AI Systems Engineering eISSN: 3093-8449

https://aise.cultechpub.com/index.php/aise

Copyright: © 2025 by the authors. This article is published by the Cultech Publishing Sdn. Bhd. under the terms of the Creative Commons Attribution4.0 International License (CC BY 4.0): https://creativecommons.org/licenses/by/4.0/

# **Intelligent IoT-Based Water Quality Monitoring and Predictive Analysis Using Machine Learning**

Sumit Kushwaha\*, Ritika Pandey

Department of Computer Applications, University Institute of Computing, Chandigarh University, Mohali-140413, Punjab, India

\*Corresponding author: Sumit Kushwaha, sumit.kushwaha1@gmail.com

#### **Abstract**

This paper proposes a scalable and intelligent system integrating Internet of Things (IoT) technologies with advanced machine learning (ML) algorithms for real-time water quality monitoring and predictive analysis. The system utilizes low-cost and reliable sensors deployed on microcontroller platforms such as ESP32 and NodeMCU to continuously collect vital water parameters, including pH, turbidity, temperature, dissolved oxygen, and total dissolved solids (TDS). Sensor data is transmitted via wireless communication protocols to cloud platforms like ThingSpeak and AWS IoT for centralized storage and preprocessing. Machine learning models, including Random Forest, Support Vector Machines (SVM), Extreme Gradient Boosting (XGBoost), and Long Short-Term Memory (LSTM) networks, are trained on historical data to forecast parameter fluctuations, detect anomalies, and compute the Water Quality Index (WQI), providing a standardized quality assessment metric. The system's automated alert and visualization layer delivers realtime dashboards and warning notifications to stakeholders, enabling timely responses to contamination events. Experimental results on self-collected IoT sensor data and benchmark public datasets demonstrate high predictive accuracy, precision, and recall, confirming the framework's suitability for continuous environmental surveillance. This approach addresses limitations of traditional costly and time-consuming laboratory tests by offering affordable, scalable, and adaptive monitoring, supporting sustainable water resource management. Future work will explore additional sensor integration, enhanced deep learning architectures, and economic feasibility studies for large-scale implementation. The proposed framework contributes significantly to environmental protection, public health safeguarding, and intelligent infrastructure development aligned with Sustainable Development Goals (SDG 6), fostering smarter and safer water ecosystems.

# **Keywords**

IoT, Machine Learning, Water Quality Monitoring, Water Quality Index (WQI), Sustainable Development Goal 6 (SDG 6), Predictive Analytics

# 1. Introduction

Water is one of the most essential natural resources on Earth and is crucial for sustaining life. Human beings, plants, and animals all rely on it as a vital requirement for survival [1]. Beyond its biological necessity, water plays a central role in agriculture, aquaculture, and industry, making it indispensable for food production, economic growth, and societal development. Recognizing this reality, the United Nations has declared access to clean and safe drinking water a fundamental human right, enshrined in Sustainable Development Goal 6 (SDG 6: Clean Water and Sanitation). Despite this recognition, the world's freshwater resources are under immense and growing pressure [2]. Rapid urbanization, industrial expansion, intensive farming practices, and the escalating impacts of climate change have combined to cause pollution, shortages, and deteriorating water quality. The World Health Organization (WHO) has highlighted the severity of this challenge, revealing that over two billion people globally still depend on unsafe water sources. This dependency exposes a large population to severe health risks from waterborne diseases such as cholera, diarrhea, and typhoid [3]. These alarming statistics underscore the urgent need for robust, reliable, and efficient water quality monitoring systems to safeguard human health, protect ecosystems, and ensure sustainable water resource management [4].

Traditionally, water quality monitoring has relied heavily on conventional laboratory-based testing. This process involves manually collecting water samples from sources such as rivers, lakes, reservoirs, treatment plants, and distribution pipelines, followed by transportation, preservation, and detailed laboratory analysis. Such traditional methods are highly regarded for their precision and reliability since they provide detailed insights into the physical, chemical, and biological properties of water [5]. However, while accurate, this methodology has multiple limitations that hinder its application in modern contexts. Firstly, it is extremely time-consuming, involving several stages of preparation, transport, and testing, all of which delay contamination detection and response. Secondly, laboratory-based water testing is costly due to the need for specialized equipment, reagents, and trained personnel, making continuous large-scale monitoring financially impractical. Thirdly, conventional assessment lacks scalability and often neglects rural, remote, and hazardous areas where water quality issues can be particularly severe. Finally, these methods fail to offer real-time updates, which makes them inadequate for detecting sudden pollutant discharges such as chemical spills or sewage leaks. With increasing pressures on global water systems, traditional water assessment methods no longer meet the requirements of timely, automated, and scalable monitoring essential for modern water management [6].

In this context, the growing field of the Internet of Things (IoT) offers groundbreaking solutions for continuous water quality monitoring. IoT refers to interconnected networks of smart devices and sensors capable of detecting, processing, and transmitting data through the internet. Within the domain of water monitoring, IoT-based systems utilize low-cost yet reliable sensors to measure parameters such as pH, turbidity, total dissolved solids (TDS), temperature, dissolved oxygen (DO), and electrical conductivity. Each parameter provides critical insights into different aspects of water health. For example, pH indicates acidity or alkalinity, turbidity reflects suspended contaminants, and TDS provides information on mineral and salt concentrations. Similarly, dissolved oxygen is vital for aquatic ecosystems, while conductivity provides information about ionic and salinity levels. Small, affordable microcontrollers such as ESP32 and NodeMCU are widely used to interface with these sensors and handle real-time communication via Wi-Fi or long-range protocols such as LoRaWAN. Data collected from the sensors is transmitted directly to cloud platforms such as ThingSpeak, AWS IoT, or Google Cloud IoT Core for further processing. These cloud services enable real-time data storage, visualization, analytics, and remote access for researchers and officials through web dashboards or mobile applications. Critically, IoT-integrated systems provide the ability to continuously monitor water quality in difficult-to-reach or hazardous areas, equipping stakeholders with near real-time information that allows rapid responses to contamination events [7,8].

While IoT systems revolutionize the way water quality data is collected by providing continuous and cost-effective monitoring, their true potential is realized when combined with machine learning (ML). Machine learning enhances IoTenabled monitoring by analyzing massive volumes of historical and real-time data to uncover hidden trends, identify anomalies, and forecast future conditions. Unlike rule-based systems, ML models can process complex, non-linear, and high-dimensional data with exceptional accuracy. This serves several crucial applications. For example, predictive analytics using ML can forecast near-term changes in parameters such as pH or turbidity, enabling proactive water management. Seasonal and long-term trend analysis helps identify recurring pollution patterns tied to agricultural runoff or industrial discharges. Anomaly detection systems recognize abnormal deviations from standard water quality levels, such as sudden chemical spills, in real time and send alerts for immediate response [9]. Among the popular ML approaches applied in water monitoring, Random Forests serve as a robust ensemble-based method for both classification and regression tasks, while Support Vector Machines (SVM) excel in anomaly detection particularly in high-dimensional parameter spaces. Gradient boosting models like XGBoost also deliver high predictive accuracy with structured data, and advanced deep learning models such as Long Short-Term Memory (LSTM) networks are especially useful because they capture temporal dependencies in time-series water data. Applications of such models include the calculation of the Water Quality Index (WQI), a simplified metric that compresses complex multi-parameter data into a single score categorized into labels such as "excellent," "good," or "unsafe." By combining IoT and ML, water monitoring systems gain the ability not only to measure and record water quality but also to predict and prevent contamination [10].

The potential of combining IoT and machine learning for water quality monitoring has already captured attention in research and practical applications. Recent studies and prototypes have shown that low-cost sensor networks connected to cloud-based systems can effectively provide real-time water data at significantly reduced expense compared to conventional methods. Deep learning approaches, particularly LSTM networks, have demonstrated strong performance in forecasting time-dependent water parameters such as turbidity and dissolved oxygen. Research has also emphasized anomaly detection approaches as a way to improve early warning systems, ensuring communities can respond to waterborne health risks before they escalate. Pilot projects funded by governments, industries, and educational institutions are increasingly exploring intelligent water management solutions, particularly within the wider context of smart city initiatives [11]. Despite these promising advances, significant challenges remain for widespread adoption. For instance, sensor calibration and maintenance are critical for accurate long-term usage, as environmental conditions often degrade sensor performance. Rural and low-connectivity areas face challenges where internet access is either limited or absent, necessitating hybrid communication technologies like LoRaWAN, GSM, or satellite IoT. Machine learning algorithms, while powerful, depend on large datasets for effective training, yet such datasets are often incomplete or unavailable in resource-constrained settings. Further concerns include data privacy, potential cyberattacks on connected infrastructure, and the affordability of sophisticated monitoring systems in developing nations. These barriers highlight

the importance of continued research, technical innovation, and supportive governance for scaling IoT-ML-based water management systems [12].

Building on these opportunities and challenges, the present research proposes a comprehensive, economical, and scalable framework for intelligent water quality monitoring by integrating IoT-based sensing with machine learning analysis. The first goal is to design and deploy a sensor-based IoT system utilizing affordable devices such as ESP32 and NodeMCU connected with cost-efficient sensors for parameters like pH, turbidity, TDS, temperature, and DO. These systems will be responsible for continuous data collection and transmission. The second objective is to integrate this IoT sensor system with cloud storage and visualization platforms such as ThingSpeak, enabling structured storage, real-time visualization, and mobile accessibility [13]. The third goal involves developing advanced machine learning models-including Random Forest, SVM, XGBoost, and LSTM-that can forecast key water parameters, compute Water Quality Index scores, and perform anomaly detection. The fourth objective is to comprehensively evaluate the system's performance against key benchmarks such as prediction accuracy, operational scalability, and cost-efficiency in comparison with time-consuming and expensive laboratory testing methods. Finally, the overall aim of this research is to demonstrate the broader societal value of combining IoT and ML for securing water safety, protecting public health, safeguarding ecosystems, and enabling smarter infrastructure in both urban and rural regions. By addressing these goals, the framework seeks to provide timely, predictive, and actionable insights into water management, thereby reducing the risks of waterborne diseases while ensuring sustainable resource utilization [14].

Water is a critical element for life, food security, and economic development, yet global water resources are under severe strain due to multiple human and environmental pressures. Conventional water quality monitoring methods, while accurate, are inadequate for today's needs in terms of scalability, speed, and cost-effectiveness. IoT-based water monitoring with real-time sensors and cloud integration provides a transformative pathway to overcome these limitations. The integration of machine learning adds a powerful analytical layer for predicting water quality changes, detecting anomalies, and computing indices such as WQI for simplified public communication. Existing literature has already demonstrated the feasibility and advantages of such combined frameworks, while also identifying current challenges that must be addressed for real-world deployment. The research proposed here builds on these insights by developing and testing a scalable IoT-ML-based system aimed at practical applications in diverse contexts. Ultimately, this approach aligns with global sustainability goals by fostering clean and safe water access, protecting communities from health risks, and ensuring a sustainable and intelligent future for water management [15].

# 2. Related Works

The use of Internet of Things (IoT) technologies combined with machine learning (ML) for water quality monitoring has rapidly evolved into an important interdisciplinary research field over the past decade. With the availability of affordable hardware, efficient communication protocols, scalable cloud storage, and intelligent data analytics frameworks, researchers and practitioners have been able to design systems that provide continuous surveillance of water bodies and early warnings of contamination events. Numerous reviews and survey papers highlight how IoT and ML integration not only advances environmental monitoring but also significantly contributes to achieving United Nations Sustainable Development Goal 6 (SDG 6), which emphasizes clean water and sanitation for all. Central to most studies is the layered pipeline model commonly used in IoT systems: sensors deployed in the environment collect raw data, this data is transmitted through wireless protocols to cloud platforms, stored and preprocessed, and finally analyzed by machine learning algorithms to derive actionable insights which are visualized through dashboards or automated alerts [16].

IoT-based sensing and hardware platforms are at the foundation of intelligent water monitoring systems. Most existing solutions rely on inexpensive electrochemical and optical sensors designed to measure critical physicochemical parameters such as pH, turbidity, dissolved oxygen (DO), total dissolved solids (TDS), temperature, and electrical conductivity (EC) [17]. These parameters are both fundamental to assessing potability and useful for ecological evaluation. For example, high turbidity often indicates the presence of suspended solids that can harbor pollutants, while dissolved oxygen illustrates the health of aquatic ecosystems. The sensors are typically coupled with low-cost and energy-efficient microcontrollers such as ESP32 and NodeMCU (ESP8266), both of which offer embedded wireless connectivity (Wi-Fi/Bluetooth) and are widely used in prototyping. Beyond prototyping, researchers often employ more powerful development boards coupled with long-range communication modules, depending on requirements for scalability or deployment in remote regions. On the software side, ThingSpeak is widely used in academia and for proofof-concept studies due to its simplicity, ease of integration with MATLAB, and suitability for quick prototyping [18]. However, when larger-scale deployments are necessary, platforms such as AWS IoT, Azure IoT, or privately managed servers are chosen, as they provide scalability, security features, and advanced processing capabilities. Despite these advancements, sensor drift, biofouling, and gradual calibration loss remain common challenges, prompting researchers to experiment with automatic drift compensation, periodic recalibration protocols, and hybrid systems that validate low-cost sensors against higher-grade instruments at regular intervals. These efforts improve the reliability of data streams in longterm deployments [19].

Beyond hardware, the issue of data quality is a critical concern and has a direct influence on the performance of machine learning models deployed in water quality analytics. Raw sensor data, especially when collected continuously in natural conditions, often contains noise spikes, abrupt missing values, or gradual drift in measurements. Preprocessing steps are therefore necessary to minimize these artifacts. Some works employ noise reduction methods such as moving averages, exponential smoothing, or median filters to stabilize measurement signals. Missing values can be imputed using interpolation techniques or more advanced machine learning methods such as k-nearest neighbors imputation [20]. Outlier detection is also important in differentiating true contamination events from faulty sensor readings. A further recommended practice is to cross-validate IoT measurements against laboratory-grade reference values wherever possible, ensuring higher consistency in data quality. A widely used metric in this context is the Water Quality Index (WQI), which consolidates multiple sensor parameters into a single score categorized into interpretable classifications such as "excellent," "good," "poor," or "unfit." The WQI has been particularly popular since it provides an easily understood way for both experts and non-experts to interpret water conditions. Some studies attempt to predict WQI directly from raw sensor readings using regression or classification models, while others first predict individual water quality parameters before calculating the final index score. In both cases, the WQI serves as a unifying framework for assessing potability and safety [21].

Machine learning is central to transforming raw IoT data into meaningful predictions. Early works in water monitoring began with classical supervised algorithms such as Random Forest (RF), Support Vector Machines (SVMs), and gradient boosting variants like XGBoost and LightGBM. These models are well suited for tabular datasets extracted from IoT streams and demonstrate robustness to noisy or partially missing data. Random Forests in particular have been used extensively for classifying water safety categories based on WQI scores, while SVMs have been applied for identifying anomalies or contamination events where data dimensionality is high. As the field matured, the availability of time-series sensor records motivated the adoption of deep learning methods [22]. Recurrent neural networks (RNNs), especially Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) architectures, became favored for predicting changes in pH, dissolved oxygen, and turbidity over short horizons ranging from hours to days. These models excel at capturing temporal dependencies that simpler ML methods cannot. Further advances incorporated Temporal Convolutional Networks (TCNs), which can model long temporal sequences with competitive efficiency compared to LSTMs. Hybrid frameworks have also become common, blending signal decomposition techniques with predictive modeling. One example is the RF-CEEMD-LSTM (Random Forest + Complete Ensemble Empirical Mode Decomposition + LSTM) framework, which achieves superior accuracy by first decomposing highly variable water signals, then applying an ensemble classifier, and finally conducting temporal sequence learning. Such hybrid approaches are reported to produce improved forecasts of river water quality indicators with minimal error margins when compared against traditional statistical models [23].

Another crucial element of water monitoring systems is anomaly detection, with the aim of developing early warning systems for pollution and contamination events. The core challenge in this research area is the limited availability of labeled contamination data. When event labels are available, supervised methods such as classification trees or ensemble classifiers can be trained to identify contamination scenarios directly. However, in most real-world contexts, labeled contamination is rare, making unsupervised anomaly detection techniques more widely applicable. Popular algorithms include autoencoders that reconstruct normal sensor behavior and detect deviations, isolation forests that detect outliers in multivariate datasets, one-class SVMs, and approaches based on prediction error, where anomalies are flagged when forecasts differ significantly from observed values. Since false positives can undermine trust in early warning systems, multiple studies recommend adopting hybrid strategies-such as combining anomaly detection outputs across multiple parameters or sensors, integrating statistical thresholds with ML-based alerts, and involving expert human validation in high-stakes decision-making. Together, these approaches yield more robust and context-aware early warning systems capable of improving public health responses [24].

The development and evaluation of these algorithms are highly dependent on the availability of data. Public datasets provide essential resources for benchmarking and reproducibility in the research community. T widely used datasets include Kaggle's water potability dataset, the UCI repository's water quality datasets, and European WaterBase collections. These resources allow cross-comparison between algorithms and provide standardized testbeds for evaluating performance. However, issues of domain shift are repeatedly reported. Models trained on one dataset often fail when deployed in different regions or under varying sensor configurations, limiting their generalizability. To overcome this limitation, researchers have explored transfer learning methods, domain adaptation techniques, and synthetic data generation or augmentation to simulate diverse deployment scenarios. Locally collected datasets continue to be the gold standard for achieving high reliability in predictive models [25].

Underlying all IoT-ML water monitoring systems are the cloud platforms that support data collection, storage, and analytics. ThingSpeak remains popular in academic environments and hobbyist projects due to its simplicity and MATLAB compatibility [26]. For large-scale or mission-critical deployments, however, cloud providers like AWS IoT Core and Microsoft Azure IoT dominate due to their scalability, integration with machine learning pipelines, and robust data security protocols. Some studies employ time-series databases such as InfluxDB or TimescaleDB to efficiently store and query large volumes of continuous water data. Case studies highlight end-to-end workflow examples, typically

following the pipeline "sensor node (ESP32/NodeMCU)  $\rightarrow$  MQTT communication  $\rightarrow$  cloud ingestion  $\rightarrow$  machine learning model  $\rightarrow$  dashboard or automated alert system." These implementations often account for practical field challenges, including protective waterproof housings for sensors, solar energy harvesting for off-grid deployment, and over-the-air firmware updates to ensure system maintainability [27].

This literature demonstrates remarkable progress in integrating IoT hardware, cloud platforms, and ML algorithms for intelligent water monitoring. The trajectory of research has shifted from proof-of-concept prototypes toward sophisticated hybrid models and scalable architectures capable of supporting city-level or regional deployments [28]. Yet, recurring challenges remain, particularly around sensor reliability, calibration, and data generalizability. Ensuring accurate long-term sensing while minimizing maintenance is one of the most pressing technical challenges. On the machine learning side, models must adapt to diverse datasets and deployment conditions without extensive retraining. Furthermore, ensuring cybersecurity, reducing false alarms, and managing operational costs are important for future adoption. Despite these challenges, the literature reflects a strong consensus on the transformative potential of IoT-ML frameworks. By providing predictive, continuous, and scalable water monitoring solutions, they not only enable more intelligent environmental governance but directly address public health concerns and sustainability goals across regions worldwide [29].

# 3. IoT-Enabled Smart Water Quality Monitoring and Machine Learning-Based Predictive Analysis System

This methodology leverages the integration of Internet of Things (IoT) technologies and Machine Learning (ML) algorithms to deliver a robust, real-time water quality monitoring and predictive analysis system. This architecture is carefully structured to address challenges associated with conventional water monitoring-namely, the lack of immediate feedback, scalability issues, and labor-intensive processes-by employing automation, wireless communication, and intelligent analytics. The workflow is organized into two primary phases: IoT-based data acquisition and machine learning-driven predictive analysis, as mentioned in figure 1. Together, these phases establish a comprehensive solution supporting both continuous surveillance and informed, rapid decision-making regarding water safety.

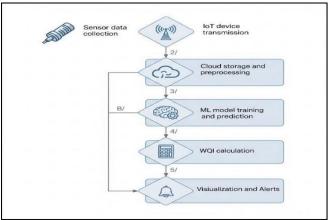


Figure 1. Workflow of the IoT and machine learning-based water quality monitoring system, illustrating data collection, transmission, processing, prediction, WQI calculation, and alert visualization.

The first phase centers around an efficient system architecture built on widely accessible and cost-effective hardware. Key sensors deployed in the system include pH sensors, TDS (Total Dissolved Solids) sensors, turbidity sensors, temperature sensors, and dissolved oxygen sensors. These devices are selected for their ability to monitor the most vital physicochemical properties of water, each providing a window into specific aspects of water quality such as acidity, mineral content, particulate contamination, temperature-driven microbial activity, and aquatic oxygenation. Sensor nodes are built around microcontrollers like the ESP32 or NodeMCU (ESP8266), platforms favored in both academic research and industrial prototyping due to their affordability, low power requirements, and built-in Wi-Fi/Bluetooth connectivity. Communication is established via Wi-Fi, enabling seamless integration with cloud platforms such as ThingSpeak or AWS IoT. These cloud services provide the backbone for centralized data storage, processing, and visualization, allowing stakeholders to access the system from remote locations using dashboards or connected devices.

The data acquisition layer forms the initial interface between the physical environment and the digital system. In this layer, the deployed sensor array continuously samples key water parameters. These readings are captured at fixed intervals-programmable depending on the use case-by the IoT microcontroller. Upon collection, the microcontroller transmits the recorded data in real time, or near-real time, to the designated cloud server. Cloud storage ensures that the time-series dataset is persistently available for synthesis, retrospective analytics, modeling, and visualization, eliminating the risk of data loss and supporting historical comparison for trend analysis.

Once data is ingested by the cloud, it undergoes thorough preprocessing to ensure reliability and analytical utility. Sensor data is commonly subject to imperfections such as transient glitches, noise spikes, missing records, or gradual sensor drift.

Preprocessing routines commence with the removal or correction of missing values and noise artifacts to prevent them from skewing downstream analyses. This may involve techniques like moving averages to smooth noise or interpolation to address gaps. Next, normalization or standardization converts sensor readings into scale-compatible values, critical for multi-parameter modeling and for reducing bias in algorithms sensitive to absolute value ranges. A central component of preprocessing is the calculation of the Water Quality Index (WQI), which aggregates multiple water quality parameters into an interpretable composite score. This index facilitates rapid classification of water samples-including labels such as "safe" or "unsafe"-and simplifies communication for both technical personnel and the general public.

The subsequent stage is the machine learning-based predictive analysis, where the intelligent core of the platform is realized. The system explores a suite of advanced ML models: Random Forest, Support Vector Machines (SVM), Extreme Gradient Boosting (XGBoost), and Long Short-Term Memory (LSTM) networks for time-series prediction. In the training phase, historical datasets are leveraged to build robust models, with each algorithm tuned for maximal predictive accuracy. Performance metrics such as RMSE (Root Mean Square Error), Accuracy, and F1-score are computed during the testing phase to rigorously evaluate model performance across prediction tasks. These tasks include forecasting future water quality values (e.g., projecting pH, turbidity, or DO levels several hours ahead), anomaly detection (e.g., alerting authorities to an abrupt increase in turbidity which may signify contamination), and classification of water as drinkable or non-drinkable. The combination of regression, classification, and anomaly detection supports both long-term trend analysis and immediate contamination warning.

The final operational layer is alerting and visualization, powered by integrated cloud dashboards provided by platforms like ThingSpeak or AWS IoT. These dashboards support the real-time display of ongoing sensor values and model predictions, supplying decision-makers and the public with current insights into water quality conditions. Moreover, the system is designed to deliver automated alerts through channels such as SMS or email whenever unsafe water conditions are forecasted, thus ensuring timely response to potential threats. The inclusion of graphical reporting tools further facilitates swift and effective interpretation for diverse stakeholder groups, including water authority officials, researchers, and community representatives.

The proposed methodology represents a multilayered approach that capitalizes on IoT's capacity for automated real-time data capture and ML's power for sophisticated prediction, classification, and anomaly detection. Through an architecture comprising robust sensing, wireless communication, scalable cloud integration, rigorous data preprocessing, model-driven analytics, and intuitive alerting, the system delivers a comprehensive toolkit for real-time water quality surveillance and public health protection. This methodology is highly adaptable, scalable to urban and rural contexts, and capable of continuous self-improvement as more data is acquired and algorithmic performance is refined. The provided workflow diagram visually encapsulates the key operational steps-sensor collection, device transmission, cloud processing, ML analytics, WQI calculation, and alerting-underscoring the end-to-end automation and intelligence of the proposed solution.

# 4. Results and Discussions

The results summarized in the provided table 1 demonstrate the effectiveness of various machine learning (ML) models in water quality monitoring applications across different datasets. Three distinct datasets were examined: self-collected sensor data from low-cost IoT devices, the publicly available Kaggle Water Potability dataset, and the UCI Water Quality dataset. Each dataset leveraged a unique combination of water quality features, ML algorithms, and evaluation metrics to assess predictive accuracy, precision, recall, and F1-score. These metrics collectively indicate that integrating ML with IoT and traditional data sources can significantly enhance real-time water quality prediction and contamination detection, supporting early intervention and improved public health outcomes.

**Table 1**. Comparison of machine learning model performance on various water quality datasets including sensor-based IoT data and publicly available datasets.

<b>Dataset Source</b>	Features Used	No. of Records	ML Model	Accuracy (%)	Precision	Recall	F1-Score
IoT Sensor Data (Self-Collected via ESP32 + pH, TDS, Turbidity, Temp, DO sensors)	pH, TDS, Turbidity, Temperature, DO	500+	Random Forest	94.6	0.93	0.92	0.92
Kaggle-Water Potability Dataset	pH, Hardness, Solids, Chloramines, Sulfates, Conductivity, Organic Carbon, Trihalomethanes, Turbidity	3276	Logistic Regression	87.2	0.85	0.84	0.84
UCI Water Quality Dataset	pH, DO, Conductivity, BOD, Nitrate, Turbidity, Iron	1980	Artificial Neural Network (ANN)	92.8	0.91	0.9	0.9

Starting with the IoT sensor data collected via ESP32 microcontrollers paired with pH, TDS, turbidity, temperature, and dissolved oxygen (DO) sensors, the Random Forest model achieved the highest performance with an accuracy of 94.6%. Additionally, the precision (0.93), recall (0.92), and F1-score (0.92) all indicate strong model reliability in correctly identifying both safe and unsafe water conditions. These results underscore the practical feasibility of low-cost, sensor-driven systems to generate high-fidelity datasets when combined with robust ensemble learning approaches. The Random Forest classifier's ability to handle noisy inputs and nonlinear relationships among multiple water parameters likely contributed to the model's superior performance.

In contrast, the Kaggle Water Potability dataset-which contains a broader list of features including pH, hardness, solids, chloramines, sulfates, conductivity, organic carbon, trihalomethanes, and turbidity-was addressed using Logistic Regression. While achieving an accuracy of 87.2%, this model demonstrated slightly lower precision (0.85), recall (0.84), and F1-score (0.84) compared to the sensor data scenario. This relative decrease may be attributed to the larger, more heterogeneous dataset and the inherently linear assumptions of Logistic Regression. Nonetheless, the findings suggest that classical statistical models can still provide useful baselines for water quality classification when the problem space is less complex or when computational simplicity is prioritized.

The UCI Water Quality dataset involved even greater diversity of parameters, including pH, DO, conductivity, biological oxygen demand (BOD), nitrate, turbidity, and iron concentration. Here, an Artificial Neural Network (ANN) was employed, which delivered accuracy of 92.8%, and high precision (0.91), recall (0.90), and F1-score (0.90). The neural network's ability to model complex, nonlinear relationships helped capture intricate patterns in the multidimensional data, outperforming the logistic model on the Kaggle dataset though slightly trailing the Random Forest's performance on IoT sensor data. This confirms ANNs' suitability for capturing complex environmental interactions, particularly when ample training data is available.

The comparative analysis, as in figure 2, highlights several key insights. Firstly, model performance depends on the nature and quality of input features as well as the dataset size and diversity. Sensor data collected through IoT systems, though smaller in volume (500+ records), delivers real-time, continuous, and highly relevant parameters for local monitoring, enabling machine learning methods like Random Forest to achieve excellent predictive results. Public datasets provide larger record counts and broader geographic coverage but are subject to heterogeneity and missing data issues, necessitating more sophisticated approaches such as neural networks to model underlying complexities. Secondly, the choice of ML algorithm significantly impacts outcomes. Ensemble methods like Random Forest are robust against overfitting and noise in relatively small datasets, while ANNs excel when complex nonlinearities dominate but require larger datasets and careful tuning. Logistic regression remains a useful, interpretable choice for baseline models but may underperform in multifaceted environments.

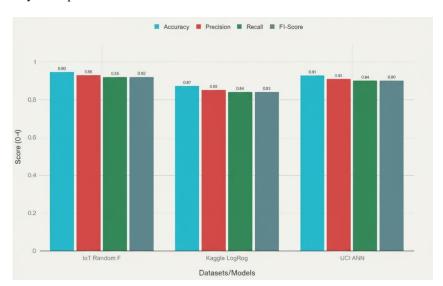


Figure 2. Performance comparison of Random Forest, Logistic Regression, and Artificial Neural Network models on IoT Sensor Data, Kaggle Water Potability Dataset, and UCI Water Quality Dataset, showing their Accuracy, Precision, Recall, and F1-Score.

From a practical standpoint, these results validate the proposed framework's core premise: integrating IoT-based sensing with advanced machine learning yields effective and scalable solutions for water quality monitoring. The high precision and recall scores for detecting unsafe water are particularly important, reflecting model sensitivity to contamination events and reducing false negatives that could jeopardize public health. Moreover, forecasting future water quality metrics and classifying water safety status can empower decision-makers with timely, actionable information. Integration with cloud platforms for visualization, alert generation, and historical trend analysis offers further operational benefits.

Challenges remain, however. The IoT dataset, while promising, is limited by its relatively small size, and larger deployments are necessary to rigorously validate model generalizability across different water bodies and seasonal variations. Additionally, pre-processing steps such as noise filtering, missing data imputation, and normalization critically affect model performance and must be automated for real-world deployment. The diversity of water quality parameters across datasets also calls for adaptable architectures that can dynamically incorporate new sensors or data types. Finally, data privacy, cybersecurity, and the cost of deploying and maintaining sensor networks must be addressed to ensure sustainability and stakeholder trust.

The comparison of ML models on IoT sensor data and established public datasets confirms that ensemble learning and neural network approaches deliver strong predictive capabilities across varying data environments. The results underscore the promise of IoT-ML integrated solutions for continuous, real-time monitoring and early detection of water quality changes. Future work should focus on increasing dataset sizes, enhancing model adaptability, and optimizing end-to-end system integration to facilitate widespread adoption in urban, rural, and developing world contexts. These efforts will aid efforts to safeguard water resources, protect public health, and promote sustainable environmental management.

#### 5. Conclusion and Future Works

The conclusion of this research emphasizes the successful development of a scalable and real-time water quality monitoring system that integrates IoT technology with machine learning algorithms. The system's multi-layer architecture effectively captures critical water quality parameters such as pH, total dissolved solids, turbidity, temperature, and dissolved oxygen through a network of sensors connected to microcontrollers. This data is transmitted to cloud platforms where it is preprocessed to maintain data integrity. The core contribution lies in the machine learning-based predictive analysis layer that leverages models like Random Forest, SVM, XGBoost, and LSTM for forecasting water quality, anomaly detection, and potability classification through an automatically computed Water Quality Index (WQI). This real-time analytical framework, complemented by user-friendly dashboards and automated alerts, equips water management authorities with actionable intelligence to promptly respond to contamination threats and maintain public safety.

This research presents a holistic and intelligent approach that transcends traditional manual monitoring by enabling proactive water resource management. It highlights how continuous data acquisition integrated with advanced machine learning fosters both immediate operational awareness and predictive foresight for safeguarding water environments. The study illustrates considerable promise for widespread application in both urban and rural contexts, addressing the pressing need for sustainable water governance and public health protection.

Future work should explore the integration of additional sensor types to expand the spectrum of detectable pollutants, including chemical contaminants and microbial indicators, enhancing the system's comprehensiveness. Further development of more sophisticated deep learning architectures and hybrid models could improve prediction accuracy and early anomaly detection capabilities, especially in complex environmental scenarios. Moreover, conducting a detailed cost-benefit analysis for large-scale deployment will be critical for real-world adoption, addressing economic feasibility alongside technical performance. Lastly, ensuring data security and privacy and improving sensor calibration using automated techniques remain important avenues for ongoing research to increase system reliability and trustworthiness. This approach sets a strong foundation for intelligent water quality monitoring, fostering sustainable water use, environmental protection, and improved public health outcomes.

# References

- [1] Singh, M., Sahoo, K. S., & Nayyar, A. (2022). Sustainable IoT solution for freshwater aquaculture management. IEEE Sensors Journal, 22, 16563-16572. https://doi.org/10.1109/JSEN.2022.3188639
- [2] Bhargavi, K., Sowmya, K. V., Ajay, P., Saketh, D., & Ravindhar, B. (2023). Water quality system for aquaculture using IoT.

  International Research Journal of Modern Engineering and Technology Sciences, 5, 21-23.

  https://doi.org/10.56726/IRJMETS35272
- [3] Baena-Navarro, R., Vergara-Villadiego, J., Carriazo-Regino, Y., Crawford-Vidal, R., & Barreiro-Pinto, F. (2024). Challenges in implementing free software in small and medium-sized enterprises in the city of Montería: A case study. Bulletin of Electrical Engineering and Informatics, 13, 586-597. https://doi.org/10.11591/eei.v13i1.6710
- [4] Carriazo-Regino, Y., Baena-Navarro, R., Torres-Hoyos, F., Vergara-Villadiego, J., & Roa-Prada, S. (2022). IoT-based drinking water quality measurement: Systematic literature review. Indonesian Journal of Electrical Engineering and Computer Science, 28, 405-418. https://doi.org/10.11591/ijeecs.v28.i1.pp405-418
- [5] Pinedo-López, J., Baena-Navarro, R., Durán-Rojas, N., Díaz-Cogollo, L., & Farak-Flórez, L. (2024). Energy transition in Colombia: An implementation proposal for SMEs. Sustainability, 16, 7263. https://doi.org/10.3390/su16177263
- [6] Vidal-Durango, J., Baena-Navarro, R., & Therán-Nieto, K. (2024). Implementation and feasibility of green hydrogen in Colombian kitchens: An analysis of innovation and sustainability. Indonesian Journal of Electrical Engineering and Computer Science, 34, 726-744. https://doi.org/10.11591/ijeecs.v34.i2.pp726-744

- [7] Kushwaha, S. (2023). A futuristic perspective on artificial intelligence. In Proceedings of the IEEE OPJU International Technology Conference on Emerging Technologies For Sustainable Development (pp. 1-6). O.P. Jindal University, Raigarh, Chhattisgarh, India.
- [8] Petkovski, A., Ajdari, J., & Zenuni, X. (2021). IoT-based solutions in aquaculture: A systematic literature review. In Proceedings of the 2021 44th International Convention on Information, Communication and Electronic Technology (MIPRO) (pp. 1358-1363). IEEE.
- [9] Abinaya, T., Ishwarya, J., & Maheswari, M. (2019). A novel methodology for monitoring and controlling of water quality in aquaculture using Internet of Things (IoT). In Proceedings of the 2019 International Conference on Computer Communication and Informatics (ICCCI) (pp. 1-4). IEEE.
- [10] Ahmed, M., Rahaman, M. O., Rahman, M., & Abul Kashem, M. (2019). Analyzing the quality of water and predicting the suitability for fish farming based on IoT in the context of Bangladesh. In Proceedings of the 2019 International Conference on Sustainable Technologies for Industry 4.0 (STI) (pp. 1-5). IEEE.
- [11] Hu, W.-C., Chen, L.-B., Huang, B.-K., & Lin, H.-M. (2022). A computer vision-based intelligent fish feeding system using deep learning techniques for aquaculture. IEEE Sensors Journal, 22, 7185-7194. https://doi.org/10.1109/JSEN.2022.3151777
- [12] Li, L., Jiang, P., Xu, H., Lin, G., Guo, D., & Wu, H. (2019). Water quality prediction based on recurrent neural network and improved evidence theory: A case study of Qiantang River, China. Environmental Science and Pollution Research, 26, 19879-19896. https://doi.org/10.1007/s11356-019-05116-y
- [13] Giao, N. T., Van Cong, N., & Nhien, H. T. H. (2021). Using remote sensing and multivariate statistics in analyzing the relationship between land use pattern and water quality in Tien Giang Province, Vietnam. Water, 13, 1093. https://doi.org/10.3390/w13081093
- [14] Rasheed Abdul Haq, K. P., & Harigovindan, V. P. (2022). Water quality prediction for smart aquaculture using hybrid deep learning models. IEEE Access, 10, 60078-60098. https://doi.org/10.1109/ACCESS.2022.3180482
- [15] Kushwaha, S. (2023). Review on artificial intelligence and human computer interaction. In Proceedings of the IEEE OPJU International Technology Conference on Emerging Technologies For Sustainable Development (pp. 1-6). O.P. Jindal University, Raigarh, Chhattisgarh, India.
- [16] Syed Taha, S. N., Abu Talip, M. S., Mohamad, M., Azizul Hasan, Z. H., & Tengku Mohmed Noor Izam, T. F. (2024). Evaluation of LoRa network performance for water quality monitoring systems. Applied Sciences, 14, 7136. https://doi.org/10.3390/app14167136
- [17] Suriasni, P. A., Faizal, F., Hermawan, W., Subhan, U., Panatarani, C., & Joni, I. M. (2024). IoT water quality monitoring and control system in moving bed biofilm reactor to reduce total ammonia nitrogen. Sensors, 24, 494. https://doi.org/10.3390/s24020494
- [18] Wang, X., Li, Y., Qiao, Q., Tavares, A., & Liang, Y. (2023). Water quality prediction based on machine learning and comprehensive weighting methods. Entropy, 25, 1186. https://doi.org/10.3390/e25081186
- [19] Dupont, C., Cousin, P., & Dupont, S. (2018). IoT for aquaculture 4.0 smart and easy-to-deploy real-time water monitoring with IoT. In Proceedings of the 2018 Global Internet of Things Summit (GIoTS) (pp. 1-5). IEEE.
- [20] Teixeira, R. R., Puccinelli, J. B., Poersch, L., Pias, M. R., Oliveira, V. M., Janati, A., & Paris, M. (2021). Towards precision aquaculture: A high performance, cost-effective IoT approach. arXiv. https://doi.org/10.48550/arXiv.2105.11493
- [21] Suhaili, W., Aziz, M., Ramlee, H., Patchmuthu, R., Shams, S., Mohamad, I., Isa, M., & Nore, B. (2023). IoT aquaculture system for sea bass and giant freshwater prawn farming in Brunei. In Proceedings of the 2023 13th International Conference on Information Technology in Asia (CITA) (pp. 60-65). IEEE.
- [22] Nayak, S., Mantri, J. K., & Swain, P. K. (2019). Design and performance analysis of rural aquaculture ponds using IoT. International Journal of Recent Technology and Engineering, 8, 3078-3081. https://doi.org/10.35940/ijrte.B2086.078219
- [23] Cheng, L., Chen, Y.-Q., Zhang, S.-X., & Zhang, S. (2024). Quantum approximate optimization via learning-based adaptive optimization. Communications Physics, 7, 83. https://doi.org/10.1038/s42005-024-01577-x
- [24] Yan, X., Zhang, T., Du, W., Meng, Q., Xu, X., & Zhao, X. (2024). A comprehensive review of machine learning for water quality prediction over the past five years. Journal of Marine Science and Engineering, 12(1), 159. https://doi.org/10.3390/jmse12010159
- [25] Kaddoura, S. (2022). Evaluation of machine learning algorithm on drinking water quality for better sustainability. Sustainability, 14(18), 11478. https://doi.org/10.3390/su141811478
- [26] Wei, T. Y., Tindik, E. S., Fui, C. F., Haviluddin, H., & Hijazi, M. H. A. (2023). Automated water quality monitoring and regression-based forecasting system for aquaculture. Bulletin of Electrical Engineering and Informatics, 12(1), 570-579. https://doi.org/10.11591/eei.v12i1.4464
- [27] Kumar, P., Tiwari, P., & Reddy, U. S. (2023). Estimating fish weight growth in aquaponic farming through machine learning techniques. In Proceedings of the 2023 3rd International Conference on Intelligent Technologies (CONIT) (pp. 1-7). IEEE.
- [28] Alashjaee, A. M., Kushwaha, S., Alamro, H., Hassan, A. A., Alanazi, F., & Mohamed, A. (2024). Optimizing 5G network performance with dynamic resource allocation, robust encryption, and Quality of Service (QoS) enhancement. PeerJ Computer Science, 10, e2567. https://doi.org/10.7717/peerj-cs.2567.
- [29] Ahmed, A. A. M., Jui, S. J. J., Chowdhury, M. A. I., Ahmed, O., & Sutradha, A. (2023). The development of dissolved oxygen forecast model using hybrid machine learning algorithm with hydro-meteorological variables. Environmental Science and Pollution Research, 30, 7851-7873. https://doi.org/10.1007/s11356-022-22601-z