in next versions
525 (omicML 2.0). omicML is presently tailored to bulk transcriptomic data and does not include
526 network-based or clinical modeling modules. In practice, many biomarker studies rely on gene co-
527 expression network analysis and survival modeling to uncover complex patterns and clinical
528 relevance, so these capabilities are absent in the current version.

529 **Conclusions**

530 omicML represents a novel end-to-end framework for biomarker discovery by integrating many
531 analytical steps into a cohesive, user-friendly platform. Its graphical interface guides users from
532 data upload to normalization, differential expression, annotation, and machine-learning evaluation,

533 therefore obviating the necessity for complex coding. This integrated pipeline unifies predictive
534 modelling and biomarker selection into a cohesive approach. omicML democratizes access to
535 complicated analyses by offering a GUI-based approach, allowing physicians and bench
536 researchers without required programming abilities to execute advanced transcriptomics
537 procedures. By reducing technical obstacles and offering a comprehensive, cohesive toolkit,
538 omicML is positioned to significantly influence translational bioinformatics and the advancement
539 of clinically pertinent molecular diagnostics.

540 Besides, omicML addressed the urgent need for mpox biomarkers and identified a six-gene model
541 (*ZNF212*, *ZNF451*, *PLAGL1*, *NFAT5*, *ICAM5*, and *RRAD*) achieving exceptional diagnostic
542 accuracy (AUROC: 0.95; AUPRC: 0.92) out of 34 clade-independent DEGs. This demonstrates
543 omicML's capacity to bridge transcriptomic insights with ML-driven validation, accelerating
544 biomarker discovery for emerging pathogens and beyond.

545 Abbreviations

546 GUI: Graphical User Interface
547 DGE: Differential gene expression
548 DEGs: Differentially expressed genes
549 LFC: Log2 fold change
550 FDR: False Discovery Rate
551 Padj: P-adjusted Value
552 PCA: Principal Component Analysis
553 UMAP: Uniform Manifold Approximation and Projection
554 t-SNE: t-distributed stochastic neighbor embedding
555 ML: Machine Learning
556 LR: Logistic Regression
557 ET: Extra Trees
558 RF: Random Forest
559 XGB: XGBoost
560 GB: Gradient Boosting
561 AB: AdaBoost
562 ACC: Accuracy
563 BACC: Balanced Accuracy
564 PREC: Precision
565 REC: Recall
566 F1: F1 Score
567 AUROC: Area Under the Receiver Operating Curve
568 AUPRC: Area Under the Precision-Recall Curve
569 MCC: Matthews Correlation Coefficient
570 KAPPA: Cohen's Kappa
571 LOGLOSS: Log Loss

572 Mpox: Monkey Pox
573 MPXV: Monkey Pox Virus
574 GEO: Gene Expression Omnibus

575 **Acknowledgments**

576 To the straightness of Kilo-Road, the sublimeness of Gazi Kalur Tila, and the symphonious rainfall
577 in SUST Campus.

578 **Author Contributions**

579 **Joy Prokash Debnath:** Methodology, Software, Formal analysis, Data Curation, Visualization,
580 Investigation, Validation, Writing – original draft
581 **Kabir Hossen:** Methodology, Software, Formal analysis, Data Curation, Visualization,
582 Investigation, Validation, Writing – original draft
583 **Md. Sayem Khandaker:** Methodology, Software, Formal analysis, Data Curation, Visualization,
584 Investigation, Validation, Writing – original draft
585 **Shawon Majid:** Software, Data Curation, Methodology, Validation, Investigation
586 **Md Mehrajul Islam:** Software, Data Curation, Methodology, Visualization, Investigation
587 **Siam Arefin:** Software, Investigation
588 **Preonath Chondrow Dev:** Conceptualization, Project administration, Resources, Supervision,
589 Writing – review and editing
590 **Saifuddin Sarkar:** Conceptualization, Project administration, Resources, Supervision, Writing –
591 review and editing
592 **Tanvir Hossain:** Conceptualization, Project administration, Resources, Supervision, Writing –
593 review and editing
594

595 **Funding**

596 This study was not supported by any funding.
597

598 **Data availability statement**

599 The datasets analyzed in the current study are available in the GEO repository.
600 GSE219036: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE219036>
601 GSE11234: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE11234>
602 GSE141932: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE141932>
603 GSE157103: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE157103>
604 GSE184320: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE184320>
605

606 **Declarations**

607 **Ethics approval and consent to participate**

608 Not applicable.

609 **Consent for publication**

610 Not applicable.

611 **Competing interests**

612 The authors declare no competing interests.

613

614 **References**

- 615 1. Ge, S. X., Son, E. W. & Yao, R. iDEP: An integrated web application for differential
616 expression and pathway analysis of RNA-Seq data. *BMC Bioinformatics* **19**, 1–24 (2018).
- 617 2. Nelson, J. W., Sklenar, J., Barnes, A. P. & Minnier, J. The START App: a web-based
618 RNAseq analysis and visualization resource. *Bioinformatics* **33**, 447–449 (2017).
- 619 3. David R. Powell. Degust: interactive RNA-seq analysis, DOI: 10.5281/zenodo.3258932.
- 620 4. ShinyNGS. <https://github.com/pinin4fjords/shinyngs>. Accessed on 25 Oct (2025).
- 621 5. Afgan, E. *et al.* The Galaxy platform for accessible, reproducible and collaborative
622 biomedical analyses: 2018 update. *Nucleic Acids Res* **46**, W537–W544 (2018).
- 623 6. Monier, B. *et al.* IRIS-EDA: An integrated RNA-Seq interpretation system for gene
624 expression data analysis. *PLoS Comput Biol* **15**, e1006792 (2019).
- 625 7. Kleverov, M. *et al.* Phantasus, a web application for visual and interactive gene expression
626 analysis. *Elife* **13**, (2024).
- 627 8. GEO2R - GEO - NCBI. <https://www.ncbi.nlm.nih.gov/geo/geo2r/>. Accessed on 25 Oct
628 (2025).
- 629 9. Debnath, J. P. *et al.* Identification of potential biomarkers for 2022 Mpox virus infection: a
630 transcriptomic network analysis and machine learning approach. *Sci Rep* **15**, 1–15 (2025).
- 631 10. Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene
632 expression and hybridization array data repository. *Nucleic Acids Res* **30**, 207–210 (2002).
- 633 11. Athar, A. *et al.* ArrayExpress update – from bulk to single-cell expression data. *Nucleic
634 Acids Res* **47**, D711–D715 (2019).
- 635 12. Tomczak, K., Czerwińska, P. & Wiznerowicz, M. The Cancer Genome Atlas (TCGA): an
636 immeasurable source of knowledge. *Contemp Oncol (Pozn)* **19**, A68–A77 (2015).
- 637 13. Leinonen, R., Sugawara, H. & Shumway, M. The sequence read archive. *Nucleic Acids
638 Res* **39**, (2011).
- 639 14. Langfelder, P. & Horvath, S. WGCNA: An R package for weighted correlation network
640 analysis. *BMC Bioinformatics* **9**, 1–13 (2008).

- 641 15. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the
642 integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat Protoc* **4**,
643 1184–1191 (2009).
- 644 16. Gao, C. H., Yu, G. & Cai, P. ggVennDiagram: An Intuitive, Easy-to-Use, and Highly
645 Customizable R Package to Generate Venn Diagram. *Front Genet* **12**, 706907 (2021).
- 646 17. Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. & Storey, J. D. The sva package for
647 removing batch effects and other unwanted variation in high-throughput experiments.
648 *Bioinformatics* **28**, 882–883 (2012).
- 649 18. Kuhn, M. Building Predictive Models in R Using the caret Package. *J Stat Softw* **28**, 1–26
650 (2008).
- 651 19. Pedregosa FABIANPEDREGOSA, F. *et al.* Scikit-learn: Machine Learning in Python. *The
652 Journal of Machine Learning Research* **12**, 2825–2830 (2011).
- 653 20. Chen, T. & Guestrin, C. XGBoost: A scalable tree boosting system. *Proceedings of the
654 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* **13-
655 17-August-2016**, 785–794 (2016).
- 656 21. Varma, S. & Simon, R. Bias in error estimation when using cross-validation for model
657 selection. *BMC Bioinformatics* **7**, 1–8 (2006).
- 658 22. Szklarczyk, D. *et al.* The STRING database in 2021: customizable protein–protein
659 networks, and functional characterization of user-uploaded gene/measurement sets.
660 *Nucleic Acids Res* **49**, D605–D612 (2021).
- 661 23. Ramírez S. FastAPI: high performance, easy to learn, fast to code, ready for production.
- 662 24. Merkel D. Docker: lightweight linux containers for consistent development and
663 deployment. *Linux journal* (2014).
- 664 25. Rubins, K. H. *et al.* Comparative Analysis of Viral Gene Expression Programs during
665 Poxvirus Infection: A Transcriptional Map of the Vaccinia and Monkeypox Genomes.
666 *PLoS One* **3**, e2628 (2008).
- 667 26. Watanabe, Y. *et al.* Virological characterization of the 2022 outbreak-causing monkeypox
668 virus using human keratinocytes and colon organoids. *J Med Virol* **95**, e28827 (2023).
- 669 27. Wu, S. *et al.* Transcriptome Analysis Reveals the Role of Cellular Calcium Disorder in
670 Varicella Zoster Virus-Induced Post-Herpetic Neuralgia. *Front Mol Neurosci* **14**, 665931
671 (2021).
- 672 28. Overmyer, K. A. *et al.* Large-Scale Multi-omic Analysis of COVID-19 Severity. *Cell Syst*
673 **12**, 23-40.e7 (2021).

- 674 29. Saluzzo, S. *et al.* Delayed antiretroviral therapy in HIV-infected individuals leads to
675 irreversible depletion of skin- and mucosa-resident memory T cells. *Immunity* **54**, 2842-
676 2858.e5 (2021).
- 677 30. Narayanan, R., DeGroat, W., Mendhe, D., Abdelhalim, H. & Ahmed, Z. IntelliGenes:
678 Interactive and user-friendly multimodal AI/ML application for biomarker discovery and
679 predictive medicine. *Biol Methods Protoc* **9**, (2024).
- 680 31. Leclercq, M. *et al.* Large-scale automatic feature selection for biomarker discovery in
681 high-dimensional omics data. *Front Genet* **10**, 449967 (2019).
- 682 32. Taverna, F. *et al.* BIOMEX: an interactive workflow for (single cell) omics data
683 interpretation and visualization. *Nucleic Acids Res* **48**, W385–W394 (2020).
- 684 33. Gahmberg, C. G., Tian, L., Ning, L. & Nyman-Huttunen, H. ICAM-5--a novel two-
685 faceted adhesion molecule in the mammalian brain. *Immunol Lett* **117**, 131–135 (2008).
- 686 34. Hsiao, B. Y., Chang, T. K., Wu, I. T. & Chen, M. Y. Rad GTPase inhibits the NF κ B
687 pathway through interacting with RelA/p65 to impede its DNA binding and target gene
688 transactivation. *Cell Signal* **26**, 1437–1444 (2014).
- 689 35. Barnes, P. J. & Karin, M. Nuclear factor-kappaB: a pivotal transcription factor in chronic
690 inflammatory diseases. *N Engl J Med* **336**, 1066–1071 (1997).
- 691 36. Zhang, C. *et al.* Tumor suppressor p53 negatively regulates glycolysis stimulated by
692 hypoxia through its target RRAD. *Oncotarget* **5**, 5535–5546 (2014).
- 693 37. Sun, Z. *et al.* Friend or Foe: Regulation, Downstream Effectors of RRAD in Cancer.
694 *Biomolecules* **13**, (2023).
- 695 38. Paetau, S., Rolova, T., Ning, L. & Gahmberg, C. G. Neuronal ICAM-5 Inhibits Microglia
696 Adhesion and Phagocytosis and Promotes an Anti-inflammatory Response in LPS
697 Stimulated Microglia. *Front Mol Neurosci* **10**, (2017).
- 698 39. Feng, Y. *et al.* Zinc finger protein 451 is a novel Smad corepressor in transforming growth
699 factor- β signaling. *J Biol Chem* **289**, 2072–2083 (2014).
- 700 40. Lee, N., Kim, D. & Kim, W. U. Role of NFAT5 in the Immune System and Pathogenesis
701 of Autoimmune Diseases. *Front Immunol* **10**, 270 (2019).
- 702 41. Vega-Benedetti, A. F. *et al.* PLAGL1: an important player in diverse pathological
703 processes. *J Appl Genet* **58**, 71–78 (2017).
- 704 42. Keck, M. K. *et al.* Amplification of the PLAG-family genes-PLAGL1 and PLAGL2-is a
705 key feature of the novel tumor type CNS embryonal tumor with PLAGL amplification.
706 *Acta Neuropathol* **145**, 49–69 (2023).