# Data-Driven Cheminformatics Models for Predicting Bioactivity of Natural Compounds in Oncology

Salvation Ifechukwude Atalor[1]

[1] Department of Pharmacy, Imanzoe Drug Store, Abuja, Nigeria

## Abstract

Advancements in data-driven cheminformatics have significantly transformed the early-stage discovery and optimization of oncology therapeutics derived from natural compounds. This review examines the integration of machine learning (ML) and quantitative structure–activity relationship phytochemicals and marine-derived agents. Emphasis is placed on the use of high-dimensional molecular descriptors, fingerprinting techniques, and graph-based neural networks for feature extraction and predictive modeling. Public bioactivity databases such as ChEMBL, PubChem BioAssay, and BindingDB are explored as primary sources for curated compound-target interaction data, which underpin supervised learning frameworks. Furthermore, the review highlights recent breakthroughs in multi-task learning, deep generative models, and transfer learning paradigms that enhance generalizability across diverse chemical scaffolds and rare oncogenic targets. Challenges such as model interpretability, data sparsity, and bioavailability prediction are discussed, with proposed strategies including explainable AI (XAI) and hybrid mechanistic-ML models. This review highlights the transformative potential of cheminformatics in accelerating oncology drug discovery by reducing reliance on labor-intensive wet-lab screening and enabling virtual prioritization of lead compounds from vast natural product libraries.

*Keywords: Cheminformatics, Bioactivity Prediction, Natural Compounds, Oncology Drug Discovery, Machine Learning Models.*

## I. INTRODUCTION

➢ *Background on Natural Compounds in Oncology*

Natural compounds have long been recognized as vital sources of chemotherapeutic agents, with extensive documentation showing their critical role in oncology drug development. Between 1981 and 2019, almost half of all approved small-molecule anticancer drugs were either natural products, semi-synthetic derivatives, or pharmacophore-inspired synthetic compounds, hilighting the unparalleled chemical diversity they offer (Newman & Cragg, 2020). These molecules, characterized by intricate stereochemistry and diverse functional groups, present opportunities to modulate complex biological targets that synthetic libraries often fail to achieve. Prominent examples include paclitaxel from Taxus brevifolia, vincristine from Catharanthus roseus, and camptothecin analogs, all of which revolutionized chemotherapy protocols.

The importance of natural product scaffolds lies not only in their bioactivity but also in their ability to inspire novel synthetic modifications for enhanced efficacy and reduced toxicity (Li & Vederas, 2009). Natural compounds exert anticancer effects through diverse mechanisms such as microtubule stabilization, topoisomerase inhibition, and induction of programmed cell death pathways. Moreover, emerging sources such as marine organisms and endophytic fungi have expanded the oncology pipeline with agents like trabectedin and salinosporamide A. Despite their successes, challenges such as supply limitations, complex isolation, and suboptimal pharmacokinetics necessitate the adoption of computational cheminformatics to accelerate virtual screening, scaffold optimization, and drug-like property prediction.

➢ *Emergence of Cheminformatics and Data-Driven Approaches*

The integration of cheminformatics into oncology research marks a paradigm shift from traditional empirical screening toward data-driven drug discovery models.

Cheminformatics, encompassing quantitative structure–activity relationship (QSAR) modeling, molecular descriptor analysis, and virtual screening, enables the systematic evaluation of large chemical libraries for potential anticancer activity (Cherkasov et al., 2014). By leveraging curated datasets and computational models, researchers can now predict the bioactivity, pharmacokinetics, and toxicity of natural compounds before costly and time-intensive experimental validation.

The emergence of machine learning and artificial intelligence (AI) further amplified cheminformatics capabilities, introducing generative algorithms capable of proposing novel chemical entities optimized for oncogenic targets (Walters & Murcko, 2020). These AI-driven platforms, employing techniques such as deep generative models and reinforcement learning, have shown promise in navigating the vast chemical space of natural products to identify bioactive analogs with improved drug-like properties. For example, deep learning frameworks can predict molecular fingerprints from limited bioactivity data, enabling virtual hit expansion for rare cancer subtypes. Additionally, cheminformatics tools now allow multitarget profiling, which is crucial for addressing the polypharmacology characteristic of many oncological diseases. Consequently, the convergence of data science and chemistry has positioned cheminformatics as an indispensable engine driving the rapid, cost-effective discovery of new oncology therapeutics sourced from natural compounds.

➢ *Scope and Objectives of the Review*

This review systematically explores how data-driven cheminformatics models are advancing the prediction of bioactivity for natural compounds in oncology drug discovery. The study covers the entire computational workflow, beginning with data acquisition from public bioactivity repositories, followed by molecular representation techniques such as descriptors, fingerprints, and graph-based models. It then examines machine learning and deep learning methods used to build predictive models capable of prioritizing promising natural compounds for anticancer activity.

The objectives of this review are threefold. First, to outline the current state of computational methodologies that enable virtual screening and bioactivity prediction in natural product libraries. Second, to evaluate case studies where cheminformatics successfully enhanced the identification of bioactive oncology candidates from natural sources. Third, to critically assess the challenges in model interpretability, data quality, and generalization across complex cancer phenotypes, while suggesting emerging solutions such as hybrid modeling and explainable AI techniques. Overall, this review aims to provide a comprehensive and technical guide for researchers and practitioners interested in accelerating oncology drug discovery using computational approaches focused on natural products.

➢ *Structure of the Paper*

The paper is organized into six major sections. Section 1 provides an introduction, covering the background of natural compounds in oncology, the emergence of cheminformatics, and the objectives of the review. Section 2 discusses data sources and molecular representations critical for modeling, including public databases and molecular encoding techniques. Section 3 examines machine learning and deep learning approaches applied to predict the bioactivity of natural products. Section 4 presents practical applications, highlighting case studies involving phytochemicals and marine-derived compounds. Section 5 addresses the current challenges in data-driven cheminformatics and reviews emerging solutions such as explainable AI and hybrid modeling. Finally, Section 6 offers future perspectives, emphasizing new computational strategies and the evolving role of cheminformatics in precision oncology.

## II. DATA SOURCES AND MOLECULAR REPRESENTATIONS

➢ *Public Bioactivity Databases: ChEMBL, PubChem BioAssay, BindingDB*

The foundation of data-driven cheminformatics models lies in access to robust and well-curated bioactivity databases. ChEMBL, PubChem BioAssay, and BindingDB are three pivotal resources that facilitate the virtual exploration of chemical-biological interaction landscapes critical for oncology research. ChEMBL is a manually curated database containing bioactivity information on over 1.9 million compounds, including data on binding, functional assays, and ADMET properties derived primarily from peer-reviewed publications (Gaulton et al., 2017). For natural products in oncology, ChEMBL provides structured annotations linking compounds to cancer-relevant targets such as kinases, nuclear receptors, and epigenetic modulators, thereby supporting predictive modeling endeavors.

Similarly, PubChem BioAssay serves as an expansive repository hosting over one million bioactivity assay records, integrating information from high-throughput screening (HTS) campaigns, particularly for oncology-related targets like p53 modulators and tyrosine kinase inhibitors (Wang et al., 2017). PubChem's integration with compound structures, assay descriptions, and experimental conditions makes it a valuable asset for constructing training datasets aimed at machine learning applications in virtual screening and drug repositioning.

BindingDB specifically focuses on binding affinities, documenting over 1.2 million binding measurements across protein-ligand complexes, with notable emphasis on kinases, GPCRs, and oncology-relevant enzymes. By providing kinetic parameters such as $IC_{50}$, $K_i$, and $K\_d$ values, BindingDB enables quantitative structure–activity relationship (QSAR) modeling with enriched data quality as represented in Table 1.

Together, these databases enable the aggregation of comprehensive datasets necessary for training supervised learning models, feature extraction, and validating bioactivity predictions for natural compounds in oncology.

The strategic utilization of these resources ensures that cheminformatics pipelines are anchored in reproducible, high-confidence biological evidence.

Table 1 Summary of Key Public Bioactivity Databases for Oncology Modeling

| Database Name | Focus Area | Key Features | Role in Oncology Modeling |
|---|---|---|---|
| ChEMBL | Curated bioactivity data from literature sources | Compound-target bioactivities, ADMET properties, cancer-relevant targets | Supports predictive modeling for oncology pathways and virtual lead prioritization |
| PubChem BioAssay | High-throughput screening assay data | Assay descriptions, compound structures, experimental conditions metadata | Enables large-scale dataset construction for machine learning in virtual screening |
| BindingDB | Experimental binding affinity data | $IC_{50}$, $K_i$, and $K\_d$ values across proteins like kinases, GPCRs, and enzymes | Provides quantitative affinity data essential for QSAR modeling and target specificity predictions |
| ChEMBL | Curated bioactivity data from literature sources | Compound-target bioactivities, ADMET properties, cancer-relevant targets | Supports predictive modeling for oncology pathways and virtual lead prioritization |

> *Molecular Descriptors, Fingerprints, and Graph Representations*

Accurately representing chemical structures is a fundamental requirement for building predictive cheminformatics models. Molecular descriptors are numerical values derived from molecular graphs, characterizing properties such as topology, geometry, electronic distribution, and atom connectivity (Todeschini & Consonni, 2009). These descriptors, which include simple counts (e.g., molecular weight, hydrogen bond donors) and complex indices (e.g., topological polar surface area, Wiener index), translate intricate molecular structures into machine-readable formats essential for supervised learning models. In oncology-focused modeling, descriptors sensitive to pharmacophoric features like planarity and lipophilicity are particularly critical for predicting bioactivity profiles.

Molecular fingerprints offer a more compact representation, encoding the presence or absence of substructures or chemical patterns as binary or hashed vectors. Popular fingerprinting methods, such as Extended Connectivity Fingerprints (ECFPs) and MACCS keys, allow for rapid similarity searches and clustering operations critical for virtual screening of natural product libraries. Fingerprints also serve as input features for classical machine learning algorithms like random forests, support vector machines, and gradient boosting frameworks in bioactivity prediction tasks.

Graph-based representations, enabled by advancements in deep learning, have further revolutionized molecular encoding. Unlike traditional descriptors or fingerprints, graph convolutional networks (GCNs) process molecules directly as graphs, where atoms are nodes and bonds are edges, learning hierarchical chemical features automatically (Duvenaud et al., 2015). This approach eliminates the need for hand-crafted features and has demonstrated superior performance in modeling natural product-derived compounds with complex fused ring systems and nonstandard functional groups. Thus, the choice of molecular representation significantly influences the performance, interpretability, and generalization capabilities of cheminformatics models targeting oncological applications.

> *Data Preprocessing and Curation Challenges*

Reliable predictive modeling in cheminformatics hinges critically on the quality and integrity of the input data. Raw chemical databases often contain redundancies, structural errors, and inconsistent annotations, necessitating rigorous preprocessing and curation protocols to ensure model robustness and generalizability (Fourches, Muratov, & Tropsha, 2016). In oncology-focused cheminformatics pipelines, particular care must be taken to correct tautomeric inconsistencies, standardize protonation states, remove salts and counterions, and verify stereochemistry. Failure to address these artifacts can propagate noise into machine learning models, leading to erroneous predictions of bioactivity for natural compounds.

An essential step in preprocessing is structure normalization, wherein variations in molecular representation, such as depiction of aromatic systems or charge assignments, are harmonized across datasets. Equally critical is the identification and removal of duplicate records, which can artificially inflate model performance metrics if inadvertently included in both training and testing sets. Williams (2012) emphasized that even high-profile databases like ChEMBL and PubChem may require secondary curation to eliminate misdrawn structures, incomplete entries, and ambiguous activity annotations.

Beyond chemical structure validation, bioactivity data must also be carefully filtered. Challenges such as ambiguous assay descriptions, inconsistent endpoint measurements (e.g., $IC_{50}$ vs. $EC_{50}$), and varying experimental conditions can undermine model fidelity if not meticulously curated. Cross-referencing experimental protocols, converting endpoint metrics into standardized units, and flagging unreliable data points are recommended

best practices. Moreover, the imbalanced distribution of active versus inactive compounds, typical in natural product oncology datasets, necessitates strategic sampling or augmentation techniques to mitigate model bias as shown in Figure 1. Consequently, robust data curation is foundational to building high-confidence cheminformatics models capable of accurately predicting anticancer bioactivities in complex natural product libraries.
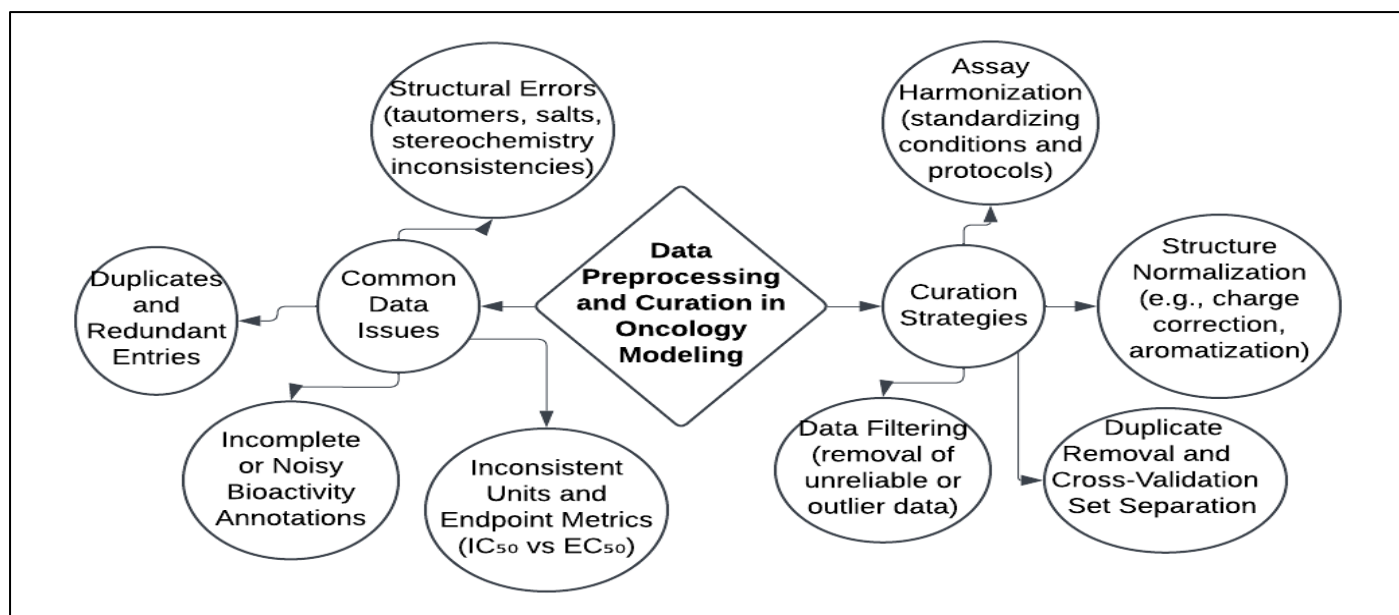


Fig 1 Diagram Showing Preprocessing Challenges and Curation Strategies in Oncology Cheminformatics.

Fig 1 presents a high-level outline of the key challenges and solutions involved in data preprocessing and curation for oncology-focused cheminformatics. At the center is the critical role of data quality in enabling reliable predictive modeling. The first branch highlights common data issues, including structural inconsistencies such as unnormalized tautomers, salts, and stereochemistry errors, as well as problems like duplicate entries, incomplete bioactivity records, and inconsistent assay units (e.g., $IC_{50}$ vs. $EC_{50}$). The second branch outlines strategic curation responses to these issues, including structure normalization to standardize molecular representations, assay harmonization to align experimental conditions, rigorous filtering to remove unreliable data, and careful duplicate removal to prevent data leakage across model training and validation. Together, these curated processes ensure that input data is clean, consistent, and suitable for high-confidence modeling in oncology drug discovery using natural products.

## III. MACHINE LEARNING AND PREDICTIVE MODELING IN CHEMINFORMATICS

➢ *Traditional QSAR Models and Regression Techniques*

Quantitative Structure–Activity Relationship (QSAR) modeling forms the foundational core of traditional cheminformatics approaches for predicting the bioactivity of chemical compounds, including natural products in oncology. QSAR models establish mathematical relationships between molecular descriptors and biological activity endpoints, providing a framework to infer the activity of untested compounds from known structural features (Cherkasov et al., 2014). Linear models, such as multiple linear regression (MLR) and partial least squares (PLS), have historically dominated early QSAR applications due to their interpretability and computational simplicity. These methods rely on the assumption that biological activity is a linear combination of selected molecular descriptors, making them particularly attractive for early-stage virtual screening.

However, real-world bioactivity data often exhibit non-linear relationships, necessitating the use of non-linear regression techniques such as support vector regression (SVR) and k-nearest neighbor (k-NN) algorithms to enhance predictive accuracy. Validation remains a critical component of QSAR modeling, with internal methods (e.g., cross-validation) and external validations (e.g., testing on independent datasets) serving as benchmarks to assess model generalizability (Gramatica, 2007). Rigorous validation ensures that models do not suffer from overfitting, a common pitfall when dealing with high-dimensional natural product datasets.

In oncology drug discovery, classical QSAR models have been used to prioritize natural compounds targeting kinases, DNA topoisomerases, and apoptotic regulators. Despite their historical success, traditional QSAR approaches face limitations when handling highly diverse, structurally complex natural products. Nevertheless, they remain valuable, especially when combined with modern ensemble strategies or used as baseline models for benchmarking more advanced machine learning architectures. Their simplicity, ease of interpretation, and relatively low computational demands make QSAR models enduring tools in cheminformatics pipelines aimed at oncology therapeutics discovery.

> *Deep Learning Architectures: Graph Neural Networks and Autoencoders*

The evolution of deep learning has introduced powerful architectures capable of automatically learning hierarchical features from molecular data, significantly enhancing bioactivity prediction in oncology cheminformatics. Graph Neural Networks (GNNs), particularly Graph Convolutional Networks (GCNs), have gained prominence for their ability to directly operate on molecular graphs without requiring manual feature engineering. In GCNs, molecules are treated as graphs where atoms are nodes and bonds are edges, and the convolutional layers iteratively aggregate information from local atomic neighborhoods to learn task-specific representations (Kipf & Welling, 2017). This approach is particularly advantageous for modeling complex natural compounds with non-canonical structures, fused ring systems, and diverse functional groups, all of which challenge traditional descriptor-based models.

Complementing GNNs, autoencoders serve as another deep learning framework with profound utility in cheminformatics. Autoencoders are neural networks trained to reconstruct input data after compressing it into a lower-dimensional latent space, effectively learning compact, information-rich representations (Hinton & Salakhutdinov, 2006). In oncology drug discovery, autoencoders are employed for unsupervised feature extraction from large chemical libraries, denoising noisy molecular descriptors, and even generating novel molecular structures through variational extensions. Their ability to capture subtle variations in chemical space makes them invaluable for navigating the highly diverse chemical scaffolds typical of natural products.

The integration of GNNs and autoencoders into cheminformatics pipelines enables end-to-end learning workflows that bypass the limitations of handcrafted features as showm in Figure 2. These deep learning models excel not only in predicting compound bioactivity but also in uncovering hidden relationships between molecular structure and anticancer efficacy. As a result, they represent a transformative leap forward in the application of artificial intelligence to oncology-focused natural product discovery.
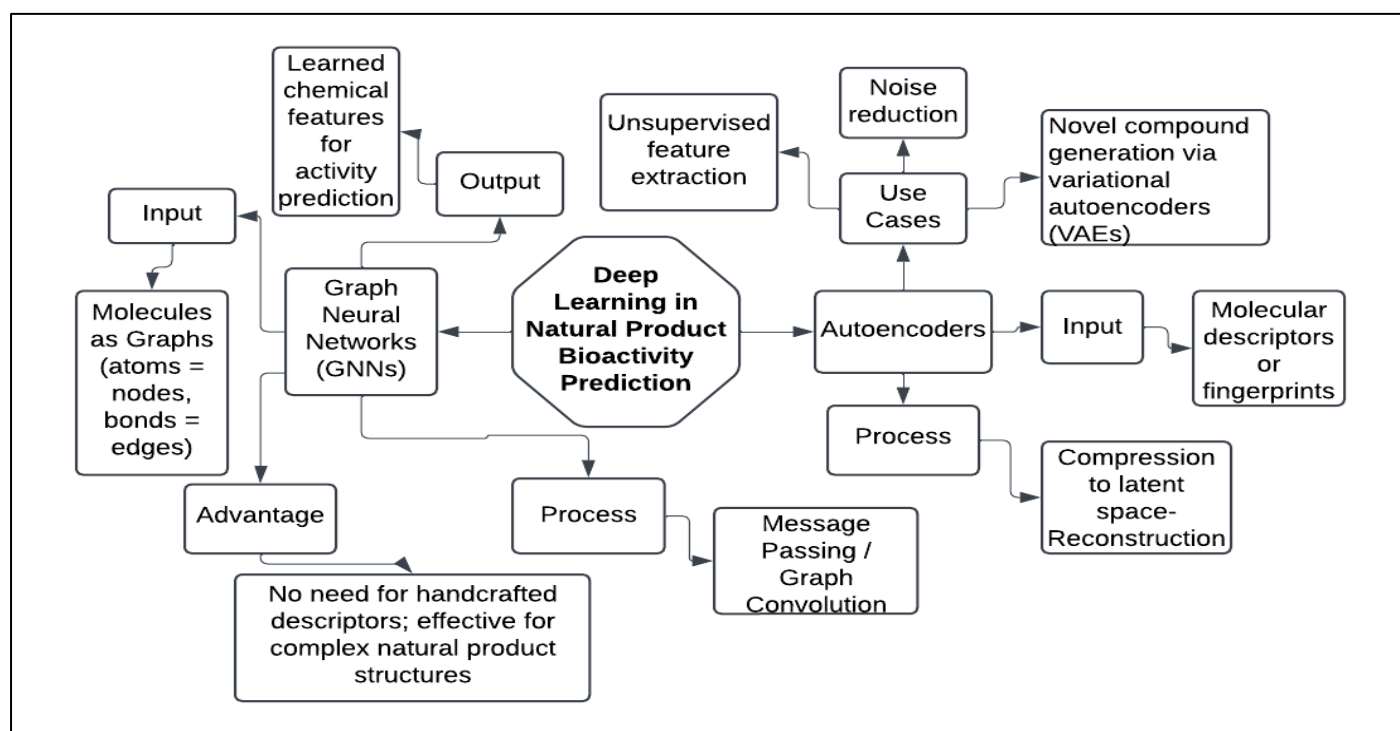


Fig 2 A Block Diagram Showing GNN and Autoencoder Architectures for Bioactivity Prediction

Figure 2 illustrates two key deep learning architectures used in oncology-focused cheminformatics: Graph Neural Networks (GNNs) and Autoencoders. At the center is a shared goal—predicting the bioactivity of natural compounds. The first branch, GNNs, processes molecular structures as graphs, treating atoms as nodes and bonds as edges, allowing the model to learn structural relationships through message passing and graph convolution. This enables accurate prediction of chemical behavior without the need for predefined descriptors, particularly useful for modeling complex natural products. The second branch, Autoencoders, begins with molecular descriptors or fingerprints and compresses them into a lower-dimensional latent space before reconstructing the input. This architecture supports unsupervised feature learning, denoising, and compound generation. Together, these architectures enhance predictive performance and flexibility in cheminformatics pipelines by extracting rich, hierarchical representations of molecular data tailored for oncology applications.

> *Multi-Task Learning and Transfer Learning in Oncology Applications*

Multi-task learning (MTL) has emerged as a powerful strategy in cheminformatics, where multiple related prediction tasks are learned simultaneously by a shared

model, leading to improved generalization and efficiency. Instead of training separate models for each oncological target or bioactivity endpoint, MTL architectures optimize shared representations that exploit correlations across tasks, such as overlapping signaling pathways or structural similarities among ligands (Ruder, 2017). For natural compound oncology modeling, MTL frameworks enable simultaneous predictions across multiple cancer biomarkers, such as kinase inhibition profiles, cytotoxicity assays, and apoptotic activity metrics, thereby reducing the risk of overfitting on sparse datasets typical of natural product libraries.

Transfer learning further enhances predictive performance, especially when available labeled data is limited, as is often the case with rare or structurally unique natural compounds. Transfer learning involves pretraining a model on a large source domain and fine-tuning it on a smaller, related target domain (Pan & Yang, 2010). In cheminformatics, models initially trained on broad chemical bioactivity datasets, such as ChEMBL or PubChem assays, can be repurposed to predict bioactivities in specialized oncology datasets with minimal additional training. This approach significantly reduces data requirements and computational costs while preserving predictive accuracy.

Recent innovations combine MTL and transfer learning to create hierarchical architectures that perform coarse-grained screening at early layers and specialized bioactivity predictions at deeper layers. Such hybrid models are particularly suited for natural product oncology research, where leveraging generalized chemical knowledge while capturing specific biological nuances is essential as presented in Table 2 . Together, multi-task and transfer learning paradigms represent critical advancements in building more versatile, efficient, and robust cheminformatics models capable of accelerating the discovery of new anticancer agents from nature's chemical repertoire.

Table 2 Multi-Task and Transfer Learning in Oncology Modeling

| Approach | Purpose | Key Techniques | Impact on Natural Product Oncology Modeling |
|---|---|---|---|
| Multi-Task Learning (MTL) | Simultaneously learn multiple related bioactivity prediction tasks | Shared neural network architectures, task-specific output heads | Improves generalization across diverse oncogenic targets and enhances learning from limited datasets |
| Transfer Learning (TL) | Transfer knowledge from broad datasets to specialized oncology targets | Pretraining on large chemical bioactivity datasets, fine-tuning on rare cancer targets | Enables accurate bioactivity prediction for rare targets with minimal labeled data |
| Hybrid MTL-TL Models | Combine MTL and TL for hierarchical task learning | Coarse-grained pretraining followed by fine-tuning on specific oncology biomarkers | Boosts model adaptability, accelerates lead discovery for underexplored cancer pathways |
| Multi-Task Learning (MTL) | Simultaneously learn multiple related bioactivity prediction tasks | Shared neural network architectures, task-specific output heads | Improves generalization across diverse oncogenic targets and enhances learning from limited datasets |

## IV. APPLICATIONS IN PREDICTING BIOACTIVITY OF NATURAL PRODUCTS

➢ *Case Studies on Phytochemicals and Marine-Derived Agents*

Natural products derived from terrestrial plants and marine organisms have yielded some of the most significant breakthroughs in oncology drug development. Phytochemicals such as paclitaxel from *Taxus brevifolia* and vinblastine from *Catharanthus roseus* exemplify the immense therapeutic potential locked within terrestrial biodiversity (Newman & Cragg, 2016). Paclitaxel operates through the stabilization of microtubules, arresting cancer cell division, and remains a cornerstone treatment for breast, ovarian, and lung cancers. Similarly, camptothecin, isolated from *Camptotheca acuminata*, led to the development of topotecan and irinotecan, critical therapies targeting DNA topoisomerase I, an essential enzyme for DNA replication in rapidly proliferating cancer cells.

Marine ecosystems have also proven to be invaluable reservoirs of anticancer agents. Compounds such as trabectedin, derived from the sea squirt *Ecteinascidia turbinata*, exhibit unique DNA minor groove binding properties that disrupt transcription processes and tumor cell survival (Molinski et al., 2009). Another notable marine-derived agent, salinosporamide A, isolated from *Salinispora tropica*, acts as a potent proteasome inhibitor and has shown efficacy in multiple myeloma models. These marine products often feature highly complex and novel chemical scaffolds that defy synthetic mimicry, making them irreplaceable sources for new drug discovery.

Data-driven cheminformatics models are increasingly applied to expedite the identification and optimization of such bioactive natural compounds. Virtual screening campaigns now leverage predictive modeling to prioritize phytochemicals and marine metabolites with high binding affinities to cancer-specific targets as presented in Table 3. Consequently, case studies in phytochemical and marine oncology reinforce the indispensable role of natural products while demonstrating how modern computational techniques are revitalizing their exploration for next-generation cancer therapeutics.

Table 3 Key Natural Compounds from Phytochemical and Marine Sources in Oncology

| Source | Example Compound | Mechanism of Action | Oncological Application |
|---|---|---|---|
| Terrestrial Plants | Paclitaxel (Taxus brevifolia) | Stabilizes microtubules, inhibits mitosis | Treatment of breast, ovarian, and lung cancers |
| Terrestrial Plants | Camptothecin (Camptotheca acuminata) | Inhibits DNA topoisomerase I | Basis for topotecan and irinotecan in solid tumor therapies |
| Marine Organisms | Trabectedin (Ecteinascidia turbinata) | Binds DNA minor groove, disrupts transcription | Approved for soft tissue sarcomas and ovarian cancer |
| Marine Organisms | Salinosporamide A (Salinispora tropica) | Inhibits proteasome function | Investigated for multiple myeloma and hematologic malignancies |

> *Target Specificity, Off-Target Effects, and Toxicity Predictions*

Predicting target specificity and minimizing off-target effects are critical components of oncology drug discovery, especially when working with structurally complex natural products. Computational approaches, including cheminformatics and machine learning, have become indispensable for early-stage prediction of both desired interactions and potential toxic liabilities. Target specificity modeling focuses on identifying compounds that interact selectively with oncogenic targets while avoiding promiscuous binding to non-cancer-related proteins, a major cause of side effects (Ekins & Williams, 2010). Techniques such as structure-based virtual screening and ligand-based similarity models are widely applied to rank natural compounds based on predicted binding affinities and selectivity indices.

Off-target profiling leverages predictive models trained on large datasets containing known bioactivity profiles across multiple protein families. By using molecular descriptors and fingerprint similarity searches, researchers can anticipate potential unintended interactions, flagging molecules likely to cause cardiotoxicity, hepatotoxicity, or neurological side effects. Toxicity prediction models also integrate chemical property thresholds (e.g., lipophilicity, molecular weight) and machine learning classifiers trained to identify structural alerts associated with toxicological outcomes (Cheng et al., 2012).

In oncology, where therapeutic windows can be narrow, early computational prediction of off-target effects helps prioritize lead compounds that balance efficacy and safety. Furthermore, advances in multi-target modeling allow simultaneous evaluation of a molecule's polypharmacological landscape, enhancing the ability to uncover synergistic interactions or harmful cross-activities. Incorporating predictive toxicity and off-target analytics into the cheminformatics workflow not only improves success rates in downstream experimental validation but also accelerates the translation of natural compounds into clinically viable oncology therapeutics.

> *Integration of Virtual Screening and Hit-to-Lead Prioritization*

Virtual screening (VS) represents a cornerstone in modern cheminformatics pipelines, particularly for oncology drug discovery using natural products. Structure-based and ligand-based virtual screening methodologies enable the rapid identification of potential bioactive compounds from large chemical libraries by evaluating their fit to specific biological targets or known active ligands (Lionta et al., 2014). In the context of natural products, VS allows researchers to navigate vast molecular diversity, prioritizing compounds with favorable binding affinity, pharmacophoric features, and drug-likeness properties, thus reducing experimental burden.

Integration with hit-to-lead prioritization strategies is critical to move beyond simple virtual hits toward compounds with optimized efficacy, selectivity, and developability. After virtual screening, candidate compounds are further evaluated using cheminformatics-driven scoring functions that predict ADMET properties, synthetic accessibility, and potential off-target effects. Techniques such as consensus scoring, where multiple predictive models are combined, increase the reliability of lead selection by mitigating bias introduced by individual algorithms (Schneider, 2010).

In natural product oncology, the hit-to-lead process often involves iterative refinement cycles, wherein top-ranked virtual hits undergo additional structure-activity relationship (SAR) modeling, docking studies, and even early-stage molecular dynamics simulations to validate binding modes. Cheminformatics models guide medicinal chemists in suggesting modifications to improve bioavailability, metabolic stability, and target engagement without compromising the structural integrity crucial to natural product-derived compounds.

Thus, the tight integration of virtual screening with systematic hit-to-lead optimization accelerates the identification of high-quality, bioactive natural compounds poised for further preclinical development as shown in Figure 3. This approach exemplifies how computational modeling transforms traditional discovery workflows into dynamic, data-driven engines for efficient oncology drug development.
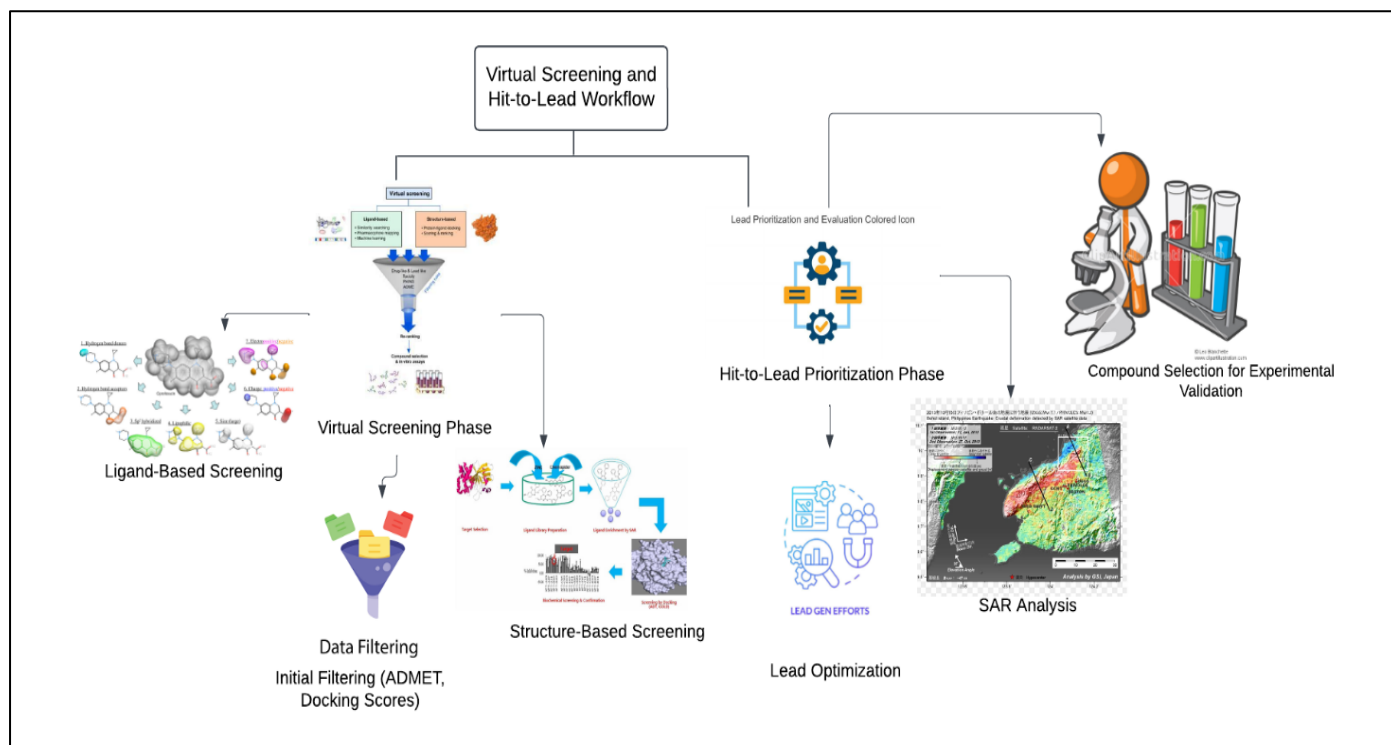
Fig 3 Workflow Diagram for Streamlined Virtual Screening and Lead Optimization

Figure 3 illustrates a streamlined workflow integrating virtual screening with hit-to-lead prioritization in oncology-focused cheminformatics. At the core is the combined process, which begins with the virtual screening phase, where both structure-based and ligand-based techniques are employed to evaluate large libraries of natural compounds. This phase includes initial filtering using docking scores and drug-likeness criteria such as ADMET properties. The process then transitions into the hit-to-lead prioritization phase, where selected virtual hits undergo structure-activity relationship (SAR) analysis and lead optimization. This step focuses on refining compound potency, selectivity, and pharmacokinetic properties to identify the most promising candidates for experimental validation. Together, these two phases create a cohesive pipeline that accelerates the discovery of effective anticancer agents from natural product sources.

## V. CHALLENGES AND EMERGING SOLUTIONS

➢ *Data Sparsity, Imbalanced Datasets, and Model Generalization*

Data sparsity and class imbalance present significant challenges in cheminformatics modeling for natural product oncology discovery. Sparsity arises due to the limited availability of experimentally validated bioactivity data for many natural compounds, leading to datasets with incomplete feature-target mappings. Machine learning models trained on sparse data often exhibit poor generalization, overfitting to the few available active compounds and failing to predict the bioactivity of novel scaffolds accurately (He & Garcia, 2009).

Compounding this issue, natural product datasets are inherently imbalanced, with a disproportionate number of inactive or weakly active compounds relative to potent bioactives. Standard classifiers tend to favor the majority class, resulting in high overall accuracy but poor sensitivity in identifying the minority (active) class — a critical failure in oncology drug discovery where identifying rare hits is paramount (Sun, Wong, & Kamel, 2009). Techniques such as Synthetic Minority Over-sampling Technique (SMOTE), adaptive synthetic sampling (ADASYN), and cost-sensitive learning have been developed to address imbalance by either enriching minority class samples or penalizing misclassifications asymmetrically during model training.

Model generalization in the presence of sparse and imbalanced data requires careful architectural and methodological choices. Ensemble methods like random forests and gradient boosting are often preferred due to their inherent robustness to noisy or skewed data. Additionally, regularization techniques such as dropout, weight decay, and early stopping help prevent overfitting in deep learning models. Data augmentation strategies, including the generation of virtual compounds via SMILES-based perturbations or molecular graph augmentations, further enhance diversity and help models capture underlying chemical-biological relationships more effectively as presented in Table 4.

Addressing sparsity and imbalance is not merely a preprocessing concern; it is integral to building cheminformatics models capable of reliably predicting novel oncology therapeutics from underexplored natural product spaces.

Table 4 Challenges and Solutions for Sparse and Imbalanced Oncology Datasets

| Challenge | Description | Key Techniques to Address It | Impact on Modeling |
|---|---|---|---|
| Data Sparsity | Limited availability of labeled bioactivity data for natural compounds | Data augmentation, transfer learning, virtual sample generation | Reduces overfitting, improves generalization to unseen compounds |
| Imbalanced Datasets | Dominance of inactive compounds over active ones in datasets | SMOTE, cost-sensitive learning, resampling techniques | Enhances model sensitivity to minority (active) class predictions |
| Model Overfitting | Models memorize training data rather than generalize | Regularization methods (dropout, weight decay), early stopping | Improves predictive robustness on new natural product scaffolds |
| Poor Generalization | Difficulty in applying models to novel or rare chemical structures | Ensemble methods, few-shot learning, graph data augmentation | Expands chemical space coverage, boosts reliability in oncology predictions |

➤ *Explainable AI (XAI) for Model Interpretability*

The increasing complexity of machine learning models in cheminformatics, particularly deep learning architectures, has led to concerns regarding their interpretability and transparency. In oncology drug discovery from natural products, where decision-making impacts critical therapeutic directions, black-box models are insufficient. Explainable AI (XAI) addresses this challenge by providing mechanisms to render machine learning predictions understandable to domain experts without compromising predictive performance (Gilpin et al., 2018).

Several XAI techniques have been tailored for cheminformatics. Feature attribution methods, such as SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations), identify which molecular substructures or descriptors contribute most significantly to a given bioactivity prediction. These methods enable researchers to discern whether a model's reasoning aligns with known pharmacophoric features or unexpected chemical motifs, facilitating trust and validation as Shown in Figure 4.

Figure 4 presents a structured sketch of how Explainable AI (XAI) enhances interpretability in cheminformatics models used for oncology drug discovery. At the core is the challenge of black-box complexity in deep learning models, particularly when high-stakes decisions depend on model outputs. The diagram branches into two major XAI techniques — SHAP and LIME — which provide molecular-level insights by identifying the contribution of specific features or substructures to prediction outcomes. Additional branches outline the benefits these techniques offer, such as aligning predictions with pharmacophoric expectations, enabling domain expert trust, and improving model validation. Finally, the strategic impact of XAI is emphasized, showing how it facilitates regulatory transparency, interpretable compound prioritization, and more cohesive collaboration between computational scientists and medicinal chemists. Together, these elements demonstrate how XAI transforms opaque machine learning models into interpretable and actionable tools in natural product oncology research.
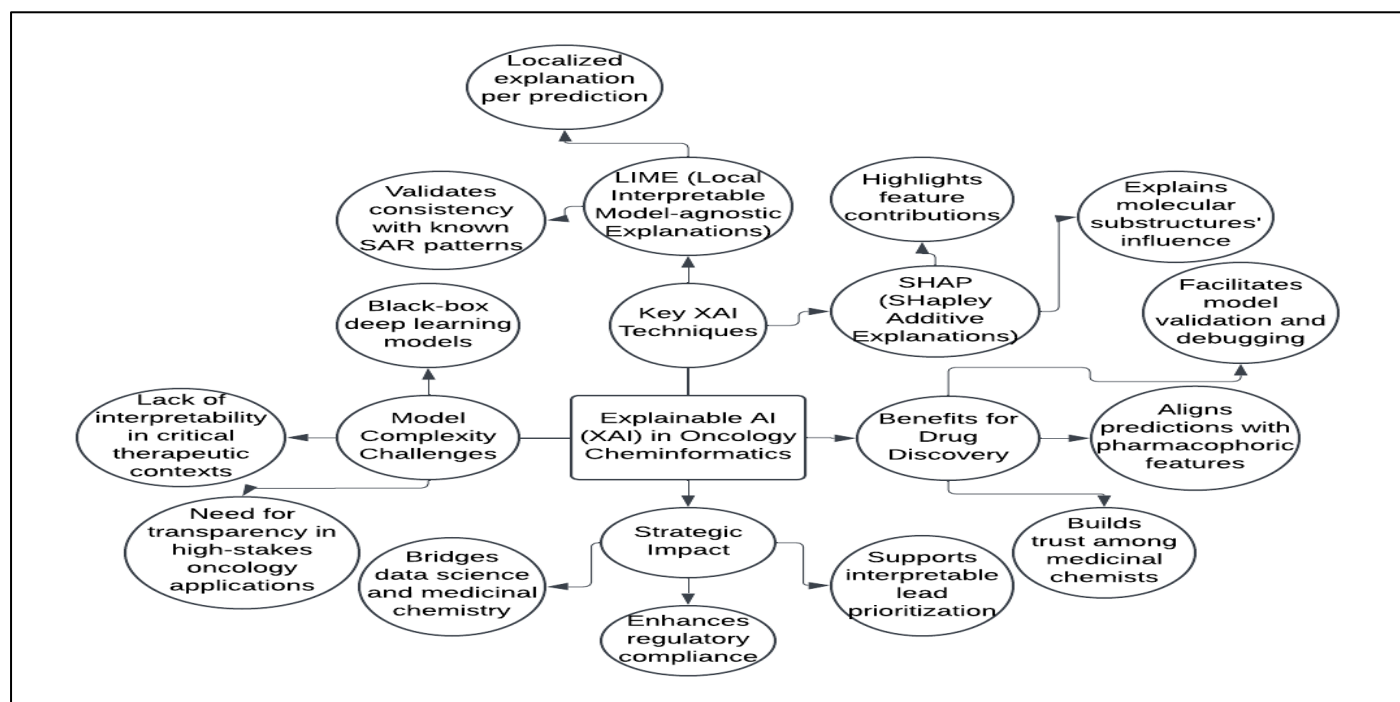


Fig 4 Diagram Showing the Role of Explainable AI in Oncology Cheminformatics

➤ *Hybrid Approaches: Mechanistic Modeling Combined with Machine Learning*

Hybrid modeling approaches that integrate mechanistic knowledge with machine learning algorithms represent a promising frontier in cheminformatics, particularly for oncology-focused natural product research. Mechanistic models encode prior biological or chemical knowledge, such as molecular binding kinetics, cellular signaling cascades, or pharmacokinetic principles, providing a structured framework that complements the data-driven flexibility of machine learning models (van der Schaar et al., 2018).

Physics-informed neural networks (PINNs), for example, embed differential equations governing chemical or biological systems directly into neural network training, ensuring that predictions adhere to known mechanistic constraints (Raissi, Perdikaris, & Karniadakis, 2019). In the context of oncology drug discovery, PINNs can model drug-target binding dynamics or simulate intracellular pathway modulations induced by natural compounds, enhancing the physiological realism of predictive outputs. By incorporating domain-specific rules into the training objective, hybrid models not only improve predictive accuracy but also enhance interpretability and generalization to unseen molecular scaffolds.

Another application of hybrid modeling involves coupling machine learning-based bioactivity prediction with mechanistic toxicity models to better anticipate adverse effects of natural products early in the discovery pipeline. Systems biology models describing tumor progression or resistance mechanisms can also be integrated with cheminformatics predictors to simulate long-term therapeutic outcomes, guiding compound prioritization beyond simple binding affinity metrics.

Hybrid approaches thus overcome some inherent limitations of purely data-driven models, such as susceptibility to overfitting or lack of biological plausibility. They enable a deeper, mechanistically coherent understanding of the interaction between natural product-derived molecules and oncological targets, ultimately facilitating more robust translation from in silico predictions to clinical candidates.

Doshi-Velez and Kim (2017) emphasized that interpretability should not be an afterthought but a core design principle in model development. For models predicting anticancer activity, integrating interpretability constraints — such as sparsity-inducing penalties or attention mechanisms focused on pharmacologically relevant atoms — enhances both model usability and regulatory compliance. Moreover, interpretable models aid in hypothesis generation, allowing researchers to propose new chemical modifications based on model-extracted structure-activity relationships (SAR).

Incorporating XAI tools into the cheminformatics pipeline transforms opaque predictive models into collaborative tools for scientific discovery, bridging computational analytics and medicinal chemistry expertise. This alignment is critical for ensuring the safe, effective, and transparent advancement of natural product-derived oncology therapeutics.

## VI. FUTURE PERSPECTIVES AND CONCLUSION

➤ *Advances in Few-Shot and Zero-Shot Learning for Rare Oncology Targets*

The challenge of predicting bioactivity for rare oncology targets with limited training data has catalyzed the adoption of few-shot and zero-shot learning techniques in cheminformatics. Few-shot learning enables models to generalize bioactivity predictions from just a handful of labeled examples, dramatically reducing the dependency on large annotated datasets. In oncology-focused natural product discovery, where unique chemical scaffolds and underexplored biological targets are common, few-shot approaches empower researchers to build predictive models even when only a small number of bioactivity measurements are available.

Zero-shot learning extends this capability further by enabling models to make predictions for entirely unseen classes or targets based solely on auxiliary information, such as target protein sequences, molecular descriptors, or ontological relationships. In natural product research, this allows predictive frameworks to infer interactions with novel cancer biomarkers without the need for explicit training examples. Embedding techniques, such as learning shared latent spaces between compounds and targets, facilitate zero-shot generalization by capturing underlying patterns between chemical structure and biological function.

Advances in meta-learning, model-agnostic meta-learning (MAML) algorithms, and transfer learning strategies have enhanced the feasibility of few-shot and zero-shot learning in real-world oncology datasets. These methods dynamically adapt model parameters to new tasks with minimal retraining, offering a pragmatic solution to the pervasive data scarcity problem in natural compound oncology pipelines. Integrating these advanced learning paradigms with traditional cheminformatics workflows ensures a more flexible and scalable approach to identifying promising anticancer leads from underrepresented regions of chemical and biological space.

➤ *Role of Cheminformatics in Precision Oncology and Personalized Therapeutics*

The integration of cheminformatics into precision oncology initiatives is transforming how natural compounds are evaluated and deployed for personalized cancer therapies. Precision oncology emphasizes tailoring therapeutic interventions based on the unique genetic, molecular, and environmental profile of each patient. Cheminformatics tools enable rapid screening and prioritization of natural compounds that align with specific

oncogenic mutations, pathway deregulations, or tumor microenvironment characteristics observed in individual patients.

By leveraging predictive bioactivity models linked to genomic and proteomic data, cheminformatics can identify compounds most likely to modulate critical disease drivers unique to a patient's cancer subtype. Moreover, structure-activity relationship (SAR) models refined through patient-specific molecular profiles allow for the customization of compound selection, optimizing efficacy while minimizing toxicity. In silico screening workflows can simulate the interaction of natural products with mutated receptors or variant enzymes, uncovering opportunities for selective targeting that conventional screening methods might overlook.

Cheminformatics also plays a vital role in optimizing combination therapies, where multiple natural compounds are selected based on their synergistic effects against heterogeneous tumor populations. Predictive modeling of drug-drug interactions, resistance mechanisms, and pharmacogenomic variations further refines therapeutic strategies, enhancing treatment durability and patient outcomes. As molecular profiling becomes standard practice in oncology care, cheminformatics-driven pipelines are poised to accelerate the development of personalized natural product-based interventions, bridging the gap between bench discovery and individualized clinical application.

➤ *Final Reflections and Recommendations*
The convergence of data-driven cheminformatics and natural product oncology research is reshaping the landscape of anticancer drug discovery. While traditional wet-lab screening remains essential, computational pipelines now enable more strategic and efficient exploration of the vast chemical diversity offered by nature. Models that accurately predict bioactivity, toxicity, and drug-likeness properties streamline the prioritization of natural compounds, reducing time and resource expenditures.

However, realizing the full potential of these approaches requires addressing critical challenges, including data sparsity, imbalance, and model interpretability. Investment in curated, high-quality datasets tailored to oncology applications is essential. Additionally, expanding the use of hybrid modeling strategies that integrate mechanistic biological knowledge with machine learning can enhance model robustness and clinical relevance. Emphasizing interpretability through explainable AI frameworks ensures that cheminformatics models serve not merely as predictive engines but as collaborative tools supporting scientific discovery and therapeutic innovation.

Future efforts should focus on expanding few-shot and zero-shot learning methodologies to enable predictions across rare and emerging cancer targets. Moreover, embedding cheminformatics workflows into precision oncology frameworks will facilitate the translation of natural product leads into personalized treatment regimens. By embracing interdisciplinary collaboration between computational scientists, medicinal chemists, and oncologists, the field can accelerate the identification of next-generation therapeutics derived from the immense, largely untapped reservoir of natural compounds.

## REFERENCES

[1]. Cheng, T., Li, Q., Zhou, Z., Wang, Y., & Bryant, S. H. (2012). Structure-based virtual screening for drug discovery: A problem-centric review. AAPS Journal, 14(1), 133–141. https://doi.org/10.1208/s12248-012-9322-0

[2]. Cherkasov, A., Muratov, E. N., Fourches, D., Varnek, A., Baskin, I. I., Cronin, M., ... & Tropsha, A. (2014). QSAR modeling: Where have you been? Where are you going to? Journal of Medicinal Chemistry, 57(12), 4977–5010. https://doi.org/10.1021/jm4004285

[3]. Cherkasov, A., Muratov, E. N., Fourches, D., Varnek, A., Baskin, I. I., Cronin, M., ... & Tropsha, A. (2014). QSAR modeling: Where have you been? Where are you going to? Journal of Medicinal Chemistry, 57(12), 4977–5010. https://doi.org/10.1021/jm4004285

[4]. Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint. https://arxiv.org/abs/1702.08608

[5]. Duvenaud, D. K., Maclaurine, D., Iparraguirre, J., Bombarell, R., Hirzel, T., Aspuru-Guzik, A., & Adams, R. P. (2015). Convolutional networks on graphs for learning molecular fingerprints. Advances in Neural Information Processing Systems, 28, 2224–2232. https://proceedings.neurips.cc/paper/2015/file/f9be311e65d81f28b116ec6fe3d9ec8b-Paper.pdf

[6]. Ekins, S., & Williams, A. J. (2010). When pharmaceutical companies publish large datasets: An abundance of riches or fool's gold? Drug Discovery Today, 15(19-20), 812–815. https://doi.org/10.1016/j.drudis.2010.08.006

[7]. Fourches, D., Muratov, E., & Tropsha, A. (2010). Trust, but verify: on the importance of chemical structure curation in cheminformatics and QSAR modeling research. Journal of chemical information and modeling, 50(7), 1189.

[8]. Fourches, D., Muratov, E., & Tropsha, A. (2016). Trust, but verify: On the importance of chemical structure curation in cheminformatics and QSAR modeling research. Journal of Chemical Information and Modeling, 56(7), 1243–1252. https://doi.org/10.1021/acs.jcim.6b00183

[9]. Gaulton, A., Hersey, A., Nowotka, M., Bento, A. P., Chambers, J., Mendez, D., ... & Leach, A. R. (2017). The ChEMBL database in 2017. Nucleic Acids Research, 45(D1), D945–D954. https://doi.org/10.1093/nar/gkw1074

[10]. Gilpin, L. H., Bau, D., Yuan, B. Z., Bajwa, A., Specter, M., & Kagal, L. (2018). Explaining

explanations: An overview of interpretability of machine learning. 2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA), 80–89. https://doi.org/10.1109/DSAA.2018.00018

[11]. 1Gramatica, P. (2007). Principles of QSAR models validation: Internal and external. QSAR & Combinatorial Science, 26(5), 694–701. https://doi.org/10.1002/qsar.200610151

[12]. He, H., & Garcia, E. A. (2009). Learning from imbalanced data. IEEE Transactions on Knowledge and Data Engineering, 21(9), 1263–1284. https://doi.org/10.1109/TKDE.2008.239

[13]. Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. Science, 313(5786), 504–507. https://doi.org/10.1126/science.1127647

[14]. Kipf, T. N., & Welling, M. (2017). Semi-supervised classification with graph convolutional networks. International Conference on Learning Representations (ICLR). https://arxiv.org/abs/1609.02907

[15]. Li, J. W.-H., & Vederas, J. C. (2009). Drug discovery and natural products: End of an era or an endless frontier? *Science, 325*(5937), 161–165. https://doi.org/10.1126/science.1168243

[16]. Li, J. W.-H., & Vederas, J. C. (2009). Drug discovery and natural products: End of an era or an endless frontier? Science, 325(5937), 161–165. https://doi.org/10.1126/science.1168243

[17]. Lionta, E., Spyrou, G., Vassilatis, D. K., & Cournia, Z. (2014). Structure-based virtual screening for drug discovery: Principles, applications, and recent advances. Current Topics in Medicinal Chemistry, 14(16), 1923–1938. https://doi.org/10.2174/1568026614666140929124445

[18]. 1Molinski, T. F., Dalisay, D. S., Lievens, S. L., & Saludes, J. P. (2009). Drug development from marine natural products. Nature Reviews Drug Discovery, 8(1), 69–85. https://doi.org/10.1038/nrd2487

[19]. Newman, D. J., & Cragg, G. M. (2016). Natural products as sources of new drugs from 1981 to 2014. Journal of Natural Products, 79(3), 629–661. https://doi.org/10.1021/acs.jnatprod.5b01055

[20]. Newman, D. J., & Cragg, G. M. (2020). Natural products as sources of new drugs over the nearly four decades from 1981 to 2019. Journal of Natural Products, 83(3), 770–803. https://doi.org/10.1021/acs.jnatprod.9b01285

[21]. Newman, D. J., & Cragg, G. M. (2020). Natural products as sources of new drugs over the nearly four decades from 1981 to 2019. Journal of Natural Products, 83(3), 770–803. https://doi.org/10.1021/acs.jnatprod.9b01285

[22]. Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. IEEE Transactions on Knowledge and Data Engineering, 22(10), 1345–1359. https://doi.org/10.1109/TKDE.2009.191

[23]. Raissi, M., Perdikaris, P., & Karniadakis, G. E. (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. Journal of Computational Physics, 378, 686–707. https://doi.org/10.1016/j.jcp.2018.10.045

[24]. 2Ruder, S. (2017). An overview of multi-task learning in deep neural networks. arXiv preprint. https://arxiv.org/abs/1706.05098

[25]. Schneider, G. (2010). Virtual screening: An endless staircase? Nature Reviews Drug Discovery, 9(4), 273–276. https://doi.org/10.1038/nrd3139

[26]. Sun, Y., Wong, A. K., & Kamel, M. S. (2009). Classification of imbalanced data: A review. International Journal of Pattern Recognition and Artificial Intelligence, 23(4), 687–719. https://doi.org/10.1142/S0218001409007326

[27]. Todeschini, R., & Consonni, V. (2009). Molecular descriptors for chemoinformatics (Vols. 1–2). Wiley-VCH Verlag GmbH & Co. KGaA. https://doi.org/10.1002/9783527628766

[28]. van der Schaar, M., Alaa, A. M., Floto, A., Gimson, A., Scholtes, S., Wood, A., & Jarrett, D. (2018). How artificial intelligence and machine learning can help healthcare systems respond to COVID-19. Machine Learning, 110(11), 1–19. https://doi.org/10.1007/s10994-020-05928-x

[29]. Walters, W. P., & Murcko, M. A. (2020). Assessing the impact of generative AI on medicinal chemistry. Nature Biotechnology, 38(2), 143–145. https://doi.org/10.1038/s41587-020-0431-6

[30]. Wang, Y., Bryant, S. H., Cheng, T., Wang, J., Gindulyte, A., Shoemaker, B. A., ... & Bolton, E. E. (2017). PubChem BioAssay: 2017 update. Nucleic Acids Research, 45(D1), D955–D963. https://doi.org/10.1093/nar/gkw1118