# Exploring Machine Learning in Lung Cancer: Predictive Modelling, Gene Associations, and Challenges

[1] K. Mary Sudha Rani[2] Dr. V. Kamakshi Prasad

**Abstract:** Lung cancer is a disease with a high mortality rate and widespread occurrence. Therefore, developing accurate prediction methods and practical gene association analyses is crucial. The utilization of high-throughput genomic data to reveal significant genetic factors has seen an increase in the application of machine-learning techniques. This document presents a thorough examination of machine learning methodologies that are presently utilized to forecast lung cancer and scrutinize gene correlations. The analysis examines different data types, such as gene expression profiles, genomic variants, and clinical data. The primary focus is on integrating multi-omics data for a more comprehensive understanding. Our study comprehensively examines a variety of machine learning algorithms, including traditional methods such as support vector machines and random forests, advanced deep learning architectures, and network-based methodologies. The following discourse explores the pragmatic utilization of the methods above in predictive modeling, biomarker identification, and drug discovery routes. The article addresses common obstacles in the field, such as interpretability and validation, and proposes potential avenues for future research, such as incorporating multi-omics data and implementing personalized medicine. This survey provides a detailed analysis of the recent advancements in machine learning techniques for lung cancer research. It aims to establish a strong basis for future improvements in diagnosis, prognosis, and treatment strategies.
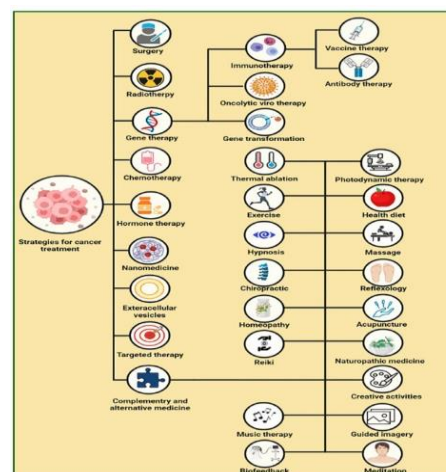
## 1. Introduction:

Cancer is a complex disease that manifests through abnormal molecular pathways. The intricate biological mechanisms and operational effects of gene sets linked to cancer are fundamental to understanding disease progression, identifying therapeutic targets, and developing personalized treatment plans. As such, contemporary cancer research leans heavily on gene-set analysis to interpret high-throughput genomic data within cancer-related biological processes [1].

Lung cancer, in particular, remains a significant global concern. With an estimated 1.8 million new cases annually, it presents a pressing public health issue [2]. It is reported that 85% of lung cancers are non-small cell lung cancers (NSCLC) [3]. Precise prediction methods and meticulous gene correlation studies are indispensable to augment lung cancer detection, diagnosis, prognosis, and treatment.

Integrating algorithms, machine learning, and multi-omics data can enhance our understanding and prediction of lung cancer. Recent advances in the field have incorporated machine-learning techniques in lung cancer research to analyze and interpret high-throughput genomic data [4]. These methodologies pave the way for identifying lung cancer-related genes through genetic and molecular data analysis.

[1]Research Scholar,CSE dept., JNTUH Hyderabad, Assistant Professor,AIML Dept., Chaitanya Bharathi Institute of Technology, Hyderabad.
kmarysudha_cseaiml@cbit.ac.in
[2] Professor,CSE dept., JNTUH Hyderabad
kamakshiprasad@jntuh.ac.in

Such techniques facilitate the analysis of varied data types, including gene expression profiles, genomic variants, and clinical data. They aid in identifying biomarkers, disease subtypes, and fundamental molecular mechanisms. With its capacity to unveil intricate patterns and relationships in genomic data, machine learning is poised to revolutionize our approach to cancer research.
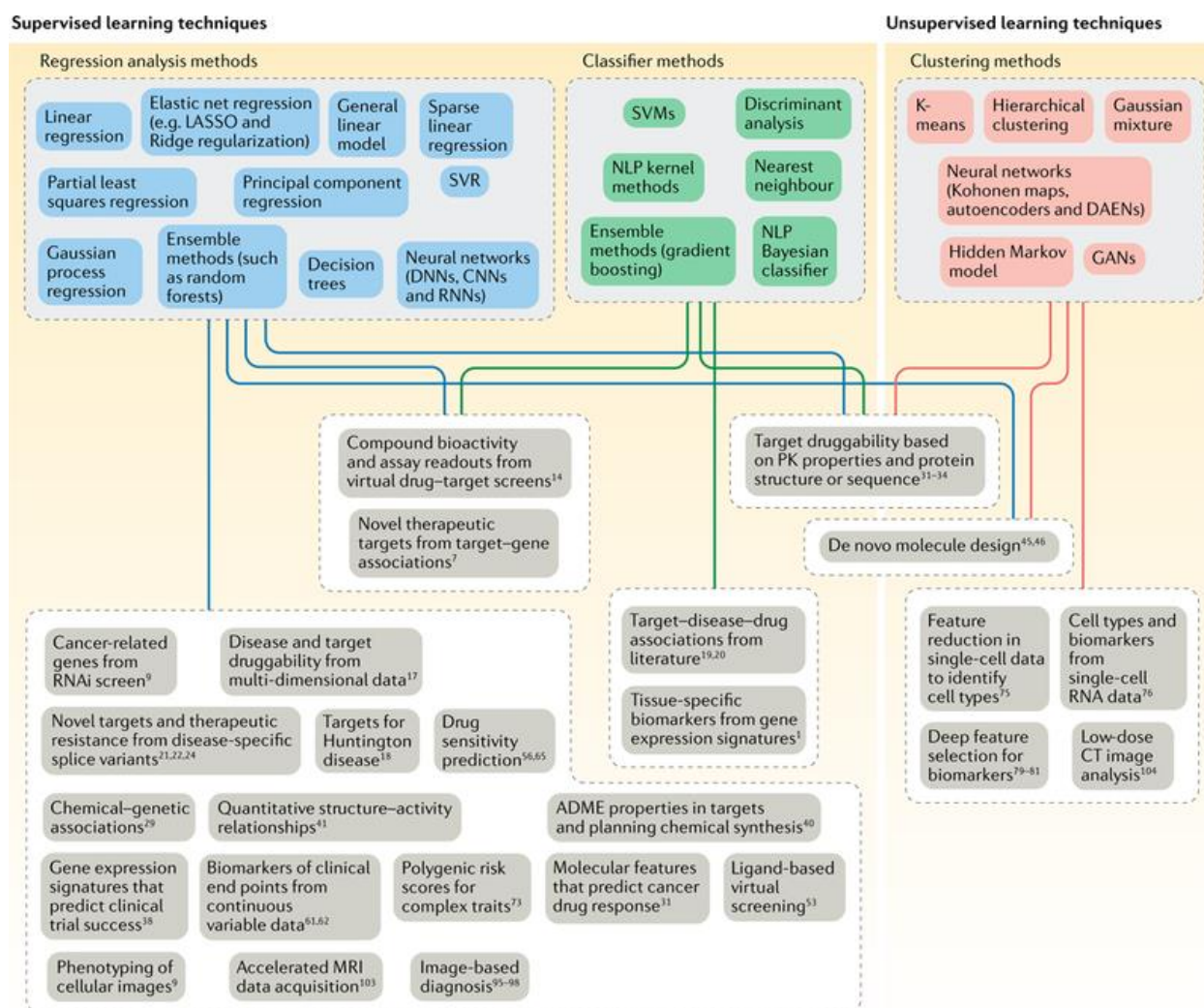
One of the promising advancements in this arena is the integration of multi-omics data, given its potential to enhance lung cancer analysis [5]. By leveraging genomic, transcriptomic, epigenomic, and other omics data, researchers are offered a more comprehensive understanding of the molecular foundation of lung cancer. This, in turn, fosters the discovery of novel diagnostic and therapeutic targets.

Machine learning methods encompass a wide range of tools, including support vector machines (SVM), random forests (RF), and logistic regression, as well as deep learning architectures like artificial neural networks (ANN), convolutional neural networks (CNN), and recurrent neural networks (RNN) [7]. These algorithms are proficient in executing lung cancer research feature selection, classification, clustering, and prediction tasks.

Gene interaction networks provide insights into gene-disease connections and may facilitate the identification of novel NSCLC genes. Graph-based algorithms such as Graph Convolutional Networks (GCN) and deep walk can analyze these networks effectively. Furthermore, network-based methodologies have emerged as promising tools for studying gene associations in lung cancer [8].
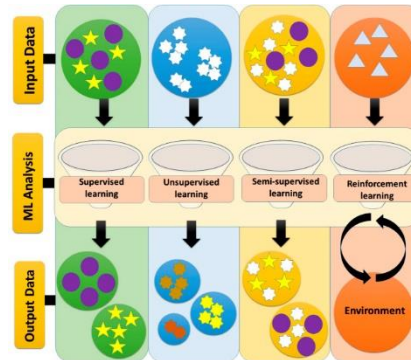
Machine learning has already led to significant strides in lung cancer research. Predictive models employing machine learning have shown promise in identifying high-risk individuals, aiding early detection and customized treatment [9]. Moreover, machine learning has unveiled prognostic biomarkers and implicated molecular pathways in the causation of lung cancer. These insights can optimize drug discovery pathways, identify therapeutic targets, and improve treatment outcomes [10].

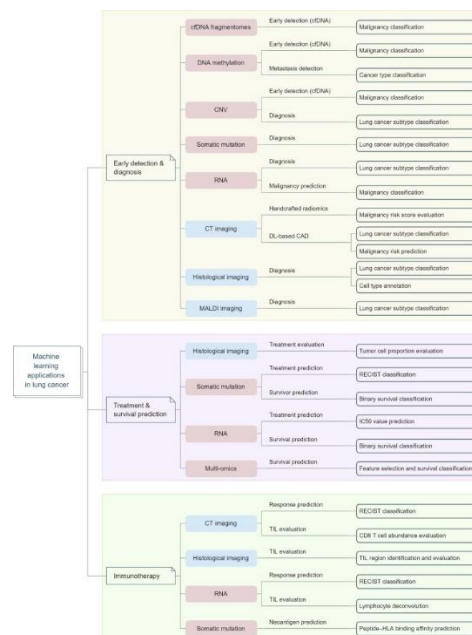**Machine learning tools and their drug discovery applications.**

Despite this progress, several challenges persist in applying machine learning for lung cancer research. The necessity for transparency and interpretability in complex algorithms presents a hurdle, especially considering the need for decision-making tools in clinical settings [11]. The heterogeneity of data sources and the requirement for large and diverse datasets for robust analysis adds further complexities to the validation and reproducibility of studies. Furthermore, harmonization, normalization, and multi-omics data integration are prerequisites for accurate analysis [12].



This paper comprehensively surveys machine-learning techniques for predicting lung cancer and analyzing gene associations. Our exploration extends across the various types of data, challenges faced in the field, and potential future directions, ultimately serving as a robust foundation for future progress in diagnosis, prognosis, and treatment strategies.



## 2.   Background:

The exponential growth and complexity of biomedical data present both challenges and opportunities. High-throughput technologies such as genomics, transcriptomics, proteomics, and metabolomics have propelled a surge in biological data production, thus providing deeper insights into health and disease states [13].

Translational bioinformatics seeks to bridge the gap between basic research and medical applications, leveraging biomedical data to understand diseases and develop personalized medicine strategies. The synthesis of computational and statistical methods with biological and clinical data yields practical insights with potential real-world impact [14].
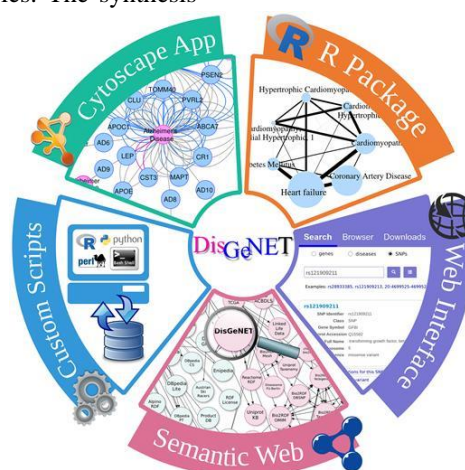


The DisGeNET platform emerges as a valuable tool for studying and interpreting human disease genes. DisGeNET serves up gene-disease associations that can be examined in the context of numerous diseases. Drawing from curated databases, scientific literature, and various online resources, it offers a comprehensive knowledge base to researchers and clinicians [15].

The platform encompasses genetic variants, gene-disease associations, disease pathways, and pharmacogenomics data. It integrates information from PubMed, ClinVar, UniProt, and GWAS Catalog. Standardized vocabularies like the Unified Medical Language System (UMLS) are used to ensure consistency and comprehensibility. What sets DisGeNET apart is its extensive catalog of gene-disease links, spanning both well-established and preliminary associations. This capability opens up new avenues of exploration and potential research subjects for investigators. DisGeNET also offers user-friendly search and filtering options, enabling users to search for genes, diseases, and associations based on criteria like evidence level, data source, or disease category. The resulting information includes supporting evidence, genetic variants, linked phenotypes, and pertinent references for gene-disease associations. DisGeNET integrates with other bioinformatics tools and workflows via APIs. This functionality facilitates the integration of DisGeNET data with other relevant datasets for comprehensive analysis.

Overall, DisGeNET significantly simplifies the process of gene-disease research and its implications for human health. It proves invaluable to professionals in genetics, genomics, precision medicine, and drug discovery, as well as clinicians and researchers. As a centralized repository of gene-disease associations, DisGeNET supports identifying therapeutic targets, discovering novel disease

mechanisms, and developing personalized treatment strategies.

**Genome-Wide Association Studies**

Genome-Wide Association Studies (GWAS) have profoundly influenced gene-disease research by identifying genetic variations associated with diseases, including lung cancer. GWAS examines numerous genetic variations across the genome to find correlations between genetic variants and disease risk. They specifically target common genetic variations, usually Single Nucleotide Polymorphisms (SNPs), that could affect the risk of specific diseases [16].

The GWAS process involves two main steps: genotyping and association analysis. Genotyping identifies genetic variations in large cases (individuals with the disease) and controls (healthy individuals), using high-throughput genotyping technologies to analyze hundreds of thousands or even millions of SNPs across the genome. The association analysis then identifies genetic variants associated with a disease phenotype, correcting for population stratification and genetic ancestry.

GWAS has made significant discoveries in lung cancer genetics. Hung et al. (2008) found a common genetic variant in the nicotinic acetylcholine receptor gene cluster (CHRNA5-CHRNA3-CHRNB4) that increased lung cancer risk. Other studies, such as those by Hu et al. (2012) and Wang et al. (2020), identified susceptibility loci for lung adenocarcinoma and linked lung cancer risk to MUC4 and PRSS8 gene variants, respectively.

While GWAS findings offer valuable insights into the genetic makeup of lung cancer and identify genes and pathways for further study, they primarily reveal correlations, not causality. To establish mechanisms and basis, functional analyses and validation are necessary.

GWAS also has some limitations. They focus on common genetic variations with moderate to large effect sizes, potentially overlooking rare or minor variations. They are also often conducted on populations with specific ancestries, limiting their applicability to others. Furthermore, while GWAS can identify genetic variations linked to disease susceptibility, they do not explain the underlying biological mechanisms.

To overcome these constraints, several methods can be employed. These include meta-analysis to increase sample size and statistical power, integrating GWAS data with functional genomics and gene expression data to rank candidate genes, and fine-mapping to reduce genomic regions of interest.

In conclusion, GWAS have significantly contributed to our understanding of the genetics of lung cancer and other complex diseases, identifying potential target genes and

pathways for further study. However, interpreting GWAS findings requires caution and additional functional studies to clarify biological mechanisms and establish causality [14].

**Next-Generation Sequencing (NGS)**

Next-Generation Sequencing (NGS) technologies have revolutionized genetic research by enabling fast and cost-effective sequencing of large amounts of DNA or RNA, providing researchers with an unprecedented volume of genomic data . Specifically, Whole Exome Sequencing (WES) and Whole Genome Sequencing (WGS) have significantly contributed to the identification of disease-related genes, including those linked to lung cancer[17].

WES focuses on sequencing the protein-coding regions of the genome, known as the exome, while WGS sequences the entire genome, including both protein-coding and non-coding regions. These NGS technologies have the advantage of detecting rare genetic variants that traditional genotyping methods may miss. These rare variants can provide insights into disease susceptibility and help uncover genes and pathways associated with lung cancer.

The application of NGS in cancer research has led to the discovery of driver mutations and novel genes associated with lung cancer. For example, using NGS, the Cancer Genome Atlas (TCGA) project has identified recurrent somatic mutations in genes like TP53, EGFR, KRAS, and ALK in lung adenocarcinoma. NGS can also detect somatic copy number alterations, structural variations, and gene fusions, revealing complex mutational patterns in lung cancer. These molecular alterations can aid in the classification of lung cancer subtypes and guide treatment decisions.

NGS technologies also present several challenges. The high volume of sequencing data produced necessitates robust computational and bioinformatics infrastructure for data storage, processing, and analysis. Advanced bioinformatics algorithms and tools are required for reading alignment, variant calling, and annotation.

Interpreting the vast number of genetic variants identified by NGS can be challenging. Distinguishing between pathogenic and benign variants is complex and often requires functional studies to confirm their effects. Additionally, the cost of WGS, which involves sequencing the entire genome, can be prohibitive, though technological advances are gradually making NGS more affordable and accessible.

In conclusion, NGS technologies, specifically WES and WGS, have substantially advanced our understanding of the genetic underpinnings of lung cancer and other diseases. They offer a powerful tool for genomic profiling,

detecting rare variants, and discovering disease-related genes and pathways. Yet, to fully harness the potential of NGS in lung cancer genetics, issues concerning data analysis, interpretation, and cost need to be addressed.

**Machine learning and deep learning**

Machine learning and deep learning have been increasingly used in bioinformatics to predict disease-gene associations. These methods leverage computational algorithms and statistical models to analyze large volumes of genomic and biomedical data, aiding the identification of disease-related genes and contributing to a better understanding of disease mechanisms [18].

Commonly utilized algorithms in disease-gene association studies include Support Vector Machines (SVM), Random Forest (RF), and logistic regression. These algorithms are capable of handling high-dimensional data and capturing complex genetic-disease relationships. They facilitate the identification of disease-related genes by performing feature selection, classification, and prediction.
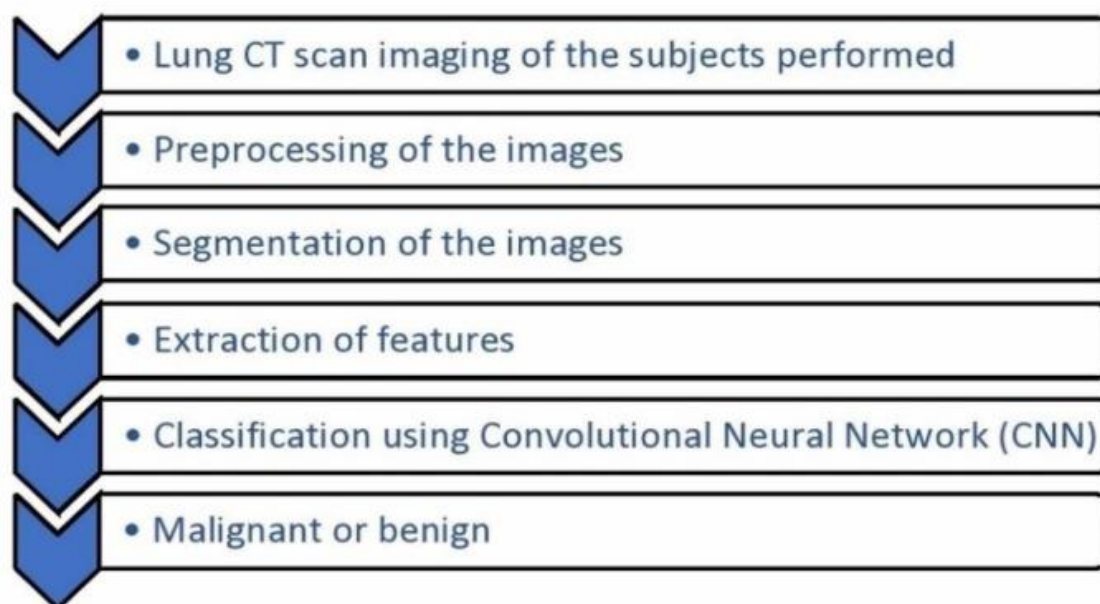
Deep learning, mainly Artificial Neural Networks (ANNs), has recently gained popularity in bioinformatics. ANNs, by mimicking the human brain's function, can learn complex patterns and extract meaningful representations from high-dimensional data. In genomics and disease-gene association studies, Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) are widely used types of ANNs.

For instance, Zhu et al. developed a novel framework for predicting disease-gene associations using embedding graph representation and Graph Convolutional Networks (GCN). Graph embedding methods represent gene interaction networks in low-dimensional vector spaces. GCN, a graph-based deep learning approach, leverages network structure to uncover gene-disease relationships. Their model outperformed other state-of-the-art methods in predicting disease-gene associations on multiple benchmark datasets, illustrating the potential of deep learning approaches in elucidating complex gene-disease relationships.

Despite their significant contributions, machine learning and deep learning methods have limitations in bioinformatics. Deep learning models, often described as 'black boxes,' can be challenging to interpret, complicating understanding the underlying biological mechanisms and drawing meaningful conclusions. Furthermore, the performance of these models heavily depends on the quality and representativeness of training data, appropriate feature selection, and careful hyperparameter tuning.

In conclusion, machine learning and deep learning methods have ushered in transformative changes in bioinformatics, notably in predicting disease-gene associations. These methods facilitate the analysis of genomic data, expose complex gene-disease relationships, and pinpoint potential disease markers. Zhu et al.'s application of graph embedding representation and GCN is a prime example of how deep learning approaches can accelerate the discovery of new disease-related genes and pathways.

- Lung CT scan imaging of the subjects performed
- Preprocessing of the images
- Segmentation of the images
- Extraction of features
- Classification using Convolutional Neural Network (CNN)
- Malignant or benign

## Steps in a deep learning model of lung nodule detection

### GCN and DeepWalk

DeepWalk and Graph Convolutional Networks (GCNs) are two well-known methods for analyzing and modeling complex biological networks and have become increasingly popular in bioinformatics [19].

DeepWalk uses a strategy of embedding nodes into low-dimensional vector spaces, capturing the network structure. It employs a Skip-Gram model to learn node embeddings from sequences of nodes generated by random walks. By mapping nodes to continuous vector representations, DeepWalk facilitates downstream analysis tasks.

On the other hand, GCN is a deep learning architecture explicitly designed for graph-structured data. GCN predicts node-level attributes by aggregating information from neighboring nodes. By harnessing spectral graph theory and convolutional operations, GCN learns node representations that capture the graph's structure locally and globally. GCN has shown exceptional performance in predicting protein-protein interactions, analyzing gene expression, and predicting disease-gene associations.

In the context of disease-gene association prediction, a combined DeepWalk-GCN framework has been proposed. In this approach, DeepWalk is used to generate gene node embeddings from the gene interaction network, capturing the structure of the gene network. The GCN model then utilizes these embeddings to predict disease-gene associations. Notably, this integrated model was able to predict network connectivity patterns accurately.

The study demonstrates the utility of DeepWalk and GCN in leveraging network structure and node embeddings to predict disease-gene associations. DeepWalk extracts informative features from the gene interaction network, which GCN then utilizes to capture complex gene-disease relationships. The synergy of these two methods improves predictive performance and deepens our understanding of gene-disease associations.

### Deep Neural Networks (DNNs)

Deep Neural Networks (DNNs) are potent machine learning models that derive their structure and function from the human brain. These models consist of multiple layers of artificial neurons, nodes, or units, which process and transform input data to generate a prediction. The strength of DNNs lies in their ability to automatically learn hierarchical representations of data, thereby capturing complex patterns and relationships[20].

DNNs have found extensive applications in fields like computer vision, natural language processing (NLP), and bioinformatics. In bioinformatics, DNNs are utilized for gene expression analysis, protein structure prediction, and disease classification, among other tasks. The deep and hierarchical structure of DNNs allows them to model complex biological systems effectively.

The abstract presents a case where a DNN was used to predict genes related to Non-Small Cell Lung Cancer (NSCLC). The DNN model was trained using a large dataset comprising gene expression profiles and clinical data from NSCLC patients. This enabled the model to learn intricate relationships between genes and NSCLC.

Employing its deep architecture, the DNN could automatically learn informative representations of the gene expression data, capturing local and global dependencies between genes and NSCLC. Further, by training on a sizable dataset, the DNN could generalize and accurately predict unseen data. The genes predicted to be related to NSCLC by the DNN provided insights into disease mechanisms and potential therapeutic targets.

In conclusion, DNNs, with their capability to process high-dimensional genomic data and identify meaningful patterns, have shown promise in predicting NSCLC-related genes. Their ability to model complex gene-disease relationships, courtesy of their deep architecture, can aid researchers in identifying new gene candidates and deepening their understanding of the genetics of NSCLC.

### DisGeNET and other Databases

DisGeNET is a comprehensive database widely used in genomics for its extensive compilation of gene-disease associations. It amalgamates data from diverse sources, such as scientific literature, curated databases, and genomic datasets, providing insights into the genetic underpinnings of various diseases[21].

Several other databases, such as OMIM (Online Mendelian Inheritance in Man), GAD (Genetic Association Database), and HGMD (Human Gene Mutation Database), also serve as valuable repositories for gene-disease association data. These databases aggregate datasets from genetic studies, clinical reports, and experimental data, enriching our understanding of gene-disease relationships.

The key advantages of gene-disease association databases include the following:

Comprehensive Coverage: They collate data from numerous sources, encompassing various diseases and genetic variants.

Centralized Repository: These databases provide easy access to a vast collection of gene-disease association data for analysis.

Curation and Quality Control: Many databases employ rigorous curation processes and quality checks to ensure data accuracy and reliability, with experts reviewing and annotating data from literature and experimental studies.

Limitations to the databases:

Data Heterogeneity: The data in these databases originate from various studies that use different methodologies and data formats. Integrating and standardizing such heterogeneous data is challenging.

Incompleteness: Publication bias, limitations of available data, and research constraints may lead to specific gene-disease associations needing to be more represented and included.

Data Quality and Reliability: Despite stringent quality checks, inaccuracies can occur. Researchers need to exercise caution and verify the original data sources.

Lack of Functional Information: While these databases highlight gene-disease associations, they often must elucidate the biological mechanisms driving these relationships.

In conclusion, databases like DisGeNET, OMIM, GAD, and HGMD provide a wealth of information on gene-disease associations. They enhance data coverage, centralization, curation, and integration from diverse sources. However, researchers should be mindful of data heterogeneity, potential gaps in the data, potential quality issues, and the limited functional information available [1,2,3,4].

**Future Directions**

Lung cancer research and gene-disease association analysis continue to evolve, promising to enhance disease comprehension and improve clinical outcomes. This evolution is primarily driven by developments in several areas, including integrating multi-omics data, exploring novel technologies and methodologies, and the advent of personalized medicine.

Integration of Multi-Omics Data: Merging data from diverse domains of the 'omics' field – genomics, transcriptomics, epigenomics, proteomics, and metabolomics – can offer a more comprehensive understanding of lung cancer at the molecular level. By combining these datasets, researchers can uncover complex molecular interactions, identify novel biomarkers, and elucidate disease mechanisms, thereby significantly enhancing gene-disease association analyses.

Exploration of Emerging Technologies and Methodologies: Rapid technological advancements profoundly impact lung cancer research. Single-cell sequencing elucidates tumor cellular heterogeneity, providing insights into the tumor microenvironment and cellular interactions. Spatial transcriptomics and imaging techniques contribute to identifying spatially regulated genes and therapeutic targets. Furthermore, AI and machine learning algorithms can detect complex patterns in large genomic datasets, underscoring their potential utility in this field.

Implementation of Personalized Medicine: Personalized medicine, with treatments tailored to an individual's genetic profile, disease characteristics, and clinical factors, is a rising paradigm. Targeted therapies and immunotherapies have already revolutionized cancer treatment. Utilizing genomic data in conjunction with patient demographics, treatment history, and lifestyle factors may aid in the formulation of personalized treatment plans. Machine learning algorithms can facilitate treatment selection, prognostic assessment, and patient stratification.

Functional Annotation and Validation: The volume of data generated by gene-disease association studies necessitates functional annotation and validation of the identified genes. Experimental models (in vitro and in vivo), functional genomics, and high-throughput screening can elucidate the biological functions of lung cancer genes. These studies can establish causality, clarify molecular mechanisms, and identify druggable targets.

Big Data and Collaborative Efforts: The integration and analysis of big data, sourced from electronic health records, population-based cohorts, and multi-center studies, can unveil novel gene-disease associations and population-specific patterns. Collaboration and data sharing foster large-scale analyses, enhancing statistical power and producing more robust findings.

In conclusion, the future of lung cancer research is bright, with advancements in multi-omics data integration, emerging technologies, personalized medicine, functional annotation, and collaboration promising to catalyze discoveries. These developments could improve diagnostics, treatments, and patient outcomes.

## 3. Conclusion

This survey paper highlights the significant contributions of gene-disease association analysis in lung cancer research, driven by machine learning and deep learning techniques. These methodologies have revolutionized the interpretation of vast genomic data, enabling the identification of biomarkers, disease subtypes, and therapeutic targets. Challenges remain in identifying causal variants and validating their functional relevance. Next-generation sequencing technologies have facilitated comprehensive genomic profiling but require robust bioinformatics pipelines and data management solutions.

Machine learning and deep learning have transformed bioinformatics, paving the way for personalized medicine predictive models. Gene-disease association databases play a crucial role, necessitating data quality, coverage, and biases considerations. Our understanding of the molecular basis of lung cancer has significantly improved through the integration of multi-omics data, advanced computational methods, and innovative technologies. These advancements hold great promise for targeted therapies, personalized medicine, and improved patient outcomes in lung cancer diagnosis, prognosis, and treatment.

## References:

[1] Knox, S. S. (2010, April 26). *From "omics" to complex disease: a systems biology approach to gene-environment interactions in cancer*. PubMed Central (PMC). https://doi.org/10.1186/1475-2867-10-11

[2] Thandra, K. C., Barsouk, A., Saginala, K., Aluru, J. S., & Barsouk, A. (2021, February 23). *Epidemiology of lung cancer*. PubMed Central (PMC). https://doi.org/10.5114/wo.2021.103829

[3] Molina, J. R., Yang, P., Cassivi, S. D., Schild, S. E., & Adjei, A. A. (n.d.). *Non–Small Cell Lung Cancer: Epidemiology, Risk Factors, Treatment, and Survivorship*. PubMed Central (PMC). https://doi.org/10.4065/83.5.584

[4] *Machine learning for multi-omics data integration in cancer*. (2022, January 22). Machine Learning for Multi-omics Data Integration in Cancer - ScienceDirect. https://doi.org/10.1016/j.isci.2022.103798

[5] Ruan, X., Ye, Y., Cheng, W., Xu, L., Huang, M., Chen, Y., Zhu, J., Lu, X., & Yan, F. (2022, May 6). *Multi-Omics Integrative Analysis of Lung Adenocarcinoma: An in silico Profiling for Precise Medicine*. Frontiers. https://doi.org/10.3389/fmed.2022.894338

[6] Cano-Gamez, E., & Trynka, G. (2020, April 6). *From GWAS to Function: Using Functional Genomics to Identify the Mechanisms Underlying Complex Diseases*. Frontiers. https://doi.org/10.3389/fgene.2020.00424

[7] Pai, A. (2020, February 17). *CNN vs. RNN vs. ANN - Analyzing 3 Types of Neural Networks in Deep Learning*. Analytics Vidhya. https://www.analyticsvidhya.com/blog/2020/02/cnn-vs-rnn-vs-mlp-analyzing-3-types-of-neural-networks-in-deep-learning/

[8] *A novel candidate disease gene prioritization method using deep graph convolutional networks and semi-supervised learning - PubMed*. (2022,

October 14). PubMed. https://doi.org/10.1186/s12859-022-04954-x

[9] Mathew, C. J., David, A. M., & Joy Mathew, C. M. (2020, December 11). *Artificial Intelligence and its future potential in lung cancer screening*. PubMed Central (PMC). https://doi.org/10.17179/excli2020-3095

[10] Vamathevan, J., Clark, D., Czodrowski, P., Dunham, I., Ferran, E., Lee, G., Li, B., Madabhushi, A., Shah, P., Spitzer, M., & Zhao, S. (n.d.). *Applications of machine learning in drug discovery and development*. PubMed Central (PMC). https://doi.org/10.1038/s41573-019-0024-5

[11] Vamathevan, J., Clark, D., Czodrowski, P., Dunham, I., Ferran, E., Lee, G., Li, B., Madabhushi, A., Shah, P., Spitzer, M., & Zhao, S. (n.d.). *Applications of machine learning in drug discovery and development*. PubMed Central (PMC). https://doi.org/10.1038/s41573-019-0024-5

[12] *Genome, transcriptome, and proteome: the rise of omics data and their integration in biomedical sciences - PubMed*. (2018, March 1). PubMed. https://doi.org/10.1093/bib/bbw114

[13] Krassowski, M., Das, V., Sahu, S. K., & Misra, B. B. (2020, November 20). *State of the Field in Multi-Omics Research: From Computational Needs to Data Mining and Sharing*. Frontiers. https://doi.org/10.3389/fgene.2020.610798

[14] *Translational bioinformatics and healthcare informatics: computational and ethical challenges - PubMed*. (2009, September 16). PubMed. https://pubmed.ncbi.nlm.nih.gov/20169020/

[15] *https://academic.oup.com/nar/article/45/D1/D833/2290909*. (n.d.). https://academic.oup.com/nar/article/45/D1/D833/2290909

[16] Witte, J. S. (n.d.). *Genome-Wide Association Studies and Beyond*. PubMed Central (PMC). https://doi.org/10.1146/annurev.publhealth.012809.103723

[17] *Next Generation Sequencing - an overview | ScienceDirect Topics*. (n.d.). Next Generation Sequencing - an Overview | ScienceDirect Topics. https://doi.org/10.1016/B978-0-12-801418-9.00002-0

[18] Abass, Y. A., & Adeshina, S. A. (n.d.). *Deep Learning Methodologies for Genomic Data Prediction: Review | Atlantis Press*. Deep Learning Methodologies for Genomic Data Prediction: Review | Atlantis Press. https://doi.org/10.2991/jaims.d.210512.001

[19] Muzio, G., O'Bray, L., & Borgwardt, K. (2021, March 1). *Biological network analysis with deep

*learning*. OUP Academic. https://doi.org/10.1093/bib/bbaa257

[20] *State-of-the-art in artificial neural network applications: A survey*. (2018, November 23). State-of-the-art in Artificial Neural Network Applications: A Survey - ScienceDirect. https://doi.org/10.1016/j.heliyon.2018.e00938

[21] Piñero, J., Bravo, L., Queralt-Rosinach, N., Gutiérrez-Sacristán, A., Deu-Pons, J., Centeno, E., García-García, J., Sanz, F., & Furlong, L. I. (2016, October 19). *DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants*. PubMed Central (PMC). https://doi.org/10.1093/nar/gkw943