

Digital Discovery

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: M. A. Ali, H. Sarker, M. Kamrun, H. Sheikh, B. Shifa, S. Ahmed, T. Islam, S. Banik and N. Kumar, *Digital Discovery*, 2025, DOI: 10.1039/D6DD00045B.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

Identification of Multi-Transcriptomic Prognostic Biomarkers to Explore Natural Therapeutics for Lung Cancer integrating Machine Learning

Md Ahad Ali^{1,2}¶*, Hrididhi Sarker^{1,3}¶, Marguba Kamrun^{1,4}, Humaira Sheikh^{2,5}, Bilkis Shifa^{1,6}, Siam Ahmed¹, Tarikul Islam⁷, Sujoy Banik⁸, Neeraj Kumar⁹

¹Computational Chemistry and Drug Design Division, Panacea Research Center, Rajshahi 6206, Bangladesh.

⁵Department of Chemistry, University of Rajshahi, Rajshahi 6205, Bangladesh.

³Department of Biochemistry and Molecular Biology, University of Rajshahi, Rajshahi 6205, Bangladesh.

⁴Department of Chemistry and Biochemistry, University of Oklahoma Norman, OK 73019, USA.

⁵Department of Chemistry, Gopalganj Science and Technology University, Gopalganj 8100, Bangladesh.

⁶Department of Biochemistry and Molecular Biology, University of Dhaka, Dhaka 1000, Bangladesh.

⁷Department of Chemistry, University of Barishal, Barishal 8254, Bangladesh.

⁸Department of Computer Science & Engineering, Rajshahi University of Engineering & Technology, Rajshahi, Bangladesh.

⁹Department of Pharmaceutical Chemistry, Bhupal Nobles' College of Pharmacy, Udaipur 313001, Rajasthan, India.

¶These authors contributed equally first author to this work.

*Correspondence: ahad.chembd@gmail.com (M. A. A.)

Abstract:

Lung cancer remains the leading cause of cancer-related mortality worldwide, underscoring the urgent need for novel therapeutic strategies. Cyclin-dependent kinase 1 (CDK1), a central cell-cycle regulator, has emerged as an oncogenic driver and potential target in lung adenocarcinoma. This study aimed to integrate transcriptomics, machine learning (ML), and advanced in silico approaches to identify natural product-derived potential inhibitors targeting CDK1. To identify robust differentially expressed genes, first we analyzed four different datasets (GSE19804, GSE10072, GSE18842, and GSE10799). Protein–protein interaction network and topological analysis highlighted CDK1 as a primary key hub gene (pKHG) enriched in cell-cycle and p53 pathways. Target validation confirmed CDK1 overexpression, prognostic significance, immune infiltration links, and mutation associations. In addition, a collected library of 9,667 natural phytochemicals was reduced through ML-based bioactivity (pIC₅₀) prediction targeting pKHG to discover potential lead molecules. Then, the selected top lead molecules were considered for further evaluation via molecular docking, molecular dynamics simulations, ADMET



analysis, and binding free-energy calculations (MM-GBSA). Among the selected phytochemicals, CID_14218027 (-7.69 kcal/mol), CID_487089 (-6.80 kcal/mol), and CID_174880 (-6.70 kcal/mol) showed the highest binding affinity score (GLIDE_XP score) and stable molecular interactions. Furthermore, MD simulations confirmed the conformational stability of ligand–protein complexes, supporting their potential as CDK1 inhibitors. This integrated omics-to-in-silico pipeline, identifies CDK1 as a robust therapeutic target and highlights natural product-derived inhibitors with favorable pharmacological and physicochemical properties. Therefore, these findings present a viable framework for accelerating precision drug discovery, with experimental validation underway. However, these findings are based solely on computational analyses and require further experimental validation to confirm CDK1 inhibitory activity, anticancer efficacy, and safety.

Keywords: Lung Cancer, CDK1 Inhibitors, Lung Adenocarcinoma, Transcriptomics and Machine Learning, Drug Discovery

1. Introduction

Lung cancer (LC) is the deadliest cancer in the world because it is the cause of more deaths annually than all breast, colon and prostate cancers put together [1,2]. In 2022, LC was responsible for about 1.82 million deaths, representing 18.7% of all cancer fatalities, far exceeding the mortality from colorectal cancer (9.3%) and breast cancer (6.9%). Approximately 2.48 million new LC diagnoses that same year, along with an age-standardized mortality rate of around 16.8 per 100,000, a pattern that tends to rise in countries with higher Human Development Index scores. If today's rates stay unchanged, population aging and growth could push the numbers to about 4.62 million new cases and roughly 3.55 million deaths by 2050, marking increases in total cases and deaths rather than in standardized mortality rates [2,3]. LC is highly lethal due to late-stage diagnosis, inherent molecular heterogeneity, and low sensitivity to the current treatment [4,5]. Though there has been a great deal of development of synthetic FDA-approved drugs-including cytotoxic chemotherapy, molecularly targeted agents, and immune checkpoint inhibitors [6–8] -the overall prognosis of the patients of LC remains challenging. However, drug-drug interactions, off-target toxicity, inadequate potency, or low



pharmacokinetics have limited the translation of these insights into effective therapies [9–11]. Though standard synthetic drugs have played a significant role in the treatment of cancer, these drugs are usually restricted by their adverse side effects, toxicities, resistance and high cost of production [12–15]. In particular, conventional chemotherapy is effective but frequently causes clinically meaningful adverse effects (notably myelosuppression and gastrointestinal toxicity), which can reduce quality of life and limit dosing intensity [16]. The number of patients who initially respond eventually develop acquired resistance (e.g., resistance after EGFR-TKI/osimertinib treatment), which drives disease progression and the need for next-line options [17,18]. Overall, these limitations justify exploration of alternative vulnerabilities such as cell-cycle control: Cyclin-Dependent Kinase is a central mitotic kinase/enzyme and is reported to be overexpressed and prognostically relevant in LC. Abnormal patterns in oncogenic signaling pathways, DNA repair mechanisms, cell death, and uncontrolled cell-cycle progression collectively contribute to tumor initiation and progression [19–22]. These molecular abnormalities suggest that cell-cycle regulators, mitotic proteins, and other important modulators represent promising points of therapeutic intervention. The transcriptomics and multi-omics study help to identify significant regulatory genes and signaling networks with clinical. Network pharmacology-based methodologies, including protein-protein interaction (PPE) networking, hub gene identification, transcriptomics factors, and pathway enrichment analysis, are some effective way to discover the key molecular regulators [23–26].

On the other hand, natural compounds, to be more specific medicinal plant derived phytochemicals have come out as an attractive alternative of those synthetic drug molecules, because their chemical diversity and evolved bioactivity can provide novel scaffolds and mechanisms that differ from current synthetic libraries [27]. The previous study shows that, more than 60 percent of the present anticancer agents are natural, such as such famous drugs as paclitaxel, camptothecin, or vincristine [28–30]. The IMPATT database is one of the large databases that contains information on 4,010 medicinal plants and 17,967 phytochemicals, along with their properties, which are often not well represented in other databases [31]. Thus, this database is a useful source of lung cancer drugs.



Nowadays, to reduce the time and cost, scientists are considering computational screening using *in silico* methodology prior to evaluating the study through experimental (*in vivo* and *in vitro*) validation. This research developed a holistic computational workflow integrating transcriptomics, machine learning (ML)-based lead screening, molecular docking and a dynamic simulation study to screen a vast set of compounds in a library and develop therapeutic drug molecules with reduced toxicity and enhanced efficacy. Previous studies have used transcriptomic and machine learning (ML) approaches in LUAD primarily to identify potential biomarkers and therapeutic targets or to predict responses to approved and investigational drugs [32–36]. Similarly, several studies have developed ML models using LUAD mRNA and mutation profiles to predict sensitivity to existing targeted and chemotherapeutic agents, achieving good predictive performance across dozens of drugs [37–40]. In contrast to these works, which mainly focus on biomarker/drug target identification, validation, and response prediction for existing drugs, this study integrates LUAD transcriptomics-based drug target identification with classical ML-based QSAR modelling to predict pIC50 values for plant-derived phytochemicals, followed by docking and molecular dynamics (MD) simulations to prioritize potent natural candidates.

Therefore, this study builds upon transcriptomic information from various publicly available GEO datasets to construct PPI networks, classical ML-based regression model to predict the bioactivity (pIC50), docking validation integrating different algorithm, evaluation of the binding strength by calculating post-docking MM-GBSA, large-scale MD simulation, and pharmacokinetics analysis to create a comprehensive approach to natural product-based drug discovery for lung cancer treatment. Compared with conventional docking-centered studies, this biologically informed workflow allows both disease-relevant target prioritization and more systematic lead selection from a large natural compound library. Thus, the present work provides a comprehensive and translationally oriented computational strategy for prioritizing CDK1-targeting phytochemicals in lung adenocarcinoma and supports future experimental validation of the identified candidate compounds. A complete guideline for this work is given below in **Figure 1**.



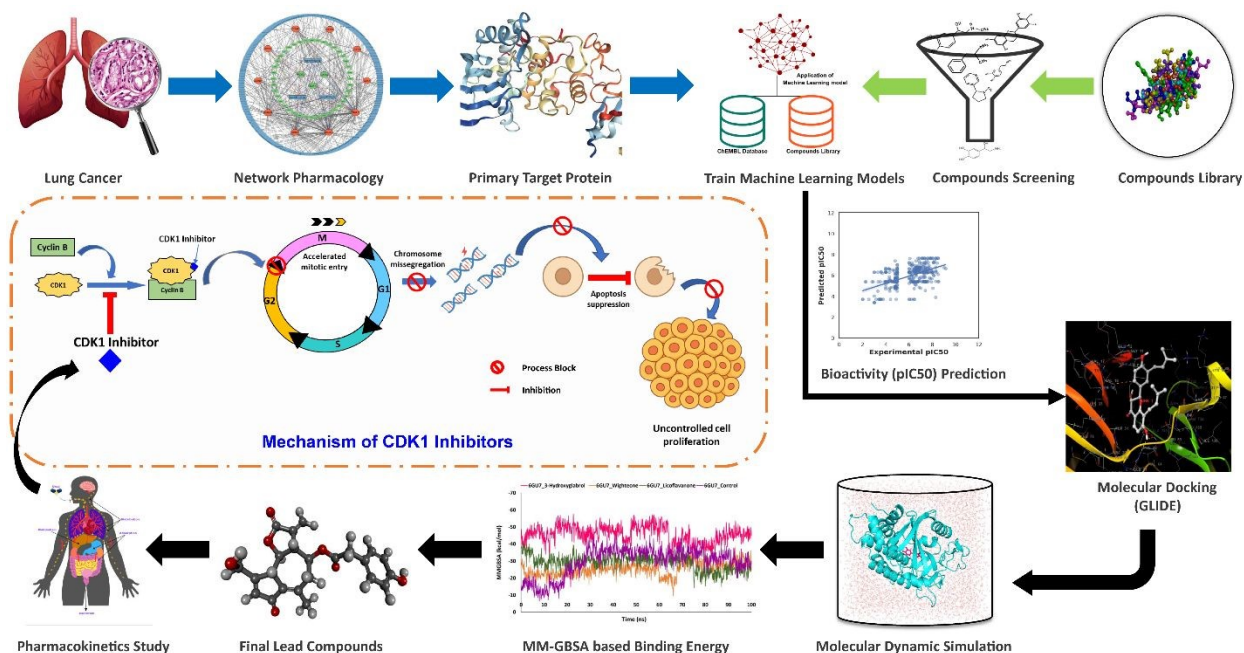


Figure 1: A Complete Graphical Representation of this Study

2. Materials and Methods

2.1 Target Identification

In this section, we present a systematic analysis of lung cancer gene expression profiles to uncover differentially expressed genes (DEGs) and highlight the primary key hub gene (pKHG), which plays a central role in understanding the underlying disease mechanisms.

2.1.1 Microarray Expression Dataset Acquisition

We used four publicly available microarray datasets from the Gene Expression Omnibus (GEO) database of National Center for Biotechnology Information (NCBI) (<https://www.ncbi.nlm.nih.gov/geo/>). The selected datasets were GSE19804, GSE10072, GSE18842, and GSE10799, which contain gene expression profiles of lung cancer and corresponding normal tissues. Details of these datasets, including platform and sample distribution, are summarized in **Table 1**.

Table 1. Summary of lung cancer GEO datasets used in this study

GEO Dataset	Number of Samples	Cancer	Control	Platform	Reference
GSE19804	120	60	60	GPL570 [HG-U133_Plus_2] Affymetrix Human Genome U133 Plus 2.0 Array	[41,42]
GSE10072	107	58	49	GPL96[HG-U133A] Affymetrix Human	[43]



GSE18842	91	46	45	Genome U133A Array GPL570 [HG-U133_Plus_2] Affymetrix Human Genome U133 Plus 2.0 Array	[44]
GSE10799	19	16	3	GPL570 [HG-U133_Plus_2] Affymetrix Human Genome U133 Plus 2.0 Array	[45]

2.1.2 DEG Analysis Based on Microarray Gene Expression Datasets

To identify differentially expressed genes (DEGs) in lung cancer, we analyzed the four selected datasets (GSE19804, GSE10072, GSE18842, and GSE10799) using GEO2R (<https://www.ncbi.nlm.nih.gov/geo/geo2r/>), an online tool provided by NCBI. In each dataset, samples were categorized into “control” and “cancer” groups, and the data were normalized using log₂ transformation. The limma package [46] was applied to detect DEGs, while the Benjamini–Hochberg false discovery rate (FDR) method was used to adjust for multiple testing. Genes with an adjusted p-value < 0.05 and logFC > +1.0 were considered significantly upregulated, whereas those with logFC < -1.0 were considered significantly downregulated.

2.1.3 Common DEGs (cDEGs) Identification

Finally, we identified the common DEGs (cDEGs) by intersecting the DEG lists from all four datasets, representing potential candidate genes associated with lung cancer. The cDEGs were extracted using the “dplyr” package [47] in R, and their overlap across datasets was visualized through Venn diagrams constructed with the R packages “ggplot2” [48] and “ggvenn” [49].

2.1.4 Identification of Key Hub Genes via PPI Network Analysis

Having determined the common DEGs (cDEGs), we investigated the interactivity of the relevant proteins with each other to reveal major molecular targets in lung cancer. The screening of a protein-protein interaction (PPI) network was performed through STRING database [50], and only interactions that were experimentally validated were considered. To concentrate the analysis on direct interactions among the cDEGs, we defined the minimum interaction score as low confidence (0.150) and did not add any additional interactors to the initial shell in order to reduce the results of the analysis to direct interactions. It was then visualized in Cytoscape (v3.10.4), with each node being a protein and each edge being the interaction between the two proteins [51]. In order to uncover the most significant proteins, we applied the CytoHubba [52] application and analyzed hub genes in eight diverse topological approaches: Degree, Maximum Neighborhood Component (MNC), Maximal Clique Centrality (MCC), Density of Maximum Neighborhood Component (DMNC), Edge Percolated Component (EPC), Bottleneck, EcCentricity, and



Closeness. Using these several topological analyses, the proteins with the largest percentage of interactions were chosen as the most important key hub genes (KHGs), which guaranteed a high value in identifying the key central target proteins in lung cancer-related networks.

2.1.5 Analysis of Transcriptional and Post-Transcriptional Regulation of KHGs

In order to understand the regulatory processes of the identified KHGs on an upstream level, we conducted a combined analysis of regulatory networks to identify the transcriptional and post-transcriptional regulators. The prediction of transcription factors was done using the TF-target interaction within the JASPAR database [53] and miRNA-target interaction within miRTarBase [54]. The regulatory networks were built with the help of the web platform NetworkAnalyst [55] (<https://www.networkanalyst.ca/>) that combines these interactions and enables topological analysis to be performed to discover essential regulators. The networks that resulted were then visualized and analyzed further using Cytoscape (v3.10.3) where nodes are regulators or target genes and edges are interactions that regulate a gene. Significant regulators were also prioritized according to their topological features of networks, thus bringing light to transcription factors and miRNAs that could be at the center of regulating lung cancer-related KHGs.

2.1.6 GO and Pathway Enrichment Analysis of KHGs

We first used the Gene Ontology (GO) and KEGG pathway enrichment analysis through the DAVID database to understand the biological functions of the identified DEGs [56]. GO grouped genes based on molecular functions, biological processes and cellular components whereas KEGG identified the relevant signaling pathways. GO terms Cutoffs GO terms A cutoff of 2 genes and p-value of less than 0.05 were used and KEGG pathways were selected with the default EASE score at DAVID. In order to guarantee the reliability of such results, DAVID-enriched GO terms and KEGG pathways were further verified with the help of Enrichr [57] and GeneCloudOmics [58]. The terms found in DAVID and those terms that we regularly found in both DAVID and the other validation platforms were retained allowing us to obtain a comprehensive and reliable set of functional categories based on which lung cancer KHGs were associated.

2.2 Target Validation

To validate the identified prime key hub gene (pKHG) in lung cancer, we focused on lung adenocarcinoma (LUAD) dataset from the available databases



2.2.1 Transcriptional and Proteomic Expression Analysis of the pKHG

The pKHG was selected from the candidate hub genes based on a combination of topological scoring and enrichment analysis, ensuring that the most biologically relevant and network-central target was prioritized for downstream validation. To strengthen the reliability of our findings, we next validated the expression pattern of the pKHG in LUAD using multiple publicly available resources. The Tumor Immune Estimation Resource (TIMER 2.0) (<http://timer.comp-genomics.org/>) [59] was first used to examine the expression of the pKHG between lung adenocarcinoma tissues and normal controls based on TCGA data. For further validation, the GEPIA2 platform (<http://gepia2.cancer-pku.cn>) [60] was employed, integrating TCGA and GTEx datasets. Boxplot parameters were set as follows: log₂FC cutoff = 1, p-value cutoff = 0.01, jitter size = 0.4, and values were log₂(TPM + 1)-transformed for visualization. The “Stage Plot” module of GEPIA2 was also utilized to determine the correlation between the expression of genes and the pathological stages (I-IV) of LUAD. In addition, two levels were applied to the UALCAN database (<http://ualcan.path.uab.edu>) [61]: (i) the study utilized TCGA RNA-seq data to validate the differences in expression of mRNA between the normal and LUAD tissues, and (ii) the study used Clinical Proteomic Tumor Analysis Consortium (CPTAC) data to test the levels of protein expression of the pKHG. This validation was consistent both at the transcriptomic and proteomic levels with this combined approach.

2.2.2 Survival Analysis of the pKHG

In order to examine the prognostic role of the pKHG in LUAD, the Kaplan-Meier survival analysis was performed on the GEPIA2 stage. The median cutoff value of 50 percent was used to classify patients in high- and low-expression groups. The overall survival (OS) and disease-free survival (DFS) were measured, and the difference between the two was tested by the use of the log-rank test. The survival plots have the time in months on the X-axis and the survival probability (%) on the Y-axis, and the dotted line indicates the 95 percent confidence intervals. This discussion presented valuable information regarding prognostic value of the pKHG in LUAD.

2.2.3 Immune Infiltration Associated with the pKHG

The correlation between the pKHG and immune cell infiltration on LUAD alone was explored using the module named “Immune-Gene” in TIMER2.0. We concentrated on CD8+ T cells, macrophages and CD4+ T cells to determine the strongest positive and negative correlations. The deconvolution algorithms were estimated to give an estimate of the immune cell infiltration and Spearman correlation with adjustment of



tumor purity was used to provide both the correlation coefficient and the significance value. The resulting associations were plotted with heatmaps, which shows in which type(s) of immune cells the hub gene in lung cancer is most tightly associated.

2.2.4 Investigation of Mutations and Alterations in the pKHG

To continue investigating the problem of cancer-related genomic changes, the cBioPortal platform (<http://cbioportal.org>) was used [62]. The patterns of genetic alteration of the identified prime key hub gene (pKHG) were systematically analyzed using this resource. The analyses of the data in the TCGA PanCancer Atlas studies (including 32 tumor types and 10,957 samples of patients) were performed with the help of the “Query by Gene” module. The resulting “Cancer Type Summary” was a summary of mutation and copy number change types of the pKHG in various cancers. Moreover, the application of the “Mutations” module was used to produce a schematic diagram of the exact mutation’s sites in the gene.

2.3 pKHG Guided *In Silico* Drug Discovery

To explore therapeutic opportunities, we performed structure-based drug discovery on the validated pKHG. This included retrieval and preparation of target protein’s 3D crystal structure, ligand collection, virtual screening of ligands, molecular docking, ADMET, Molecular dynamics (MD) simulation, MMGBSA, PCA (Principal Component Analysis) and 3D-FEL (Free Energy Landscape).

2.3.1 Retrieval and Preparation of Target Protein

The pKHG crystallographic structure was derived in the RCSB protein database [63]. To enable the subsequent analysis of the protein, all the existing heteroatoms, ligands and water molecules were eliminated with the help of BIOVIA Discovery Studio 2021 [64]. The protein structure was then energy-minimized using Swiss-PDB Viewer (spdbv) 4.1.0 [65].

2.3.2 Compound Library Construction

In order to detect possible inhibitors, we built a universal phytochemical library with the aid of the IMPPAT 3.0 database [31]. We took a total of 33 traditionally used medicinal plants depending on their ethnopharmacological significance, and extracted phytochemicals in various parts of the plants, such as roots, leaves, and seeds. Out of these chosen plants, 9,667 phytochemicals were obtained, which were



collected together to form a comprehensive compound library to be used in future virtual screening and molecular docking studies.

2.3.3 Physicochemical Property-Based Ligand Screening

Virtual Screening (VS) is an essential part of contemporary drug discovery, which also provides a computational method of finding the possible bioactive compounds in large chemical libraries [66]. Pharmacokinetic and toxicity-related parameters can help in this screening process with the help of ADMETlab 2.0 [67], a web-based platform. In order to narrow the selection down, the Rule of Five (RO5) by Lipinski [68] is often used as a criterion in determining drug-likeness. Under this rule, the compound has higher chances of being better oral bioavailable; therefore, it has no more than five hydrogen bond donors, no more than ten hydrogen bond acceptors, a molecular weight less than 500 Da, and the logP less than five. Phytochemicals that fulfil this requirement are said to be good leads to further computational and experimental studies. Duplicate compounds were eliminated and those that survived were those which had 3D structure available and would be used in further studies.

2.3.3 ML-based Bioactivity Prediction (pIC50) of the Selected Compounds

The potent inhibitory concentration (pIC50) value serves as an important indicator of a drug's potency, showing the concentration needed to reduce the activity of a biological target by half [69]. Over time, many studies have worked to improve these prediction strategies, making it easier to focus on the most promising drug candidates [70–72]. In this context, the present study applied machine learning (ML) techniques to forecast the bioactivity of the selected compounds against pKHG.

a. Dataset curation and preparation

For dataset construction, bioactivity records and chemical structures were obtained from the ChEMBL database, a well-established platform for QSAR investigations. Compounds were then curated by retaining only those entries that reported IC50 values against pKHG. For selecting the final ChEMBL database we mainly focused on the types of protein, organism, availability of the compounds activities, including standard type (IC50) and standard unit (nm). This careful refinement produced a dataset of reliable and biologically relevant molecules, which served as the basis for regression modeling and subsequent analyses.

b. Molecular Descriptor Calculations



Molecular descriptors play a central role in QSAR modeling because they translate the structural and physicochemical features of compounds into measurable values, allowing meaningful patterns to be recognized [73]. In this study, RDKit was used to calculate Morgan fingerprints and ECFP4 descriptors for each compound. These descriptors offered a detailed numerical profile of the compounds' structural characteristics, which served as essential inputs for building predictive models to estimate bioactivity in QSAR modeling.

c. Data Splitting into Train and Test

The final dataset of 987 compounds was split into training and test subsets using the `train_test_split` function from scikit-learn with stratification on the class labels to preserve the original class distribution in both subsets. It was split into 789 molecules in the training set and 198 molecules in the test set, typical among cheminformatics and other machine learning tasks to use 20-30% of the data to test the algorithm. Random or stratified random splits are widely used as a baseline in cheminformatics modeling, although recent work has highlighted that more stringent splitting strategies (e.g., scaffold or clustering-based splits) may provide more realistic assessments of external generalization for chemical datasets [74,75].

d. ML Model Development and Validation

For predictive modeling, we selected the best regression model, with a high-performance machine learning algorithm valued for its speed, accuracy, interpretability, and the coefficient of determination (R^2) values of test and train dataset. To enhance the model's effectiveness, Recursive Feature Elimination (RFE) was applied to remove less informative features, reducing noise and making the model more interpretable. Here, we used `optuna` to tune the hyperparameter. The predictive performance of the model was assessed using widely adopted metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and coefficient of determination (R^2). These metrics were calculated according to the previously published formula.

Mean Absolute Error (MAE):

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Squared Error (MSE):



$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE):

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

R-squared (R²):

$$1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Here, y_i refers to the experimentally observed IC₅₀ values, \hat{y}_i represents the predicted values generated by the model, and \bar{y} indicates the mean of the observed IC₅₀ values.

e. Model Implementation for Bioactivity (pIC₅₀) Predictions

After model development we applied our trained and fine-tuned model to predict the pIC₅₀ values of the collected compounds against pKHG. Compounds showing the highest predicted pIC₅₀ values (>6.5), reflecting greater potential potency, were considered for further investigation. These selected compounds were then analyzed through molecular docking study to evaluate their binding affinities and interactions with the pKHG.

2.3.4 Molecular docking via AutoDock vina.

Molecular docking is important in the discovery and optimization of possible drug candidates. It can be used to estimate the possibility of two molecules, such as between a protein and a ligand, interacting with each other, and can provide information on their binding behavior and possible efficacy [76]. This is why in this study, PyRx was used [77] in molecular docking of the chosen compounds with target protein. PyRx offers an easy-to-use graphical interface with AutoDock Vina to make docking a lot easier. It further combines AutoDock Vina, and Open Babel, which provides a complete package of molecular docking and analysis [78]. Lastly, we analyzed the interaction between the ligands and the target receptor through PyMOL and BIOVIA Discovery Studio 2021, which enabled us to see the interaction and to have a better perspective about their binding patterns.

2.3.5 Molecular docking validation via Schrödinger software.

Following the first docking in AutoDock Vina, the shortlisted compounds were further filtered based on their highest binding affinity and were selected for further docking re-scoring using the GLIDE module.



Firstly, these potential hit compounds were generated and optimized with LigPrep module of Schrodinger suites [79]. The protein structure was preprocessed, and the grid was generated through the protein preparation wizard and receptor grid generation module, respectively. This involved the correction of missing hydrogen atoms, the closure of side-chain gaps and loops using Prime, the elimination of water molecules that are distant to the active site and the creation of suitable protonation states using Epik at physiological pH of 7.0 ± 2.0 . Lastly, the validated ligands were re-docked into the active sites of the target proteins using the GLIDE module of the Maestro (*Version* 11.8.012) software in Extra Precision (XP) mode which made it such that only reliable targets were studied in more detail in terms of their binding interactions.

2.3.6 Post docking MMGBSA

To gain a better idea of ligand-protein interactions, post-docking refinement methods, especially, Molecular Mechanics Generalized Born Surface Area (MM-GBSA) were used. Although docking scores provide an initial value of binding affinity, they do not include all of the solvation effects or entropic contributions, which may affect the binding stability and strength of ligand binding [80]. MM-GBSA overcomes these drawbacks by offering a more accurate estimate of the binding free energy (ΔG), which provides a more accurate image of the affinity between the ligand and receptor [81]. In this research, the MM-GBSA analysis was only applied to the ligands that yielded good results in extra precision (XP) mode, using OPLS4 force field in Schrodinger suite, PRIME module. This selective approach was necessary because not all ligands produce reliable or meaningful scores during XP docking due to differences in binding modes and structural flexibility [82]. The binding free energy for each ligand was calculated using the formula:

$$\Delta G_{\text{bind}} = \Delta G_{\text{complex}} - (\Delta G_{\text{receptor}} + \Delta G_{\text{ligand}})$$

Here, $\Delta G_{\text{complex}}$ represents the free energy of the ligand-protein complex, $\Delta G_{\text{receptor}}$ is the receptor's free energy, and ΔG_{ligand} is the ligand's free energy. More negative ΔG values correspond to stronger binding interactions [83].

2.3.7 Molecular Dynamics Simulation and post simulation MM-GBSA Calculation

The molecular dynamics (MD) simulations were performed using the Desmond module of Schrödinger to examine the stability and conformational dynamics of the protein–ligand complexes [84]. System



preparation was performed using the System Builder wizard. Each complex was solvated in a simple point charge (SPC) water model inside an orthorhombic simulation box, ensuring a minimum distance of 10 Å between the protein surface and the box edges. To emulate physiological ionic conditions, the systems were neutralized and supplemented with 0.15 M Na⁺ and Cl⁻ ions. Energy minimization was conducted for 100 ps to remove unfavorable contacts, followed by equilibration under NVT and NPT ensembles at 300 K and 1 atm. The production MD simulations were conducted for 100 ns with a 2-fs timestep. Trajectories were saved at 20-ps intervals, yielding 5000 frames in total. We used the OPLS3e force field [85], which delivers improved parameterization for biomolecules as well as drug-like compounds. The simulation trajectories were analyzed for RMSD, RMSF, SASA, radius of gyration, hydrogen bonding, and principal component dynamics to assess conformational stability.

Then the MM-GBSA-based binding free energy (ΔG_{bind}) of the selected complex was calculated with the gmx_MMPBSA package [86,87]. This approach integrates van der Waals, electrostatic, and solvation energy components, along with solvent-accessible surface area (SASA) contributions, providing a thermodynamic estimate of ligand affinity. The Desmond trajectory files were converted to GROMACS-compatible formats using Schrödinger utilities, while the corresponding topology files were generated through InterMol conversion of *.cms files to *.gro and *.top formats [88]. The free energy associated with binding (ΔG_{bind}) was determined based on the following relationship:

$$\Delta G_{\text{bind}} = \langle G_{\text{PL}} \rangle - \langle G_{\text{P}} \rangle - \langle G_{\text{L}} \rangle$$

In this equation, $\langle G_{\text{PL}} \rangle$, $\langle G_{\text{P}} \rangle$, and $\langle G_{\text{L}} \rangle$ correspond to the average free energies of the complex, the unbound protein, and the free ligand, respectively.

Accordingly, the overall binding energy is given by:

$$\Delta G_{\text{bind}} = \Delta E_{\text{MM}} + \Delta G_{\text{SOLV}} - T\Delta S$$

where ΔE_{MM} represents the gas-phase molecular mechanics energy (including van der Waals and electrostatic components), ΔG_{SOLV} is the change in solvation free energy, and $T\Delta S$ reflects the contribution from entropy.

2.3.8 Pharmacokinetics (ADME & Toxicity) Evaluation

After completing molecular docking and simulation studies, it became essential to examine whether the top-performing compounds possessed properties suitable for real-world drug development. The pharmacokinetic analysis has been done based on the SwissADME web tool (<http://www.swissadme.ch/>), which predicts the important ADME parameters (absorption, distribution, metabolism, and excretion)



and the physicochemical properties including solubility, lipophilicity, and molecular flexibility [89]. In this analysis, short listing was done on compounds that had an optimal balance of potency and pharmacokinetic feasibility.

In order to supplement these results, the ProTox-III server (https://tox-new.charite.de/protox_III) was employed to make predictions regarding different types of toxicity, such as hepatotoxicity, nephrotoxicity, cardiotoxicity, and neurotoxicity [90]. It uses deep machine learning algorithms to train on experimental data to give confidence in estimating toxicity.

3. Result

3.1 Target Identification

3.1.1 Differential Gene Expression Analysis and DEG Selection

In order to examine the transcriptional changes in lung cancer, four independent GEO microarray datasets were analyzed (GSE19804, GSE10072, GSE18842 and GSE10799). We found extensive transcriptional changes in lung cancer versus normal tissues, as many of the genes were significantly upregulated or downregulated in all the four datasets. **Table S1** contained lists of these genes of each dataset. DEGs visualization was performed as a Volcano plot (**Figure 2A**), revealing obvious differences in significantly upregulated and downregulated genes in the state of lung cancer and normal tissues. As a step to find strong molecular signatures, we then identified common DEGs among the four datasets. The intersection analysis showed that there are common DEGs (cDEGs), which were integrated in a Venn diagram (**Figure 2B**). **Table S2** gives the detailed list of these cDEGs. These shared genes are the possible candidates that can be important to the pathogenesis of lung cancer and they were prioritized to be further analyzed with respect to their functions.



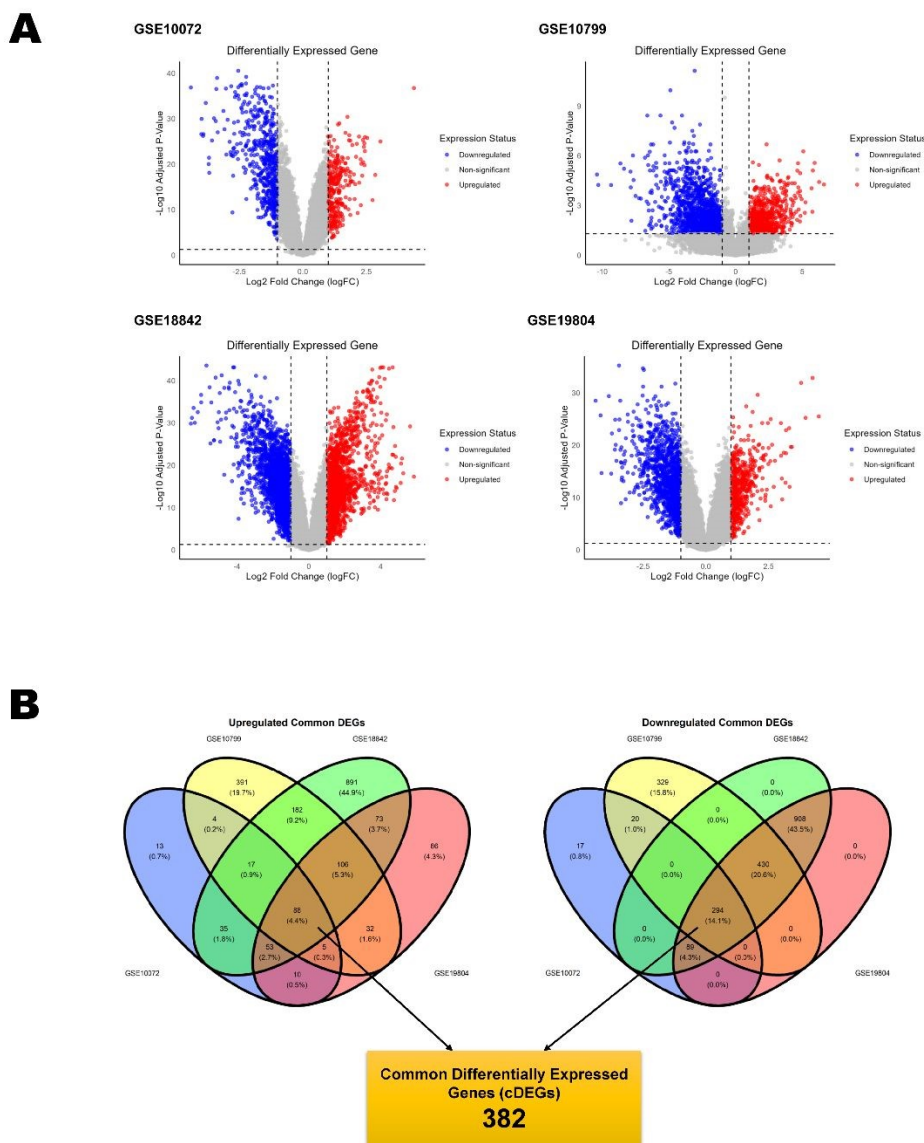


Figure 2. (A) Volcano plots illustrating the distribution of significantly upregulated and downregulated genes across each dataset. (B) Venn diagram showing the overlap of DEGs among the four datasets and the common DEGs (cDEGs).

3.1.2 Identification of Key Hub Proteins via PPI Network Analysis

Having built the protein-protein interaction (PPI) network using the common DEGs (cDEGs), we obtained the network having 375 nodes and 349 edges with the average node degree of 1.86 and the average local clustering coefficient of 0.299. They were supposed to have 222 edges and the PPI



enrichment p-value was 2.55×10^{-15} , which meant that the interactions observed were considerably more than should have occurred by chance. All of these PPI networks are depicted in **Figure 3A** in which nodes correspond to individual proteins and edges are experimentally confirmed interactions.

To identify the most influential proteins within the network, we applied eight topological methods using the CytoHubba plugin. For each method, the top 10 hub genes were determined. The results of these analyses are summarized in **Figure 3B**, highlighting the proteins that consistently appear as central hubs across multiple metrics. These key hub genes (KHGs) are likely to play critical roles in lung cancer biology and represent promising targets for further investigation. Among all the identified KHGs, CDK1 exhibited the highest number of protein-protein connections, and therefore its interaction network is highlighted in **Figure 3A**.



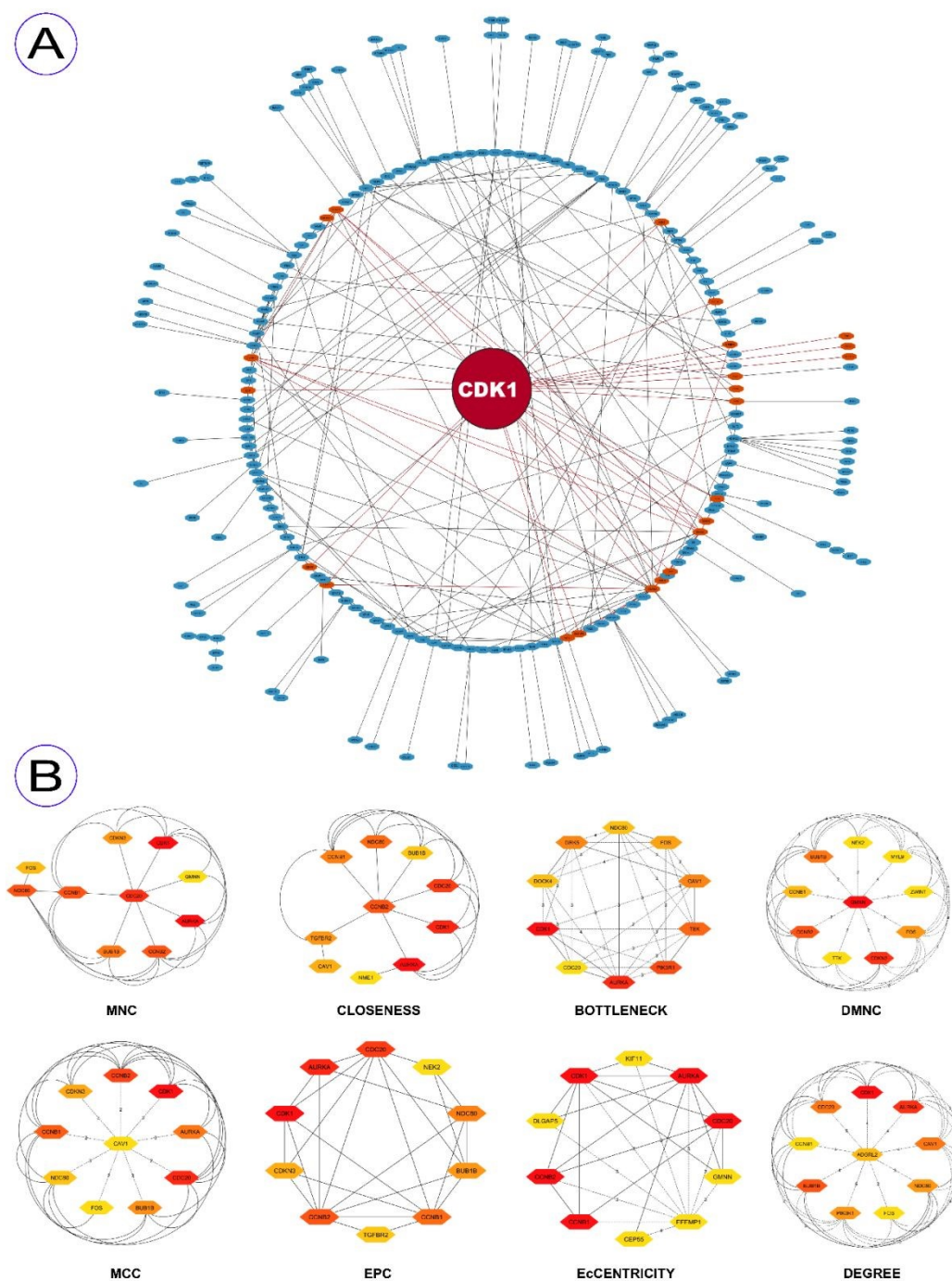


Figure 3. Protein–protein interaction (PPI) network of common DEGs (cDEGs) in lung cancer. (A) Overall PPI network showing nodes as proteins and edges as experimentally validated interactions. (B) Top 10 hub genes identified using eight different topological methods (Degree, MNC, MCC, DMNC, EPC, Bottleneck, EcCentricity, and Closeness).



3.1.3 Analysis of Transcriptional and Post-Transcriptional Regulation of KHGs

After identifying hub genes using multiple topological algorithms in the PPI network (eight different methods), we observed that although several hub genes were common across most methods, some unique genes were identified by individual algorithms. Each centrality method emphasizes different network properties and therefore captures complementary aspects of network biology [91,92]. To minimize the bias associated with any single algorithm and to ensure that potentially relevant candidates were not excluded, we considered the union set of hub genes ($n=26$) obtained from all methods. To better understand how the identified KHGs are regulated, we carried out an integrated analysis focusing on their interactions with transcription factors (TFs) and microRNAs (miRNAs). As illustrated in **Figure 4A** and **Figure 4B**, the networks highlight kHG–TF and kHG–miRNA interactions, respectively. From the network's topological evaluation, we pinpointed five key TFs (FOXC1, GATA2, YY1, E2F1, and HINFP) along with five major miRNAs (hsa-miR-192-5p, hsa-miR-92a-3p, hsa-miR-193b-3p, hsa-miR-215-5p, and hsa-miR-155-5p) as central regulators. These molecules appear to play crucial roles in controlling gene expression at both the transcriptional and post-transcriptional levels.

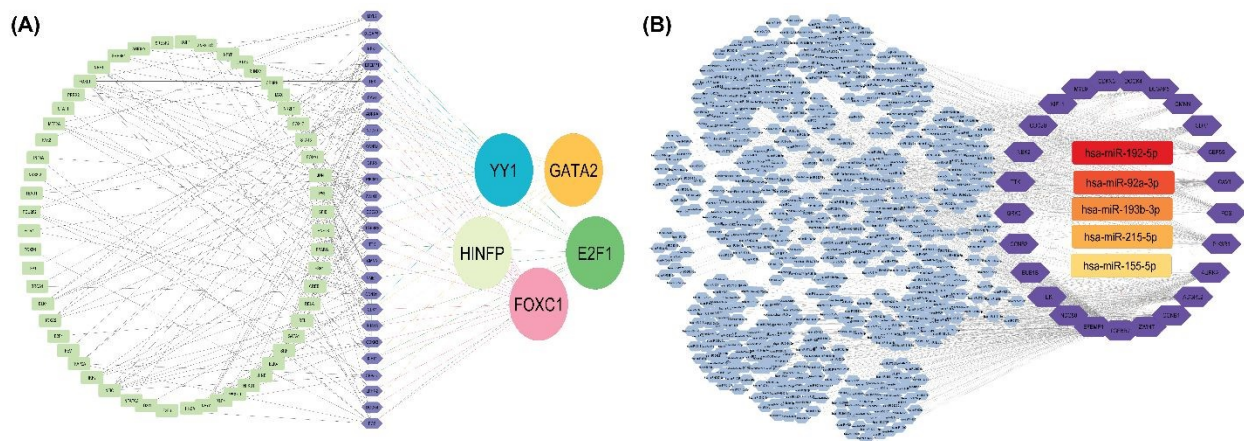


Figure 4. Integrated regulatory network of DEGs with transcription factors (A) and microRNAs (B).

3.1.4 GO and Pathway Enrichment Analysis of KHGs

In order to shed light on the biological significance of the identified key hub genes (KHGs) in lung cancer, the enrichment analysis of Gene Ontology (GO) and KEGG pathways were conducted in DAVID, EnrichR, and GeneCloudOmics. All of the terms detected in DAVID and those that were present throughout the validation platforms were kept, which led to a full range of functional categories of the KHGs of lung cancer. The enriched categories were categorized into biological processes (BP), cellular



components (CC), molecular functions (MF), and signaling pathways, and the lists of corresponding detail are presented in the Supplementary Section (Table S3). In GO categories, enrichment of functions related to the cell cycle was observed to be strong. Mitotic cell cycle progression, spindle organization, chromosome segregation, and checkpoint signalling terms were highly enriched in biological processes (BP), and AURKA, CDK1, CCNB1, NDC80 and CDC20 among others have been repeatedly involved. These results indicate that the destabilization of mitotic fidelity is one of the key markers in lung cancer progression. KHGs in the cellular components (CC) group were enriched to known mitotic structures, e.g. kinetochore, centrosome, spindle pole, and mitotic spindle, and KHGs genes like NDC80, BUB1B and AURKA were seen as major drivers of these enrichments. This indicates that there is structural dysregulation of chromosome segregation machinery in the lung cancer cells. In molecular functions (MF), the representation of kinase activities (protein kinase, serine/threonine kinase) and binding ATP functions were most enriched with the central position of AURKA, CDK1, NEK2, and TTK. Multiple histone kinase activities also arose, and these hub genes might play a role in mitosis-related epigenetic regulation. In accordance with the GO results, KEGG pathway analysis showed 14 highly enriched pathways with the strongest association made by the cell cycle pathway ($p = 3.83E-07$) with multiple KHGs (CDK1, CCNB1, CCNB2, BUB1B, CDC20, NDC80, TTK). Ways connected to tumor suppressive functions and cellular stressful reactions, including p53 signaling and FoxO signaling were additionally greatly enriched and include these genes as part of pathways that regulate apoptosis, senescence, and genomic stability. Interestingly, a number of infection-related pathways (e.g., HTLV-1, HIV-1 and Hepatitis B) were also enriched, implying that there are common molecular pathways between virus infection and oncogenesis which could play a role in the pathology of lung cancer.

3.2 Target Validation

Based on the integrated evidence from hub gene ranking within the PPI network and the enrichment analysis highlighting its strong involvement in significant pathways, CDK1 was prioritized as the prime key hub gene (pKHG) for downstream validation.

3.2.1 Evaluation of Transcriptional and Proteomic Expression Levels of the CDK1 in Lung Cancer

To validate CDK1 as the prime key hub gene, we first assessed its expression profile across pan-cancer datasets using TIMER2.0. As shown in Figure 5A, CDK1 expression was significantly elevated in multiple tumor types, including lung adenocarcinoma (LUAD), with $p < 0.001$ indicating strong statistical significance. Further validation in LUAD using GEPIA2 confirmed the upregulation of CDK1. The



boxplot analysis (**Figure 5B-i**) demonstrated a significant increase in CDK1 transcript levels in LUAD tissues compared to normal controls, while the stage plot (**Figure 5B-ii**) revealed that CDK1 expression remained significantly high across pathological stages (I–IV) ($p = 0.000579$). These findings suggest that CDK1 overexpression is maintained throughout disease progression. In addition, UALCAN-based analyses provided consistent evidence at both transcriptomic and proteomic levels. TCGA RNA-seq data showed markedly higher CDK1 expression in LUAD tumors compared with normal samples (**Figure 5C-i**), and CPTAC proteomic data similarly confirmed elevated protein expression of CDK1 in tumor tissues relative to controls (**Figure 5C-ii**). Collectively, these results demonstrate that CDK1 is robustly upregulated in LUAD at both mRNA and protein levels, reinforcing its candidacy as a biologically relevant prime key hub gene for downstream validation.



3.2.2 Survival Analysis of the CDK1 in Lung Cancer

To explore the clinical relevance of CDK1 expression in LUAD, we performed Kaplan–Meier survival analysis using the GEPIA2 platform. Patients were stratified into high- and low-expression groups based on the median cutoff value. In the overall survival (OS) analysis, patients with high CDK1 expression showed a markedly poorer outcome compared to those with low expression (**Figure 6A**). The difference was statistically significant (log-rank $p = 2.6 \times 10^{-5}$), with a hazard ratio (HR) of 1.9, suggesting that elevated CDK1 expression nearly doubled the risk of death in LUAD. Consistent with this, the disease-free survival (DFS) analysis also indicated that patients with higher CDK1 levels experienced earlier recurrence and shorter DFS than those in the low-expression group (**Figure 6B**). The association was significant (log-rank $p = 0.027$; HR = 1.4), reinforcing the unfavorable role of CDK1 overexpression in disease progression. Together, these findings highlight CDK1 as a potential prognostic marker in LUAD, where its overexpression is linked to both reduced survival and increased risk of relapse.

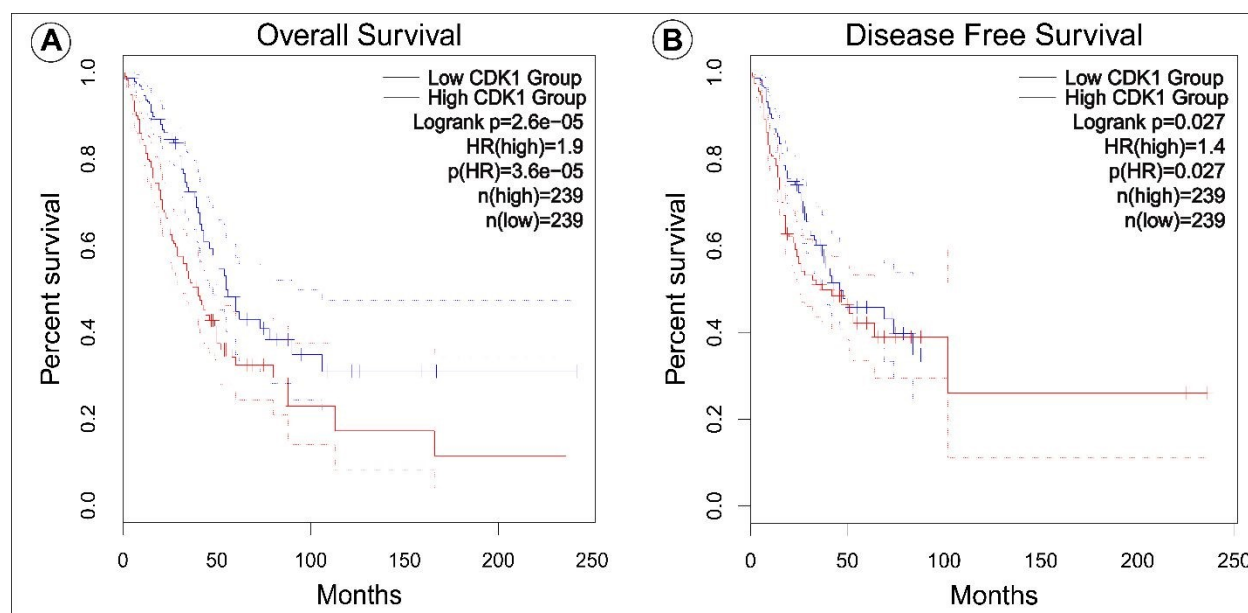


Figure 6. Kaplan–Meier survival analysis of LUAD patients stratified by CDK1 expression. (A) Overall survival (OS). (B) Disease-free survival (DFS).

3.2.3 Evaluation of Immune Infiltration Associated with the CDK1

Using TCGA information, we examined the connection between the expression of CDK1 and the immune cell infiltration in LUAD. Data of CD8+ T cell infiltration analysis in LUAD indicated that the levels of CDK1 expression and T cell CD8+ infiltration had a significant positive correlation as determined by



various deconvolution algorithms. A clear band in the heatmap (**Figure S1A**) showed the LUAD and this band showed a significant association. The best positive relationship was found with T cell CD8+ (naïve_XCELL) infiltration with a Rho of 0.235 and p-value of 1.35e-07. On the other hand, the most significant negative relationship was observed with T cell CD8+ (EPIC) with a Rho of -0.207 a p-value of 3.46e-06. In CD4+ T cell infiltration in LUAD, the heatmap (**Figure S1B**) once again was showing an orange band that was representative of LUAD-specific associations. The most significant correlation was found with T cell CD4+ (Th2XCELL) with the Rho value of 0.813 and p-value of 2.96e-117. The most negative correlation was found with T cell CD4+ (central memory_XCELL) with a Rho value of -0.322 and a p-value of 2.27e-13. When it comes to the macrophage infiltration in LUAD, the heatmap (**Figure S1C**) was once again represented by an orange band, pointing to the relevant associations. The best positive correlation was on Macrophage (M1_CIBERSORT) infiltration at a Rho of 0.367 at a p-value of 3.39e-17. It was the strongest negative correlation, but with Macrophage (M2_TIDE), which has a Rho value of -0.346 and a p-value of 2.45e-15. Combined with the analysis, we have found that there is a strong association between the expression of CDK1 and immune cell infiltration in LUAD. In particular, there were positive and negative correlations between CD8 + and CD4 + T cell subsets, and a significant correlation between CD4 + Th2 cells (a subtype that is a humoral immunity). On the other hand, the negative correlation with the CD4+ central memory cells (long-term immune memory) showed a condition-dependent effect. In addition, the expression of CDK1 was positively associated with M1 macrophage infiltration (pro-inflammatory/anti-tumor phenotype) and negatively correlated with M2 macrophages (immunosuppressive/pro-tumor phenotype). This movement between M1 and M2, also known as macrophage polarization, is the functional re-programming of tumor microenvironment macrophages. Taken together, these results indicate that CDK1 can have an effect on the immune microenvironment of LUAD, and it could have an impact on tumor progression and sensitivity to therapy.

3.2.4 Investigation of Mutations and Alterations in the CDK1

To investigate the genomic landscape of CDK1 in human cancers, we queried the cBioPortal database using data from the TCGA Pan Cancer Atlas. The highest alteration frequency was observed in uterine corpus endometrial carcinoma (UCEC), approaching 8%, predominantly driven by amplification. This was followed by skin cutaneous melanoma at approximately 3%, Cholangiocarcinoma and Uterine corpus endometrial carcinoma at 2.5-2.8%, and lower rates in breast invasive carcinoma (**Figure S2A**). Overall, mutations were the most prevalent alteration type across cohorts, with structural variants (purple), amplifications (red), deep deletions (blue), and multiple alterations (black) occurring less



frequently but notably in cancers such as LUAD and BRCA. The presence of these alterations, particularly amplification, positions CDK1 as a significant oncogenic driver in a subset of LUAD patients. This suggests that targeted therapies against CDK1 (such as selective CDK inhibitors) could be a viable therapeutic strategy for those patients whose tumors harbor these specific genetic alterations. Further analysis via the "Mutations" module generated a schematic lollipop diagram of CDK1 mutations (**Figure S2B**) mapping alterations onto the 297-amino acid protein sequence. Missense mutations (green) dominate, while truncating (black), splice (brown), and fusion (purple) variants were less common. Detailed mutation distributions were presented in **Table S4**. This data shows that CDK1 can be mutated in LUAD, specifically through missense mutations that are predicted to be functionally damaging. This analysis identified and characterizes three non-synonymous somatic mutations in the CDK1 gene (S178L, I136N, E57V) within a cohort of lung adenocarcinoma (LUAD) patients from The Cancer Genome Atlas (TCGA). While none were currently annotated as known drivers in major clinical databases (OncoKB, CIViC), in silico functional prediction tools suggest that two of these mutations (I136N and E57V) are likely pathogenic. These mutations may represent a candidate biomarker for sensitivity to CDK inhibitors in a subset of LUAD patients.

3.3 Structure-Based Drug Discovery Targeting CDK1

3.3.1 Structural Preparation and Optimization of the Target Protein

The three-dimensional crystal structure of the CDK1 protein (PDB ID: 6GU7) was obtained from the RCSB Protein Data Bank [93]. Prior to its use in computational experiments, the structure underwent a refinement process as outlined in the methodology section. This preparation involved the removal of heteroatoms and water molecules, followed by energy minimization to ensure a stable and optimized conformation suitable for subsequent molecular modeling analyses.

3.3.2 Compound Filtering and Selection Results

In this study we selected a total of 33 natural medicinal plants as a source of compounds to identify the potential inhibitors of CDK1 from IMPPAT database (**Table S5**). Initially we retrieved almost 9,667 compounds from the database with their IMPPAT IDs and SMILES. The Lipinski's Rule of Five (RO5) results reveal that 4,236 were found to have zero (0) violations out of 5,786 phytochemicals, indicating favorable drug-likeness profiles. Following the removal of duplicates, 2,113 unique compounds were retained from all of that selected medicinal plants (**Table S6**), these compounds were then selected for further computational analyses.



3.3.3 Bioactivity prediction (pIC₅₀) using ML approach

a. Model development, Optimization and Validation

To calculate the bioactivity of our selected compounds targeting CDK1 we choose ChEMBL dataset (ChEMBL308) containing the 3,348 experimentally validated inhibitors and 4,536 activities against CDK1 with their IC₅₀ information (Access on November 20, 2025). The aim was to develop a classical regression model for CDK1 inhibitors, so that we can calculate the IC₅₀ and pIC₅₀ of our compound library based on trained model. According to our study, we get random forest (RF), lighGBM, XGBoost, and Stacking as the best performed model. The evaluation metrics indicates that the difference of R² value for RF, LightGBM, XGBoost, and Stacking were 0.177, 0.117, 0.241, and 0.214, respectively (**Figure S3A**). Thus, based on this finding we can observed that, the LightGBM shows the lowest difference in R² values between test and train, indicating best fit for prediction outcomes.

The LightGBM regression among other machine-learning methods has been receiving considerable attention due to its excellent predictive accuracy, computational efficiency, and the capacity to handle large data volumes [94,95]. As an ensemble of decision trees, LightGBM is well suited to capture nonlinear structure–activity relationships while maintaining fast training on large molecular descriptor and fingerprint sets. The LightGBM model was trained on the training set using optimized hyperparameters, including `n_estimators`, `num_leaves`, `max_depth`, `learning_rate`, `subsample`, `colsample_bytree`, `min_child_samples`, `reg_alpha`, and `reg_lambda`. Additionally, LightGBM model tuned to improve generalization and reduce overfitting. Model performance was then evaluated on the test set using standard QSAR regression metrics: Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared (R²). These parameters offered a comprehensive assessment of how well the model's predictions aligned with the actual experimental data. The test model demonstrated impressive performance, with MAE, MSE, RMSE, and R² values of 0.520, 0.544, 0.737, and 0.758, respectively, shown in **Figure S3A**. The overall results suggest that the accuracy is comparatively higher than the others. Therefore, finally we selected the LightGBM model for predicting the bioactivity (pIC₅₀) value of our selected natural compounds.

b. Application to Bioactivity Calculation and Compound Screening

Once the model's predictive performance was validated, it was used to virtually screen a curated phytocompound library of 2,113 phytocompounds by predicting their pIC₅₀ values using the optimized ML model. A total of 380 compounds showed predicted pIC₅₀ > 6.5 (~IC₅₀ < 1 μM), suggesting potential inhibitory activity against the target and prioritizing them for follow-up validation (see **Table S7**) [96]. As



the experimental and predicted values lie close to each other in the scatter plot (**Figure S3B**), it indicates that the QSAR model is strong and has excellent predictive accuracy, as most data points closely align along the regression line. These results indicate the extent to which the model is useful in identifying potentially promising molecules that can be subjected to further investigation either with the use of molecular docking or through experimental studies to confirm their utility as a possible drug target.

3.3.4 Molecular docking analysis via AutoDock vina

Out of 380 previously selected phytochemicals, a total of 358 with available 3D SDF structures were extracted from the PubChem database and subjected to molecular docking against the CDK1 protein (6GU7) using AutoDock Vina implemented in PyRx. The docking was performed using a grid box centered at the protein's active site with coordinates $X = 22.0351$, $Y = 12.6323$, and $Z = 4.7655$. Additionally, we also compare the docking binding affinity score with the known inhibitor of CDK1 (FB8/AZD5438) and we set this binding affinity of AZD5438 as our cut-off value -8.8 kcal/mol. Finally, we select 29 compounds for further analysis and evaluation of their binding strength (**Table S8**).

3.3.5 Molecular docking validation via Schrödinger software

For further accuracy assessment, extra precision (XP) docking was carried out using the GLIDE module within the Schrödinger's Maestro environment for the selected 29 compounds. From the initially docked compounds, top 5 ligands (CID_115196, CID_5317764, CID_14218027, CID_487089, and CID_174880) exhibiting binding affinities below -9.2 kcal/mol in AutoDock Vina were shortlisted for this validation step, along with the control compound. Their XP docking results reveals their binding affinity ranged from -6.80 to -7.85 kcal/mol (**Table 2**). Particularly, CID_115196, CID_5317764, and CID_14218027 showed the highest Maestro XP docking scores of -7.85 , -7.70 , and -7.69 kcal/mol respectively, indicating the strongest predicted interactions with the target protein among these ligands. Also, CID_487089 (-6.80 kcal/mol) and CID_174880 (-6.70 kcal/mol) indicated better binding than the reference compound AZD5438 (-5.65 kcal/mol). Thus, these docking results suggested to select these top 5 compounds as the candidate drug molecules for further computational study.

Table 2. Docking scores (Binding Affinity) in kcal/mol of the selected ligands and positive control compound using AutoDock Vina and GLIDE_XP mode.

Compound	Compounds Name	Source	Binding Affinity	Binding Affinity
----------	----------------	--------	------------------	------------------



CID			by AutoDock	by GLIDE_XP
115196	Brassinolide	<i>Camellia sinensis</i>	-9.2	-7.85
5317764	Glycyrrhisoflavon	<i>Glycyrrhiza glabra</i>	-9.5	-7.70
14218027	Licoflavanone	<i>Glycyrrhiza glabra</i>	-9.5	-7.69
487089	3-(3,4-Dihydroxyphenyl)-5,7-dihydroxy-6,8-bis(3-methylbut-2-enyl)chroman-4-one	<i>Glycyrrhiza glabra</i>	-9.2	-6.80
174880	Lactupicrin	<i>Cichorium intybus</i>	-9.4	-6.70
16747683	AZD5438		-8.8	-5.15

3.3.6 Post docking MMGBSA analysis

To evaluate the docking-score-based screening results, we further estimated the binding free energies (ΔG_{bind}) of the selected protein–ligand complexes using the MM-GBSA approach. The analysis revealed that CID_174880 exhibited the most favorable binding energy (-50.44 kcal/mol), indicating the highest predicted binding stability among all evaluated compounds. Similarly, CID_487089 also demonstrated strong binding affinity with a ΔG_{bind} value of -46.78 kcal/mol. CID_14218027 (-36.28 kcal/mol) and CID_5317764 (-36.15 kcal/mol) showed moderate but stable binding energies, comparable to each other. In contrast, CID_115196 exhibited the weakest binding among the tested phytochemicals (-27.81 kcal/mol). Notably, the reference inhibitor AZD5438 displayed a strong binding free energy of -44.58 kcal/mol, consistent with its reported inhibitory activity against the target protein. Importantly, several candidate compounds, particularly CID_174880 and CID_487089 showed more favorable MM-GBSA energies than the control, while others such as CID_14218027 and CID_5317764 demonstrated comparable binding profiles. Overall, these findings suggest that the top-ranked compounds possess binding affinities equal to or stronger than the reference inhibitor, supporting their potential as promising candidates for further stability and dynamic analyses (Table 3). Lastly, based on the above findings we selected top 3 compounds for further analysis to validate their binding strength and stability.

Table 3. Post Docking MM-GBSA Binding Free Energies of Top Docked Compounds Compared to Control (AZD5438)

Compound CID	Compounds Name	Post-Docking MMGBSA (ΔG_{Bind})
174880	Lactupicrin	-50.44



487089	3-(3,4-Dihydroxyphenyl)-5,7-dihydroxy-6,8-bis(3-methylbut-2-enyl)chroman-4-one	-46.78
14218027	Licoflavanone	-36.28
5317764	Glycyrrhisoflavon	-36.15
115196	Brassinolide	-27.81
16747683	AZD5438	-44.58

3.3.7 Assessment of Molecular Interactions between Protein and Ligands

All docked ligands exhibited stable binding within the CDK1 active site through a combination of hydrogen bonding and hydrophobic interactions, closely resembling the interaction behavior of the control compound AZD5438. The control formed four well-oriented hydrogen bonds (THR14, LYS130, ASN133, LEU83) along with diverse hydrophobic contacts, establishing a strong benchmark profile. Among the phytochemicals, 6GU7_487089 showed the most competitive interaction pattern, forming multiple hydrogen bonds with key residues such as LYS130, ASP86, LEU83, and ASN133, along with dense hydrophobic interactions involving ILE10, ALA31, VAL64, and LEU135, indicating a highly stable and well-packed binding mode. Similarly, 6GU7_174880 demonstrated strong anchoring through four hydrogen bonds (LEU83, LYS89, ASP86, GLU81) supported by π -alkyl and alkyl interactions with ALA31, ILE10, and LEU135, although its hydrophobic network was comparatively less extensive than 6GU7_487089. In contrast, 6GU7_14218027 exhibited a balanced but slightly weaker profile, with hydrogen bonds primarily centered on LYS33, LYS89, and GLN132 and hydrophobic interactions involving ALA31, VAL18, and PHE80. Overall, while AZD5438 maintains the most diverse interaction pattern, 6GU7_487089 emerges as the closest competitor in terms of binding stability and interaction density, followed by 6GU7_174880, whereas 6GU7_14218027 shows comparatively moderate interaction strength. These findings suggest that certain phytochemicals, particularly 6GU7_487089, have the potential to mimic the binding efficiency of the control ligand within the CDK1 active site. Full interaction details are provided in **Table S9**, and the corresponding 3D and 2D interaction maps are illustrated in **Figure 7**.



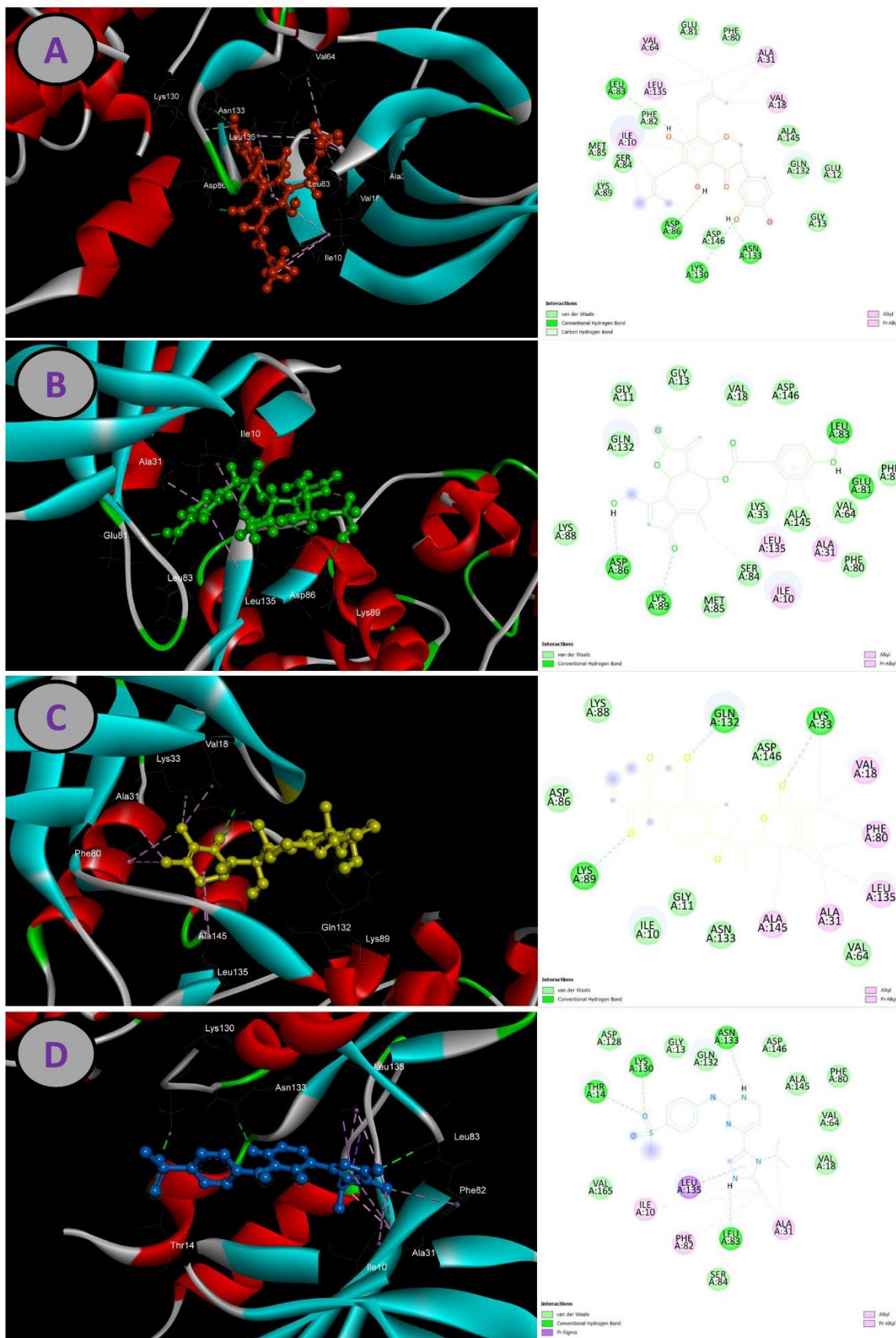


Figure 7. 3D (Left) and 2D (Right) interaction representations of CDK1 complexes: (A) 6GU7_487089, (B) 6GU7_174880, (C) 6GU7_14218027, and (D) Control (AZD5438), highlighting key hydrogen-bonding and hydrophobic interactions within the binding pocket.

3.3.8 Docking Protocol Validation by Re-Docking

Docking protocol validation is typically carried out by re-docking the co-crystallized ligand back into the active site of the target protein to check how reliable the docking method is. The accuracy of this process is measured by calculating the root mean square deviation (RMSD) between the experimentally determined ligand position and the re-docked pose. In this study, the co-crystal ligand AZD5438 was re-docked into the binding pocket of CDK1 (PDB ID: 6GU7) using the Maestro platform. The best docking pose showed an RMSD value of 1.360 Å when superimposed with the original ligand conformation (**Figure 8**). Since this RMSD value is well below the commonly accepted threshold of ≤ 2.0 Å, it indicates that the docking protocol can accurately reproduce the experimentally observed binding mode. Overall, this validation confirms that the docking setup is reliable and suitable for further analysis of binding affinities and molecular interactions.

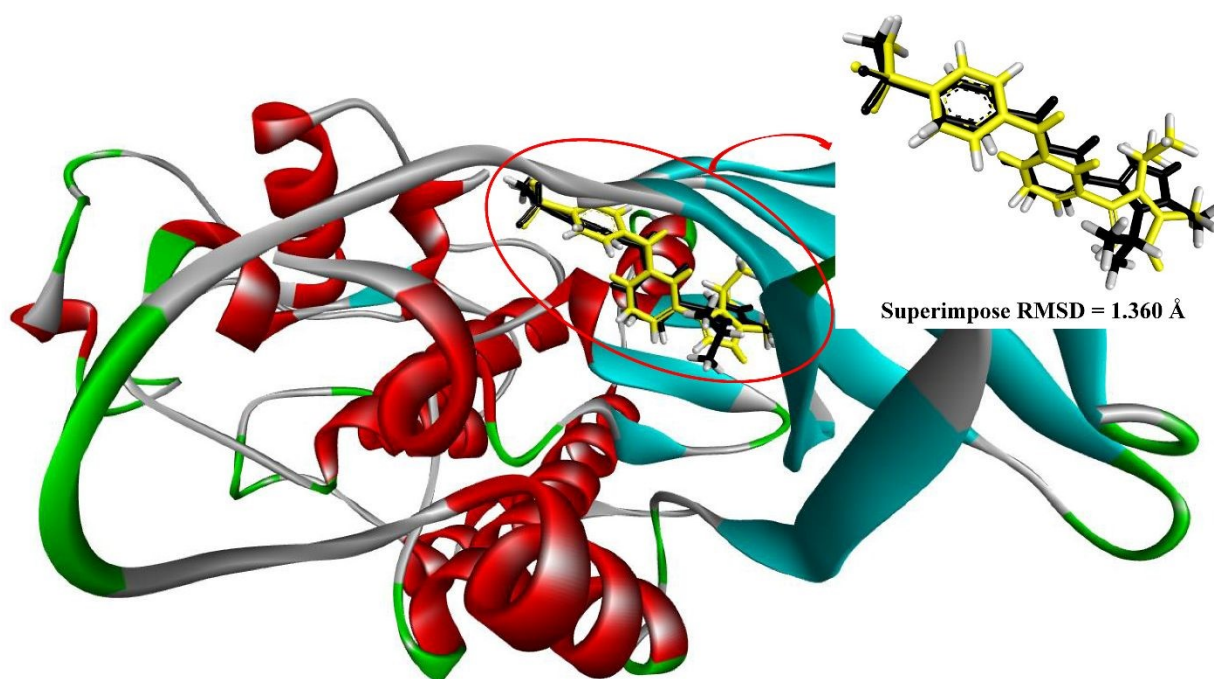


Figure 8. Graphical illustration of the re-docking procedure and superimposition between co-crystallized ligand (green) and the re-docked co-crystal ligand (blue) using GLIDE docking.



3.3.9 Molecular Dynamic Simulation

To make sense of the behavior of all four CDK1 complexes of 6GU7_487089, 6GU7_174880, 6GU7_14218027, and 6GU7_AZD5438 (Control), a 100 ns molecular dynamics (MD) simulation to each of the four complexes was conducted to understand their behavior at a condition similar to the human body. With this simulation we could track the evolution in the structures of the structures with time, and how the ligands remained bound to the protein. In order to have a clear understanding of their stability, we analysed some of the important indicators such as RMSD, RMSF, radius of gyration (Rg), and solvent-accessible surface area (SASA). Besides that, principal component analysis (PCA) and free energy landscape (FEL) mapping were conducted to describe the large-scale motions and found the most stable conformational states that each complex visited throughout the trajectory of 100 ns.

3.3.9.1 Root-mean-square deviation (RMSD)

The root mean square deviation (RMSD) evaluates the structural integrity and conformational shifts in the docked complex [70], with elevated values signaling reduced stability and lower values reflecting stable system performance [97]. The RMSD analysis showed that the apo protein (6GU7) presented the highest average structural deviation (3.245 Å), which is the highest conformational flexibility that is shown in the absence of a ligand. In contrast, all ligand-bound complexes exhibited markedly lower average RMSD values, indicating enhanced structural stability upon ligand binding. Among the screened compounds, CID_487089 was the compound with the lowest average RMSD (2.446 Å), which was then followed by CID_14218027 (2.470 Å), indicating that the compounds have a strong stabilizing interaction in the binding pocket. CID_174880 showed comparable stability to CID_14218027, with an average RMSD of 2.476 Å. The control ligand showed an average RMSD of 3.002 Å, comparable to the apo form, implying minimal stabilizing influence. Overall, the average RMSD results confirm that CID_487089 forms the most stable complex with 6GU7 during the MD simulation (**Figure 9A**).

3.3.9.2 Root mean square fluctuation (RMSF)

The residue-level dynamics of the system were examined through RMSF analysis over the equilibrated phase of the simulation. As expected, the apo protein showed the highest average fluctuation (1.571 Å), reflecting its greater structural freedom in the absence of a bound ligand. In contrast, all ligand-bound complexes exhibited noticeably lower average RMSF values, suggesting that ligand interaction contributes to a more stable and ordered residue environment. Among the screened compounds, CID_174880 exhibited the strongest stabilizing effect, as reflected by its lowest average RMSF value (1.181



Å), suggesting tight binding and effective restriction of residue mobility. The control compound showed a moderate decrease in flexibility (1.390 Å), indicating a reasonable but less pronounced stabilizing interaction. Meanwhile, CID_487089 and CID_14218027 displayed intermediate stabilization, with average RMSF values of 1.253 Å and 1.375 Å, respectively, implying balanced binding interactions that reduce flexibility without overly restricting protein dynamics (**Figure 9B**).

3.3.9.3 Solvent-accessible surface area (SASA)

SASA reflects the solvent-exposed surface of a protein, with lower values generally indicating greater ligand burial and potentially stronger binding [98]. Analysis of the 6GU7 protein complexes revealed notable differences in solvent-accessible surface area (SASA) among the ligands. The 6GU7_487089 complex exhibited the highest average SASA value (148.17 Å²), indicating a relatively larger solvent-exposed surface and suggesting a less compact binding conformation. In contrast, the 6GU7_174880 complex showed the lowest average SASA value (126.15 Å²), implying that the ligand is more deeply buried within the binding pocket and forms a more compact and stable protein–ligand interface. The 6GU7_14218027 (130.14 Å²) and control complex (131.81 Å²) displayed intermediate SASA values, reflecting moderately compact structural arrangements with balanced solvent exposure. Overall, 6GU7_174880 demonstrated the most favorable SASA profile in terms of structural compactness, highlighting its potential as the most effective stabilizing ligand among the tested compounds (**Figure 9C**).

3.3.9.4 The radius of gyration (Rg)

The radius of gyration (Rg) is a key indicator of protein structural compactness and stability in molecular dynamics simulations. In this study, the control complex (6GU7_Control) exhibited a relatively high average Rg value (4.60 Å), indicating a less compact conformation. Similarly, the 6GU7_487089 complex showed the highest Rg value among the ligand-bound systems (4.69 Å), suggesting comparatively lower compactness and weaker structural stabilization. In contrast, 6GU7_14218027 (4.40 Å) and 6GU7_174880 (4.44 Å) displayed lower average Rg values, reflecting more compact and stable conformations. Among these, 6GU7_14218027 exhibited the lowest Rg value, indicating the highest degree of structural compactness and suggesting a stronger stabilizing effect on the protein. Overall, ligand binding generally enhanced protein compactness compared to the control, with 6GU7_14218027 emerging as the most effective stabilizer based on Rg analysis, followed closely by 6GU7_174880 (**Figure 9D**).



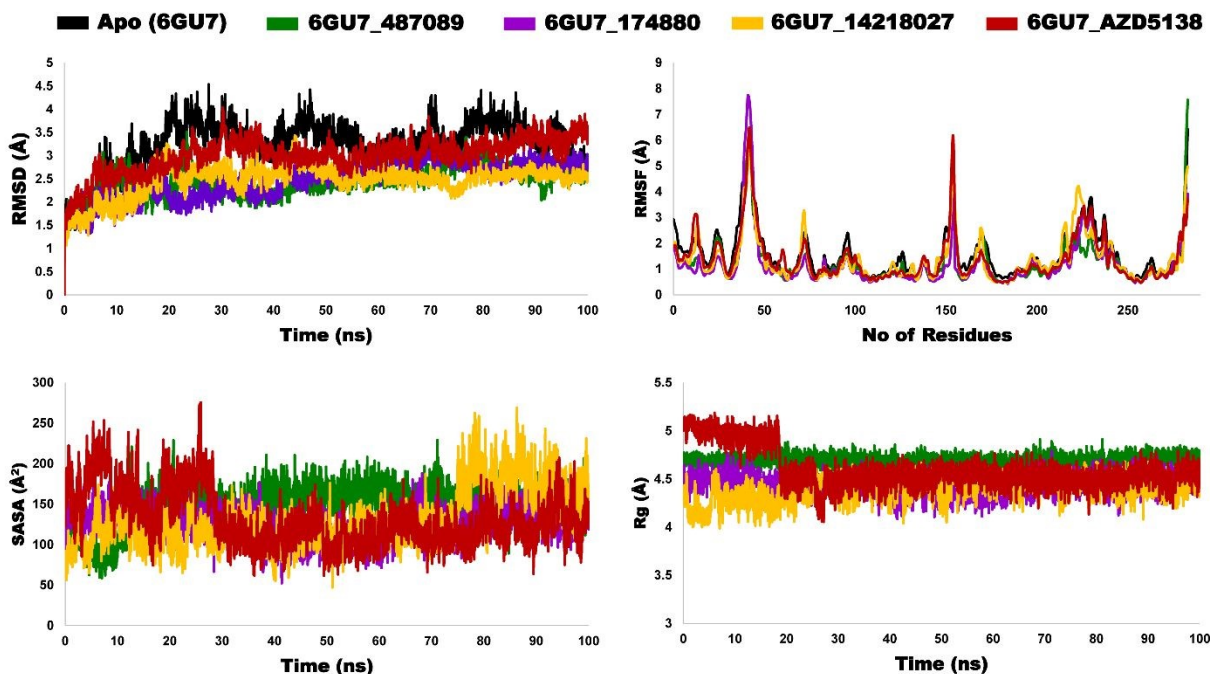


Figure 9. Molecular dynamics analysis of CDK1 (PDB ID: 6GU7) in complex with three test ligands (6GU7_487089, 6GU7_174880, and 6GU7_14218027) and the reference ligand (6GU7_Control) over 100 ns. (A) RMSD, (B) RMSF, (C) SASA, and (D) Rg profiles.

3.3.10 Post simulation MMGBSA analysis

The post simulation MMGBSA analysis of 6GU7_487089, 6GU7_174880, 6GU7_14218027, and the reference ligand complex (6GU7_AZD5438) against CDK1 (PDB ID: 6GU7) were analyzed over the 100 ns MD simulation (**Figure 10**). The 6GU7_487089 complex exhibited the most favorable binding free energy (-40.29 kcal mol⁻¹), indicating the strongest and most stable interaction with the target protein. This was followed by 6GU7_174880 (-36.06 kcal mol⁻¹), which also demonstrated a relatively strong binding affinity. In contrast, 6GU7_14218027 (-29.13 kcal mol⁻¹) and the control complex (-29.53 kcal mol⁻¹) showed comparatively higher ΔG values, suggesting weaker binding interactions. Notably, the binding affinity of 6GU7_14218027 was slightly lower than that of the control, indicating limited improvement over the reference ligand. Overall, 6GU7_487089 emerged as the most promising candidate with the highest binding affinity, followed by 6GU7_174880. These findings are consistent with their favorable stability profiles observed in molecular dynamics simulations, further supporting their potential as effective inhibitors.



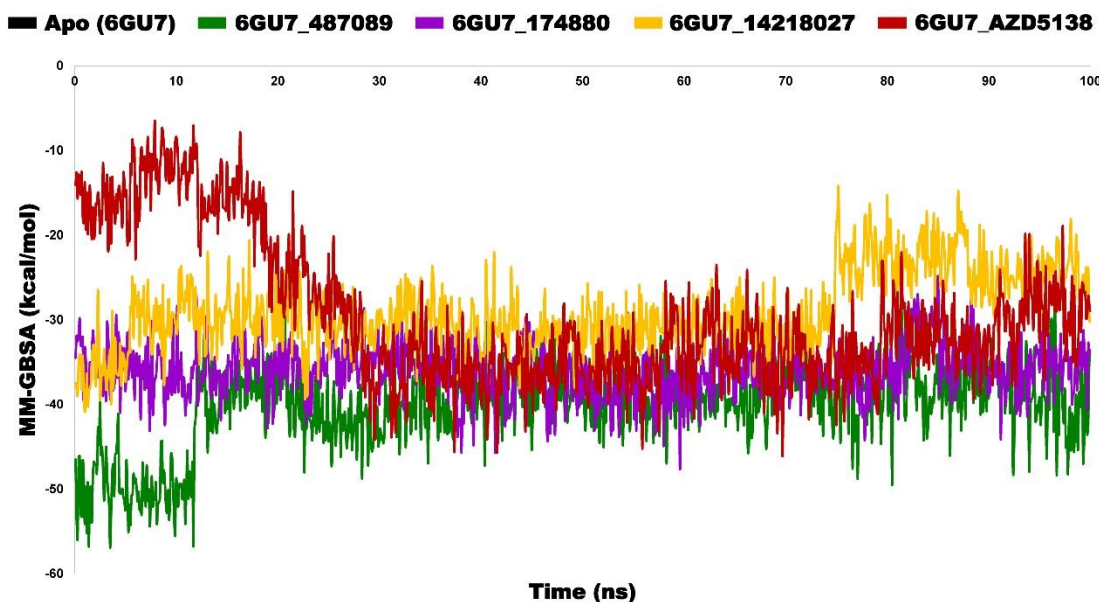


Figure 10. Comparative MM/GBSA Free Energy Analysis of Ligand-Bound and Control CDK1 Systems

3.3.11 Principal component analysis

To understand how each ligand influenced the overall motion of the protein, we performed a Principal Component Analysis (PCA). Among the studied complexes, 6GU7_174880 showed the highest contribution in PC1 (44.97%), which was very close to the control complex, 6GU7_AZD5438 (44.18%). This suggests that most of the protein's motion in the presence of ligand CID_174880 is mainly concentrated in one dominant direction. Such a pattern usually indicates a more stable and well-defined conformational movement, similar to that observed for the reference inhibitor. On the other hand, 6GU7_487089 and 6GU7_14218027 showed lower PC1 values, with 31.37% and 30.43%, respectively. At the same time, these two complexes had relatively higher contributions in PC2 and PC3, especially PC2 values of 16.13% and 17.88%. This may indicate that the protein motion is spread across more than one direction, suggesting comparatively higher flexibility and the possibility of multiple conformational changes during the simulation period. For the control complex, the contributions of the first three principal components were 44.18% (PC1), 10.85% (PC2), and 7.21% (PC3), which reflects a stable dynamic behavior of the protein–ligand complex. Interestingly, the PCA profile of 6GU7_174880 was found to be quite similar to that of the control, which may support its potential as a promising ligand with favorable binding stability. In contrast, the higher PC2 and PC3 values observed for 6GU7_487089 and 6GU7_14218027 suggest that these complexes may have relatively more flexible interaction patterns with the target protein (Figure S4).



3.3.12 Gibbs free energy landscape (FEL) analysis

The Gibbs free energy landscape (FEL) analysis was carried out to explore the conformational stability and energy minima of the 6GU7 protein in complex with the selected ligands and the reference control. In the FEL plots, the dark blue regions represent the lowest energy states, indicating the most stable conformations sampled during the simulation. Among the studied complexes, 6GU7_174880 (**Figure S5B**) and the control complex, 6GU7_AZD5438 (**Figure S5D**), showed a well-defined and deep energy basin with a comparatively compact distribution. This suggests that both complexes remained in a more stable conformational state throughout the simulation and experienced less conformational fluctuation. The similarity between the ligand CID_174880 and the control indicates that this compound may induce a stable binding conformation comparable to the reference inhibitor. In contrast, 6GU7_487089 (**Figure S5A**) displayed a broader low-energy basin with a relatively wider spread across the conformational space. This may indicate that the complex explored multiple conformational states during the simulation, suggesting a comparatively more flexible dynamic behavior. Similarly, 6GU7_14218027 (**Figure S5C**) also showed an extended energy basin with noticeable spreading of the low-energy region. Such a pattern may reflect the presence of multiple metastable conformations and greater structural flexibility compared with the control complex. Overall, the FEL results suggest that 6GU7_174880 exhibited a more stable conformational landscape that closely resembles the control, whereas 6GU7_487089 and 6GU7_14218027 showed relatively broader energy minima, indicating more flexible binding-associated motions (**Figure S5**).

3.4 Pharmacokinetics Study

3.4.1 Physicochemical and ADME properties prediction

The physicochemical and ADME properties of the selected phytochemicals were evaluated and compared with the reference control compound, AZD5438, to assess their drug-likeness and pharmacokinetic suitability. The molecular weight of all selected compounds was found within an acceptable range for oral drug candidates, ranging from 340.37 to 424.49 g/mol, which is comparable to the control (371.46 g/mol). Similarly, the topological polar surface area (TPSA) values of the phytochemicals (86.99–110.13 Å²) were also within the favorable range, suggesting good membrane permeability and absorption potential. The consensus LogP values varied among the compounds, where CID_487089 (4.35) showed relatively higher lipophilicity compared to the control (2.51), while CID_174880 (1.87) and CID_14218027 (3.33) remained within an acceptable range. All selected



phytochemicals showed high gastrointestinal (GI) absorption, similar to the control compound, indicating their potential suitability for oral administration. In addition, CID_487089 and CID_14218027 were predicted as non-substrates of P-glycoprotein (P-gp), similar to the control, whereas CID_174880 was identified as a P-gp substrate, which may slightly influence its cellular efflux behavior. None violated Lipinski rule of five, Vebers rule, Egan rule, or Ghoses filter (0 violations of all compounds) and this once again indicates that they had good drug-likeness properties [99–102]. None of the compounds, including the control, were predicted to be BBB permeant, which may reduce the possibility of central nervous system-related side effects [103]. The bioavailability score was 0.55 for all selected compounds, which is comparable to the control and suggests favorable oral bioavailability. Regarding structural alerts, CID_14218027 demonstrated a more favorable profile with no PAINS alert and only one Brenk alert, whereas CID_487089 and CID_174880 showed relatively higher alert counts. Among the evaluated compounds, CID_14218027 exhibited broad-spectrum inhibition across all major CYP isoforms (CYP1A2, CYP2C19, CYP2C9, CYP2D6, and CYP3A4), indicating a high likelihood of metabolic interference and an increased risk of drug–drug interactions. In contrast, CID_487089 showed selective inhibition of CYP2C9 and CYP3A4, suggesting a comparatively moderate interaction risk. Notably, CID_174880 demonstrated no inhibitory activity against any of the assessed CYP enzymes, reflecting a more favorable metabolic profile with minimal risk of CYP-mediated interactions. The control compound also inhibited CYP2C9 and CYP3A4, further emphasizing that compounds with minimal or no CYP inhibition, such as CID_174880, are more desirable from a pharmacokinetic and safety perspective. The detailed physicochemical and ADME parameters of all compounds are summarized in **Table S10**, and their overall ADME profiles are visualized in **Figure 11**.



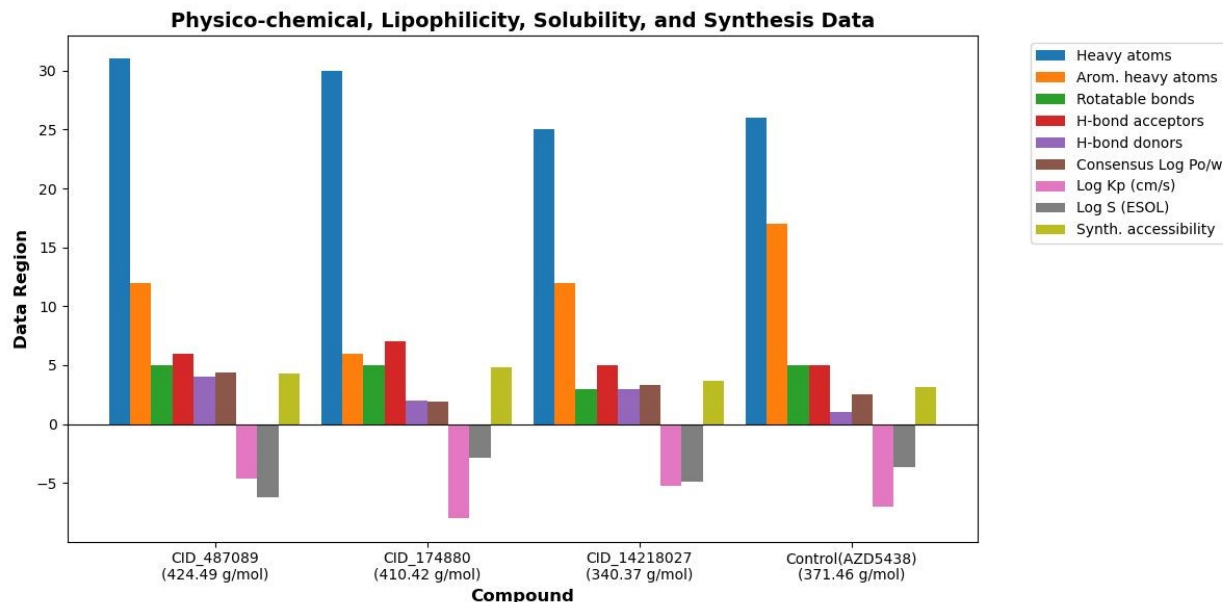


Figure 11. Physicochemical and ADME-Related Property Comparison of CID_487089, CID_174880, CID_14218027 and Control AZD5438.

3.4.2 Toxicity Analysis

The toxicity assessment revealed that all three selected phytochemicals exhibited a favorable safety profile with no predicted hepatotoxicity, neurotoxicity, cardiotoxicity, mutagenicity, or cytotoxicity. However, all compounds showed potential respiratory toxicity, similar to the reference drug AZD5438, suggesting a possible target-related adverse effect. Notably, the phytochemicals demonstrated immunotoxicity, whereas the control compound was inactive in this regard, indicating a limitation of the selected ligands. Importantly, none of the compounds were predicted to be carcinogenic, in contrast to the control drug, highlighting a significant safety advantage. Based on LD₅₀ values and toxicity classification, CID_14218027 emerged as the safest candidate with low acute toxicity (Class 4), whereas CID_487089 showed high toxicity (Class 2), making it less suitable for further development. (Table 4). Therefore, comprehensive experimental validation is essential to confirm the safety and therapeutic applicability of these compounds.

Table 4. Predicted Toxicity Profiles of Selected Lung Cancer Drug Candidates (CID_487089, CID_174880, CID_14218027) and Control (AZD5438)



Target	CID_487089	CID_174880	CID_14218027	Control (AZD5438)
Hepatotoxicity	Inactive	Inactive	Inactive	Inactive
Neurotoxicity	Inactive	Inactive	Inactive	Inactive
Cardiotoxicity	Inactive	Inactive	Inactive	Inactive
Respiratory toxicity	Active	Active	Active	Active
Immunotoxicity	Active	Active	Active	Inactive
Carcinogenicity	Inactive	Inactive	Inactive	Active
Mutagenicity	Inactive	Inactive	Inactive	Inactive
Cytotoxicity	Inactive	Inactive	Inactive	Inactive
LD ₅₀ mg/kg	10	300	2000	500
Toxicity Class	2	3	4	4

4. Discussion

Lung cancer is one of the most detrimental malignancies in the entire world. It grows fast, is often detected late, and does not respond well to available treatments [2]. These challenges highlight the need for new treatments guided by a clear strategy for finding targets and developing drugs. Using transcriptomic profiling, network analysis, pathway studies, and computer-based drug discovery, this research systematically identified a key molecular target in lung cancer and proposed a drug to target it. Transcriptomics is the study of all RNA molecules produced in a cell, helping us understand which genes are active and how their activity changes in different conditions [104]. In cancer research, it's often used to find important targets because cancer cells tend to show unusual patterns of gene activity compared to normal cells [105]. However, results from a single dataset can be affected by technical differences, experimental platforms, or variations among populations. To address this, four independent GEO microarray datasets (GSE19804, GSE10072, GSE18842, and GSE10799) were analyzed across different patient groups (case and control). From this combined analysis, 88 genes were consistently upregulated and 294 genes downregulated. Repeated changes in these certain genes suggest a real connection to this cancer. After identifying these common differentially expressed genes, we studied how their proteins



interact. The strong enrichment of the protein-protein interaction (PPI) network indicates that these connections are not random, showing a real biological relationship between these genes in lung cancer. Then, based on this PPI network we identified the major key hub genes (kHGs) using several analytical methods, which ensured that our results were robust and not dependent on a single approach. Among these interactions, the cyclin dependent kinase 1 (CDK1) showed the highest number of interactions, suggesting it may play a central role in regulating important cancer-related processes. To clarify kHGs regulation, we analyzed their interactions with transcription factors (TFs) and microRNAs (miRNAs). Network analysis (**Figure 4**) identified FOXC1, GATA2, YY1, E2F1, and HINFP as major TFs, while hsa-miR-192-5p, hsa-miR-92a-3p, hsa-miR-193b-3p, hsa-miR-215-5p, and hsa-miR-155-5p were identified as central miRNA regulators. The study of TFs and miRNAs revealed important hints on the regulation of kHGs at the different levels. Transcription factors are known to activate the expression of genes [106], whereas the miRNAs are known to regulate the activity of genes once they are transcribed [107]. This combined method enabled us to determine vital regulatory factors which could contribute to dysregulated gene expression in lung cancer and gives a clear understanding to further functional and therapeutic studies. Moreover, we performed functional enrichment analyses (GO and KEGG) using multiple pathway databases to understand the role of these hub genes in lung cancer cells. These enrichment results indicate that lung cancer progression is driven by widespread disruption of cell cycle regulation, mitotic structure, and signaling control rather than by isolated gene alterations. The overall involvement of CDK1 in key mitotic events, strong kinase activity, and extensive interaction with other cell cycle regulators suggest that it functions as a major coordinator of uncontrolled cell division in lung cancer. Thus, this central and recurring role makes CDK1 the most promising candidate for further functional validation and therapeutic exploration.

The reason why LUAD was chosen as a validation is because it is the most prevalent type of lung cancer and it is a substantial percentage of the cases of lung cancer, as opposed to other types [108]. Multiple analyses of expressions in different independent platforms indicated that not only is CDK1 transcriptionally upregulated, but also overexpressed at the protein level in LUAD. In addition, the patients whose CDK1 level was higher had poorer survival and earlier recurrence, which is an obvious clinical effect. In addition to expression, the correlation of CDK1 and immune cell infiltration indicates that its activity moves beyond the expression to influence the tumor microenvironment, especially at the T-cell subsets and macrophage polarization. This implies that CDK1 can regulate tumor growth as well as immune response. Lastly, the somatic mutations and copy number changes observed on CDK1 contribute to the oncogenic significance of this protein in a group of LUAD patients. From a biological perspective,



CDK1 is more than a statistically prioritized hub gene. As a key regulator of the G2/M transition and mitotic entry, CDK1 promotes sustained tumor cell proliferation when aberrantly activated. Prior studies in lung cancer have also linked elevated CDK1 expression with poor survival, enhanced cell-cycle and DNA-repair signaling, and altered immune-associated pathways in LUAD [109,110]. These observations strengthen the view that CDK1 may act as a functionally important driver of lung cancer progression and a rational candidate for therapeutic targeting.

In this study, we used a structure-based drug discovery approach to efficiently screen a large library of phytochemicals [111]. In this field, the integration of machine learning (ML), molecular docking, and molecular dynamics (MD) simulation studies has transformed the identification and optimization of novel drug candidates [112,113]. We performed cheminformatics-based screening procedure to screen physicochemical and drug-likeness properties of 9,577 phytochemicals. After removing duplicate entries and unwanted sub-structure finally 1802 from 2,113 phytochemicals that have no violation according RO5 and Veber's rule. ML is revolutionizing the computational drug discovery approach to reduce the traditional experimental time and cost. ML is the best option to calculate the bioactivity of a large dataset with high accuracy within a very short time [114]. To identify the most promising candidates, we applied a ML-based pIC50 prediction model, which helped us identify compounds with higher chances of biological activity. This step reduced the dataset to identify most potential 380 phytochemicals with pIC50 >6.5 for further computational analysis including molecular docking, molecular dynamic simulation, and pharmacokinetics analysis. The docking results and post docking MM-GBSA binding energy (ΔG) highlights top-ranked three phytochemicals (CID_487089, CID_174880, and CID_14218027) as candidate drug molecules.

Furthermore, the molecular dynamics (MD) simulations were used to understand how the selected phytochemicals interact with CDK1 over time. While the apo protein and AZD5438-bound complex showed noticeable structural fluctuations, all phytochemical-bound systems exhibited improved stability throughout the 100 ns simulation. Notably, CID_174880 consistently outperformed the control compound by maintaining lower RMSD and RMSF values (**Figure 9**), indicating a stronger and more stable binding mode. Structural compactness analyses further highlighted this difference. The CID_174880 complex showed the lowest SASA and a tightly maintained Rg, reflecting deeper ligand burial and a more compact protein structure than AZD5438 [115]. In contrast, the control complex remained relatively solvent-exposed and flexible, suggesting weaker stabilization of CDK1. Post-simulation MM-GBSA results further supported these findings, with CID_487089 showing the most favorable binding free energy (-40.29 kcal mol⁻¹), followed by CID_174880 (-36.06 kcal mol⁻¹), both outperforming the control



AZD5438 (-29.53 kcal mol⁻¹). Dynamic motion analysis using PCA showed clear differences between the control compound (AZD5438) and the phytochemicals [116], which suggests a better fit between the ligand and the binding site. Also, the FEL analysis further supported these findings by demonstrating the stability and conformational states of the ligand–protein interactions. Among the studied compounds, CID_174880 showed a dynamic and energy profile closely comparable to the control, suggesting favorable binding stability. The control compound was predicted to be carcinogenic, while none of the phytochemicals showed this risk, indicating better overall safety. However, all compounds exhibited some potential respiratory toxicity, and the phytochemicals also showed immunotoxic effects, which may limit their development. Based on LD₅₀ values and toxicity classes, CID_14218027 was the safest (Class 4), whereas CID_487089 showed higher toxicity (Class 2). CID_174880 displayed a toxicity profile close to the control but with slightly improved safety. Overall, CID_174880 demonstrated promising inhibitory potential with comparatively better safety profiles than the control compound. Therefore, this study could be a useful resource to identify natural CDK1 inhibitors after validating through the experimental (*in vivo* & *in vitro*) study.

5. Conclusion

This paper displays CDK1 as an important lung cancer driver, disrupting cell-cycle, p53 pathway, and immune microenvironment changes, such as augmented M1 macrophages, diminished M2 polarization, and unfavorable prognosis owing to overexpression and mutations, including I136N and E57V. These discoveries indicate a shared vulnerability in lung cancer, which puts mitotic dysregulation as one of the prospective treatment targets. Among the phytochemicals evaluated, Lactupicrin (CID_174880), derived from *Cichorium intybus*, emerged as the most promising *in silico* prioritized candidate compared with the reference compound AZD5438, based on its favorable docking performance, structural stability, dynamic behavior, and predicted pharmacokinetic and toxicity profiles. Although minor fluctuations were observed in some metrics, all CDK1 complexes maintained overall stability and compatibility. Docking results, MD simulations, MM-GBSA calculations and ADMET analyses consistently support Lactupicrin as a computationally prioritized candidate for further development. Further experimental validation, including *in vitro* biochemical assays, cell-based functional studies, and *in vivo* investigations, is required to confirm CDK1 inhibitory activity, anticancer efficacy, and safety. Overall, this research provides a clear workflow for identifying effective CDK1-targeting compounds and offers a roadmap for developing broad-spectrum, mechanism-based therapies for lung cancer, bridging computational predictions with potential real-world applications.



6. Limitations of this study

Although this study applied a comprehensive multi-omics and in silico approach, several limitations should be considered. First, the transcriptomic analysis was based solely on publicly available microarray datasets, which may be affected by batch effects, platform-related biases, and limited clinical information. While using four independent cohorts improves reliability, the absence of RNA-seq data may limit the ability to capture finer gene expression variability. Second, the drug discovery pipeline, including machine learning predictions, molecular docking, MM-GBSA analysis, MD simulations, and ADEMT, relies entirely on computational models. Although these methods are useful for prioritizing candidates, they cannot fully represent the complexity of biological systems, tumor heterogeneity, or real pharmacodynamic behavior. Finally, the identified drug candidates, especially Lactupicrin, have not yet been validated in in vitro or in vivo models. Therefore, the therapeutic potential, optimal dosage, and safety profiles of the candidate drugs still need to be confirmed through experimental validation.

7. Biographical Note

Md. Ahad Ali is a computational chemist and bioinformatics researcher affiliated with Panacea Research Center and the University of Rajshahi, where he works on transcriptomics, machine learning, and structure-based drug discovery.

Acknowledgements

Not applicable

Funding

This research has not received any funding

Institutional Review Board Statement

Not applicable.

Informed Consent Statement

Not applicable.

Competing Interest

All authors declare no conflict of interests

Data Availability Statement

The original data and contributions presented in this study are included in the article and its Supplementary Information (Table S1-S10 and Figure S1-S5). The training set of our selected compounds,



the test set, and the chEMBL datasets (containing the consensus predictions) and related python code for running the qsar models to predict the bioactivity of the selected compounds, can be found on our GitHub repository (https://github.com/ahad004/LUAD_ML_QSAR_Modeling), or can be accessed using the following <https://doi.org/10.5281/zenodo.19841200>.

Author's contribution

Conceptualization – Md. Ahad Ali and Hridhhi Sarker; **Methodology** –Md. Ahad Ali, Hridhhi Sarker, and Humaira Sheikh; **Data curation** – Hridhhi Sarker, Marguba Kamrun, Bilkis Shifa, and Sujoy Banik; **Formal analysis** – Md. Ahad Ali, Hridhhi Sarker, Humaira Sheikh, and Siam Ahmed; **Visualization** – Hridhhi Sarker, Md. Ahad Ali, Marguba Kamrun, and Tarikul Islam; **Validation** – Md. Ahad Ali, Hridhhi Sarker, and Tarikul Islam; **Project administration** – Md. Ahad Ali; **Software & Resources** – Neeraj Kumar, Sujoy Banik, and Bilkis Shifa; **Supervision** – Md. Ahad Ali; **Writing – original draft**, Hridhhi Sarker, Md. Ahad Ali, and Humaira Sheikh; **Writing – review & editing**, Md. Ahad Ali, Marguba Kamrun, and Neeraj Kumar.

8. References

1. Ji, Y.; Zhang, Y.; Liu, S.; Li, J.; Jin, Q.; Wu, J.; Duan, H.; Liu, X.; Yang, L.; Huang, Y. The Epidemiological Landscape of Lung Cancer: Current Status, Temporal Trend and Future Projections Based on the Latest Estimates from GLOBOCAN 2022. *J. Natl. Cancer Cent.* **2025**, *5*, 278–286, doi:10.1016/j.jncc.2025.01.003.
2. Bray, F.; Laversanne, M.; Sung, H.; Ferlay, J.; Siegel, R.L.; Soerjomataram, I.; Jemal, A. Global Cancer Statistics 2022: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries. *CA. Cancer J. Clin.* **2024**, *74*, 229–263, doi:10.3322/caac.21834.
3. Zhou, J.; Xu, Y.; Liu, J.; Feng, L.; Yu, J.; Chen, D. Global Burden of Lung Cancer in 2022 and Projections to 2050: Incidence and Mortality Estimates from GLOBOCAN. *Cancer Epidemiol.* **2024**, *93*, 102693, doi:10.1016/j.canep.2024.102693.
4. Bharadwaj, S.; Mierzwicka, J.M.; Vaňková, L.; Malý, P. Unraveling the Molecular-Pathological Characteristics and Cellular Complexity of the Tumor Immune Microenvironment in Metastatic Non-Small Cell Lung Cancer. *Cell Commun. Signal.* **2025**, *23*, 400, doi:10.1186/s12964-025-02410-w.
5. Gu, Z.; Heng, Y.; Fan, R.; Luo, J.; Ju, L. Single-Cell RNA Sequencing Reveals Cellular and Molecular Heterogeneity in Extensive-Stage Small Cell Lung Cancer with Different Chemotherapy Responses. *Cancer Cell Int.* **2025**, *25*, 157, doi:10.1186/s12935-025-03785-z.
6. Jachowski, A.; Marcinkowski, M.; Szydłowski, J.; Grabarczyk, O.; Nogaj, Z.; Marcin, Ł.; Pławski, A.; Jagodziński, P.P.; Słowikowski, B.K. Modern Therapies of Nonsmall Cell Lung Cancer. *J. Appl. Genet.* **2023**, *64*, 695–711, doi:10.1007/s13353-023-00786-4.
7. Araghi, M.; Mannani, R.; Heidarnejad maleki, A.; Hamidi, A.; Rostami, S.; Safa, S.H.; Faramarzi, F.; Khorasani, S.; Alimohammadi, M.; Tahmasebi, S.; et al. Recent Advances in Non-Small Cell Lung Cancer Targeted Therapy; an Update Review. *Cancer Cell Int.* **2023**, *23*, 162, doi:10.1186/s12935-023-02990-y.



8. Su, P.-L.; Furuya, N.; Asrar, A.; Rolfo, C.; Li, Z.; Carbone, D.P.; He, K. Recent Advances in Therapeutic Strategies for Non-Small Cell Lung Cancer. *J. Hematol. Oncol.* **2025**, *18*, 35, doi:10.1186/s13045-025-01679-1.
9. Valdez Capuccino, L.; Kleitke, T.; Szokol, B.; Svajda, L.; Martin, F.; Bonechi, F.; Krekó, M.; Azami, S.; Montinaro, A.; Wang, Y.; et al. CDK9 Inhibition as an Effective Therapy for Small Cell Lung Cancer. *Cell Death Dis.* **2024**, *15*, 345, doi:10.1038/s41419-024-06724-4.
10. Osoegawa, A.; Takumi, Y.; Hashimoto, T.; Nakatsuji, S.; Hori, M.; Sakai, M.; Karashima, T.; Abe, M.; Miyawaki, M.; Sugio, K. Cyclin-Dependent Kinase (CDK) 4/6 Inhibition in Non-Small Cell Lung Cancer with Epidermal Growth Factor Receptor (EGFR) Mutations. *Invest. New Drugs* **2023**, *41*, 183–192, doi:10.1007/s10637-023-01337-8.
11. Panagiotou, E.; Gomatou, G.; Trontzas, I.P.; Syrigos, N.; Kotteas, E. Cyclin-Dependent Kinase (CDK) Inhibitors in Solid Tumors: A Review of Clinical Trials. *Clin. Transl. Oncol.* **2022**, *24*, 161–192, doi:10.1007/s12094-021-02688-5.
12. Huang, X.; Yin, Y.; Saha, G.; Francis, I.; Saha, S.C. A Comprehensive Numerical Study on the Transport and Deposition of Nasal Sprayed Pharmaceutical Aerosols in a Nasal-To-Lung Respiratory Tract Model. *Part. Part. Syst. Character.* **2025**, *42*, doi:10.1002/ppsc.202400004.
13. Li, X.-Q.; Cheng, X.-J.; Wu, J.; Wu, K.-F.; Liu, T. Targeted Inhibition of the PI3K/AKT/MTOR Pathway by (+)-Anthrabenoxocinone Induces Cell Cycle Arrest, Apoptosis, and Autophagy in Non-Small Cell Lung Cancer. *Cell. Mol. Biol. Lett.* **2024**, *29*, 58, doi:10.1186/s11658-024-00578-6.
14. Kitai, H.; Choi, P.H.; Yang, Y.C.; Boyer, J.A.; Whaley, A.; Pancholi, P.; Thant, C.; Reiter, J.; Chen, K.; Markov, V.; et al. Combined Inhibition of KRASG12C and MTORC1 Kinase Is Synergistic in Non-Small Cell Lung Cancer. *Nat. Commun.* **2024**, *15*, 6076, doi:10.1038/s41467-024-50063-z.
15. Huang, J.-L.; Wu, L.-M.; Wu, S.-Q.; Yuan, F.-Y.; Weng, H.-Z.; Huang, D.; Gan, L.; Chen, S.-B.; Tang, G.-H.; Yin, S. A Small Molecule Targets LIC1 to Suppress Lung Tumor Growth by Inducing Autophagy. *Nat. Chem. Biol.* **2025**, doi:10.1038/s41589-025-02040-w.
16. Crawford, J.; Herndon, D.; Gmitter, K.; Weiss, J. The Impact of Myelosuppression on Quality of Life of Patients Treated with Chemotherapy. *Future Oncol.* **2024**, *20*, 1515–1530, doi:10.2217/fon-2023-0513.
17. Lazzari, C.; Gregorc, V.; Karachaliou, N.; Rosell, R.; Santarpia, M. Mechanisms of Resistance to Osimertinib. *J. Thorac. Dis.* **2020**, *12*, 2851–2858, doi:10.21037/jtd.2019.08.30.
18. Astolfi, L.; Ghiselli, S.; Guaran, V.; Chicca, M.; Simoni, E.; Olivetto, E.; Lelli, G.; Martini, A. Correlation of Adverse Effects of Cisplatin Administration in Patients Affected by Solid Tumours: A Retrospective Evaluation. *Oncol Rep* **2013**, *29*, 1285–1292, doi:10.3892/or.2013.2279.
19. Alsatari, E.S.; Smith, K.R.; Galappaththi, S.P.L.; Turbat-Herrera, E.A.; Dasgupta, S. The Current Roadmap of Lung Cancer Biology, Genomics and Racial Disparity. *Int. J. Mol. Sci.* **2025**, *26*, 3818, doi:10.3390/ijms26083818.
20. Gupta, G.; Samuel, V.P.; M., R.M.; Rani, B.; Sasikumar, Y.; Nayak, P.P.; Sudan, P.; Goyal, K.; Oliver, B.G.; Chakraborty, A.; et al. Caspase-Independent Cell Death in Lung Cancer: From Mechanisms to Clinical Applications. *Naunyn. Schmiedebergs. Arch. Pharmacol.* **2025**, *398*, 13031–13048, doi:10.1007/s00210-025-04149-0.
21. Zhang, J.; Zeng, X.; Guo, Q.; Sheng, Z.; Chen, Y.; Wan, S.; Zhang, L.; Zhang, P. Small Cell Lung Cancer: Emerging Subtypes, Signaling Pathways, and Therapeutic Vulnerabilities. *Exp. Hematol. Oncol.* **2024**, *13*, 78, doi:10.1186/s40164-024-00548-w.
22. MISHRA, N.; SONI, A.; KUMARI, M.; SINGH, G.; SHARMA, S.K.; SINGH, S.K. Targeting Cell Cycle Regulators: A



- New Paradigm in Cancer Therapeutics. *BIOCELL* **2024**, *48*, 1639–1666, doi:10.32604/biocell.2024.056503.
23. Alibakhshi, A.; Alagheband Bahrami, A.; Mohammadi, E.; Ahangarzadeh, S.; Mobasheri, M. In-Silico Design of a New Multi-Epitope Vaccine Candidate against SARS-CoV-2. *Acta Virol.* **2024**, *67*, doi:10.3389/av.2023.12481.
24. Abdel Razek, F.S.; Ibrahim, S.D.; Megahed, A.A.; Sadik, A.S.; El-Masry, S.S.A. In Silico Molecular Docking Analysis of Green Tea Bioactive Compounds Targeting Banana Bunchy Top Virus Proteins. *Discov. Appl. Sci.* **2025**, *7*, 1232, doi:10.1007/s42452-025-07466-4.
25. Jamal, Q.M.S.; Khan, S.; Khan, M.; Ansai, A.A.; Ashraf, J.M.; Habibullah, M.; Farasani, A.; Madkhali, A.M.; Lohani, M. Smoking May Increase the Risk of COVID-19 Infection: Evidence from In Silico Analysis. *J. Pharm. Res. Int.* **2021**, 12–21, doi:10.9734/jpri/2021/v33i22B31394.
26. Sajid, Z.; Akhtar, T.; Ahmad, K.; Haroon, M. Molecular Docking Simulation and ADMET/Pharmacokinetic Screening of Newly Designed 2-(2-(Aryl)-4-oxo-4,5-dihydrothiazol-5-yl)Acetohydrazides as Potential Antitubercular Agents. *ChemistrySelect* **2024**, *9*, doi:10.1002/slct.202403715.
27. Chunarkar-Patil, P.; Kaleem, M.; Mishra, R.; Ray, S.; Ahmad, A.; Verma, D.; Bhayye, S.; Dubey, R.; Singh, H.N.; Kumar, S. Anticancer Drug Discovery Based on Natural Products: From Computational Approaches to Clinical Studies. *Biomedicines* **2024**, *12*, doi:10.3390/biomedicines12010201.
28. Huang, M.; Lu, J.-J.; Ding, J. Natural Products in Cancer Therapy: Past, Present and Future. *Nat. Products Bioprospect.* **2021**, *11*, 5–13, doi:10.1007/s13659-020-00293-7.
29. Asma, S.T.; Acaroz, U.; Imre, K.; Morar, A.; Shah, S.R.A.; Hussain, S.Z.; Arslan-Acaroz, D.; Demirbas, H.; Hajrulai-Musliu, Z.; Istanbulgul, F.R.; et al. Natural Products/Bioactive Compounds as a Source of Anticancer Drugs. *Cancers (Basel)*. **2022**, *14*, 6203, doi:10.3390/cancers14246203.
30. Cragg, G.M.; Newman, D.J. Natural Products: A Continuing Source of Novel Drug Leads. *Biochim. Biophys. Acta - Gen. Subj.* **2013**, *1830*, 3670–3695, doi:10.1016/j.bbagen.2013.02.008.
31. Mohanraj, K.; Karthikeyan, B.S.; Vivek-Ananth, R.P.; Chand, R.P.B.; Aparna, S.R.; Mangalapandi, P.; Samal, A. IMPPAT: A Curated Database of Indian Medicinal Plants, Phytochemistry and Therapeutics. *Sci. Rep.* **2018**, *8*, doi:10.1038/s41598-018-22631-z.
32. Soyer, S.M.; Ozbek, P.; Kasavi, C. Lung Adenocarcinoma Systems Biomarker and Drug Candidates Identified by Machine Learning, Gene Expression Data, and Integrative Bioinformatics Pipeline. *Omi. A J. Integr. Biol.* **2024**, *28*, 408–420, doi:10.1089/omi.2024.0121.
33. Li, Y.; Cai, Y.; Ji, L.; Wang, B.; Shi, D.; Li, X. Machine Learning and Bioinformatics Analysis of Diagnostic Biomarkers Associated with the Occurrence and Development of Lung Adenocarcinoma. *PeerJ* **2024**, *12*, doi:10.7717/peerj.17746.
34. Roh, W.; Geffen, Y.; Cha, H.; Miller, M.; Anand, S.; Kim, J.; Heiman, D.I.; Gainor, J.F.; Laird, P.W.; Cherniack, A.D.; et al. High-Resolution Profiling of Lung Adenocarcinoma Identifies Expression Subtypes with Specific Biomarkers and Clinically Relevant Vulnerabilities. *Cancer Res.* **2022**, *82*, 3917–3931, doi:10.1158/0008-5472.CAN-22-0432.
35. Huang, R.X.; Siriwan, D.; Cho, W.C.; Wan, T.K.; Du, Y.R.; Bennett, A.N.; He, Q.E.; Liu, J.D.; Huang, X.T.; Chan, K.H.K. Lung Adenocarcinoma-Related Target Gene Prediction and Drug Repositioning. *Front. Pharmacol.* **2022**, *13*, doi:10.3389/fphar.2022.936758.
36. Li, Y.; Ma, K.; Wang, H.; Liu, Z.; Li, Z. Identification of Therapeutic Targets in Lung Adenocarcinoma Using Mendelian Randomization and Multi-Omics. *Discov. Oncol.* **2025**, *16*, doi:10.1007/s12672-025-02835-2.
37. Jia, K.; Wang, Y.; Cao, Q.; Wang, Y. Extensive Prediction of Drug Response in Mutation-Subtype-Specific LUAD with Machine Learning Approach. *Oncol. Res.* **2024**, *32*, 409–419, doi:10.32604/or.2023.042863.



38. Qureshi, R.; Basit, S.A.; Shamsi, J.A.; Fan, X.; Nawaz, M.; Yan, H.; Alam, T. Machine Learning Based Personalized Drug Response Prediction for Lung Cancer Patients. *Sci. Rep.* **2022**, *12*, doi:10.1038/s41598-022-23649-0.
39. Sobhan, M.; Mondal, A.M. Explainable Machine Learning to Identify Patient-Specific Biomarkers for Lung Cancer. *Proc. - 2022 IEEE Int. Conf. Bioinforma. Biomed. BIBM 2022* **2022**, 3152–3159, doi:10.1109/BIBM55620.2022.9995516.
40. Okyay, T.M.; Yilmaz, I.; Koldas, M. Machine Learning-Based Bioactivity Prediction of Porphyrin Derivatives: Molecular Descriptors, Clustering, and Model Evaluation. *Photochem. Photobiol. Sci.* **2025**, *24*, 923–937, doi:10.1007/s43630-025-00733-8.
41. Lu, T.-P.; Tsai, M.-H.; Lee, J.-M.; Hsu, C.-P.; Chen, P.-C.; Lin, C.-W.; Shih, J.-Y.; Yang, P.-C.; Hsiao, C.K.; Lai, L.-C.; et al. Identification of a Novel Biomarker, SEMA5A, for Non-Small Cell Lung Carcinoma in Nonsmoking Women. *Cancer Epidemiol. Biomarkers Prev.* **2010**, *19*, 2590–2597, doi:10.1158/1055-9965.EPI-10-0332.
42. Lu, T.-P.; Hsiao, C.K.; Lai, L.-C.; Tsai, M.-H.; Hsu, C.-P.; Lee, J.-M.; Chuang, E.Y. Identification of Regulatory SNPs Associated with Genetic Modifications in Lung Adenocarcinoma. *BMC Res. Notes* **2015**, *8*, 92, doi:10.1186/s13104-015-1053-8.
43. Landi, M.T.; Dracheva, T.; Rotunno, M.; Figueroa, J.D.; Liu, H.; Dasgupta, A.; Mann, F.E.; Fukuoka, J.; Hames, M.; Bergen, A.W.; et al. Gene Expression Signature of Cigarette Smoking and Its Role in Lung Adenocarcinoma Development and Survival. *PLoS One* **2008**, *3*, e1651, doi:10.1371/journal.pone.0001651.
44. Sanchez-Palencia, A.; Gomez-Morales, M.; Gomez-Capilla, J.A.; Pedraza, V.; Boyero, L.; Rosell, R.; Fárez-Vidal, M.E. Gene Expression Profiling Reveals Novel Biomarkers in Nonsmall Cell Lung Cancer. *Int. J. cancer* **2011**, *129*, 355–364, doi:10.1002/ijc.25704.
45. Wrage, M.; Ruosaari, S.; Eijk, P.P.; Kaifi, J.T.; Hollmén, J.; Yekebas, E.F.; Izbicki, J.R.; Brakenhoff, R.H.; Streichert, T.; Riethdorf, S.; et al. Genomic Profiles Associated with Early Micrometastasis in Lung Cancer: Relevance of 4q Deletion. *Clin. Cancer Res.* **2009**, *15*, 1566–1574, doi:10.1158/1078-0432.CCR-08-2188.
46. Smyth, G.K. Linear Models and Empirical Bayes Methods for Assessing Differential Expression in Microarray Experiments. *Stat. Appl. Genet. Mol. Biol.* **2004**, *3*, doi:10.2202/1544-6115.1027.
47. Singh, G.; Soman, B. Data Transformation Using Dplyr Package in R. **2019**.
48. Valero-Mora, P.M. Ggplot2: Elegant Graphics for Data Analysis. *J. Stat. Softw.* **2010**, *35*, doi:10.18637/jss.v035.b01.
49. Gao, C.H.; Yu, G.; Cai, P. GgVennDiagram: An Intuitive, Easy-to-Use, and Highly Customizable R Package to Generate Venn Diagram. *Front. Genet.* **2021**, *12*, doi:10.3389/fgene.2021.706907.
50. Szklarczyk, D.; Morris, J.H.; Cook, H.; Kuhn, M.; Wyder, S.; Simonovic, M.; Santos, A.; Doncheva, N.T.; Roth, A.; Bork, P.; et al. The STRING Database in 2017: Quality-Controlled Protein–Protein Association Networks, Made Broadly Accessible. *Nucleic Acids Res.* **2017**, *45*, D362–D368, doi:10.1093/nar/gkw937.
51. Shannon, P.; Markiel, A.; Ozier, O.; Baliga, N.S.; Wang, J.T.; Ramage, D.; Amin, N.; Schwikowski, B.; Ideker, T. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks. *Genome Res.* **2003**, *13*, 2498–2504, doi:10.1101/gr.1239303.
52. Chin, C.-H.; Chen, S.-H.; Wu, H.-H.; Ho, C.-W.; Ko, M.-T.; Lin, C.-Y. CytoHubba: Identifying Hub Objects and Sub-Networks from Complex Interactome. *BMC Syst. Biol.* **2014**, *8*, S11, doi:10.1186/1752-0509-8-S4-S11.
53. Khan, A.; Fornes, O.; Stigliani, A.; Gheorghe, M.; Castro-Mondragon, J.A.; Van Der Lee, R.; Bessy, A.; Chèneby, J.; Kulkarni, S.R.; Tan, G.; et al. JASPAR 2018: Update of the Open-Access Database of Transcription Factor Binding Profiles and Its Web Framework. *Nucleic Acids Res.* **2018**, *46*, D260–D266, doi:10.1093/nar/gkx1126.



54. Hsu, S.-D.; Lin, F.-M.; Wu, W.-Y.; Liang, C.; Huang, W.-C.; Chan, W.-L.; Tsai, W.-T.; Chen, G.-Z.; Lee, C.-J.; Chiu, C.-M.; et al. MiRTarBase: A Database Curates Experimentally Validated MicroRNA–Target Interactions. *Nucleic Acids Res.* **2011**, *39*, D163–D169, doi:10.1093/nar/gkq1107.
55. Xia, J.; Gill, E.E.; Hancock, R.E.W. NetworkAnalyst for Statistical, Visual and Network-Based Meta-Analysis of Gene Expression Data. *Nat. Protoc.* **2015**, *10*, 823–844, doi:10.1038/nprot.2015.052.
56. Sherman, B.T.; Hao, M.; Qiu, J.; Jiao, X.; Baseler, M.W.; Lane, H.C.; Imamichi, T.; Chang, W. DAVID: A Web Server for Functional Enrichment Analysis and Functional Annotation of Gene Lists (2021 Update). *Nucleic Acids Res.* **2022**, *50*, W216–W221, doi:10.1093/nar/gkac194.
57. Kuleshov, M. V.; Jones, M.R.; Rouillard, A.D.; Fernandez, N.F.; Duan, Q.; Wang, Z.; Koplev, S.; Jenkins, S.L.; Jagodnik, K.M.; Lachmann, A.; et al. Enrichr: A Comprehensive Gene Set Enrichment Analysis Web Server 2016 Update. *Nucleic Acids Res.* **2016**, *44*, doi:10.1093/nar/gkw377.
58. Helmy, M.; Agrawal, R.; Ali, J.; Soudy, M.; Bui, T.T.; Selvarajoo, K. GeneCloudOmics: A Data Analytic Cloud Platform for High-Throughput Gene Expression Analysis. *Front. Bioinforma.* **2021**, *1*, doi:10.3389/fbinf.2021.693836.
59. Li, T.; Fu, J.; Zeng, Z.; Cohen, D.; Li, J.; Chen, Q.; Li, B.; Liu, X.S. TIMER2.0 for Analysis of Tumor-Infiltrating Immune Cells. *Nucleic Acids Res.* **2020**, *48*, doi:10.1093/NAR/GKAA407.
60. Tang, Z.; Kang, B.; Li, C.; Chen, T.; Zhang, Z. GEPIA2: An Enhanced Web Server for Large-Scale Expression Profiling and Interactive Analysis. *Nucleic Acids Res.* **2019**, *47*, W556–W560, doi:10.1093/nar/gkz430.
61. Chandrashekar, D.S.; Karthikeyan, S.K.; Korla, P.K.; Patel, H.; Shovon, A.R.; Athar, M.; Netto, G.J.; Qin, Z.S.; Kumar, S.; Manne, U.; et al. UALCAN: An Update to the Integrated Cancer Data Analysis Platform. *Neoplasia (United States)* **2022**, *25*, doi:10.1016/j.neo.2022.01.001.
62. Gao, J.; Aksoy, B.A.; Dogrusoz, U.; Dresdner, G.; Gross, B.; Sumer, S.O.; Sun, Y.; Jacobsen, A.; Sinha, R.; Larsson, E.; et al. Integrative Analysis of Complex Cancer Genomics and Clinical Profiles Using the CBioPortal. *Sci. Signal.* **2013**, *6*, doi:10.1126/scisignal.2004088.
63. Burley, S.K.; Bhikadiya, C.; Bi, C.; Bittrich, S.; Chen, L.; Crichlow, G. V.; Christie, C.H.; Dalenberg, K.; Di Costanzo, L.; Duarte, J.M.; et al. RCSB Protein Data Bank: Powerful New Tools for Exploring 3D Structures of Biological Macromolecules for Basic and Applied Research and Education in Fundamental Biology, Biomedicine, Biotechnology, Bioengineering and Energy Sciences. *Nucleic Acids Res.* **2021**, *49*, D437–D451, doi:10.1093/nar/gkaa1038.
64. BIOVIA, D.S. Discovery Studio Visualizer V21.1.0.20298. *BIOVIA, Dassault Systèmes* **2005**.
65. Guex, N.; Peitsch, M.C. SWISS-MODEL and the Swiss-Pdb Viewer: An Environment for Comparative Protein Modeling. *Electrophoresis* **1997**, *18*, 2714–2723, doi:10.1002/elps.1150181505.
66. Puellas, A.A.; Bastos, L.L.; Paixão, V.M.; Araujo, S.C.; de Melo Minardi, R.C. Virtual Screening. In: 2024; pp. 209–236.
67. Xiong, G.; Wu, Z.; Yi, J.; Fu, L.; Yang, Z.; Hsieh, C.; Yin, M.; Zeng, X.; Wu, C.; Lu, A.; et al. ADMETlab 2.0: An Integrated Online Platform for Accurate and Comprehensive Predictions of ADMET Properties. *Nucleic Acids Res.* **2021**, *49*, W5–W14, doi:10.1093/nar/gkab255.
68. Lipinski, C.A. Lead- and Drug-like Compounds: The Rule-of-Five Revolution. *Drug Discov. Today Technol.* **2004**, *1*, 337–341, doi:10.1016/j.ddtec.2004.11.007.
69. W. Caldwell, G.; Yan, Z.; Lang, W.; A. Masucci, J. The IC50 Concept Revisited. *Curr. Top. Med. Chem.* **2012**, *12*, 1282–1290, doi:10.2174/156802612800672844.
70. Ali, M.A.; Sarker, H.; Khan, T.; Sheikh, H.; Saif, A.; Farid, F. Bin; Afrin, S.; Khatun, M.A.; Kumar, N. Multi-



- Omics Pan-Cancer Profiling of CDK2 and in Silico Identification of Plant-Derived Inhibitors Using Machine Learning Approaches. *RSC Adv.* **2025**, *15*, 36938–36968, doi:10.1039/D5RA05535K.
71. Saif, A.; Islam, M.T.; Raihan, M.O.; Yousefi, N.; Rahman, M.A.; Faridi, H.; Hasan, A.R.; Hossain, M.M.; Saleem, R.M.; Albadrani, G.M.; et al. Pan-Cancer Analysis of CDC7 in Human Tumors: Integrative Multi-Omics Insights and Discovery of Novel Marine-Based Inhibitors through Machine Learning and Computational Approaches. *Comput. Biol. Med.* **2025**, *190*, doi:10.1016/j.combiomed.2025.110044.
 72. Gubler, H.; Schopfer, U.; Jacoby, E. Theoretical and Experimental Relationships between Percent Inhibition and IC50 Data Observed in High-Throughput Screening. *J. Biomol. Screen.* **2013**, *18*, 1–13, doi:10.1177/1087057112455219.
 73. Danishuddin; Khan, A.U. Descriptors and Their Selection Methods in QSAR Analysis: Paradigm for Drug Design. *Drug Discov. Today* **2016**, *21*, 1291–1302, doi:10.1016/j.drudis.2016.06.013.
 74. Guo, Q.; Hernandez-Hernandez, S.; Ballester, P.J. UMAP-Based Clustering Split for Rigorous Evaluation of AI Models for Virtual Screening on Cancer Cell Lines*. *J. Cheminform.* **2025**, *17*, doi:10.1186/s13321-025-01039-8.
 75. Majumdar, S.; Basak, S.C. Beware of External Validation! - A Comparative Study of Several Validation Techniques Used in QSAR Modelling. *Curr. Comput. Aided. Drug Des.* **2018**, *14*, 284–291, doi:10.2174/1573409914666180426144304.
 76. Meng, X.-Y.; Zhang, H.-X.; Mezei, M.; Cui, M. Molecular Docking: A Powerful Approach for Structure-Based Drug Discovery. *Curr. Comput. Aided-Drug Des.* **2012**, *7*, 146–157, doi:10.2174/157340911795677602.
 77. Kondapuram, S.K.; Sarvagalla, S.; Coumar, M.S. Docking-Based Virtual Screening Using PyRx Tool: Autophagy Target Vps34 as a Case Study. In *Molecular Docking for Computer-Aided Drug Design*; Elsevier, 2021; pp. 463–477.
 78. Danish Ahmad, A.V.; Khan, S.W.; Yasar, Q.; Shaikh, M.S.; Khan, M.M. Computational Biology Approach to Predict Molecular Mechanism in Cancer. *Oral Oncol. Reports* **2024**, *12*, 100651, doi:10.1016/j.oor.2024.100651.
 79. Owoloye, A.J.; Ligali, F.C.; Enejoh, O.A.; Musa, A.Z.; Aina, O.; Idowu, E.T.; Oyebola, K.M. Molecular Docking, Simulation and Binding Free Energy Analysis of Small Molecules as PfHT1 Inhibitors. *PLoS One* **2022**, *17*, e0268269, doi:10.1371/journal.pone.0268269.
 80. Ogbodo, U.C.; Enejoh, O.A.; Okonkwo, C.H.; Gnanasekar, P.; Gachanja, P.W.; Osata, S.; Atanda, H.C.; Iwuchukwu, E.A.; Achilonu, I.; Awe, O.I. Computational Identification of Potential Inhibitors Targeting Cdk1 in Colorectal Cancer. *Front. Chem.* **2023**, *11*, doi:10.3389/fchem.2023.1264808.
 81. Kollman, P.A.; Massova, I.; Reyes, C.; Kuhn, B.; Huo, S.; Chong, L.; Lee, M.; Lee, T.; Duan, Y.; Wang, W.; et al. Calculating Structures and Free Energies of Complex Molecules: Combining Molecular Mechanics and Continuum Models. *Acc. Chem. Res.* **2000**, *33*, 889–897, doi:10.1021/ar000033j.
 82. Friesner, R.A.; Murphy, R.B.; Repasky, M.P.; Frye, L.L.; Greenwood, J.R.; Halgren, T.A.; Sanschagrin, P.C.; Mainz, D.T. Extra Precision Glide: Docking and Scoring Incorporating a Model of Hydrophobic Enclosure for Protein–Ligand Complexes. *J. Med. Chem.* **2006**, *49*, 6177–6196, doi:10.1021/jm051256o.
 83. Gilson, M.K.; Given, J.A.; Bush, B.L.; McCammon, J.A. The Statistical-Thermodynamic Basis for Computation of Binding Affinities: A Critical Review. *Biophys. J.* **1997**, *72*, 1047–1069, doi:10.1016/S0006-3495(97)78756-3.
 84. Arya, H.; Bhatt, T.K. Molecular Dynamics Simulations. In *The Design & Development of Novel Drugs and Vaccines*; Elsevier, 2021; pp. 65–81.
 85. Roos, K.; Wu, C.; Damm, W.; Reboul, M.; Stevenson, J.M.; Lu, C.; Dahlgren, M.K.; Mondal, S.; Chen, W.;



- Wang, L.; et al. OPLS3e: Extending Force Field Coverage for Drug-Like Small Molecules. *J. Chem. Theory Comput.* **2019**, *15*, 1863–1874, doi:10.1021/acs.jctc.8b01026.
86. Valdés-Tresanco, M.S.; Valdés-Tresanco, M.E.; Valiente, P.A.; Moreno, E. Gmx_MMPBSA: A New Tool to Perform End-State Free Energy Calculations with GROMACS. *J. Chem. Theory Comput.* **2021**, *17*, 6281–6291, doi:10.1021/acs.jctc.1c00645.
87. Miller, B.R.; McGee, T.D.; Swails, J.M.; Homeyer, N.; Gohlke, H.; Roitberg, A.E. MMPBSA.Py: An Efficient Program for End-State Free Energy Calculations. *J. Chem. Theory Comput.* **2012**, *8*, 3314–3321, doi:10.1021/ct300418h.
88. Shirts, M.R.; Klein, C.; Swails, J.M.; Yin, J.; Gilson, M.K.; Mobley, D.L.; Case, D.A.; Zhong, E.D. Lessons Learned from Comparing Molecular Dynamics Engines on the SAMPL5 Dataset. *J. Comput. Aided. Mol. Des.* **2017**, *31*, 147–161, doi:10.1007/s10822-016-9977-1.
89. Daina, A.; Michielin, O.; Zoete, V. SwissADME: A Free Web Tool to Evaluate Pharmacokinetics, Drug-Likeness and Medicinal Chemistry Friendliness of Small Molecules. *Sci. Rep.* **2017**, *7*, 42717, doi:10.1038/srep42717.
90. Banerjee, P.; Kemmler, E.; Dunkel, M.; Preissner, R. ProTox 3.0: A Webserver for the Prediction of Toxicity of Chemicals. *Nucleic Acids Res.* **2024**, gkae303, doi:10.1093/nar/gkae303.
91. Ashtiani, M.; Salehzadeh-Yazdi, A.; Razaghi-Moghadam, Z.; Hennig, H.; Wolkenhauer, O.; Mirzaie, M.; Jafari, M. A Systematic Survey of Centrality Measures for Protein-Protein Interaction Networks. *BMC Syst. Biol.* **2018**, *12*, 80, doi:10.1186/s12918-018-0598-2.
92. Wang, M.; Wang, H.; Zheng, H. A Mini Review of Node Centrality Metrics in Biological Networks. *Int. J. Netw. Dyn. Intell.* **2022**, 99–110, doi:10.53941/ijndi0101009.
93. Wood, D.J.; Korolchuk, S.; Tatum, N.J.; Wang, L.-Z.; Endicott, J.A.; Noble, M.E.M.; Martin, M.P. Differences in the Conformational Energy Landscape of CDK1 and CDK2 Suggest a Mechanism for Achieving Selective CDK Inhibition. *Cell Chem. Biol.* **2019**, *26*, 121–130.e5, doi:10.1016/j.chembiol.2018.10.015.
94. Ke, G.; Meng, Q.; Finley, T.; Wang, T.; Chen, W.; Ma, W.; Ye, Q.; Liu, T.-Y. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In Proceedings of the Advances in Neural Information Processing Systems; Guyon, I., Luxburg, U. Von, Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R., Eds.; Curran Associates, Inc., 2017; Vol. 30.
95. Zhao, X.; Liu, Y.; Zhao, Q. Improved LightGBM for Extremely Imbalanced Data and Application to Credit Card Fraud Detection. *IEEE Access* **2024**, *12*, 159316–159335, doi:10.1109/ACCESS.2024.3487212.
96. Bento, A.P.; Gaulton, A.; Hersey, A.; Bellis, L.J.; Chambers, J.; Davies, M.; Krüger, F.A.; Light, Y.; Mak, L.; McGlinchey, S.; et al. The ChEMBL Bioactivity Database: An Update. *Nucleic Acids Res.* **2014**, *42*, D1083–90, doi:10.1093/nar/gkt1031.
97. Zhao, L.; Jiang, W.; Zhu, Z.; Pan, F.; Xing, X.; Zhou, F.; Zhao, L. Rosemarinic Acid-Induced Destabilization of A β Peptides: Insights from Molecular Dynamics Simulations. *Foods (Basel, Switzerland)* **2024**, *13*, doi:10.3390/foods13244170.
98. Abouzied, A.S.; Alqarni, S.; Younes, K.M.; Alanazi, S.M.; Alrashed, D.M.; Alhathal, R.K.; Huwaimel, B.; Elkashlan, A.M. Structural and Free Energy Landscape Analysis for the Discovery of Antiviral Compounds Targeting the Cap-Binding Domain of Influenza Polymerase PB2. *Sci. Rep.* **2024**, *14*, 25441, doi:10.1038/s41598-024-69816-3.
99. Veber, D.F.; Johnson, S.R.; Cheng, H.Y.; Smith, B.R.; Ward, K.W.; Kopple, K.D. Molecular Properties That Influence the Oral Bioavailability of Drug Candidates. *J. Med. Chem.* **2002**, *45*, 2615–2623, doi:10.1021/jm020017n.



100. Ndombera, F.T. Revisiting Cheminformatics and Mechanisms of Action of Chloroquine and Hydroxychloroquine in Targeting Covid-19. *JBCG* **2020**, *3*, doi:10.17303/jbcg.2020.3.101.
101. Rashid, M.; Afzal, O.; Altamimi, A.S.A. BENZIMIDAZOLE MOLECULE HYBRID WITH OXADIAZOLE RING AS ANTIPROLIFERATIVE AGENTS: IN-SILICO ANALYSIS, SYNTHESIS AND BIOLOGICAL EVALUATION. *J. Chil. Chem. Soc.* **2021**, *66*, 5164–5182, doi:10.4067/S0717-97072021000205164.
102. Benet, L.Z.; Hosey, C.M.; Ursu, O.; Oprea, T.I. BDDCS, the Rule of 5 and Drugability. *Adv. Drug Deliv. Rev.* **2016**, *101*, 89–98, doi:10.1016/j.addr.2016.05.007.
103. Sun, Y.; Zabihi, M.; Li, Q.; Li, X.; Kim, B.J.; Ubogu, E.E.; Raja, S.N.; Wesselmann, U.; Zhao, C. Drug Permeability: From the Blood–Brain Barrier to the Peripheral Nerve Barriers. *Adv. Ther.* **2023**, *6*, doi:10.1002/adtp.202200150.
104. Heavey, M.K.; Durmusoglu, D.; Crook, N.; Anselmo, A.C. Discovery and Delivery Strategies for Engineered Live Biotherapeutic Products. *Trends Biotechnol.* **2022**, *40*, 354–369, doi:10.1016/j.tibtech.2021.08.002.
105. Sibai, M.; Cervilla, S.; Grases, D.; Musulen, E.; Lazcano, R.; Mo, C.-K.; Davalos, V.; Fortian, A.; Bernat, A.; Romeo, M.; et al. The Spatial Landscape of Cancer Hallmarks Reveals Patterns of Tumor Ecological Dynamics and Drug Sensitivity. *Cell Rep.* **2025**, *44*, 115229, doi:10.1016/j.celrep.2024.115229.
106. Boija, A.; Klein, I.A.; Sabari, B.R.; Dall’Agnese, A.; Coffey, E.L.; Zamudio, A. V.; Li, C.H.; Shrinivas, K.; Manteiga, J.C.; Hannett, N.M.; et al. Transcription Factors Activate Genes through the Phase-Separation Capacity of Their Activation Domains. *Cell* **2018**, *175*, 1842–1855.e16, doi:10.1016/j.cell.2018.10.042.
107. Catalanotto, C.; Cogoni, C.; Zardo, G. MicroRNA in Control of Gene Expression: An Overview of Nuclear Functions. *Int. J. Mol. Sci.* **2016**, *17*, doi:10.3390/ijms17101712.
108. Chen, Y.; Jin, L.; Jiang, Z.; Liu, S.; Feng, W. Identifying and Validating Potential Biomarkers of Early Stage Lung Adenocarcinoma Diagnosis and Prognosis. *Front. Oncol.* **2021**, *11*, doi:10.3389/fonc.2021.644426.
109. Du, Q.; Liu, W.; Mei, T.; Wang, J.; Qin, T.; Huang, D. Prognostic and Immunological Characteristics of CDK1 in Lung Adenocarcinoma: A Systematic Analysis. *Front. Oncol.* **2023**, *13*, doi:10.3389/fonc.2023.1128443.
110. Li, S.; Li, H.; Cao, Y.; Geng, H.; Ren, F.; Li, K.; Dai, C.; Li, N. Integrated Bioinformatics Analysis Reveals CDK1 and PLK1 as Potential Therapeutic Targets of Lung Adenocarcinoma. *Medicine (Baltimore)*. **2021**, *100*, e26474, doi:10.1097/MD.00000000000026474.
111. Batool, M.; Ahmad, B.; Choi, S. A Structure-Based Drug Discovery Paradigm. *Int. J. Mol. Sci.* **2019**, *20*, doi:10.3390/ijms20112783.
112. Enejoh, O.A.; Okonkwo, C.H.; Nortey, H.; Kemiki, O.A.; Moses, A.; Mbaoji, F.N.; Yusuf, A.S.; Awe, O.I. Machine Learning and Molecular Dynamics Simulations Predict Potential TGR5 Agonists for Type 2 Diabetes Treatment. *Front. Chem.* **2024**, *12*, doi:10.3389/fchem.2024.1503593.
113. Di Stefano, M.; Galati, S.; Ortore, G.; Caligiuri, I.; Rizzolio, F.; Ceni, C.; Bertini, S.; Bononi, G.; Granchi, C.; Macchia, M.; et al. Machine Learning-Based Virtual Screening for the Identification of Cdk5 Inhibitors. *Int. J. Mol. Sci.* **2022**, *23*, doi:10.3390/ijms231810653.
114. Bhimanwar, R.S.; Lokhande, K.B.; Shrivastava, A.; Singh, A.; Chitlange, S.S.; Mittal, A. Identification of Potential Drug Candidates as TGR5 Agonist to Combat Type II Diabetes Using in Silico Docking and Molecular Dynamics Simulation Studies. *J. Biomol. Struct. Dyn.* **2023**, *41*, 13314–13331, doi:10.1080/07391102.2023.2173654.
115. Hurmach, V.; Karaushu, V.; Prylutska, S.; Klestova, Z.; Vyzhva, S.; Prylutsky, Y.; Ritter, U.; Garamus, V. In Silico Analysis of C60 Fullerene Interaction with TMPRSS2: Toward Novel COVID-19 Prevention Approaches. *Molecules* **2025**, *30*, doi:10.3390/molecules30234586.



116. David, C.C.; Jacobs, D.J. Principal Component Analysis: A Method for Determining the Essential Dynamics of Proteins. In; 2014; pp. 193–226.



Data Availability Statement (DAS)

The original data and contributions presented in this study are included in the article and its Supplementary Information (Table S1-S10 and Figure S1-S5). The training set of our selected compounds, the test set, and the chEMBL datasets (containing the consensus predictions) and related python code for running the qsar models to predict the bioactivity of the selected compounds, can be found on our GitHub repository (https://github.com/ahad004/LUAD_ML_QSAR_Modeling), or can be accessed using the following <https://doi.org/10.5281/zenodo.19841200>.

