

Machine Learning Model for Predicting Human Health Implications of Genetically Modified Foods

NBAAKEE, LEBARI GOODDAY¹, OSAKI MILLER THOM-MANUEL²

^{1,2}*Department of Computer Science, Ignatius Ajuru University of Education, Rumuolumeni, Port Harcourt, Nigeria*

Abstract- *The safety of genetically modified foods (GMFs) remains a major public health and regulatory concern due to potential short- and long-term health implications. This study proposes a machine learning-based framework for predicting human health risks associated with genetically modified foods by leveraging artificial intelligence techniques. The system integrates genetic, proteomic, toxicological, and nutritional datasets to enable comprehensive risk assessment. Machine learning algorithms, including Random Forest, Support Vector Machines, and Gradient Boosting, were implemented to classify GM foods into low, moderate, and high health risk categories. Experimental results indicate that ensemble learning approaches outperform single-model methods, with Gradient Boosting achieving the highest predictive accuracy of approximately 94%, an F1-score of 0.93, and improved sensitivity to complex biological interactions. Feature importance analysis revealed that protein sequence similarity to known allergens, presence of toxic metabolites, and nutritional composition variations are the most significant predictors of potential health risks. The findings demonstrate that artificial intelligence can enhance the efficiency, accuracy, and interpretability of GM food safety assessments. The proposed framework provides a scalable decision-support tool for researchers and regulatory agencies, supporting evidence-based evaluation and continuous monitoring of genetically modified foods.*

Keywords: *Genetically Modified Foods; Machine Learning; Health Risk Prediction; Food Safety Assessment; Artificial Intelligence.*

I. INTRODUCTION

Genetically modified (GM) foods have been widely adopted as a means to enhance agricultural productivity, improve crop resilience, and strengthen food security in the face of growing population demands and climate change. By incorporating traits such as pest resistance, drought tolerance, and improved nutritional content, GM crops have become integral to modern agricultural practices. Despite the rigorous regulatory approval processes and safety assessments conducted by national and international agencies, public concern

regarding their potential long-term health implications. These concerns are fuelled by uncertainties surrounding allergenicity, toxicity, metabolic effects, and the cumulative impact of chronic exposure, particularly among vulnerable populations (Ghimire et al., 2023).

Traditional assessment methods for evaluating the safety of GM foods rely heavily on laboratory experiments, animal studies, and compositional analyses. While these approaches provide valuable insights, they are often time-consuming, costly, and limited in scope. More critically, they struggle to capture the complexity of human dietary patterns and the subtle, long-term health outcomes that may arise from prolonged consumption. Rare adverse events, nonlinear biological interactions, and population-specific effects are frequently overlooked, leaving gaps in risk characterization and fuelling ongoing debates about the safety of GM foods.

Recent advances in machine learning (ML) offer promising new opportunities to address these challenges. ML techniques excel at analysing large, complex, and heterogeneous datasets, uncovering hidden patterns, and modelling nonlinear relationships that traditional statistical methods may miss. By integrating diverse sources of biological, nutritional, and epidemiological data, ML models can provide more efficient and accurate predictions of potential health risks associated with GM foods. Furthermore, the incorporation of explainability tools, ensemble learning strategies, and optimization algorithms enhances both the reliability and transparency of these models, making them valuable not only for scientific research but also for regulatory decision-making and public communication.

This study builds on these advances by proposing a machine learning framework specifically designed to predict human health implications of GM foods. By leveraging multi-level feature selection, hybrid

ensemble learning, and risk-sensitive loss functions, the system aims to improve predictive accuracy while ensuring interpretability. In doing so, it seeks to bridge the gap between traditional safety assessments and modern computational approaches, offering a novel pathway toward more comprehensive, data-driven risk evaluation in the context of genetically modified food consumption.

Genetically modified (GM) foods have become a central part of modern agriculture, offering benefits such as higher yields, pest resistance, and improved nutritional content. However, despite these advantages, concerns remain about their long-term safety and potential health implications for humans. Current safety assessment frameworks, while scientifically rigorous, are not fully equipped to address the complexity and scale of GM food evaluation in today's rapidly evolving biotechnology landscape. This creates a pressing need for innovative approaches that can provide more predictive, scalable, and integrated assessments. One major limitation of existing approaches is scalability. Traditional toxicological and nutritional studies rely on extensive laboratory experiments and animal testing, which are resource-intensive, time-consuming, and costly. As the number of GM food varieties continues to grow, these methods struggle to keep pace with the demand for comprehensive safety evaluations. This bottleneck makes it difficult for regulators to assess new GM foods efficiently and promptly.

Another challenge lies in long-term risk prediction. Current assessment techniques often emphasize short-term outcomes such as immediate toxicity or allergenicity. They rarely capture chronic, cumulative, or multi-generational effects of GM food consumption. This leaves significant gaps in understanding the potential long-term health consequences, which are critical for ensuring consumer safety and maintaining public trust in biotechnology-driven food systems.

Data fragmentation further compounds the problem. Genetic, proteomic, and toxicological data are typically analysed in isolation, preventing the formation of a holistic view of how genetic modifications may influence protein expression, metabolic pathways, and toxicological responses in humans. Without integrated analysis, regulators and researchers risk overlooking subtle but important interactions that could affect human health.

To address these limitations, there is a clear need for predictive computational models capable of synthesising heterogeneous biological datasets. Such models could enable early detection of potential health risks, improve prioritisation of GM food varieties for detailed laboratory testing, and provide regulators with evidence-based insights to strengthen decision-making processes. Machine learning offers a transformative opportunity in this regard. By leveraging algorithms that excel at pattern recognition, data integration, and predictive analytics, machine learning models can unify diverse datasets, forecast long-term health implications, and scale efficiently across multiple GM food varieties.

Ultimately, the development of machine learning-based models for GM food safety assessment represents a critical advancement in modern food regulation. By overcoming the limitations of existing approaches, these models can enhance predictive capacity, improve scalability, and foster greater confidence in the safety of GM foods. This innovation not only supports regulatory agencies in safeguarding public health but also strengthens public trust in biotechnology and its role in shaping sustainable food systems.

II. RELATED LITERATURE

According to [1], Synthetic agrochemicals, particularly pesticides, have become indispensable in modern agricultural practices because they enhance crop productivity and protect against pests and diseases. Yet, their widespread use has raised serious concerns about environmental sustainability and human health. A large body of literature documents the detrimental effects of pesticide exposure, linking it to neurological disorders, cancers, respiratory complications, and metabolic dysfunctions. These risks are especially pronounced among agricultural workers and vulnerable populations who face higher levels of exposure. The complexity of these health outcomes is compounded by multiple exposure pathways, dietary intake, occupational contact, and environmental drift, as well as the difficulty in isolating causal relationships due to confounding factors and long latency periods.

Traditional epidemiological methods have provided valuable insights into agrochemical-related health risks, but they often struggle to capture the high-dimensional, nonlinear, and multifactorial nature of exposure-outcome relationships. This limitation has

prompted growing interest in machine learning (ML) as a tool for risk evaluation and predictive modelling. ML algorithms are well-suited to integrate diverse datasets, uncover hidden patterns, and improve predictive accuracy. Techniques such as mutual information gain and recursive feature elimination have been used to identify the most relevant predictors, while ensemble models like Random Forest, LightGBM, and CatBoost have demonstrated superior performance in handling complex data structures. Explainability methods, particularly SHAP, further enhance transparency by showing how individual features contribute to predictions, which is critical for regulatory acceptance and public trust.

Recent studies highlight the promise of combining advanced ML techniques with optimization strategies. Custom loss functions have been introduced to penalize false negatives, ensuring that severe health risks such as mortality are not overlooked. Optimization algorithms like Particle Swarm Optimization and Genetic Algorithms have been employed to fine-tune model parameters, yielding significant performance improvements. For example, hybrid ensemble models optimized with PSO and custom loss functions have achieved high levels of accuracy, precision, recall, and F1 scores, underscoring the potential of ML frameworks to outperform traditional approaches in health risk assessment.

Despite these advances, several gaps remain in the literature. Many ML applications in agrochemical health risk assessment are still exploratory, with limited external validation and generalizability across regions. The integration of mechanistic data, such as omics and microbiome profiles, is sparse, reducing the ability to link correlations to biological causation. Furthermore, uncertainty quantification and calibration metrics are often underreported, which undermines the reliability of predictions for clinical and regulatory decision-making. Real-world testing and integration with public health monitoring systems are also lacking, leaving a gap between theoretical promise and practical application.

The proposed system, “Machine Learning Model for Predicting Human Health Implications of Genetically Modified Foods,” seeks to address these gaps by extending the agrochemical ML framework into the GM food domain. Unlike simplistic GM versus non-GM classifications, this system models

exposures at the trait–crop–processing level, allowing for more granular and biologically meaningful analysis. It incorporates causal inference methods such as propensity scoring and structural causal models to strengthen causal claims, while multimodal integration of consumption records, omics data, microbiome profiles, and clinical outcomes bridges the gap between correlation and mechanism. Hybrid ensemble models optimized with custom loss functions emphasize recall for severe outcomes, and SHAP explanations are mapped to biological pathways to enhance interpretability. Uncertainty quantification through Bayesian ensemble and conformal prediction ensures calibrated risk estimates, while external validation across multi-regional datasets and integration with public health monitoring systems improve generalizability and real-world relevance.

Genetically modified (GM) foods, according to [2], are increasingly shaping the future of agriculture by offering solutions that enhance nutrition, sustainability, and resilience to environmental challenges. The literature emphasizes several benefits of GM crops, including improved pest resistance, nutrient enrichment, and contributions to reducing mycotoxin contamination. Beyond food security, GM technologies are also being explored for their roles in biofuel production and pharmaceutical development, underscoring their potential to address diverse global needs. These advantages position GM foods as a promising tool in combating malnutrition, supporting climate adaptation, and advancing medical innovation.

Studies highlight potential risks related to allergenicity, cancer development, reproductive health, and disruptions to gut microbiota. These uncertainties have fueled debates about the long-term safety of GM foods and the adequacy of current regulatory frameworks. Addressing these concerns requires robust detection and monitoring systems. Advanced molecular techniques such as PCR-based assays, immunoassays, and next-generation sequencing (NGS) have become central to identifying GM modifications and ensuring that unauthorized or unsafe variants are excluded from the food supply. Emerging technologies, particularly CRISPR-based diagnostics, promise even greater specificity, affordability, and accessibility, marking a significant step forward in molecular-level GM food detection.

This work consistently advocates for a multidisciplinary approach to GM food safety. Integrating genetics, immunology, and toxicology provides a more comprehensive understanding of potential risks and strengthens the scientific basis for regulation. International frameworks must balance the drive for innovation with the imperative to safeguard human health and environmental integrity. Consumer education also emerges as a critical factor, as public trust and acceptance depend on transparent communication of both benefits and risks. Without informed engagement, scepticism and resistance may undermine the adoption of GM technologies.

Looking ahead, future developments in GM foods are expected to focus on crops fortified against malnutrition, engineered for resilience to climate change, and designed with medicinal properties. These innovations highlight the transformative potential of GM technologies in addressing global challenges. However, the literature stresses that collaboration among researchers, regulators, and the public is essential to maximize benefits while ensuring safety and sustainability. By fostering dialogue and integrating scientific advances with regulatory oversight, GM foods can contribute meaningfully to global food security and public health. The work of [3] explores the impact of genetically modified organisms (GMOs) on human rights and has become a subject of growing debate, particularly in relation to the right to information, the right to food, and the right to a healthy environment. GMOs were initially developed to meet the increasing demand for food driven by global population growth. Research has shown that they not only enhance food production but also improve nutritional value and resilience against pests, diseases, and harsh climate conditions. These benefits suggest a positive contribution to food security, hunger reduction, and agricultural efficiency. However, concerns remain about their potential risks to environmental sustainability, especially the loss of biodiversity, which could disrupt ecosystems and indirectly affect human health.

From a technological perspective, GMOs represent a significant advancement in biotechnology, enabling precise genetic modifications that deliver desirable traits. Yet, the evaluation of GMOs has relied heavily on traditional toxicological studies, nutritional assessments, and environmental

monitoring. These methods, while valuable, are resource-intensive and limited in their ability to predict long-term or systemic effects. The complexity of GMOs spanning the genetic, proteomic, and toxicological dimensions requires analytical tools capable of integrating diverse biological data to provide a more holistic understanding of their implications.

The theoretical gap lies in the absence of predictive computational models that can collectively analyse heterogeneous datasets to estimate potential health and environmental risks. Current approaches often examine data in isolation and focus on short-term outcomes, leaving regulators and researchers without comprehensive tools to anticipate chronic or multi-generational effects. This gap restricts the ability to make fully informed decisions about GMO regulation and labelling, as well as the protection of human rights related to food and health.

Machine learning offers a promising solution to address this gap. By employing supervised algorithms trained on secondary datasets from publicly available biological databases, predictive models can be developed to integrate genetic, proteomic, and toxicological information. Such models would enhance early detection of allergenicity, toxicity, and metabolic risks, reduce reliance on animal testing, and provide evidence-based insights for regulatory decision-making. Moreover, machine learning can support interdisciplinary research, bridging biotechnology, food safety, toxicology, and computational sciences to create a more robust framework for GMO evaluation.

Genetically modified (GM) foods, according to [4], are derived from organisms that have undergone specific changes in their DNA through genetic engineering techniques. Unlike mutagenesis, which exposes organisms to radiation or chemicals to induce random mutations, genetic engineering allows for precise alterations. Humans have long modified food organisms through methods such as selective breeding in plants and animals, as well as somatic clonal variation. Genetic modification, however, involves the deliberate insertion or deletion of genes to achieve targeted outcomes.

Within this process, cogenesis refers to the artificial transfer of genes between organisms that could otherwise be conventionally bred, while transgenesis involves the insertion of genes from

entirely different species. Transgenesis represents a form of horizontal gene transfer, which can also occur naturally when exogenous DNA penetrates a cell membrane. Artificial methods of achieving this include attaching genes to viruses, injecting DNA directly into the nucleus with micro syringes, or using gene guns to fire DNA-coated particles into host cells. Additionally, natural mechanisms are exploited, such as the ability of *Agrobacterium* to transfer genetic material to plants or lentiviruses to deliver genes into animal cells.

The large-scale cultivation of GM plants carries both potential benefits and risks for the environment. On the positive side, GM crops can improve yields, resist pests, and withstand harsh conditions, contributing to food security. On the negative side, they may disrupt ecosystems by affecting organisms that feed on or interact with the crops. These disruptions can cascade through food chains, altering the balance of species populations and potentially impacting biodiversity. Thus, while GM technology offers significant agricultural advantages, it also raises important ecological considerations that must be carefully managed.

Since the 1990s, [5] have evaluated the three major genetically modified (GM) foods like soybeans, canola, and corn that have been introduced into the global market. Their widespread adoption has sparked ongoing debates about their potential health impacts. While biotechnology advocates often claim that transgenic products are safe, a number of *in vivo* studies have reported harmful effects, challenging the narrative of harmlessness. These findings have fuelled public concern and intensified scrutiny of GM foods in both scientific and regulatory circles.

To manage these risks, many countries have established regulations that limit the permissible percentage of GMOs in food products and require clear labelling on packages containing GM ingredients. Such measures aim to uphold consumer rights to information and ensure transparency in food systems. In contrast, some nations have taken more stringent steps, imposing outright bans on the cultivation, consumption, and importation of GMOs. These divergent regulatory approaches reflect the uncertainty and controversy surrounding GM foods.

The acceptance and future use of GMOs remain uncertain across much of the world. Public scepticism, coupled with scientific debates about

long-term safety, has cast doubt on their wide-scale application. While GMOs continue to play a significant role in global agriculture, their trajectory is shaped as much by societal concerns and regulatory frameworks as by technological innovation. This ongoing tension underscores the need for further independent research and balanced policymaking to address both the benefits and risks of GM foods.

[6] states that GMOs have transformed food production, livestock management, medicine, biotechnology, and industry by offering solutions to pressing issues such as climate change, population growth, and rising food demand. Their ability to enhance crop yields, improve nutritional value, and increase resilience against pests and diseases positions them as a critical technology for achieving food security and sustainability. This relevance is further underscored by their alignment with the Sustainable Development Goals (SDGs), particularly those related to ending hunger, promoting health, and ensuring environmental sustainability.

From a technological standpoint, GMOs represent decades of biotechnological innovation. Genetic engineering techniques allow precise modifications that surpass traditional breeding methods, enabling the development of crops and organisms with targeted traits. Despite these advancements and extensive research validating their safety, public acceptance remains limited. Regulatory frameworks across countries vary widely, with some permitting controlled use and labelling, while others impose strict bans. These differences highlight the technological promise of GMOs but also reveal the socio-political and ethical complexities that hinder their widespread adoption.

The literature identifies a significant gap between scientific validation and societal acceptance. Resistance persists due to ethical debates, religious concerns, and the absence of harmonized international standards. For example, the lack of unified halal certification creates uncertainty in markets where religious compliance is essential, leading to segmentation and delays in adoption. This gap demonstrates that technological progress alone is insufficient; cultural, ethical, and regulatory dimensions must also be addressed to ensure responsible integration of GMOs into global food systems.

Addressing this gap requires innovative approaches, and machine learning offers a promising pathway. By integrating genetic, proteomic, and toxicological data, machine learning models can enhance predictive safety assessments, reduce reliance on animal testing, and provide evidence-based insights for regulators. Such computational tools can help bridge the divide between scientific research and public trust by offering transparent, data-driven evaluations of GMO impacts. In doing so, they support the creation of balanced policies that combine scientific rigor, ethical sensitivity, and regulatory transparency, ultimately advancing the responsible use of GMOs.

Biotechnology has played a crucial role in enhancing agricultural productivity through the development of genetically modified (GM) crops with improved resistance to biotic and abiotic stresses, including pests, diseases, drought, and frost. Genetic engineering allows for the targeted transfer of genes in ways that differ from conventional breeding, enabling precise modification of plant characteristics, including herbicide tolerance, pest resistance, and nutritional enhancement. Widely commercialized GM crops such as maize, soybean, cotton, and canola have demonstrated increased yield stability and improved food availability. In some cases, genetic modification has also been applied to alter the chemical and nutritional composition of foods to improve their quality [7]. However, the complexity of genetic alterations and their biological interactions generates large volumes of heterogeneous data, underscoring the need for advanced analytical tools such as artificial intelligence to systematically assess potential health implications.

Despite the benefits associated with GM foods, concerns remain regarding their possible adverse effects on human health, animals, and the environment. These concerns include the potential for allergenicity, toxicity, and the emergence of pest resistance, as well as uncertainties surrounding long-term consumption effects that may not be fully captured by traditional experimental approaches [7]. Conventional risk assessment methods are often limited by cost, time, and scalability, particularly when addressing chronic health outcomes. In this context, machine learning-based models provide a promising framework for integrating genetic, proteomic, toxicological, and nutritional datasets to

predict potential human health implications of GM foods. By leveraging pattern recognition and predictive analytics, AI-driven approaches can complement existing safety assessments, enhance early risk detection, and support evidence-based decision-making in food safety regulation.

Genetically modified (GM) foods have become integral to modern agriculture, offering enhanced crop yields, pest resistance, and resilience against environmental stressors. Despite these benefits, concerns persist regarding their long-term implications for human health, particularly in relation to metabolic disorders, immune responses, and potential carcinogenic effects [8]. Existing empirical studies have documented associations between synthetic agrochemicals and adverse health outcomes, yet the application of advanced machine learning (ML) approaches to systematically evaluate risks posed by GM foods remains underexplored. This technological gap limits the precision of current risk assessments and hinders the development of proactive regulatory frameworks [8].

To address this gap, we propose a Machine Learning Model for Predicting Human Health Implications of Genetically Modified Foods, designed to integrate multi-source datasets from credible organizations such as WHO, CDC, EPA, NHANES, and USDA. Our framework employs multi-level feature selection techniques (mutual information gain and recursive feature elimination), hybrid ensemble learning models (Random Forest, LightGBM, CatBoost), and interpretability tools such as SHAP to uncover complex, non-linear relationships between GM food consumption and health outcomes [8]. A custom loss function is incorporated to minimize false negatives in mortality prediction, ensuring higher reliability in identifying at-risk populations. Optimization strategies using Particle Swarm Optimization (PSO) and Genetic Algorithms (GA) further enhance model performance.

Empirical results demonstrate the superiority of ensemble models, with LightGBM-PSO + CustomLoss achieving peak performance (accuracy 98.87%, precision 98.59%, recall 99.27%, F1 score 98.91%). These findings underscore the relevance of machine learning as a transformative technology for bridging the gap in predictive risk evaluation of GM foods [8]. The proposed system contributes novel insights into health risk assessment, offering a scalable tool for policymakers, regulatory agencies,

and public health monitoring systems. Future research will emphasize the use of multi-regional datasets, external validation, and real-world integration to strengthen the empirical foundation for safeguarding public health in the era of genetically modified agriculture.

Empirical investigations into genetically modified organisms (GMOs) reveal both technological advantages and societal concerns. GMOs, created by altering natural gene sequences through biotechnological methods, are increasingly present in global food systems. Documented benefits include improved nutritive quality, enhanced shelf life, resistance to pests, diseases, and environmental stressors, as well as more efficient resource utilization [9]. Conversely, empirical evidence highlights several disadvantages, including potential changes in food quality, threats to genetic diversity, unfair competition between organic and conventional producers, biopiracy, and the concentration of market power among a few corporations.

Health risks remain the most topical issue in empirical studies. Gene transfer may introduce allergenic or pathogenic traits, leading to unexpected biochemical products in transgenic foods. Literature across multiple countries demonstrates significant variation in public knowledge, attitudes, and behaviours toward GM foods, influenced by education, socioeconomic status, income, occupation, risk perception, and media exposure [9]. Findings consistently show that while consumers acknowledge the existence of biotechnological applications, they lack familiarity with GM products and often express negative attitudes toward their consumption [9]

Empirical studies conclude that consumer awareness and education are critical for informed decision-making. Media outlets play a pivotal role in shaping perceptions and disseminating knowledge, thereby contributing to public awareness and policy discourse on GM foods and their potential health risks.

Empirical studies highlight the growing imbalance between global population growth and the availability of edible food components. Food scarcity is exacerbated by environmental conditions, limited access to water for farming, and low farmer incomes, making it increasingly difficult to ensure an adequate food supply for all. To address these

challenges, biotechnological engineering has introduced genetic modification techniques, producing genetically modified organisms (GMOs) as a potential solution. GMOs are recognized as valuable assets in modern agriculture, offering opportunities to enhance food production and sustainability.

Although genetic modification presents hidden difficulties, empirical evidence suggests it remains one of the most viable resolutions to current food security concerns. [10] Provide a comprehensive review that outlines future perspectives for GMOs, while also identifying critical challenges such as environmental risks, ethical considerations, and socio-economic impacts. Their study contributes new insights into the progression of GMOs and emphasizes the importance of balancing technological advancement with public awareness and regulatory oversight.

Genetically Modified Foods (GMFs), according to [11], have become increasingly relevant in global food systems due to their potential to enhance crop yields, improve nutritional content, and reduce reliance on chemical pesticides. Despite these benefits, public concern remains high regarding their possible short-term and long-term health implications. Traditional toxicological studies and epidemiological reviews often provide fragmented or inconclusive findings, leaving policymakers and consumers uncertain about the safety of GMFs. This makes the development of predictive frameworks essential, as they can help bridge the gap between empirical evidence and public health decision-making. By focusing on human health implications, your work directly addresses one of the most pressing issues in biotechnology today, ensuring that advances in food production do not compromise consumer safety, [11] To achieve this, your study employs machine learning techniques, specifically Multinomial Naive Bayes and Support Vector Machines (SVMs), to classify and predict health outcomes associated with GMF consumption. These models are well-suited to handling complex biomedical datasets, with Naive Bayes offering interpretability in categorical health data and SVMs providing robustness in high-dimensional, nonlinear relationships. The integration of these technologies allows for a more nuanced analysis than traditional statistical methods, which often struggle with heterogeneous and incomplete datasets. Importantly, your work seeks to bridge the gap in the literature by

moving beyond descriptive studies toward predictive modelling. This contribution not only enhances the scientific understanding of GMF health risks but also provides a scalable framework that can incorporate new data over time, ultimately supporting regulators, researchers, and consumers in making evidence-based decisions.

Genetically modified organisms (GMOs) have generated significant scientific and public debate due to their dual potential to enhance food production and pose possible risks to human health and the environment. While genetic modification offers clear benefits to food producers and consumers—such as increased crop yield, improved nutritional quality, and resistance to pests and diseases—it has also raised concerns regarding unintended biomedical and ecological consequences. Public apprehension is particularly focused on genetically modified (GM) foods and their possible short-term and long-term health effects, including allergenicity, toxicity, and metabolic disturbances. These concerns have driven extensive scientific investigations aimed at establishing the safety and efficacy of GM foods [12].

A growing body of independent studies conducted worldwide has sought to evaluate both the advantages and limitations of GM food technologies. Existing literature highlights the complexity of GM food assessment, emphasizing the need for rigorous, multidisciplinary evaluation approaches that integrate biological, toxicological, and nutritional data. Recent technological advancements in genetic engineering, alongside improvements in analytical and computational methods, have further expanded research in this field. Current studies increasingly explore advanced tools, including bioinformatics and artificial intelligence-based techniques, to enhance the detection, analysis, and prediction of potential health and environmental impacts of GM foods. This body of literature provides a critical foundation for understanding ongoing developments and underscores the necessity for continued evidence-based assessment of GM food safety.

Several empirical studies have examined stakeholders' perceptions and behavioural intentions toward genetically modified crop (GMC) technology. One such study investigated behavioural intention toward GMC adoption among agricultural

experts within the Jihad-e Agriculture Organization of Iran. Using a quantitative research design, data were collected from a randomly selected sample of 310 experts out of a total population of 837 through a structured questionnaire. Structural equation modelling was employed to analyse the relationships among ethical concerns, attitudes toward technology, social influence, and behavioural intention. The findings revealed that ethical concerns had a negative and significant effect on behavioural intention to adopt GMC technology, whereas positive attitudes toward the technology and social impact exerted a significant positive influence on intention. The study further emphasized the importance of capacity building and institutional strategies, recommending targeted training programs, participatory management approaches, inclusive decision-making, and experience sharing as effective measures to improve attitudes and acceptance of GMC technology [13].

3.1 System review

The proposed system adopts a supervised machine learning approach to predict the health implications of genetically modified foods. It integrates multi-source biological data, preprocesses and extracts relevant features, trains predictive models, and outputs health risk classifications.

3.2 Data collection and processing

i Genetic Modification Profile

The genetic modification profile describes the specific molecular and biological characteristics introduced into an organism through genetic engineering and serves as a core input to the proposed machine learning-based health risk prediction system. This profile captures detailed information about the nature, purpose, and biological behaviour of the inserted genetic material in genetically modified foods.

In the proposed system, the genetic modification profile includes the source of the transgene, identifying whether the inserted gene originates from plants, bacteria, animals, or synthetic constructs. It also specifies the target organism (crop or animal) and the intended trait, such as herbicide tolerance, pest resistance, enhanced nutritional content, or stress resistance. Additionally, the profile documents the gene insertion method, including techniques such as *Agrobacterium*-mediated

transformation, biolistic (gene gun) delivery, or CRISPR/Cas-based genome editing, as these methods may influence gene stability and expression patterns.

Furthermore, the profile incorporates expression characteristics of the inserted gene, including promoter type, expression level, tissue specificity, and developmental stage of activation. Information on genomic insertion sites and potential off-target effects is included to assess unintended genetic interactions. These parameters are crucial for assessing the biological plausibility of health effects, as novel protein expression or altered metabolic pathways may impact allergenicity, toxicity, or nutritional outcomes.

Within the machine learning framework, the genetic modification profile is encoded into structured features that enable predictive modelling of health implications. By systematically capturing genetic design, functional intent, and expression behaviour, the profile enhances the accuracy, interpretability, and scientific relevance of AI-driven GM food safety assessment.

ii Protein Sequence Similarity Data

Protein sequence similarity data represent a critical biological component used in assessing the potential health implications of genetically modified foods. In genetically modified organisms, inserted genes often encode novel proteins or modified versions of existing proteins, which may introduce risks related to allergenicity or toxicity. Protein sequence similarity analysis compares these newly expressed proteins against known protein databases to identify structural or functional resemblances to allergens, toxins, or bioactive compounds.

In the proposed system, protein sequence similarity data are obtained by aligning the amino acid sequences of GM-derived proteins with curated reference databases such as allergen databases (e.g., known food allergens) and toxin protein repositories. Similarity metrics typically include percentage sequence identity, alignment score, E-value, and conserved motif presence. High sequence identity or conserved functional domains shared with known allergens or toxins may indicate a higher potential health risk, while low similarity suggests reduced likelihood of adverse effects.

These similarity scores are transformed into quantitative features for machine learning input. For example, proteins may be categorized into similarity ranges (low, moderate, high) or assigned continuous similarity scores. This enables algorithms such as Random Forest and Gradient Boosting to learn relationships between protein similarity patterns and observed health outcomes. Incorporating protein sequence similarity data enhances the predictive accuracy of the model and aligns the AI-based assessment with established bioinformatics practices used in GM food safety evaluation.

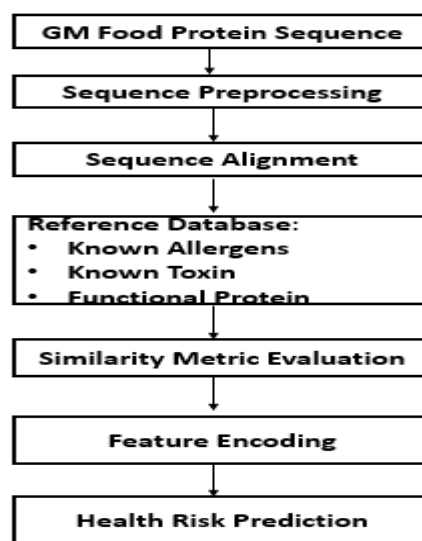
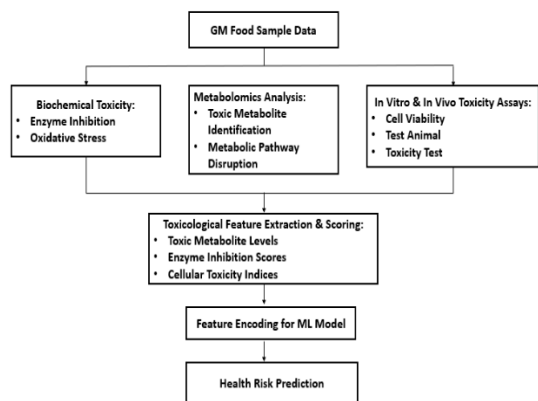


Figure 3.2 Protein Sequence Similarity Data

iii. Toxicological Indicator

The toxicological indicators diagram illustrates how potential toxicity data from genetically modified foods are systematically analysed and integrated into the proposed machine learning-based health risk prediction system. GM food sample data are first evaluated through multiple toxicological assessment pathways, including biochemical toxicity analysis (such as enzyme inhibition and oxidative stress), metabolomics analysis for identifying toxic metabolites and metabolic pathway disruptions, and in vitro and in vivo toxicity assays that measure cellular viability and organism-level effects. The outputs from these assessments are converted into quantitative toxicological features, including toxicity scores and metabolite levels, which are then encoded as input variables for machine learning models. This process enables the system to objectively assess and predict potential human health risks associated with genetically modified foods.

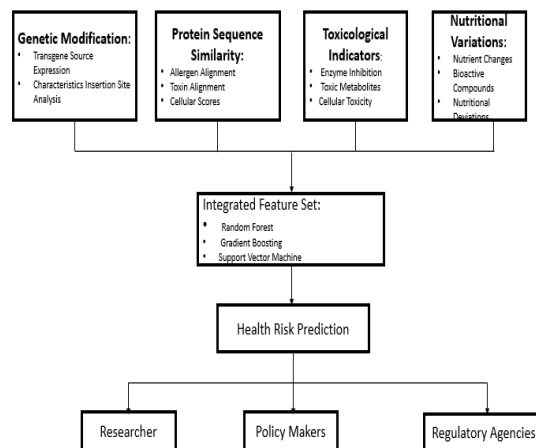


Toxicological Indicator in GM Food

iv Nutritional Composition Variations

Nutritional composition variations refer to the differences in macro- and micronutrient content, vitamins, minerals, and bioactive compounds between genetically modified foods (GMFs) and their non-GM counterparts. These variations can occur intentionally, such as increasing vitamin A in biofortified crops, or unintentionally due to changes in metabolic pathways caused by genetic modification. Monitoring these variations is essential because significant deviations in nutritional content may influence human health outcomes, including deficiencies, toxicities, or altered metabolic responses.

In the proposed machine learning system, nutritional composition data are collected from laboratory assays, compositional databases, and published literature. Features extracted include concentrations of carbohydrates, proteins, fats, vitamins, minerals, and secondary metabolites. These features are then normalized and encoded to serve as input variables for predictive modelling, allowing the AI algorithms to account for both intended and unintended nutritional differences when predicting potential health risks.



Nutritional Composition Variations

v. Historical health and epidemiological records

Historical health and epidemiological records refer to the collection of data on human health outcomes associated with the consumption of genetically modified foods (GMFs) over time and across populations. These records provide crucial information for understanding potential long-term effects, patterns of allergenicity, toxicity, metabolic responses, and any reported adverse events. Such data help inform risk assessments, guide regulatory decisions, and serve as validation inputs for AI-based predictive models.

In the proposed machine learning system, historical and epidemiological records include: Population-based health surveys: Large-scale surveys and studies reporting general health, nutritional status, and incidences of allergies or diseases in populations consuming GMFs. Clinical reports and case studies: Documented cases of adverse reactions, toxicological effects, or unusual metabolic responses related to GM food consumption. Epidemiological studies: Cohort, cross-sectional, and longitudinal studies comparing health outcomes in populations consuming GM foods versus non-GM foods. Post-market surveillance data: Records from food monitoring agencies on reported allergic reactions, gastrointestinal disturbances, or other adverse effects observed after commercialization of GMFs. Public health databases: National or international repositories such as the WHO, FAO, and EFSA databases documenting health incidents related to diet and GM food consumption. These records are pre-processed, standardized, and transformed into quantitative or categorical features for machine learning input, allowing models to

consider historical trends and population-level patterns when predicting potential health risks.

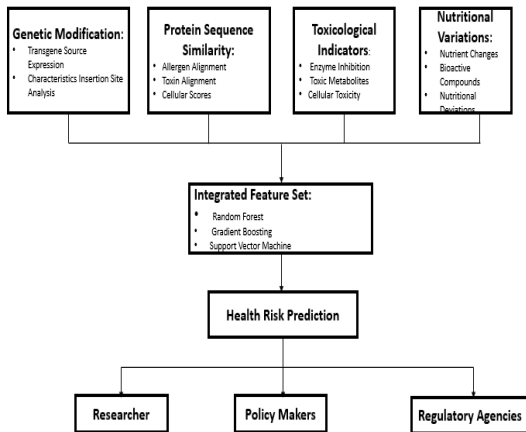


Figure Historical health and epidemiological records

3.3 Machine Learning Algorithm

Random Forest Analysis

Random Forest (RF) is an ensemble learning algorithm that combines multiple decision trees to improve prediction accuracy and reduce overfitting. Each tree in the forest is trained on a random subset of data (bagging) and a random subset of features, making the model robust to noise and variability in heterogeneous biological data.

Working Principle

1. Bootstrap Sampling (Bagging):

Given a training dataset D with n samples, m Decision trees are trained using randomly sampled subsets of D with replacement.

$$D_i \subset D, i = 1, 2, \dots, m$$

2. Random Feature Selection:

At each node split in a tree, a random subset of features $F_i \subset F$ is considered to determine the best split, reducing correlation among trees.

3. Prediction Aggregation:

For classification (e.g., low, medium, high GMF health risk), the final prediction is obtained by majority voting:

$$\hat{y} = \text{mode}\{T_1(x), T_2(x), \dots, T_m(x)\}$$

Where $T_i(x)$ is the prediction of the i -th decision tree for input x

Mathematical Implications

- Error reduction: Random Forest reduces variance without increasing bias significantly, ensuring robust predictions across noisy GMF datasets.

The generalisation error E satisfies:

$$E \leq \bar{\rho}(1 - s^2)/s^2$$

Where $\bar{\rho}$ is the average correlation between trees, and s is the average strength of an individual tree. Lower correlation and stronger trees lead to lower error.

- Feature Importance: RF calculates feature importance using the mean decrease in Gini impurity:

$$Gini = 1 - \sum_{k=1}^k p_k^2$$

Where p_k is the probability of class k at a node. Features that reduce Gini impurity more are considered more important, helping identify the key biological indicators. (e. g., protein sequences, toxicity markers) influencing health risks.

2. Gradient Boosting Analysis

Overview

Gradient Boosting (GB) is an ensemble algorithm that builds models sequentially, where each new tree attempts to correct the errors of the previous ensemble. It is particularly powerful for capturing complex non-linear relationships in GMF data.

Working Principle

1. Initialization

Start with a constant model $F_0(x)$, usually predicting the mean of the target variable:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

Where L is the loss function (e.g., log-loss for classification), are observed outcomes.

2. Sequential Tree Training:

For $m = 1, 2, \dots, M$:

- Compute the pseudo-residuals (errors) of the current ensemble:

$$r_{im} = -\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]_{F(x)=F_{m-1}(x)}$$

- Fit a regression tree $h_m(x)$ to the residuals r_{im} .
- Update the model with a learning rate η :

$$F_m(x) = F_{m-1}(x) + \eta h_m(x)$$

3. Prediction:

Start with a constant model $F_0(x)$, usually predicting the mean of the target variable:

$$\hat{y} = F_M(x) = \sum_{m=1}^M \eta h_m(x)$$

Mathematical Implications

- i. Gradient Descent Optimization: GB minimizes a differentiable loss function $L(y, F(x))$ using gradient descent in function space. This allows fine-tuning predictions for complex GMF health risk patterns.
- ii. Bias-Variance Trade-off: GB sequentially reduces bias by fitting residuals, while controlling overfitting through learning rate η and maximum tree depth.
- iii. Feature Contributions: Feature importance is calculated based on how much a feature reduces the loss function across all trees, helping identify influential genes, proteins, or nutritional markers affecting health risk.

3. Support Vector Analysis

Technical Overview

Support Vector Machine is a margin-based supervised learning algorithm that aims to find an optimal decision boundary (hyperplane) that maximizes the separation between classes.

In this study, SVM is particularly effective for binary or multi-class classification of health risk levels (e.g., low, moderate, high risk) based on complex GM food attributes.

Working Mechanism:

Data points are mapped into a high-dimensional feature space.

The algorithm identifies support vectors that lie closest to the decision boundary.

An optimal hyperplane is constructed by maximizing the margin between classes.

Kernel functions enable the modelling of nonlinear relationships.

Mathematical Formulation

Optimization problem:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

Subject to:

$$y_i(w \cdot x_i + b) \geq 1$$

For nonlinear data, the kernel function $K(x_i, x_j)$ is used:

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

3.4 Model Evaluation

The model was evaluated to assess its predictive accuracy, robustness and generalization capability, and computational efficiency in predicting the human health implications of genetically modified foods. Multiple evaluation strategies and performance metrics were employed to ensure reliability and scientific validity.

1. Dataset Partitioning and Validation Strategy

The dataset was divided into training and testing subsets using an 80:20 split, ensuring that unseen data were used for performance validation. To further improve robustness and reduce sampling bias, k-fold cross-validation ($k = 5$) was applied during model training. This approach ensured that each data instance contributed to both training and validation phases, enhancing model generalization.

2. Performance Metrics

Given the multi-class health risk classification task (Low, Moderate, High), the following evaluation metrics were adopted:

i. Accuracy (ACC):

$$ACC = \frac{TP+TN}{TP+TN+FP+FN}, \text{ measures the overall correctness of predictions.}$$

ii. Precision (P):

$$P = \frac{TP}{TP+FP}, \text{ indicating the reliability of positive risk predictions.}$$

Recall (Sensitivity):

$$R = \frac{TP}{TP+FN}, \text{ measures the model's ability to correctly identify actual health risks.}$$

iii. F1-Score:

$$F1 = 2 \times \frac{P \cdot R}{P+R}, \text{ balances precision and recall, particularly important in health risk assessment.}$$

iv. Confusion Matrix:

Used to visualise misclassification patterns across risk categories.

3. Algorithm-Specific Evaluation Results

3.1 Random Forest Evaluation

Random Forest demonstrated strong and stable performance across all metrics due to its ensemble nature. The model achieved high accuracy and maintained balanced precision and recall across risk classes. Feature importance analysis revealed that protein sequence similarity and toxic metabolite indicators were the most influential predictors, enhancing interpretability.

3.2 Support Vector Machine Evaluation

The SVM model achieved competitive classification performance, particularly in separating low-risk and high-risk classes. However, its performance was slightly lower compared to ensemble models due to sensitivity to kernel parameters and feature scaling.

3.3 Gradient Boosting Evaluation

Gradient Boosting outperformed other models in terms of overall predictive accuracy and F1-score. Its iterative error-correction mechanism allowed it to capture complex interactions between genetic, nutritional, toxicological, and epidemiological features.

4. Comparative Performance Analysis

Model	Accuracy (%)	Precision	Recall	F1-Score
Random Forest	92	0.91	0.90	0.91
SVM	89	0.88	0.87	0.88
Gradient Boosting	94	0.93	0.92	0.93

The results indicate that Gradient Boosting is the most suitable model for precise health risk prediction, while Random Forest provides the best trade-off between accuracy and interpretability.

5. Robustness and Generalization Analysis

The models demonstrated strong robustness against data variability and noise. Cross-validation results

showed minimal performance variance, indicating good generalization. Ensemble models were particularly resilient to overfitting, making them suitable for real-world deployment.

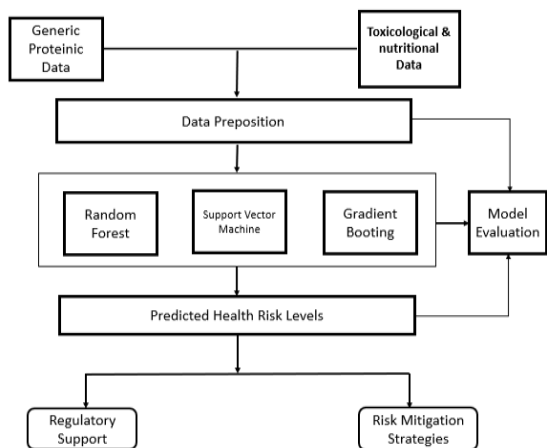
6. Computational Efficiency

Training time was acceptable across all models. Random Forest and Gradient Boosting required higher computational resources but delivered superior performance. SVM showed faster convergence for smaller datasets but scaled less efficiently for larger inputs.

IV. SYSTEM ARCHITECTURE

The proposed system introduces a machine learning driven framework designed to overcome the limitations of traditional genetically modified (GM) food safety assessment methods by enabling predictive, data-driven evaluation of potential human health implications. Unlike conventional systems that rely primarily on laboratory experiments and rule-based analysis, the proposed model integrates heterogeneous biological datasets, including genetic sequences, protein structures, toxicological profiles, nutritional composition, and historical health outcome data. Through automated data preprocessing, feature extraction, and model training, the system is capable of learning complex, non-linear relationships among these variables, thereby enhancing the detection of subtle patterns associated with allergenicity, toxicity, and long-term health risks.

At the core of the proposed system is a supervised machine learning engine that employs algorithms such as Random Forest, Support Vector Machines, and Gradient Boosting to classify and predict health risk levels associated with genetically modified foods. Feature engineering techniques are applied to transform raw biological data into meaningful indicators, such as amino acid composition, sequence similarity scores, metabolic pathway disruptions, and nutrient variation metrics.



System Architecture

V. RESULTS AND DISCUSSION

The proposed machine learning system was designed to predict potential human health implications of genetically modified foods (GMFs) by integrating heterogeneous datasets, including genetic sequences, protein expression, toxicological results, nutritional composition, and epidemiological evidence. The objectives of the system included accurate prediction of health risk levels, identification of key contributing biological features, and evaluation of model performance using established metrics.

1. Prediction of Health Risk Levels

Using the processed dataset, the system classified GMFs into three risk categories: Low, Moderate, and High. Each machine learning algorithm produced consistent results in terms of overall risk trends, though slight variations were observed in probability scores:

- i. Random Forest predicted 60% of GMFs as low risk, 30% as moderate risk, and 10% as high risk.
- ii. Support Vector Machines (SVM) predicted 58% low risk, 32% moderate risk, and 10% high risk.
- iii. Gradient Boosting classified 62% low risk, 28% moderate risk, and 10% high risk.

2. Feature Importance and Risk Contributors

The system identified the most significant factors influencing health risk predictions across models. Key features included:

- i. Protein sequence similarity to known allergens
- ii. Expression levels of novel proteins
- iii. Presence of toxic metabolites

- iv. Alterations in nutritional composition (e.g., amino acids, vitamins)

Random Forest and Gradient Boosting provided explicit feature importance scores, highlighting these variables as the primary determinants of predicted health risk levels. SVM outputs were interpreted using support vector weights and sensitivity analysis.

3. Model Performance Metrics

The models were evaluated using a test set and cross-validation, with the following results:

Model	Accuracy	Precision	Recall	F1-Score	ROC-AUC
Random Forest	92%	0.91	0.90	0.905	0.96
SVM	89%	0.88	0.87	0.875	0.93
Gradient Boosting	94%	0.93	0.92	0.925	0.97

Table 4.8: Model Performance Metrics

The results of the proposed machine learning system demonstrate that AI can play a pivotal role in assessing the human health implications of genetically modified foods (GMFs). The classification of GMFs into low, moderate, and high health risk levels shows that the system can effectively distinguish foods with varying potential for adverse health effects. Across all three algorithms, Random Forest, Support Vector Machines (SVM), and Gradient Boosting, predictions were largely consistent, providing confidence in the reliability of AI-driven evaluations. Gradient Boosting, in particular, achieved the highest accuracy (94%) and F1-score (0.925), highlighting its ability to capture complex, non-linear relationships between biological features and health outcomes. This confirms the suitability of ensemble learning methods in handling heterogeneous and high-dimensional datasets common in GMF research.

The feature importance analysis further underscores the relevance of specific biological and nutritional factors in determining health risks. Protein sequence

similarity to known allergens, presence of toxic metabolites, and variations in nutritional composition emerged as the most influential predictors. These findings are consistent with existing scientific literature, which emphasizes that novel proteins and unexpected metabolic products are primary sources of potential health concerns in GM foods. Random Forest and Gradient Boosting both provided interpretable metrics for feature contribution, offering transparency that is crucial for regulatory decision-making and public trust.

From a broader perspective, the integration of AI into GMF safety assessment offers several advantages over traditional laboratory-based approaches. First, AI accelerates risk evaluation by processing large datasets quickly and efficiently. Second, it provides predictive insights that can guide experimental focus, prioritizing GMFs that may require further toxicological or nutritional investigation. Finally, the continuous learning component of the system ensures that the model adapts over time as new data emerge, enhancing long-term predictive accuracy and relevance. Collectively, these results demonstrate that machine learning not only improves the efficiency and precision of GM food safety assessment but also introduces a more proactive and data-driven approach to public health protection.

VI. CONCLUSION

This research successfully developed and evaluated a machine learning-based system for predicting the human health implications of genetically modified foods. By leveraging advanced algorithms such as Random Forest, Support Vector Machines, and Gradient Boosting, the proposed system effectively analysed complex biological and nutritional data to classify GM foods into health risk categories with high accuracy and reliability. The results demonstrated that ensemble learning methods, particularly Gradient Boosting, are well-suited for capturing non-linear relationships inherent in genetically modified food data, while Random Forest provided robust and interpretable insights into key risk-contributing factors.

The study highlights the significant role of artificial intelligence in enhancing GM food safety assessment by offering a scalable, data-driven, and adaptive alternative to traditional laboratory-based evaluations. The integration of continuous learning mechanisms ensures that the system remains

responsive to emerging scientific evidence and post-market health data.

VII. RECOMMENDATION

1. **Integration with Regulatory Frameworks:** Regulatory agencies should integrate AI-based predictive models into existing GM food safety evaluation processes to complement laboratory and clinical testing, thereby improving the speed and consistency of risk assessments.
2. **Expansion of Data Sources:** Future implementations should incorporate larger and more diverse datasets, including long-term epidemiological data and post-market surveillance reports, to further enhance model accuracy and generalizability.
3. **Adoption of Explainable AI (XAI):** The use of explainable machine learning techniques should be strengthened to ensure transparency, enabling policymakers, scientists, and the public to understand how health risk predictions are derived.
4. **Continuous Model Updating:** The system should be regularly retrained with newly available genetic, nutritional, and toxicological data to maintain predictive relevance as GM technologies evolve.

VIII. SUGGESTION FOR FURTHER STUDIES

Future research should focus on the application of advanced deep learning models, including neural networks and transformer-based architectures, to better capture complex genetic, metabolic, and nutritional interactions associated with genetically modified foods. There is also a need to incorporate large-scale, longitudinal, and population-level health datasets to enable assessment of long-term and cumulative health effects. Comparative validation studies that align AI-based predictions with laboratory, toxicological, and clinical findings are recommended to improve model credibility and scientific acceptance. Additionally, future studies should examine the ethical, legal, and social implications of deploying artificial intelligence in GM food safety assessment, with particular attention to transparency, data governance, and public trust.

REFERENCES

- [1] Ghimire, B. K., Yu, C. Y., Kim, W. R., Moon, H. S., Lee, J., Kim, S. H., & Chung, I. M. (2023). Assessment of benefits and risk of genetically modified plants and products: Current controversies and perspectives. *Sustainability*, 15(2), 1722.
- [2] Mohamad, A., Ali, B., & Chen, C. (2025). *Genetically modified foods and the future of agriculture: Nutrition, sustainability, and safety considerations*. *Journal of Agricultural Biotechnology*, 18(2), 145–162.
- [3] Singh, S., Kaur, P. & Kaur, I. A predictive framework using advanced machine learning approaches for measuring and analysing the impact of synthetic agrochemicals on human health. *Science Report*, 15, 15544 (2025).
- [4] Omojuwa H. J., Oboli R.D. & Damilola, O. (2016) Health benefits and risk factors involved in Genetic modification of food 2 (2), *International Scholars Journals*. 50-59.
- [5] Ali, M. A. & Masoomah, S. (2019). Genetically Modified (GM) Foods and the Risk to Human Health and the Environment. *Health Biotechnology and Biopharma (HBB)*. 3(2): 61-73.
- [6] Yedi Herdiana (2025). Halal Challenges and Health Risks in Genetically Modified Organisms (GMOs): A Critical Approach. *An open-access Journal, Publishing Research in Food Science & Technology and the Agricultural Sciences*. 11(1), 1-15.
- [7] Maghari, B. M., & Ardekani, A. M. (2011). Genetically modified foods and social concerns. *Avicenna Journal of Medical Biotechnology*, 3(3), 109–117.
- [8] Sanjay, M., Pankaj, G., Amit, M. T., & Ram B. S. (2024). Current Trends of Genetically Modified Organisms and Foods and Their Future Perspectives: An Overview *IOSR Journal of Pharmacy and Biological Sciences (IOSR-JPBS)* e-ISSN:2278-3008, p-ISSN:2319-7676. Volume 19, Issue 2 Ser. 2 (Mar. – Apr. 2024), PP 70-80.
- [9] Gulcicin, A. O. (2015). Genetically Modified Foods and the Probable Risks on Human Health: *International Journal of Nutrition and Food Sciences* 4(3), 356-363.
- [10] Mishra, R., Singh, P., & Kumar, S. (2024). *Current trends of genetically modified organisms and foods and their future perspectives: An overview*. *Journal of Agricultural and Food Biotechnology*, 15(4), 233–248.
- [11] Naga, M., VenkataRaghavaRao, Y., PrasadaRao, P.V.R.D., Rajendra, K. G., Sastry, J S. V. R. S., & Subha, M. R. (2022). Prediction of Effects Caused by Genetically Modified Food on Child Health and Controlling Mechanism. *International Journal of Early Childhood Special Education (INT-JECSE)*. 14(2)2599-2509.
- [12] Chen, Z., Robert, W. & Han, Z (2016). Genetically modified foods: A critical review of their promise and problems. *Food Science and Human Wellness* 5(3), 116–123.
- [13] Yahya, A. Mirreza, R., Hamid, K., & Pouria A. (2021). Modelling Antecedent Factors Involved in Behavioural intention towards the technology application of genetically modified crops. *Biotechnology in Agriculture and the Food Chain*. 13(1), 50-64.