

Environmental Science Advances

Accepted Manuscript

View Article Online
View Journal

This article can be cited before page numbers have been issued, to do this please use: K. R. Siddique, D. Bose, R. Bhattacharya, R. Villamarin Rodriguez and A. Ray, *Environ. Sci.: Adv.*, 2025, DOI: 10.1039/D5VA00240K.



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

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

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

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

The integration of artificial intelligence (AI) and bioinformatics in microbial bioremediation, as presented in this work, offers a transformative approach to addressing environmental pollution caused by heavy metals and untreated wastewater. By leveraging advanced computational models such as Random Forest, Artificial Neural Networks, and Support Vector Machines, alongside bioinformatics tools like AlphaFold2, QIIME, and MG-RAST, this research enhances the efficiency, precision, and scalability of bioremediation processes. These technologies enable real-time monitoring, accurate prediction of microbial behavior, and optimization of pollutant degradation pathways, overcoming limitations of traditional methods like bioaugmentation and bio-stimulation. The application of AI-driven biosensors, metagenomics, and gene-editing techniques, such as CRISPR, facilitates sustainable and cost-effective solutions for ecological restoration. This interdisciplinary approach not only mitigates the adverse impacts of contaminants on ecosystems and human health but also paves the way for innovative, data-driven strategies to achieve long-term environmental sustainability and resilience.



Artificial Intelligence Driven Bioinformatics for Sustainable Bioremediation: Integrating Computational Intelligence with Ecological Restoration

Kashif R. Siddique¹, Debajyoti Bose^{2*}, Riya Bhattacharya², Raul Villamarin Rodriguez², Aritra Ray³

¹Department of Analytics, School of Business, Woxsen University, Hyderabad, Telangana, India

²Centre of Excellence in Health Technology, Ecosystems & Biodiversity, Woxsen University, Hyderabad, Telangana, India

³College of Engineering, School of Mechanical Engineering, Purdue University, West Lafayette, Indiana

*Corresponding author email: debajyoti1024@gmail.com

Abstract: Environmental pollution from heavy metals and untreated wastewater poses significant risks to ecosystems and human health, highlighting the urgent need for innovative remediation strategies. Bioremediation employs microorganisms to break down contaminants, presents a sustainable and economical solution. However, conventional techniques such as bioaugmentation and bio-stimulation face challenges due to inefficiencies and the absence of real-time monitoring. This narrative review consolidates the latest developments in AI-driven bioinformatics aimed at enhancing microbial bioremediation, with an emphasis on the degradation of heavy metals and wastewater pollutants. Advanced computational models such as Random Forest, Artificial Neural Networks, and Support Vector Machines demonstrate high predictive accuracy ($R^2 > 0.99$) in analysing microbial behaviour, pollutant dynamics, and optimizing processes. Bioinformatics tools such as AlphaFold2, I-TASSER, and metagenomic platforms such as QIIME and MG-RAST facilitate accurate identification of microbial communities, genes, and degradation pathways. AI-powered biosensors and advanced deep learning enable the continuous observation of enzymatic activity and the effectiveness of treatments. The combination of AI, metagenomics, and gene editing techniques, such as CRISPR, presents scalable approaches for achieving sustainable bioremediation. Present work emphasizes innovative tools, practical applications such as ANN-RF hybrid models, and prospective pathways, highlighting the significant impact of computational intelligence on ecological restoration.

Keywords: Bioremediation; AI; Bioinformatics; Metagenomics; Machine Learning.



1. Introduction

Bioremediation is a process of breaking down contaminants using microorganisms (1). Pollution in any form is a major global issue which adversely affects the environment, human, animals and plant health (2). Innovative pollution-reduction strategies are required to solve environmental problems and their harmful impact on human health. Microbial degradation of pollutants and toxins is a promising strategy for environmental remediation (3,4). Optimization of bioremediation processes is a major task for environmental research. Artificial intelligence or AI is one of the optimum tools to change environmental conditions, AI algorithms have gained popularity in environmental research due to its capability to handle big and complex data, feature extraction, discovering patterns and ability to generate timely ideas to tackle environmental issues (5,6). Although the potential of AI can only be evaluated by interdisciplinary research elements. This has been developed as a cost effective and viable method for repairing polluted ecosystems. The contaminated effects of heavy metals can lead to variety of health challenges like neurological disorders, reproductive problems, kidney damage and cancer (7).

According to World Health Organizations in 2021 and United Nations Environment Program, around 80% of the wastewater from worldwide is discharged into the untreated environment from economically developing countries without any proper treatment (8–10). AI-driven optimization of bioremediation using algorithms can be used to evaluate ecological data, estimate pollutant behavior and optimize the process and increase efficiency. Traditional approaches such as bioaugmentation and bio-stimulation have been used for years, still they have certain limits. Bioaugmentation is a process of infusing microorganisms to break down the contaminant and enhance the metabolic activity of the existing bacterial community (11,12). Similarly, bio-stimulation is a process of adding nutrients, oxygen and other chemicals that can improve the pollution degrading capability (13,14). These traditional bioremediation processes are time-consuming and lack real-time monitoring capabilities, consideration of site heterogeneity, and prediction of microbial behavior (15). To overcome the limitations of the traditional bioremediation approaches bioinformatics and computational models are being used to predict the microorganism's behavior in response to the environmental conditions and pollutions.

Machine learning or ML models like random forests can predict microbial survival based on site parameters also accelerates the identification of best parameters. Bioinformatics approaches, such as metagenomics, have shown potential to replace traditional culture methods by accelerating microbial identification for degradation processes and enabling visualization of bacterial communities. Recently computational and bioinformatics techniques have shown potential in various areas of basic and applied sciences. This narrative review synthesizes recent advances in AI- driven bioinformatics for optimizing microbial bioremediation of heavy metals, highlighting key tools, challenges and future directions.



2. AI in Microbial Selection and Optimization

AI leverages digital algorithms to perform complex tasks, advanced machine learning algorithms like Artificial Neural Network (ANN), Support Vector Machine (SVM), Random Forest, are some promising algorithms (16,17). Techniques such as clustering and classification help in identifying and characterizing diverse microbial communities for efficient degradation process, ANN and ensemble methods can predict microbial community dynamics and explore the relation between environmental conditions and microbial community for effective bioremediation strategies as represented in **Figure 1**. In a recent study by one group observed Random Forest has played an important role in prediction of bacterial microbiota changes in various contaminated environments. Authors have utilized web of science to compile data and collected 6 variables to analyze, three algorithms were used to predict the microplastics effects on antibiotic degrading bacteria, random forest model has performed the best with AUC of 85% and 88% (18). AI-driven sensors in association of deep learning algorithms are used for real time monitoring and provide insights on catalytic activities of enzymes (19). Another investigation has shown that biosensors and learning tools can aid to the prediction of long-term treatment outcomes. Researchers have investigated the interaction of four commercial dyes by the *Bacillus megaterium* H2 azoreductases using in silico analysis (ligand binding site modeling, molecular docking, and molecular dynamics simulation), the obtained binding results suggested the stabilization between the complexes (20). A group of researchers also presented a hypothesis on how nano-biosensors can be used to detect pesticides in the field which can be modified and exploited in the bioremediation process(21). Random Forest models offer high interpretability by revealing feature importance, aiding in understanding key microbial and environmental factors driving bioremediation outcomes. In contrast, Artificial Neural Networks and their hybrids, while highly accurate, are less interpretable due to their complex architectures, necessitating advanced techniques like SHAP or LIME for practical deployment.



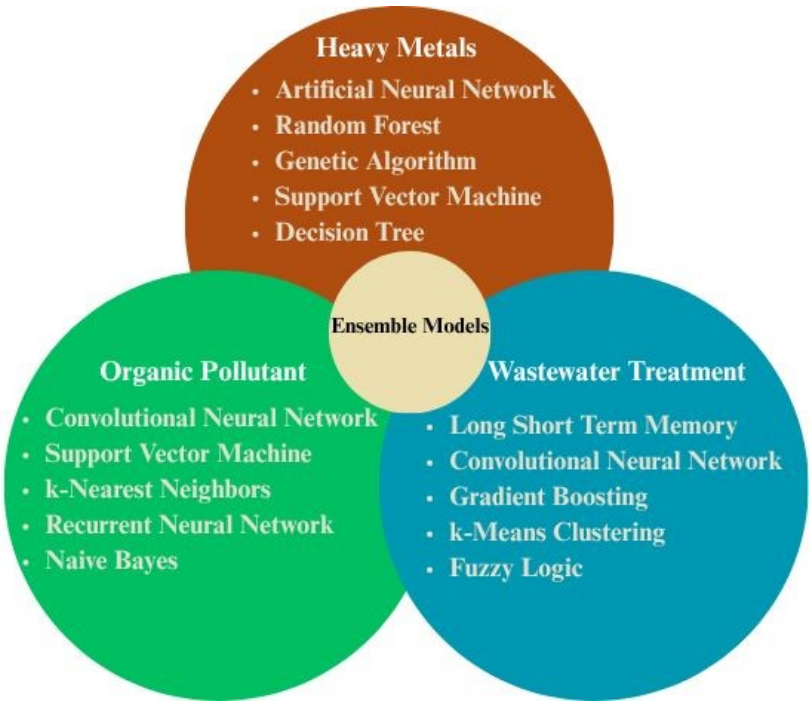


Figure 1: Machine learning models used in Heavy Metals, Organic Pollutants, and Wastewater Treatment, with Ensemble Models at the intersection.

One evaluation has utilized 22 metagenomic and genomic data of microbial community integrated with AI and ML algorithms to enhance degradation of environmental contaminants and toxins while providing insights of genetic and functional potentials of these communities (22). Some studies have shown that ML can predict pathways that are involved in degradation by analyzing sequences while functional profiling using AI enables identification of key enzymes and taxa responsible for degradation which helps optimization of bioremediation process (1,23). Biomarkers like genes, metabolites and proteins can indicate the activity and presence of specific microbial communities capable of degrading pollutants. Algorithms like SVM, ANN and RF can identify potential biomarkers by recognizing correlation between the microbial activities and pollutant concentrations as indicated by two independent research groups working at different timelines (24,25). These integration strategies of ML with microbiome analysis are shown in **Figure 2**.



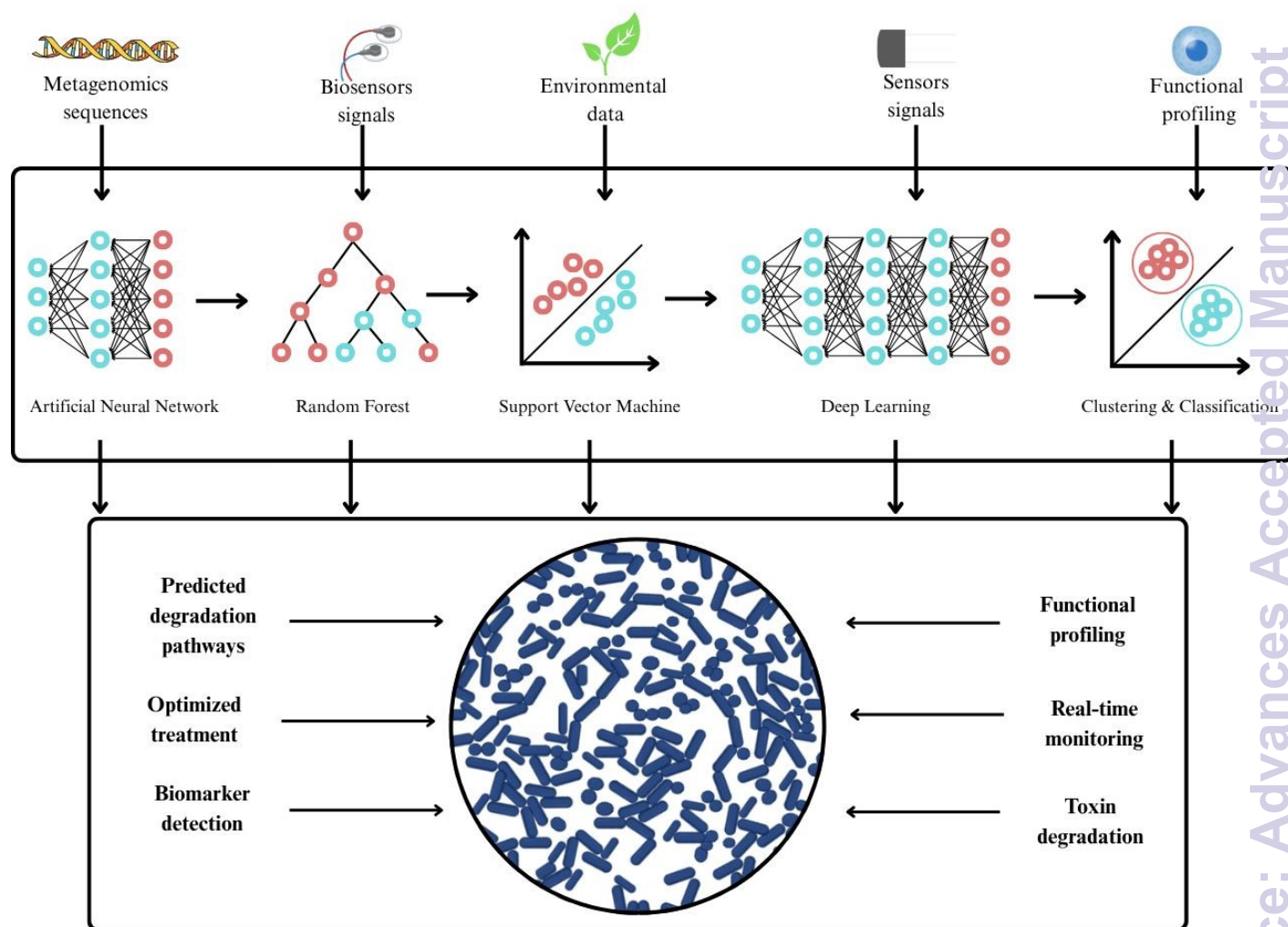


Figure 2: Integration of Machine Learning Techniques in Microbiome Analysis for Environmental and Clinical Applications.

3. Bioinformatics and Metagenomics

Bioinformatics aids in using microarray data by enhancing the structural characterization of proteins (26). In a recent study (27) have constructed a comprehensive proteins library that supports heavy metal bio-removal which were modeled through Alphafold2. Bioinformatics can help in Insilco studies and analyzation of the data and can also enhance the bioremediation using the databases for gene identification and microbial degradation pathways of compounds Table.2. Tools such as I-TASSER, Phyre2 and SWISS-MODEL can also be used for protein structural prediction which is often the initial step for detection of active sites to determine enzymatic function. Other tools like CASTp are used for automated detection of active sites. Implementation of bioinformatics tools such as PathPred and University of Minnesota Pathway Prediction System (UMPPS) are freely available tools that provides users with variety of biochemical reactions



ultimately leads to modification of pathways (4). Tools like BLAST is an alignment tool applied for the identification of resemblances among sequences of both protein and protein which is based on the hypothesis that homologous sequences are expected to function similarly. Other tools like genome scale metabolic model (GSMM), Constraint based reconstruction and analysis (COBRA) uses genetic information from databases for construction of metabolic pathways, while KEGG, Biocyc and others provide all the probable metabolic pathways.

Complete wide range analysis of bacterial communities will play important role in identification of new genes and metabolic pathways for bioremediation. Bioinformatics approaches have increased our capability to detect pollution sources and screen the fluctuations in the microorganisms throughout the process. An overview of successful studies where AI can be successfully integrated into bioremediation models is shown in **Table 1**. Additionally, metagenomics is a critical tool in the bioremediation provides an insightful grasp of the structural and functional properties of the bacterial communities that are engage in bioremediation process. as represented in **Table 2**. This comprises the direst isolation of genome from the sources to identify composition without any process of isolation and cultivation in the laboratory. The variability in sequencing technologies and biased datasets can limit model generalizability across diverse contaminated sites. Models trained on specific datasets may not perform well on unseen environmental conditions due to site-specific factors like pH or microbial diversity.

Table 1: AI application in different bioremediation approaches, leveraging the capabilities of AI algorithms, to improve efficiency, accuracy, and sustainability.

AI application	Bioremediation approach	References
Heavy Metal Removal	Constructed Wetlands	(16). (28)
Microbial Selection	Bioaugmentation	(16) (7)
Pollution Monitoring	Data Driven Monitoring	(29), (30)
Nutrient Supplementation	Bio Stimulation	(31), (32)
Real Time Adaptive Control	Dynamic Bioremediation	(33)
Wastewater Treatment	Optimized Constructed Wetlands	(31)
Nanotechnology Integration	Enhanced Bioremediation with Nanoparticles	(7)

Table 2: Tools used in metagenomic analysis enabling the study of microbial communities in polluted environments to identify and understand the microorganisms involved in pollutant degradation (34,35).

Name of the Software	Application Areas in Bioremediation
Metagenome Seq	Evaluation of the abundance of 16S rRNA genes in meta-profiling
UCLUST	A clustering tool, which utilizes USEARCH to allocate sequences to Clusters
Mothur	Used in the quality analysis of reads for taxonomic classification
NGSQC toolkit	Method of performing quality control analysis in a direct environment
RDP (Ribosomal Database Project)	Biodiversity analysis, sequence arrangement, alignment, trimming, and taxonomic classification of sequences
Pfam	A large collection of families and domains expressed by profile HMMs and multiple sequence alignments
Prodigal	Identification of translation initiation sites in prokaryotic genes
CAMERA	A server for a metagenomic database containing sequences from environmental samples collected during the GOS
envDB	Prokaryotic taxa environmental distribution database and tool server
myPhyloDB	A tool used for the purpose of storage and metagenomic analysis
FUNGI Path	A database used for metagenomic and pathological studies of fungi
PyNASt	Aligned sequences of representative OTUs
Meta MIS	Analysis of microbial interaction
FOAM	Created to screen environmental metagenomic sequencing datasets and to offer a novel functional ontology specialized in categorizing gene functions pertinent to environmental microbes using HMMs



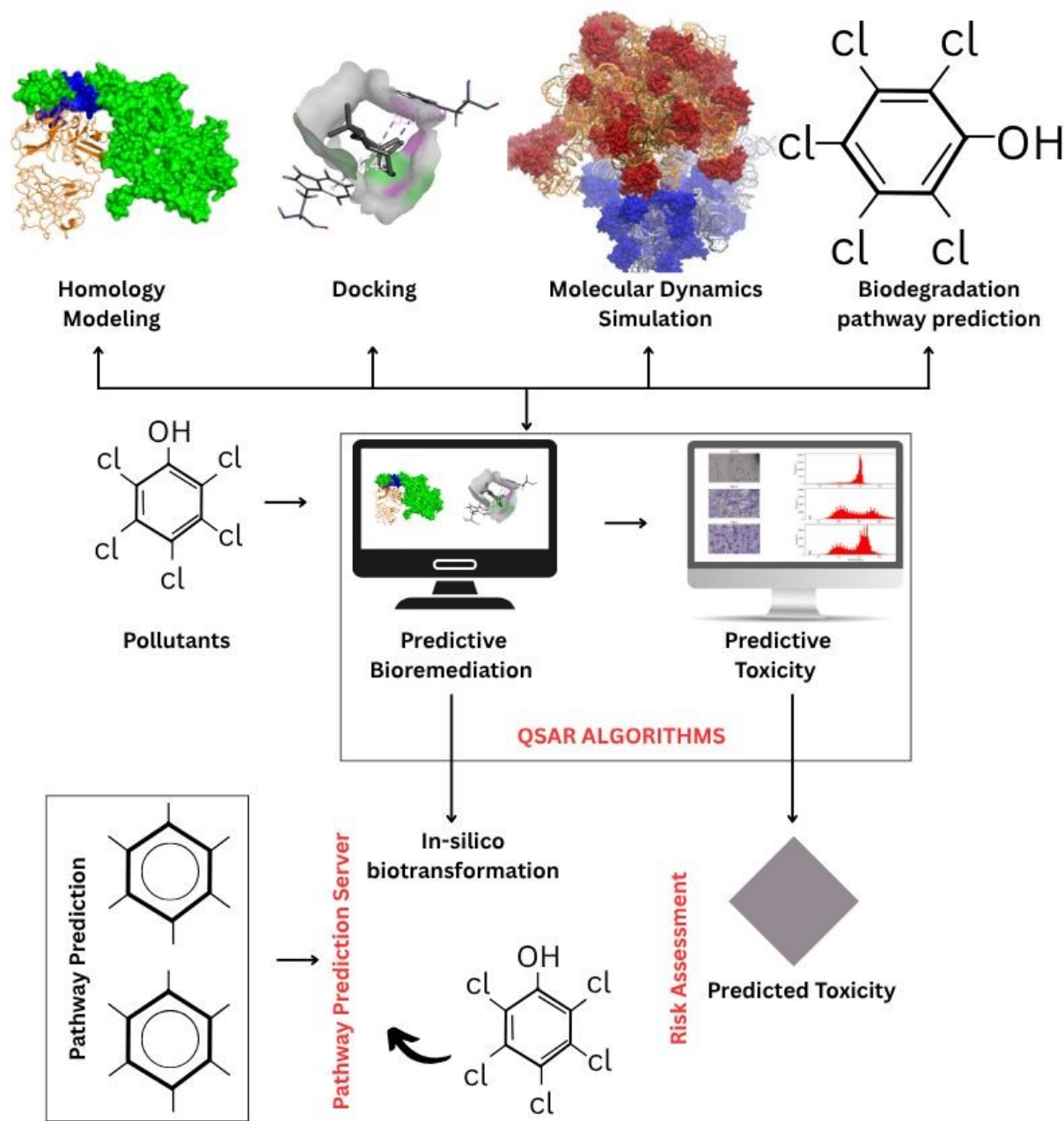


Figure 3: Graphical illustration of different computational techniques used in predictive bioremediation.

The workflow as shown in **Figure 3** of metagenomics during the extraction of the DNA from the sources and construction of genomic library, analysis and sequencing of the data to target genes for further application. The online accessibility of several Metagenomic Analysis Shell for Microbial Communities (SmashCommunity), Meta Genome Analyzer (MEGAN), Metagenomic Rapid Annotation using Subsystems Technology (MG-RAST) and Community



Cyberinfrastructure for Advanced Microbial Ecology Research and Analysis (CAMERA) equips the scientists with cutting-edge approaches for the metagenomic based studies. Metagenomics is a strong tool for exposing diverse bacterial species and their importance, but it also has some limitations such as lack of ideal guideline for each specific bacterial species, there are several issues with the quality of metagenomics because of longer reads, large data and numerous error models. Similar datasets can be interpreted differently by distinct sequencing technologies. Amplicon based sequencing suffers biased amplification and is unable to amplify unknown regions. High quality datasets are crucial for training AI and ML algorithms; incomplete and biased data can lead to inaccurate predictions.

4. Real World Applications

Real time applications in evaluations have involved hybrid techniques by merging two ML techniques including ANN-RF and SVM-RF (9,10). Four-layer ANN, four-layer RNN, typical adaptive neuro FIS and typical CNN are the four most common neural networks structures used in water treatment and monitoring as represented in **Figure 4**. ANN is the most preferred in comparison with other algorithms such as Genetic Algorithms (GA) and SVM. ML models with output percentage of absorbate, absorption capacity, effluent concentrations, water quality parameters.

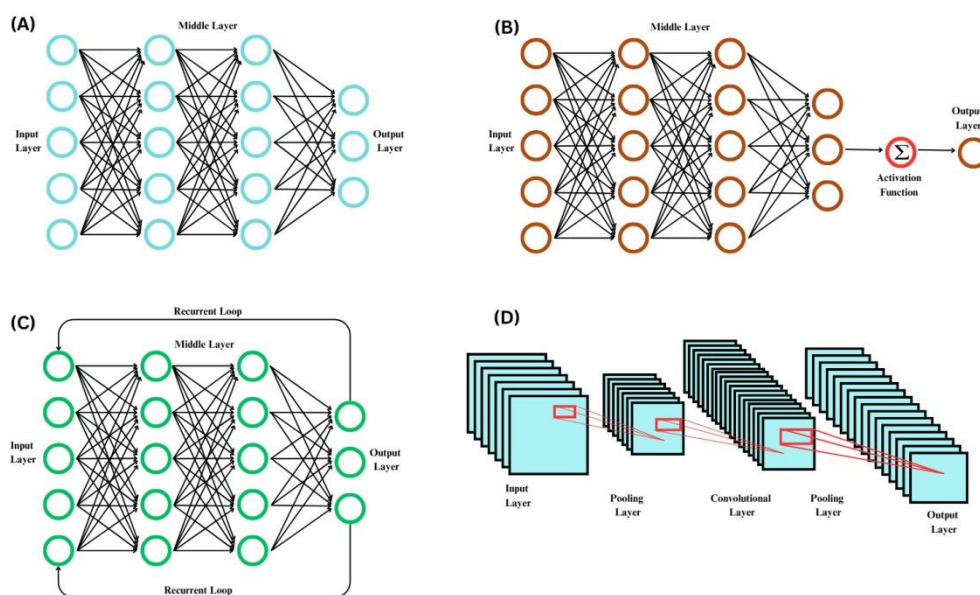


Figure 4: Illustration of Different Neural Network Architectures Used in Machine Learning and Deep Learning Applications for Bioremediation studies. (A) Four-layered ANN (B) Four-layered ANFIS (C) Four layered RNN (D) Four Layered CNN.



Other models with notable success are ANFIS, SVM and RF generally achieving R^2 values of more than 0.9 and in some cases greater than 0.99. Another process known as membrane filtration which refers to separation of contaminants using barriers or filter can be optimized using ML techniques. Like the previous process, ANNs are the most dominant model among RNN, SVM and ANN. ML techniques like Recurrent Neural Network (RNN), Random Forest (RF), Support Vector Machines (SVM), Artificial Neural Network (ANN) and Fuzzy Interference Systems (FIS) as shown in **Table 3**.

Table 3: Instances of successful implementation of ML models used in bioremediation, water treatment and monitoring.

Techniques	Applications	Water Treatment Applications	References
Random Forest	Regression, Classification and SVM	Adsorption percent removal, dissolved oxygen modeling	(36)
Support Vector Machines/ Regressions	SVM, Regression, Classification/ Pattern Analysis	Membrane process parameter, Biological Oxygen demand (BOD) and Chemical Oxygen Demand (COD) modelling	(37) (38)
Fuzzy Inference System	Regression, Classification and stochastic algorithm	Disinfection by product modeling	(39)
Artificial Neural Network	Regression, Classification and SVM	Adsorption process modeling, Dissolved oxygen concentration modeling.	(40)
k-Nearest Neighbor	Classification and SVM	Aquaponics growth stage classification	(41) (42)
Recurrent Neural Network / Long Short-Term Memory	Regression, Classification and SVM	Membrane process parameter, Dissolved oxygen concentration modeling	(43)
Adaptive Neuro-Fuzzy Inference Systems	Regression, Classification and SVM	Membrane process parameter, Biological Oxygen demand (BOD) and Chemical Oxygen Demand (COD) modelling	(44)
Convolutional Neural Network	Regression, Classification and SVM	Adsorption process modeling	(45)



One study reported the use of three different models for the prediction of copper removal in an absorption process using clay as the primary adsorbent (2). These three ML models include ANN, SVM and RF developed using open-source software for programming language R, the dataset was divided into two parts for training and testing with ratio of 8:2 providing 80% for the training and 20% for the testing. The ANN consisted of a four layered model having one input and output layer and two hidden layers, each neuron relying on linear output function. RF model utilizes 76 samples to develop decision trees and SVM model was developed using linear kernel. RF and ANN models showed the best performance in terms of accuracy achieving greater than 0.99, while SVM achieved 0.93.

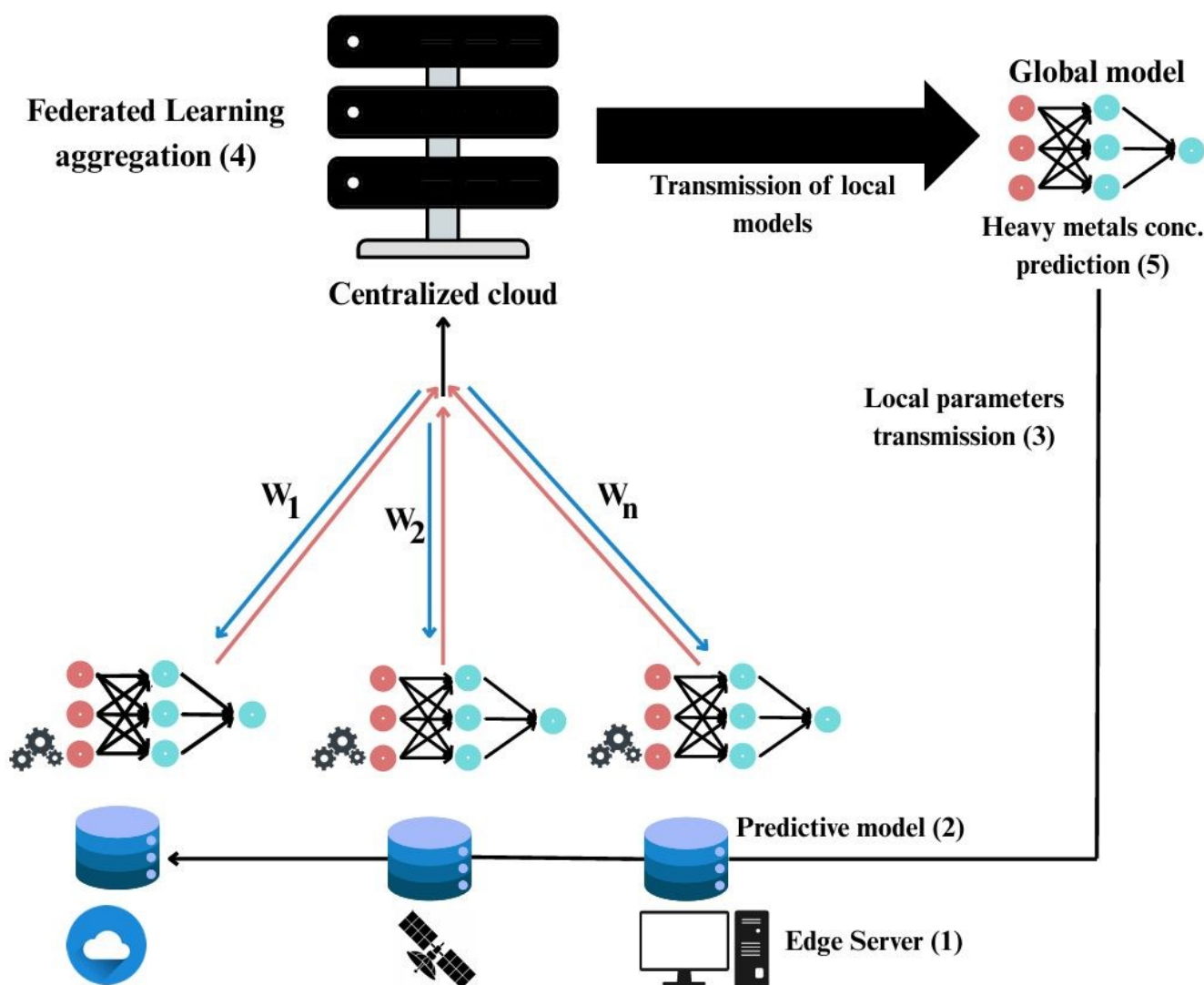


Figure 5: Federated learning system for predicting heavy metal concentrations using decentralized edge models aggregated into a global model.



Studies on water quality management use many models such as ANN, ANFIS, RNN, RF. ANFIS has outperformed typical ANN and SVM and in some cases outperformed by hybrid models, ANN and ANFIS models achieved R^2 values greater than 0.999 with both models while forecasting water level (46,47). Although the studies compare different models but it lacks comparisons to simpler baseline models (e.g., linear regression or traditional statistical methods) to contextualize AI model performance.

Baseline comparisons would clarify whether complex AI models provide significant improvements over simpler approaches, justifying their computational cost. In other studies authors have compared both ANFIS and ANN models for estimation of Water Quality Index (WQI), both the models showed relative success by achieving accuracy greater than 0.99 (48,49). Federated learning technology and Edge cloud-based server is also used to develop automated system for prediction and monitoring as represented in **Figure 5** (50). High R^2 values (e.g., >0.99) suggest potential overfitting, especially for complex models like ANN and hybrid ANN-RF, overfitting occurs when models fit training data too closely, failing to generalize to new data. The use of an 80:20 train-test split is a basic validation approach, but it may not detect overfitting in small or non-diverse datasets which is common in bioremediation studies. Techniques like k-fold cross-validation or regularization could mitigate overfitting.

To illustrate the practical impact of AI-driven bioremediation, two case studies demonstrate successful implementations at meaningful scales. First, a municipal wastewater treatment facility in Hyderabad, India, employed an ANN-RF hybrid model to optimize microbial degradation of heavy metals (e.g., cadmium, lead) in effluent from industrial discharges. Using metagenomic data analyzed via QIIME (Table 2), the model identified key microbial taxa and predicted optimal bioaugmentation strategies, achieving a 95% reduction in heavy metal concentrations across 10,000 m³ of wastewater daily. Real-time biosensors, integrated with deep learning algorithms (Figure 1), enabled continuous monitoring, ensuring stable performance over six months. Second, a field-scale bioremediation project at a chromium-contaminated mining site in Odisha, India, utilized RF models to predict microbial survival under varying pH and temperature conditions. By integrating metagenomic insights from MG-RAST (Table 2), the project enhanced bio-stimulation, reducing chromium levels by 90% across a 5-hectare site over one year. These cases highlight the scalability of AI-driven approaches, leveraging models like ANN-RF and RF (Table 3) to address large-scale environmental challenges, though site-specific recalibration remains critical for sustained success.

In a study authors have utilized a conventional response surface methodology and ANN model for enhancing fluoranthene degradation by mycobacterium litorale. The study uses optimized CaCl₂, KH₂PO₄ and NH₄NO₃. The designed model maximized fluoranthene degradation. It was obtained that the model could efficiently simulate the degradation process and the output received from the ANN model were reliable, precise and reproducible. The authors have claimed that the degradation



on the 3rd day was better with 51.28% in comparison to an un optimized degradation method which was only 26.37% after 7 days (51).

5. Challenges and Prospects

To fully realize the promises of multi-omics integration in systems biology researchers will need to address some challenges which include documentation and integration of data collection and analytical methods, establishment of databases on metabolites and pathways. CRISPR and gene editing are some novel technologies that offers the possibility of optimizing bacterial metabolism with high accuracy. The overall goal is moving towards a more efficient, cost effective, sustainable solution for bioremediation. For training AI and ML, there is a requirement of high quality and comprehensive datasets as incomplete and biased data will lead to incorrect predictions, requirement of significant. Databases availability is another major concern for AI and ML; existing databases may not cover full diversity of microbial communities and their metabolites.

Practical Limitations of Deploying AI-Driven Bioremediation Systems

Deployment of AI driven bioremediation systems especially involvement of gene editing tools like CRISPR must comply with environmental regulations. For instance, regulations like EPA guidelines or the European Union's Genetic Modification Directives have strict controls on releasing genetically modified organisms (GMOs) to the environment to prevent ecological imbalances. AI systems relying on the metagenomic data involve handling large datasets which may raise concerns, Regulatory frameworks like the General Data Protection Regulation (GDPR) in Europe could complicate data management. Also, we noticed lack of ideal guidelines for metagenomic analysis (e.g., variability in sequencing technologies). This absence of standardized protocols for AI model validation and deployment in bioremediation can hinder regulatory approval, as agencies may require consistent, reproducible methodologies. Use of advanced models like ANN, RF and SVM which generally requires significant computational power but in developing countries access to high performance computational infrastructure or cloud-based servers may be limited increasing the costs and reduces the scalability. Implementing AI-driven systems, including biosensors and real-time monitoring tools requires substantial investment in hardware, software, and maintenance. Small-scale or rural bioremediation projects may lack funding, as the manuscript notes no funding was received for this work, reflecting a broader challenge in resource allocation. Apart from these there are other challenges that are related to data availability and quality, the need for high-quality datasets for training AI models. In real-world settings, collecting comprehensive, unbiased data from diverse contaminated sites is resource-intensive, requiring specialized equipment and trained personnel, which may not be feasible in low-resource environments. Operating AI-driven systems, such as those integrating metagenomics or machine learning models demands expertise in bioinformatics, data science, and environmental microbiology. The manuscript's reliance on complex tools like AlphaFold2 and I-TASSER suggests a steep learning curve, which may not be practical in regions with limited access to trained professionals. Real-time monitoring systems, such as AI-powered biosensors require ongoing



maintenance and troubleshooting. Without onsite expertise, system failures or misinterpretations of AI outputs could compromise bioremediation efficiency. Deploying these systems in non-ideal settings necessitates training local personnel, which can be time-consuming and costly. The manuscript's call for interdisciplinary collaboration underscores the need for knowledge transfer, but this is challenging in areas with limited educational infrastructure

Fundamental limitations of AI approaches in environmental contexts pose significant challenges to their widespread adoption in bioremediation. Data scarcity, due to the high cost and complexity of collecting diverse metagenomic and environmental datasets, restricts the training of robust AI models like ANN and RF, particularly for rare pollutants or understudied microbial communities. Environmental heterogeneity, characterized by variations in site-specific factors such as pH, temperature, and pollutant profiles, limits the generalizability of models trained on specific datasets, as seen in applications like heavy metal removal (Table 1). Furthermore, scaling laboratory results, such as those from AlphaFold2 or QIIME analyses, to field applications is hindered by uncontrolled environmental variables and the need for real-time model recalibration, necessitating robust validation strategies and adaptive algorithms to ensure reliable performance across diverse ecological settings.

Expansion of coverage is required for more complete understanding of microbial metabolism (52). A further major difficulty is that distinct sequencing technologies interpret same datasets differently, missing environmental context limits application of supervised ML, there is still need of wet-lab validation because of the limited biodegradation database availability. Researchers can combine more sensitive tools to be used for these challenging tasks, such as advanced structural elucidation and optimization. By enhancing multi-omics and genetic engineering tools can help to develop a more sustainable bioremediation strategy. Further refinements in metagenomics and bioinformatics will help in precise ecological interpretations. These omics approaches can anticipate microbial metabolism at contaminated site therefore by utilizing multi-omics approaches can lead to new hypothesis and theories. Another area of investigation would be an experimental approach from a multidisciplinary perspective to achieve better prediction and application.

Generated data might not be utilized into different ML techniques to test effectiveness. Another major challenge is management of several hazardous organic compounds which can be reduced by the development of new monitoring techniques. Most of the literature on metagenomic focuses on limited range of enzymes, mostly esterases and oxidases while the other also play an essential role in biodegradation of pollutants and toxins (6). Closed loops systems should emphasize the development of closed systems where AI/ ML models' predictions are iteratively validated through experimental testing and refinement creating a feedback cycle that will enhance both predictive accuracy and practical applicability.

Outlook



In hindsight, bioinformatics tools like AlphaFold, I-TASSER and SWISS-MODEL are used for protein modeling, which is an important step for the metagenomics, recent development in sequencing technologies have removed the obstacles and opened the door for the metagenomics revealing novel information about microbial diversity. Metagenomics can analyze genetic material of microbial communities to solve problems and discover new enzymes, genes and metabolic pathways. In metagenomics tools like QIIME, UPARSE, MOTHUR facilitate bioinformatics analysis with MG-RAST and MetaPhlAn2 providing phylogenetic analysis and functional insights, these web-based tools are crucial for advancing omics. AI and ML have shown potential to solve complex issues facing during bioremediation process, ML algorithms like ANN, CNN, RNN and RF have achieved accuracy of more than 0.99. ML models have optimized, predicted, modeled and automated some applications of bioremediation and its strategies. To advance AI driven bioremediation researchers should prioritize standardized protocols that integrate metagenomics with ML frameworks for accurate and reproducible analysis. There is a need for robust publicly available datasets representing diverse contaminated environments data to enhance generalizability. Enhancements of bioinformatics pipelines such as QIIME and MG -RAST to manage noise in the datasets and integration of protein structure prediction tools to accelerate novel discoveries.

ML algorithms have shown better results in prediction, modelling and optimization of water treatment process to degrade pollutant and toxins. Though many studies have reviewed with success but there are sets of challenges and limitations which include data management, reproducibility and transparency that must be addressed. Interdisciplinary collaborations integrating bioinformatics, metagenomics and machine learning will be pivotal in overcoming these challenges. By bridging the gaps between these interdisciplinary collaborations future research can be transformative in bioremediation ultimately contributing to more effective bioremediation strategies and environmental solutions.

DECLARATIONS

Ethics approval and consent to participate: Not Applicable

Consent for publication: All authors agree to this submission.

Availability of data and materials: Data will be made available from the corresponding author upon reasonable request.

CONFLICT OF INTEREST STATEMENT

The authors declare no known competing interests.

FUNDING DECLARATION



No funding was received for this work.

REFERENCES

1. Haque S, Srivastava N, Pal DB, Alkhanani MF, Almalki AH, Areeshi MY, et al. Functional microbiome strategies for the bioremediation of petroleum-hydrocarbon and heavy metal contaminated soils: A review. *Science of The Total Environment*. 2022 Aug;833:155222.
2. Bhagat SK, Pyrgaki K, Salih SQ, Tiyaasha T, Beyaztas U, Shahid S, et al. Prediction of copper ions adsorption by attapulgitite adsorbent using tuned-artificial intelligence model. *Chemosphere*. 2021 Aug;276:130162.
3. Bhattacharya R, Bose D, Ganti P, Rizvi A, Halder G, Sarkar A. Bioelectricity production and bioremediation potential of *Withania somnifera* in plant microbial fuel cells with food wastes as enrichment. 2023.
4. Heirendt L, Arreckx S, Pfau T, Mendoza SN, Richelle A, Heinken A, et al. Creation and analysis of biochemical constraint-based models using the COBRA Toolbox v.3.0. *Nat Protoc*. 2019 Mar 20;14(3):639–702.
5. Gu J, Wang Z, Kuen J, Ma L, Shahroudy A, Shuai B, et al. Recent advances in convolutional neural networks. *Pattern Recognit*. 2018 May;77:354–77.
6. Firincă C, Zamfir LG, Constantin M, Răut I, Jecu ML, Doni M, et al. Innovative Approaches and Evolving Strategies in Heavy Metal Bioremediation: Current Limitations and Future Opportunities. *J Xenobiot*. 2025 Apr 26;15(3):63.
7. Patowary R, Devi A, Mukherjee AK. Advanced bioremediation by an amalgamation of nanotechnology and modern artificial intelligence for efficient restoration of crude petroleum oil-contaminated sites: a prospective study. *Environmental Science and Pollution Research*. 2023 May 23;30(30):74459–84.
8. World Health Organization: WHO. (2023, November 21). Antimicrobial resistance. <https://www.who.int/news-room/fact-sheets/detail/antimicrobial-resistance> [Internet]. World Health Organization: WHO, 2023.
9. Alam G, Ihsanullah I, Naushad Mu, Sillanpää M. Applications of artificial intelligence in water treatment for optimization and automation of adsorption processes: Recent advances and prospects. *Chemical Engineering Journal*. 2022 Jan;427:130011.



10. Xie Y, Chen Y, Lian Q, Yin H, Peng J, Sheng M, et al. Enhancing Real-Time Prediction of Effluent Water Quality of Wastewater Treatment Plant Based on Improved Feedforward Neural Network Coupled with Optimization Algorithm. *Water (Basel)*. 2022 Mar 27;14(7):1053.
11. Bose D, Bhattacharya R, Gopinath M, Vijay P, Krishnakumar B. Bioelectricity production and bioremediation from sugarcane industry wastewater using microbial fuel cells with activated carbon cathodes. *Results in Engineering*. 2023 Jun;18:101052.
12. Raper E, Stephenson T, Anderson DR, Fisher R, Soares A. Industrial wastewater treatment through bioaugmentation. *Process Safety and Environmental Protection*. 2018 Aug;118:178–87.
13. Bhattacharya R, Bose D, Yadav J, Sharma B, Sangli E, Patel A, et al. Bioremediation and bioelectricity from Himalayan rock soil in sediment-microbial fuel cell using carbon rich substrates. *Fuel*. 2023 Jun;341:127019.
14. Nivetha N, Srivarshine B, Sowmya B, Rajendiran M, Saravanan P, Rajeshkannan R, et al. A comprehensive review on bio-stimulation and bio-enhancement towards remediation of heavy metals degeneration. *Chemosphere*. 2023 Jan;312:137099.
15. Bose D, Santra M, Sanka RVSP, Krishnakumar B. Bioremediation analysis of <sc>sediment-</sc> microbial fuel cells for energy recovery from microbial activity in soil. *Int J Energy Res*. 2021 Mar 25;45(4):6436–45.
16. Gupta PK, Yadav B, Kumar A, Himanshu SK. Machine learning and artificial intelligence application in constructed wetlands for industrial effluent treatment: advances and challenges in assessment and bioremediation modeling. In: *Bioremediation for Environmental Sustainability*. Elsevier; 2021. p. 403–14.
17. Gupta PK, Yadav B, Kumar A, Himanshu SK. Machine learning and artificial intelligence application in constructed wetlands for industrial effluent treatment: advances and challenges in assessment and bioremediation modeling. In: *Bioremediation for Environmental Sustainability*. Elsevier; 2021. p. 403–14.
18. Wang J, Peng C, Liu X. Prediction of bacterial microbiota changes in various microplastics-contaminated environments based on machine learning. *J Environ Chem Eng*. 2025 Oct;13(5):117461.
19. Yang Y, Jerger A, Feng S, Wang Z, Brasfield C, Cheung MS, et al. Improved enzyme functional annotation prediction using contrastive learning with structural inference. *Commun Biol*. 2024 Dec 23;7(1):1690.



20. Oyewusi HA, Wahab RA, Akinyede KA, Albadrani GM, Al-Ghadi MQ, Abdel-Daim MM, et al. Bioinformatics analysis and molecular dynamics simulations of azoreductases (AzrBmH2) from *Bacillus megaterium* H2 for the decolorization of commercial dyes. *Environ Sci Eur*. 2024 Feb 18;36(1):31.
21. Srinivasan S, Raajasubramanian D, Ashokkumar N, Vinothkumar V, Paramaguru N, Selvaraj P, et al. Nanobiosensors based on on-site detection approaches for rapid pesticide sensing in the agricultural arena: A systematic review of the current status and perspectives. *Biotechnol Bioeng*. 2024 Sep 10;121(9):2585–603.
22. Ayilara MS, Babalola OO. Bioremediation of environmental wastes: the role of microorganisms. *Frontiers in Agronomy*. 2023 May 30;5.
23. Wijaya J, Park J, Yang Y, Siddiqui SI, Oh S. A metagenome-derived artificial intelligence modeling framework advances the predictive diagnosis and interpretation of petroleum-polluted groundwater. *J Hazard Mater*. 2024 Jul;472:134513.
24. Vikram S, Guerrero LD, Makhalanyane TP, Le PT, Seely M, Cowan DA. Metagenomic analysis provides insights into functional capacity in a hyperarid desert soil niche community. *Environ Microbiol*. 2016 Jun 10;18(6):1875–88.
25. Zhong S, Zhang K, Bagheri M, Burken JG, Gu A, Li B, et al. Machine Learning: New Ideas and Tools in Environmental Science and Engineering. *Environ Sci Technol*. 2021 Aug 17;acs.est.1c01339.
26. Chettri D, Verma AK, Chirania M, Verma AK. Metagenomic approaches in bioremediation of environmental pollutants. *Environmental Pollution*. 2024 Dec;363:125297.
27. Uttarotai T, Mukjang N, Chaisoung N, Pathom-Aree W, Pekkoh J, Pumas C, et al. Putative Protein Discovery from Microalgal Genomes as a Synthetic Biology Protein Library for Heavy Metal Bio-Removal. *Biology (Basel)*. 2022 Aug 17;11(8):1226.
28. Biswas PP, Chen WH, Lam SS, Park YK, Chang JS, Hoang AT. A comprehensive study of artificial neural network for sensitivity analysis and hazardous elements sorption predictions via bone char for wastewater treatment. *J Hazard Mater*. 2024 Mar;465:133154.
29. Wani AK, Rahayu F, Ben Amor I, Quadir M, Murianingrum M, Parnidi P, et al. Environmental resilience through artificial intelligence: innovations in monitoring and management. *Environmental Science and Pollution Research*. 2024 Feb 15;31(12):18379–95.



30. Tao H, Al-Khafaji ZS, Qi C, Zounemat-Kermani M, Kisi O, Tiyyasha T, et al. Artificial intelligence models for suspended river sediment prediction: state-of-the art, modeling framework appraisal, and proposed future research directions. *Engineering Applications of Computational Fluid Mechanics*. 2021 Jan 1;15(1):1585–612.
31. Sahu S, Kaur A, Singh G, Kumar Arya S. Harnessing the potential of microalgae-bacteria interaction for eco-friendly wastewater treatment: A review on new strategies involving machine learning and artificial intelligence. *J Environ Manage*. 2023 Nov;346:119004.
32. Bibri SE, Krogstie J, Kaboli A, Alahi A. Smarter eco-cities and their leading-edge artificial intelligence of things solutions for environmental sustainability: A comprehensive systematic review. *Environmental Science and Ecotechnology*. 2024 May;19:100330.
33. Chang NB, Mohiuddin G, Crawford AJ, Bai K, Jin KR. Diagnosis of the artificial intelligence-based predictions of flow regime in a constructed wetland for stormwater pollution control. *Ecol Inform*. 2015 Jul;28:42–60.
34. Bharagava RN, Purchase D, Saxena G, Mulla SI. Applications of Metagenomics in Microbial Bioremediation of Pollutants. In: *Microbial Diversity in the Genomic Era*. Elsevier; 2019. p. 459–77.
35. Arya P, Ravindra. Metagenomics based approach to reveal the secrets of unculturable microbial diversity from aquatic environment. In: *Recent Advancements in Microbial Diversity*. Elsevier; 2020. p. 537–59.
36. Liu Y, Wang Y, Zhang J. New Machine Learning Algorithm: Random Forest. In 2012. p. 246–52.
37. Cortes C, Vapnik V. Support-vector networks. *Mach Learn*. 1995 Sep;20(3):273–97.
38. Chua KS. Efficient computations for large least square support vector machine classifiers. *Pattern Recognit Lett*. 2003 Jan;24(1–3):75–80.
39. Moraga C, Trillas E, Guadarrama S. Multiple-Valued Logic and Artificial Intelligence Fundamentals of Fuzzy Control Revisited. In: *Artificial Intelligence in Logic Design*. Dordrecht: Springer Netherlands; 2004. p. 9–37.
40. Uhrig RE. Introduction to artificial neural networks. In: *Proceedings of IECON '95 - 21st Annual Conference on IEEE Industrial Electronics*. IEEE; p. 33–7.
41. Samet H. K-Nearest Neighbor Finding Using MaxNearestDist. *IEEE Trans Pattern Anal Mach Intell*. 2008 Feb;30(2):243–52.



42. Jiang L, Cai Z, Wang D, Jiang S. Survey of Improving K-Nearest-Neighbor for Classification. In: Fourth International Conference on Fuzzy Systems and Knowledge Discovery (FSKD 2007). IEEE; 2007. p. 679–83.
43. DiPietro R, Hager GD. Deep learning: RNNs and LSTM. In: Handbook of Medical Image Computing and Computer Assisted Intervention. Elsevier; 2020. p. 503–19.
44. Farhoudi J, Hosseini SM, Sedghi-Asl M. Application of neuro-fuzzy model to estimate the characteristics of local scour downstream of stilling basins. *Journal of Hydroinformatics*. 2010 Mar 1;12(2):201–11.
45. Gu J, Wang Z, Kuen J, Ma L, Shahroudy A, Shuai B, et al. Recent advances in convolutional neural networks. *Pattern Recognit*. 2018 May;77:354–77.
46. Al-Yaari M, Aldhyani THH, Rushd S. Prediction of Arsenic Removal from Contaminated Water Using Artificial Neural Network Model. *Applied Sciences*. 2022 Jan 19;12(3):999.
47. Mazaheri H, Ghaedi M, Ahmadi Azqhandi MH, Asfaram A. Application of machine/statistical learning, artificial intelligence and statistical experimental design for the modeling and optimization of methylene blue and Cd(<scp>ii</scp>) removal from a binary aqueous solution by natural walnut carbon. *Physical Chemistry Chemical Physics*. 2017;19(18):11299–317.
48. Mazloom MS, Rezaei F, Hemmati-Sarapardeh A, Husein MM, Zendehboudi S, Bemani A. Artificial Intelligence Based Methods for Asphaltenes Adsorption by Nanocomposites: Application of Group Method of Data Handling, Least Squares Support Vector Machine, and Artificial Neural Networks. *Nanomaterials*. 2020 May 6;10(5):890.
49. Mesellem Y, Hadj AA El, Laidi M, Hanini S, Hentabli M. Computational intelligence techniques for modeling of dynamic adsorption of organic pollutants on activated carbon. *Neural Comput Appl*. 2021 Oct 4;33(19):12493–512.
50. Yaseen ZM. The next generation of soil and water bodies heavy metals prediction and detection: New expert system based Edge Cloud Server and Federated Learning technology. *Environmental Pollution*. 2022 Nov;313:120081.
51. Dudhagara DR, Rajpara RK, Bhatt JK, Gosai HB, Dave BP. Bioengineering for polycyclic aromatic hydrocarbon degradation by *Mycobacterium litorale*: Statistical and artificial neural network (ANN) approach. *Chemometrics and Intelligent Laboratory Systems*. 2016 Dec;159:155–63.



52. Khanna S, Kumar A. Bioinformatics Toward Improving Bioremediation. In: Biotechnological Innovations for Environmental Bioremediation. Singapore: Springer Nature Singapore; 2022. p. 631–69.



Availability of data and materials: Data will be made available from the corresponding author upon reasonable request.

[View Article Online](#)
DOI: 10.1039/D5VA00240K

