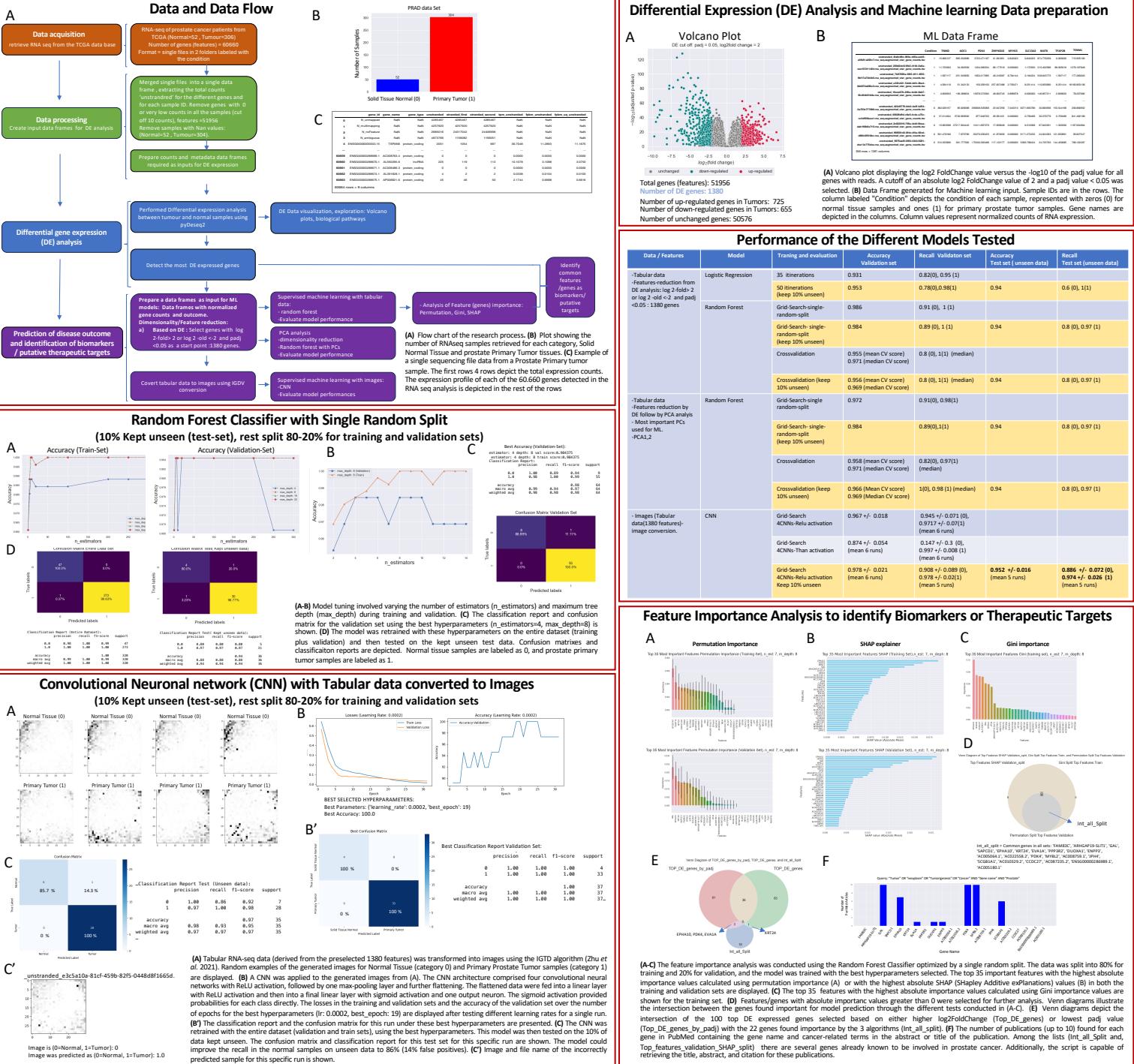


# Prostate Cancer Prediction and Biomarker Identification Using Machine Learning and Deep Learning Algorithms on Transcriptome Data from The Cancer Genome Atlas (TCGA) Database

## ABSTRACT

The search for novel RNA biomarkers and innovative methods to identify cancerous tissues can significantly advance the development of RNA-based diagnostic and therapeutic strategies, leading to more effective and personalized approaches for cancer treatment and management. In this project, we investigated the feasibility of predicting or diagnosing prostate cancer, which ranks among the most prevalent cancers in the male population, by applying machine learning (ML) and convolutional neural network (CNN) algorithms to gene expression data of normal and primary tumor prostate gland samples. Genes/features used as input for ML were reduced by preselecting the most differentially expressed (DE) genes between cancer and normal samples. Machine learning algorithms (logistic regression, random forest, random forest on the most important principal components (PCs)) were applied to predict cancer outcomes using RNA expression data on the selected genes. A CNN was also tested on the same tabular data converted to images. Moreover, through an examination of the disturbed gene expression patterns in prostate cancer samples and the genes important for predicting cancer versus normal tissue outcomes by machine learning, we also set up to discover putative novel RNA biomarkers for prostate cancer.

## RESULTS



## CONCLUSIONS AND OUTLOOK

- Machine learning applied to RNAseq data has successfully predicted prostate cancer outcomes.
- Random forest outperformed logistic regression, enhancing recall for under-represented normal tissue.
- PCA feature reduction was effective; 2 PCs matched RF performance with 1,380 features.
- Transforming tabular data into images for a CNN improved model performance, particularly recall for the underrepresented category; visualization provided insights not easily discernible from 1,380 tabular features.
- Main issues: unbalanced, limited data and no accessible independent dataset for final validation. While models showed high accuracy, they struggled with underrepresented normal samples but excelled in classifying tumor samples.

- Stratified splitting improved Random Forest performance on underrepresented samples. Further enhancement of CNN could be achieved with stratification and cross-validation.
- Generating synthetic RNA-seq data and utilizing independent datasets is recommended.
- Optimizing Random Forest by adjusting hyperparameters (`min_samples_split`, `min_samples_leaf`, `max_features`) is advised to boost stability, reduce overfitting, and enhance performance..

(A-C) The feature importance analysis was conducted using the Random Forest Classifier optimized by a single random split. The data was split into 80% for training and 20% for validation, and the model was trained with the best hyperparameters selected. The top 35 important features with the highest absolute importance values calculated using permutation importance (A) and the highest absolute SHAP (Shapley Additive exPlanation) values (B) in both the training and validation sets are shown. (D) The top 25 features with the highest absolute importance values calculated using Gini importance (C) are shown for the training set. (E) Features/genes with absolute importance values greater than 0 were selected for further analysis. Venn diagrams illustrate the intersection of the 100 top DE expressed genes selected based on either higher log2FoldChange (`Top_DE_genes`) or lower padj value (`Top_DE_genes_by_padj`) with the 22 genes found importance by the 3 algorithms (`Int_all_Split`). (F) The number of publications (up to 10) found for each gene in PubMed containing the gene name and cancer-related terms in the abstract or title of the publication. Among the lists (`Int_all_Split` and `Top_features_validation_SHAP_split`) there are several genes already known to be involved in prostate cancer. Additionally, the script is capable of retrieving the title, abstract, and citation for these publications.