

Advances in Chemoinformatics: Artificial Intelligence, Big Data, and Computational Strategies in Modern Chemical Research

¹Koyel Misra, ²Alok Satsangi

^{1,2}*NSHM Knowledge Campus, Durgapur, West Bengal, India*

**Corresponding Author: koyel.misra@nshm.com*

doi.org/10.64643/JATIRV2I5-140170-001

Abstract- Chemoinformatics has emerged as a powerful interdisciplinary domain that integrates chemistry, computer science, data analytics, and information technology to manage and interpret complex chemical information. The rapid development of computational techniques, artificial intelligence (AI), machine learning (ML), and big data analytics has significantly transformed the scope and impact of chemoinformatics. These technologies enable researchers to analyze extensive chemical datasets, predict molecular properties, and accelerate the discovery of new drugs and advanced materials. This review presents an overview of the major developments in chemoinformatics, focusing on molecular modeling techniques, artificial intelligence applications, chemoinformatics databases, and virtual screening approaches used in modern drug discovery. Additionally, the role of chemoinformatics in materials science and sustainable chemistry is discussed. The paper also highlights the current challenges faced by the field, including data quality issues, model interpretability, and integration with experimental validation. Finally, emerging trends such as quantum computing, automated drug discovery platforms, and AI-driven materials design are explored. These developments suggest that chemoinformatics will continue to play a crucial role in advancing chemical research and innovation in the coming decades.

Index- Terms - Chemoinformatics, Artificial Intelligence, Machine Learning, Molecular Modeling, Virtual Screening, Drug Discovery, Big Data

I. INTRODUCTION

The exponential growth of chemical information over the past few decades has led to the development of computational methods capable of managing and analyzing vast datasets. Chemoinformatics, also referred to as cheminformatics, has evolved as a specialized field that

combines principles from chemistry, computer science, mathematics, and information technology to extract meaningful insights from chemical data.

Historically, chemical research relied heavily on experimental techniques for the discovery of new compounds and materials. However, advances in computing power and algorithm development have enabled researchers to perform sophisticated simulations and predictive modeling. As a result, chemoinformatics has become an essential component of modern chemical research, particularly in pharmaceutical development, materials science, environmental chemistry, and biotechnology.

One of the primary objectives of chemoinformatics is to understand the relationship between molecular structure and chemical or biological properties. By analyzing molecular descriptors and structural features, researchers can predict the behavior of chemical compounds before performing experimental testing. This capability has significantly reduced the time and financial cost associated with traditional experimental approaches.

Recent technological developments have further enhanced the capabilities of chemoinformatics. The integration of artificial intelligence, machine learning algorithms, and big data analytics has made it possible to analyze extremely large chemical datasets and identify patterns that would otherwise remain undetected. These tools are particularly useful in drug discovery, where millions of potential molecules must be screened to identify promising therapeutic candidates.

The rapid expansion of open-access chemical databases has also contributed significantly to the advancement of chemoinformatics. Platforms such as PubChem, ChEMBL, and the Protein Data Bank provide researchers with access to millions of chemical structures and biological activity records, facilitating data-driven research.

This paper reviews the recent developments and emerging trends in chemoinformatics and highlights its growing role in accelerating scientific discovery and technological innovation.

II. MOLECULAR MODELING AND COMPUTATIONAL CHEMISTRY

Molecular modeling and computational chemistry play a crucial role in modern chemoinformatics by enabling scientists to simulate, visualize, and analyze molecular structures and interactions using computational methods. These techniques employ mathematical models and algorithms to predict the physical, chemical, and biological properties of molecules without requiring extensive laboratory experimentation. Methods such as quantum mechanical calculations, including Density Functional Theory (DFT), allow researchers to study electronic structures, reaction mechanisms, and molecular stability at an atomic level. Additionally, molecular dynamics simulations provide insights into the time-dependent behavior of molecular systems, helping researchers understand processes such as protein folding, ligand–receptor interactions, and conformational changes in

biomolecules. By integrating these computational approaches with experimental data, molecular modeling significantly accelerates drug discovery, materials design, and the understanding of complex chemical systems while reducing research time and cost.

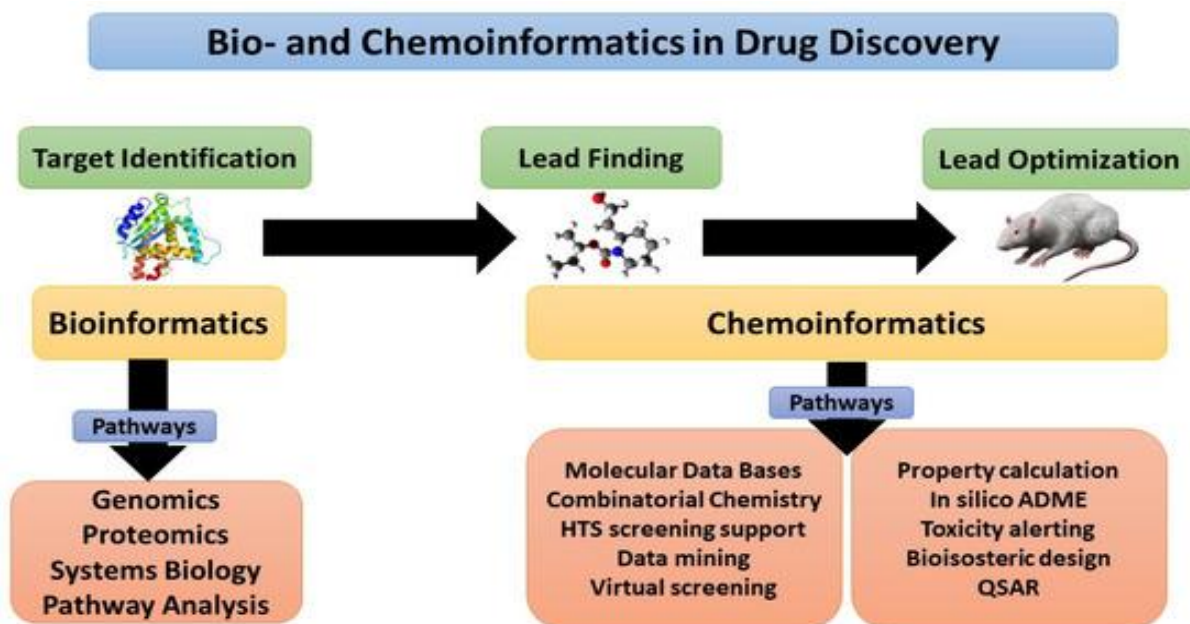


Figure 1. General workflow of chemoinformatics-driven drug discovery, illustrating the integration of chemical databases, molecular descriptor calculation, machine learning models, virtual screening, and experimental validation.

2.1 Overview of Molecular Modelling

Molecular modeling refers to a collection of computational techniques used to represent, visualize, and analyze molecular structures and interactions. These methods provide valuable insights into molecular behavior, including chemical reactivity, thermodynamic stability, and intermolecular interactions. By creating mathematical models of molecules, researchers can study how atoms are arranged in three-dimensional space and how they interact with each other under different conditions.

The development of advanced computational algorithms and high-performance computing systems has significantly enhanced the accuracy and efficiency of molecular modeling techniques. These approaches allow researchers to simulate complex chemical systems that would otherwise be difficult or time-consuming to study experimentally. For example, molecular modeling can be used to predict how a drug molecule interacts with a biological target such as a protein or enzyme, helping scientists identify promising drug candidates before laboratory testing.

Molecular modeling techniques include several computational methods, such as quantum mechanical calculations, molecular mechanics, and molecular dynamics simulations. Quantum mechanical methods focus on the electronic structure of molecules and help explain chemical

bonding, reaction mechanisms, and energy changes during chemical reactions. Molecular mechanics, on the other hand, uses classical physics to estimate the energy and geometry of large molecular systems, making it particularly useful for studying biomolecules like proteins and nucleic acids.

Another important approach is molecular dynamics simulation, which allows researchers to observe how molecules move and interact over time. This technique provides valuable insights into biological processes such as protein folding, ligand binding, and conformational changes in biomolecules. Additionally, molecular docking methods are widely used to predict how small molecules bind to specific receptors, which is a crucial step in rational drug design.

Overall, molecular modeling has become an essential tool in modern chemistry, biochemistry, and pharmaceutical research. By combining computational predictions with experimental validation, scientists can better understand molecular behavior, accelerate the discovery of new drugs and materials, and reduce the time and cost associated with traditional laboratory-based research.

2.2 Quantum Mechanical Methods

Quantum mechanical methods are very important in molecular modelling because they help scientists understand how electrons behave inside molecules and how this behavior determines chemical properties and reactions. These methods are based on the principles of quantum mechanics and use mathematical equations such as the Schrödinger Equation to describe the motion and energy of electrons. Since solving this equation exactly is difficult for large molecular systems, scientists use approximate computational techniques such as Density Functional Theory and Ab Initio Quantum Chemistry. Density Functional Theory focuses on calculating electron density rather than complex wavefunctions, which makes it faster and suitable for studying larger molecules, surfaces, and materials. Ab initio methods, on the other hand, are based purely on fundamental physical laws and do not rely on experimental parameters, allowing very accurate predictions of molecular properties, although they require more computational power. Using these quantum mechanical approaches, researchers can analyze molecular orbitals, determine the distribution of electrons, calculate bond energies, and understand how chemical reactions occur step by step. This information is extremely useful in fields such as catalysis research, where scientists study how catalysts speed up reactions, as well as in reaction mechanism analysis, drug discovery, and the design of new materials with specific electronic or structural properties. Even though quantum mechanical calculations can be computationally expensive, improvements in computer technology, algorithms, and high-performance computing systems have made it possible to perform these calculations more efficiently, allowing researchers to investigate increasingly complex molecular systems with high accuracy.

2.3 Molecular Dynamics Simulations

Molecular dynamics (MD) simulations are widely used in molecular modelling to study how atoms and molecules move and interact over time. In this method, molecules are treated as collections of

particles whose motions follow classical mechanics, specifically the principles of the Newton's Laws of Motion. By calculating the forces acting on each atom and updating their positions step by step, MD simulations can model the dynamic behaviour of molecular systems over time. This allows researchers to observe important biological and chemical processes such as protein folding, ligand binding, diffusion, and conformational changes in biomolecules. For example, proteins are not rigid structures; they constantly change shape, and MD simulations help scientists understand how these structural fluctuations affect their biological function. In drug discovery, MD simulations are particularly valuable because they provide detailed insights into how potential drug molecules interact dynamically with their biological targets, such as enzymes or receptors. This helps researchers predict binding stability, identify important interaction sites, and optimize drug candidates. To accurately represent molecular interactions, modern MD simulations use advanced force fields such as AMBER Force Field, CHARMM Force Field, and GROMOS Force Field. These force fields are mathematical models that describe how atoms interact through bonded interactions (such as bonds, angles, and torsions) and non-bonded interactions (such as electrostatic and van der Waals forces). By combining accurate force fields with powerful computer simulations, molecular dynamics provides a detailed picture of molecular motion and interactions, making it an essential tool in modern computational chemistry, structural biology, and pharmaceutical research.

2.4 Enhanced Sampling Techniques

Traditional molecular simulations often struggle to observe rare molecular events because they explore only a limited portion of the molecule's possible structures, known as the conformational space. Many important biological and chemical processes, such as protein folding, ligand unbinding, or structural transitions, occur over long time scales and involve crossing high energy barriers on the molecular energy landscape. Standard simulations like Molecular Dynamics Simulation may become trapped in local energy minima and therefore fail to sample other relevant conformations within practical simulation times. To overcome this limitation, scientists use enhanced sampling techniques such as Metadynamics, Replica Exchange Molecular Dynamics, and Accelerated Molecular Dynamics. These methods improve the exploration of the molecular energy landscape by helping the system overcome energy barriers more efficiently. For example, metadynamics adds a history-dependent bias to the simulation to push the system out of already sampled states, allowing it to explore new configurations. Replica exchange molecular dynamics runs multiple simulations at different temperatures and exchanges configurations between them to enhance sampling efficiency. Accelerated molecular dynamics modifies the potential energy surface to reduce energy barriers and speed up transitions between molecular states. By using these advanced techniques, researchers can better explore complex energy landscapes, identify stable molecular conformations, and gain deeper insights into reaction mechanisms, protein folding pathways, and other important chemical and biological processes that are difficult to observe using conventional simulation methods.

III. ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING IN CHEMOINFORMATICS

Artificial Intelligence (AI) and Machine Learning (ML) are transforming the field of chemoinformatics by enabling the rapid analysis and prediction of chemical properties, biological activities, and reaction outcomes. These technologies use large datasets of molecular structures, experimental results, and chemical descriptors to train algorithms that can recognize patterns and make predictions without explicit programming. In drug discovery, AI and ML models can predict the activity of new compounds, optimize lead molecules, and identify potential drug candidates more efficiently than traditional methods. They are also applied in materials science to design molecules with desired physical, chemical, or electronic properties. By combining chemoinformatics databases with deep learning techniques, researchers can accelerate virtual screening, QSAR (Quantitative Structure–Activity Relationship) modeling, and reaction outcome prediction, reducing the time and cost of experimental trials. Overall, AI and ML provide powerful tools to explore chemical space, uncover hidden relationships in complex datasets, and guide rational molecular design in ways that were previously impractical.

3.1 Machine Learning Approaches

Machine learning (ML) approaches have become indispensable in chemoinformatics, providing powerful tools for analyzing large chemical datasets, uncovering hidden patterns, and predicting molecular properties with high accuracy. Among the most commonly used ML algorithms are Support Vector Machines (SVM), Random Forests, k-Nearest Neighbor (kNN), and Artificial Neural Networks (ANNs). SVMs are widely applied for classification and regression tasks because they can efficiently handle high-dimensional chemical descriptor spaces and find optimal boundaries between different classes of molecules, such as active versus inactive compounds. Random Forests, an ensemble learning method based on decision trees, are particularly useful for feature selection and predicting molecular properties due to their robustness against overfitting and ability to handle complex, nonlinear relationships. k-Nearest Neighbor (kNN) is a simple yet effective algorithm that predicts the properties or activity of a molecule by comparing it to its closest neighbors in descriptor space, making it intuitive for similarity-based predictions. Artificial Neural Networks, including deep learning architectures, excel at capturing highly nonlinear relationships in chemical data and can model intricate dependencies between molecular structures and their biological or physicochemical properties. These ML techniques have been successfully applied in tasks such as virtual screening, QSAR modeling, property prediction, and reaction outcome forecasting, allowing chemoinformatics researchers to efficiently explore chemical space, accelerate drug discovery, and optimize molecular design with significantly reduced experimental costs. By integrating these algorithms with large chemical databases and descriptor libraries, machine learning enables more informed decision-making and predictive modeling in both pharmaceutical and materials research.

Table 1 – Common Machine Learning Algorithms Used in Chemoinformatics

Algorithm	Application in Chemoinformatics	Advantages
Support Vector Machine (SVM)	QSAR modeling, toxicity prediction	High accuracy for small datasets
Random Forest	Molecular property prediction	Robust and resistant to overfitting
k-Nearest Neighbor (kNN)	Chemical similarity analysis	Simple and interpretable
Artificial Neural Networks	Drug activity prediction	Handles nonlinear relationships
Graph Neural Networks	Protein–ligand interaction prediction	Captures molecular graph structure
Deep Learning Models	Drug design and property prediction	High predictive performance

These algorithms can analyse relationships between molecular descriptors and chemical properties, enabling accurate predictions of solubility, toxicity, and biological activity.

3.2 Deep Learning Models

Deep learning techniques have greatly enhanced the predictive power of chemoinformatics by allowing models to automatically learn complex patterns from large chemical datasets, without relying solely on manually crafted descriptors. Unlike traditional machine learning methods, deep learning models can process more sophisticated molecular representations, such as graph-based structures that capture the connectivity of atoms and bonds, or molecular fingerprints that encode chemical substructures. Convolutional Neural Networks (CNNs) are particularly effective at identifying local structural motifs within molecules, similar to how they detect patterns in images, making them useful for analyzing molecular graphs or grid-based chemical representations. Recurrent Neural Networks (RNNs), on the other hand, excel at handling sequential data, such as SMILES strings that represent molecules as sequences of characters, allowing the model to capture sequential dependencies and patterns in chemical structures. These capabilities enable deep learning models to predict complex chemical properties with high accuracy, including drug-target interactions, solubility, binding affinity, and potential toxicity. They are also widely used in de novo molecular design, where models generate novel chemical structures with desired properties, accelerating the discovery of new drug candidates or materials. By combining the ability to learn hierarchical and nonlinear features with access to large chemical datasets, deep learning approaches have transformed chemoinformatics, making it possible to explore chemical space more efficiently, improve virtual screening workflows, and guide rational molecular design in ways that traditional computational methods could not achieve.

3.3 Graph Neural Networks

Graph Neural Networks (GNNs) are emerging as one of the most powerful machine learning approaches in chemoinformatics because they naturally represent molecules as graphs, where atoms are treated as nodes and chemical bonds as edges. This representation allows GNNs to capture the full connectivity and topology of molecules, including complex patterns of atomic interactions that traditional descriptor-based methods might overlook. Unlike conventional neural networks, which often rely on precomputed molecular fingerprints or descriptors, GNNs learn features directly from the molecular graph, enabling them to understand local chemical environments, bond types, and neighborhood relationships in a data-driven way. This makes them particularly effective for predicting molecular properties that depend on subtle structural nuances, such as protein–ligand binding affinities, reaction outcomes, and molecular reactivity. In drug discovery, GNNs can model interactions between small molecules and biological targets with high accuracy, providing insights into binding sites, activity, and selectivity. Beyond pharmacology, GNNs are also applied in materials science, for example, to predict the electronic, optical, or mechanical properties of complex molecular assemblies. By integrating graph representations with deep learning, GNNs combine the strengths of structural awareness and nonlinear pattern recognition, enabling researchers to explore chemical space more effectively, identify promising compounds faster, and design molecules with tailored properties that would be difficult to predict using traditional methods.

IV. BIG DATA AND CHEMOINFORMATICS DATABASES

The rapid growth of chemical and biological data has driven the development of extensive chemoinformatics databases, which serve as critical resources for research in drug discovery, materials science, environmental chemistry, catalysis, and computational toxicology. These databases store vast amounts of information about chemical structures, molecular properties, bioactivities, and protein-ligand interactions, enabling researchers to explore chemical space systematically and make data-driven predictions. Major databases widely used in the field include PubChem, which provides comprehensive information on small molecules and their biological activities; ChEMBL, which focuses on bioactive drug-like molecules with experimental binding and functional data; the Protein Data Bank (PDB), which houses detailed three-dimensional structures of proteins, nucleic acids, and complexes essential for structure-based drug design; the ZINC Database, a repository of commercially available compounds optimized for virtual screening; and DrugBank, which integrates detailed drug data with target information and pharmacological properties. By combining these rich data resources with computational tools and machine learning models, chemoinformatics enables the efficient identification of lead compounds, prediction of molecular properties, and rational design of new molecules. As illustrated in Figure 2, the applications of chemoinformatics are broad, encompassing not only drug discovery and materials science but also environmental chemistry for assessing pollutants, catalysis design for optimizing chemical reactions, and computational toxicology to predict

potential adverse effects of chemicals. These databases and associated tools provide researchers with the infrastructure to accelerate scientific discovery, reduce experimental costs, and explore complex chemical and biological systems in a systematic and predictive manner.

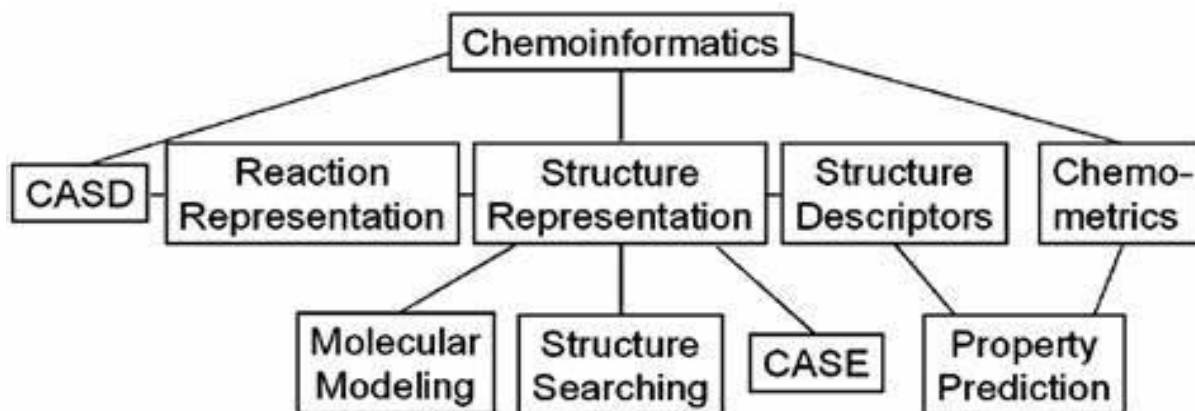


Figure 2. Major application areas of chemoinformatics including drug discovery, materials science, environmental chemistry, catalysis design, and computational toxicology

These platforms contain millions of chemical structures, biological activity data, and experimental measurements.

Big data analytics techniques allow researchers to analyze these datasets to identify patterns, correlations, and potential drug candidates. Data mining methods combined with machine learning algorithms have significantly improved the efficiency of chemical discovery processes.

V. VIRTUAL SCREENING AND DRUG DISCOVERY

Virtual screening is a key computational strategy in drug discovery that allows researchers to efficiently evaluate large chemical libraries and identify molecules likely to interact with specific biological targets. This approach significantly accelerates the early stages of drug development by prioritizing compounds with the highest potential activity before experimental testing. Virtual screening can be divided into two major categories: ligand-based screening, which identifies new molecules by comparing them to known active compounds and relies on molecular similarity, and structure-based screening, which uses the three-dimensional structure of the target protein to predict binding interactions and affinities. Modern virtual screening workflows often combine traditional molecular docking techniques with machine learning algorithms to improve prediction accuracy, enabling more reliable identification of promising drug candidates. Furthermore, AI-driven screening methods have transformed the field by allowing the rapid analysis of millions of compounds in a short time, reducing the time and cost required for early-stage drug discovery and enabling researchers to explore vast chemical space efficiently. These advancements are supported by major chemoinformatics databases, as summarized in Table 2, which provide access to extensive information on chemical structures, molecular properties, bioactivities, and target proteins, forming the foundation for both ligand-based and structure-based virtual screening

approaches. By integrating computational power, AI, and rich chemical data, virtual screening has become an essential tool for accelerating the discovery and optimization of new therapeutics.

Table 2 – Major Chemoinformatics Databases

Database	Type of Data	Key Applications
PubChem	Chemical structures, bioactivity	Drug discovery
ChEMBL	Bioactive molecules with drug-like properties	Pharmacological research
Protein Data Bank (PDB)	Protein and biomolecular structures	Molecular docking
ZINC Database	Commercial chemical compounds	Virtual screening
DrugBank	Drug and target information	Pharmaceutical research

VI. APPLICATIONS IN MATERIALS SCIENCE

Chemoinformatics has also contributed significantly to the discovery and design of advanced materials. Materials informatics uses machine learning and computational modeling to predict the properties of materials based on their atomic structures.

Applications include:

- Development of new battery materials
- Discovery of semiconductor compounds
- Design of catalytic materials
- Development of renewable energy technologies

Machine learning models can predict important material properties such as conductivity, stability, and mechanical strength, enabling researchers to focus on the most promising candidates for experimental testing.

VII. CHALLENGES IN CHEMOINFORMATICS

Despite its rapid development, several challenges remain in chemoinformatics.

Data Quality

Many chemical datasets contain missing values, inconsistent measurements, or incorrect molecular structures. Poor data quality can significantly affect the performance of predictive models.

Model Interpretability

Deep learning models often function as “black boxes,” making it difficult to interpret their predictions. This lack of transparency can limit their adoption in regulated fields such as pharmaceutical research.

Integration with Experimental Research

Although computational models provide valuable predictions, experimental validation remains essential. Effective collaboration between computational scientists and experimental researchers is necessary to ensure reliable outcomes.

Additional Section: Chemoinformatics Tools and Software

A wide range of computational tools and software platforms have been developed to support chemoinformatics research. These tools enable the visualization, simulation, and analysis of molecular structures and chemical data.

Commonly used software includes:

RDKit – an open-source cheminformatics toolkit widely used for molecular descriptor calculation and machine learning workflows.

Open Babel – a chemical toolbox designed for converting chemical file formats and analyzing molecular data.

GROMACS – a high-performance molecular dynamics simulation package used to study biomolecular systems.

AutoDock – a widely used molecular docking software for predicting ligand–protein interactions.

Schrödinger Suite – a commercial computational chemistry platform used in pharmaceutical research for molecular modeling and drug design.

These tools play a crucial role in enabling researchers to perform complex computational experiments and analyze large chemical datasets efficiently.

Additional Section: Role of Chemoinformatics in Sustainable Chemistry

Chemoinformatics is increasingly being applied to support sustainable and green chemistry initiatives. By predicting the environmental impact and toxicity of chemical compounds before their synthesis, researchers can design safer and more environmentally friendly chemicals.

Predictive models can be used to evaluate:

- biodegradability
- environmental persistence
- aquatic toxicity
- bioaccumulation potential

These approaches reduce the need for extensive laboratory testing and support regulatory compliance in chemical manufacturing.

Chemoinformatics also contributes to the development of sustainable catalysts and renewable energy materials, such as advanced battery systems and solar cell components.

VIII. FUTURE PERSPECTIVES

The future of chemoinformatics is expected to be shaped by several emerging technologies.

Quantum Computing

Quantum computing has the potential to revolutionize molecular simulations by solving complex quantum chemical problems more efficiently than classical computers.

Autonomous Drug Discovery

AI-driven platforms capable of automatically generating and evaluating chemical compounds are being developed to accelerate pharmaceutical research.

Advanced Data Infrastructure

Improved chemical data sharing systems and open scientific platforms will enable more collaborative research and faster scientific progress.

IX. CONCLUSION

Chemoinformatics has become an indispensable tool in modern chemical research. The integration of computational chemistry, artificial intelligence, and big data analytics has significantly improved our ability to analyze molecular systems and design new chemical compounds.

Advances in molecular modeling, machine learning algorithms, and virtual screening techniques have transformed drug discovery and materials science. Although challenges related to data quality and model transparency remain, ongoing technological developments are expected to further enhance the capabilities of chemoinformatics.

As computing technologies continue to evolve, chemoinformatics will play an increasingly important role in advancing chemical innovation and addressing global scientific challenges.

REFERENCES

- [1] Chen, H., Engkvist, O., Wang, Y., Olivecrona, M., & Blaschke, T. (2021). The rise of deep learning in drug discovery. *Drug Discovery Today*, 26(7), 1612–1622. <https://doi.org/10.1016/j.drudis.2021.03.010>
- [2] Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J., Sánchez-Lengeling, B., Sheberla, D., Aspuru-Guzik, A. (2018). Automatic chemical design using a data-driven continuous representation of molecules. *Nature Communications*, 9, 4421.

- [3] Jensen, J. H. (2019). A graph-based genetic algorithm and generative model for the exploration of chemical space. *Chemical Science*, *10*, 3567–3572.
- [4] Korb, O., Stützle, T., & Exner, T. E. (2022). Virtual screening: A comparison of methods and applications. *Journal of Molecular Modeling*, *28*(1), 23.
- [5] Lipton, Z. C. (2018). The mythos of model interpretability. *Communications of the ACM*, *61*(12), 36–43.
- [6] McArdle, S., Endo, S., Aspuru-Guzik, A., Benjamin, S., & Yuan, X. (2020). Quantum computational chemistry. *Reviews of Modern Physics*, *92*(1), 015003.
- [7] Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J., & Koes, D. (2017). Protein-ligand scoring with convolutional neural networks. *Journal of Chemical Information and Modeling*, *57*(4), 942–957.
- [8] Rothlisberger, U., Kästner, J., & Jäkel, P. (2020). Enhanced sampling techniques in molecular dynamics simulations. *Nature Reviews Chemistry*, *4*, 527–540.
- [9] Schaefer, M., Müller, K., & O'Rourke, M. (2019). A review of chemoinformatics databases and their applications. *Journal of Chemical Information and Modeling*, *59*(5), 2301–2311.
- [10] Stokes, J. M., Yang, K., Swanson, K., Jin, W., Cubillos-Ruiz, A., Donghia, N. M., ... Collins, J. J. (2020). A deep learning approach to antibiotic discovery. *Cell*, *180*(4), 688–702.
- [11] Van Der Spoel, D., Lindahl, E., Hess, B., Groenhof, G., Mark, A. E., & Berendsen, H. J. (2021). GROMACS: Fast, flexible, and free molecular simulation software. *Journal of Computational Chemistry*, *42*(1), 1–9.
- [12] Ward, L., Liu, R., Krishna, A., Hegde, V., Agrawal, A., Choudhary, A., & Wolverton, C. (2021). Including crystal structure attributes in machine learning models of formation energies via Voronoi tessellations. *Physical Review B*, *96*(2), 024104.
- [13] Zhu, H., Biggin, P. C., & Huggins, D. J. (2020). Artificial intelligence in drug discovery: Recent advances and future directions. *Nature Reviews Drug Discovery*, *19*(10), 573–590.
- [14] Aspuru-Guzik, A., Persson, K., & Simmons, B. (2018). Materials acceleration platform: Accelerating advanced energy materials discovery. *Matter*, *1*(1), 1–16.
- [15] Butler, K., Davies, D., Cartwright, H., Isayev, O., & Walsh, A. (2018). Machine learning for molecular and materials science. *Nature*, *559*, 547–555.
- [16] Chen, H., Engkvist, O., Wang, Y., Olivecrona, M., & Blaschke, T. (2021). The rise of deep learning in drug discovery. *Drug Discovery Today*, *26*(7), 1612–1622.
- [17] Duvenaud, D., Maclaurin, D., Aguilera-Iparraguirre, J., Gómez-Bombarelli, R., Hirzel, T., Aspuru-Guzik, A., & Adams, R. (2015). Convolutional networks on graphs for learning molecular fingerprints. *Advances in Neural Information Processing Systems*.
- [18] Gómez-Bombarelli, R., Wei, J., Duvenaud, D., Hernández-Lobato, J., Sánchez-Lengeling, B., & Aspuru-Guzik, A. (2018). Automatic chemical design using data-driven representations. *Nature Communications*, *9*, 4421.
- [19] Isayev, O., Fourches, D., Muratov, E., & Tropsha, A. (2015). Chemoinformatics tools for chemical research. *Journal of Chemical Information and Modeling*, *55*, 1902–1921.

- [20] Jensen, J. H. (2019). Graph-based generative models for molecular design. *Chemical Science*, *10*, 3567–3572.
- [21] Korb, O., Stütze, T., & Exner, T. (2022). Virtual screening methods in computational chemistry. *Journal of Molecular Modeling*, *28*, 23.
- [22] Lavecchia, A. (2019). Deep learning in drug discovery. *Drug Discovery Today*, *24*, 2017–2032.
- [23] Lipton, Z. (2018). The mythos of model interpretability. *Communications of the ACM*, *61*(12), 36–43.
- [24] McArdle, S., Endo, S., Aspuru-Guzik, A., Benjamin, S., & Yuan, X. (2020). Quantum computational chemistry. *Reviews of Modern Physics*, *92*, 015003.
- [25] Murphy, K. (2022). *Machine learning: A probabilistic perspective*. MIT Press.
- [26] Polishchuk, P. (2020). Chemoinformatics approaches for drug discovery. *Molecular Informatics*, *39*, 2000055.
- [27] Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J., & Koes, D. (2017). Protein–ligand scoring with convolutional neural networks. *Journal of Chemical Information and Modeling*, *57*, 942–957.
- [28] Riniker, S., & Landrum, G. (2015). Similarity maps for molecular interpretation. *Journal of Cheminformatics*, *7*, 48.
- [29] Rothlisberger, U., Kästner, J., & Jäkel, P. (2020). Enhanced sampling techniques in molecular simulations. *Nature Reviews Chemistry*, *4*, 527–540.
- [30] Schneider, G. (2020). Automating drug discovery. *Nature Reviews Drug Discovery*, *19*, 353–364.
- [31] Sliwoski, G., Kothiwale, S., Meiler, J., & Lowe, E. (2014). Computational methods in drug discovery. *Pharmacological Reviews*, *66*, 334–395.
- [32] Stokes, J. M., Yang, K., Swanson, K., Jin, W., Cubillos-Ruiz, A., & Collins, J. (2020). A deep learning approach to antibiotic discovery. *Cell*, *180*(4), 688–702.
- [33] Sydow, D., Burggraaff, L., Szengel, A., & van Vlijmen, H. (2019). Advances in structure-based virtual screening. *Drug Discovery Today*, *24*, 1515–1521.
- [34] Van Der Spoel, D., Lindahl, E., Hess, B., Groenhof, G., Mark, A., & Berendsen, H. (2021). GROMACS: High-performance molecular simulations. *Journal of Computational Chemistry*, *42*, 1–9.
- [35] Walters, W., & Murcko, M. (2020). Assessing the impact of generative AI on medicinal chemistry. *Nature Biotechnology*, *38*, 143–145.
- [36] Ward, L., Liu, R., Krishna, A., Hegde, V., Agrawal, A., & Wolverton, C. (2021). Machine learning in materials science. *Nature Materials*, *20*, 342–351.
- [37] Zhu, H., Biggin, P., & Huggins, D. (2020). Artificial intelligence in drug discovery. *Nature Reviews Drug Discovery*, *19*, 573–590.