

A Discussion on Ethics, Trust and Security Governance in the Evolution of Artificial Intelligence Technology

Shufeng Wang^{1,*}

¹ School of Computer Science and Technology, Taiyuan Normal University, Shanxi 030619, China

*** Correspondence:**

Shufeng Wang

18235230439@163.com

Received: 10 March 2026 / Accepted: 27 March 2026 / Published online: 30 March 2026

Abstract

With the rapid advancement and widespread deployment of artificial intelligence, issues related to ethics, trust, and security have become increasingly prominent, posing significant challenges to its sustainable development. This paper aims to systematically investigate these intertwined challenges by constructing a structured analytical framework encompassing three core dimensions: ethics, trust, and security. Methodologically, the study analyzes the evolutionary trajectory and application contexts of AI, and synthesizes key risk patterns and governance concerns. It examines major ethical issues, including algorithmic bias, data privacy, responsibility attribution, and value alignment, and further explores the underlying tension between technological rationality and social values. In addition, the paper proposes a multi-level trust formation mechanism based on technological reliability, institutional assurance, and public cognition, emphasizing that trust emerges from the joint interaction of technical systems and governance structures. At the security level, it highlights the complexity and cross-domain nature of AI risks, and advocates for a full life-cycle governance framework integrating risk assessment, technical auditing, and continuous supervision. The study's key contribution lies in integrating ethical, trust, and security perspectives into a unified analytical framework, and in proposing a collaborative governance approach involving governments, enterprises, and international stakeholders.

Keywords: Security Risk Management; AI Ethics; Algorithmic Fairness; Explainable AI

1. Introduction

Over the past few decades, artificial intelligence technology has undergone leapfrog development and, through technological breakthroughs such as deep learning, big data and computing power improvement, has gradually moved from laboratory research to large-scale application, and has penetrated into many key fields such as medical care, finance, transportation,

education and entertainment. The spread of AI technology has not only significantly improved the efficiency of social operation and resource allocation capabilities, but also profoundly reshaped the human-computer interaction mode and decision-making mechanism. However, while the widespread application of artificial intelligence has brought significant economic and social benefits, it has also triggered a series of complex and deep-seated social problems (Russell et al., 2021). Ethical issues have become a core topic of widespread discussion globally. From the perspective of value norms and social governance, the decision-making process of AI systems often lacks transparency, and structural biases or discrimination mechanisms may be embedded in the data collection, processing and algorithm modeling process. Such technical biases may not only damage individual privacy and personal dignity, but may also exacerbate social inequality and resource allocation imbalance, thus posing a serious challenge to the principle of social fairness. In this context, building a systematic, operable and universally binding artificial intelligence ethical framework has become an important issue in the global governance system. At the same time, the rapid iteration and widespread application of artificial intelligence technology are constantly impacting the existing social norm system and ethical boundary structure. As the participation of intelligent systems in the decision-making process continues to deepen, the traditional responsibility allocation mechanism and rights and obligations framework built with human subjects as the core are facing structural challenges (Batoool et al., 2025). These deep-seated theoretical propositions urgently need to be systematically responded to and institutionalized within the framework of ethical philosophy, jurisprudence and public governance.

From the perspective of social trust mechanism, trust constitutes the basic condition for the widespread application of artificial intelligence in social embedding. Its formation and maintenance directly affect the social acceptance and sustainable development capability of AI system. Although artificial intelligence has shown significant advantages in efficiency optimization and decision support, due to the complexity of algorithm models and the opacity of decision paths, the public often lacks a full understanding of its operating logic and risk control mechanism, which in turn affects the degree of trust in the system. Therefore, how to build a trustworthy artificial intelligence system with verifiability and accountability by enhancing the interpretability of algorithms, improving the transparency of operation and clarifying the responsibility allocation mechanism has become an important direction for current technology governance and institutional innovation. At the level of security governance, the risks faced by artificial intelligence systems also show a trend of increasing complexity and diversification. With the continuous expansion of algorithm capabilities and application scenarios, problems such as malicious attacks, adversarial sample interference, data tampering and model loss of control are constantly emerging, and their potential impact may affect public safety, economic stability and even national security (McCormack et al., 2024). Especially under the dual effect of "black box" characteristics and adaptive learning mechanism, the unpredictability of some system behaviors has been significantly enhanced, increasing the difficulty of risk assessment and responsibility definition. Therefore, while continuously promoting technological innovation and industrial upgrading, constructing a systematic safety assessment system and a multi-level risk prevention and control mechanism to ensure the controllability and reliability of artificial intelligence systems has become a critical issue that society must address. This study aims to systematically

explore the ethical, trust, and security issues in artificial intelligence, analyze how these issues intertwine and influence each other, and provide useful insights and references for academic research and practical applications in related fields.

2. Basic Concepts and Theoretical Foundations of Artificial Intelligence

2.1. Definition of Artificial Intelligence

Artificial Intelligence, as an important branch of computer science, has been evolving and expanding with technological development since its proposal in the mid-20th century. Generally speaking, artificial intelligence refers to a series of theories, methods and technologies that simulate, extend or expand human intelligent activities through computer systems. Its core goal is to enable machines to have the ability to perceive the environment, understand information, learn from experience, reason and make decisions, and act autonomously (Li et al., 2024). The formal proposal of the concept of artificial intelligence is usually traced back to the Dartmouth Conference in 1956. At this conference, John McCarthy first proposed the term "Artificial Intelligence" and defined it as "the science and engineering of how to make machines exhibit human-like intelligent behavior", emphasizing that through algorithm design and program implementation, machines can exhibit human-like cognitive abilities in specific tasks. At the same time, Alan Turing, in the earlier Turing Test, explored from a behaviorist perspective whether machines can exhibit intelligent reactions that are indistinguishable from humans, providing an important testable theoretical basis for artificial intelligence research.

From the perspective of the degree of intelligence realization, artificial intelligence is usually divided into weak artificial intelligence (Weak AI) and strong artificial intelligence (Strong AI). Weak artificial intelligence, also known as narrow-domain artificial intelligence, mainly focuses on achieving high efficiency in specific tasks or application scenarios, such as image recognition, speech recognition and natural language processing. Most practical application systems currently belong to this category. In contrast, strong artificial intelligence refers to intelligent systems that possess general cognitive abilities comparable to or even surpassing those of humans, capable of transfer learning, autonomous understanding and creative thinking between different tasks. Although strong artificial intelligence is still in the theoretical exploration stage, its concept is of great significance in the research on artificial intelligence ethics, philosophy and future technological development. Overall, artificial intelligence is not a single technology or method, but a typical interdisciplinary research field that integrates the theoretical foundations of multiple disciplines such as computer science, mathematics, cognitive science, neuroscience and philosophy (Kattnig et al., 2024). Therefore, in academic research, the definition of artificial intelligence should not only focus on the algorithm and system implementation at the technical level, but also systematically understand and analyze it from multiple dimensions such as the essence of intelligence, behavioral performance and decision-making mechanism, so as to more comprehensively reveal the development connotation and research value of artificial intelligence.

2.2. Main research content of artificial intelligence

As a highly comprehensive interdisciplinary field, artificial intelligence covers multiple levels

from basic theory to applied technology. Around the core goal of building an intelligent system with perception, learning, reasoning and decision-making capabilities, artificial intelligence has gradually formed several interrelated research directions, mainly including machine learning, deep learning, natural language processing, computer vision, knowledge representation and reasoning, and reinforcement learning. These research directions are intertwined in theory and method, and together constitute the main framework of the artificial intelligence technology system (Sistla, 2024). Among them, machine learning, as an important theoretical foundation of artificial intelligence, takes data-driven as its core paradigm. Its basic idea is to enable computer systems to automatically identify potential patterns and complete prediction or decision-making tasks through training and modeling of sample data. Compared with the traditional method of relying on explicit rule programming, machine learning emphasizes the optimization of parameters and the abstraction of patterns in the process of data input and model training, thereby significantly improving the system's adaptive ability and generalization performance. Based on the differences in learning mechanisms and data structures, machine learning can usually be divided into supervised learning, unsupervised learning and reinforcement learning. Through its high flexibility and adaptability in data modeling, pattern recognition, and decision optimization, machine learning not only provides systematic algorithmic tools and modeling methods for artificial intelligence research, but also gradually forms a sustainable and expandable technical framework and methodological foundation, making it an indispensable core technical support in the design and implementation of current intelligent systems.

Based on the continuous development of the theoretical system of machine learning, deep learning, as an important branch, has significantly improved the ability to represent and process complex data by constructing multi-layer neural network structures. The core advantage of deep learning is that it can automatically learn hierarchical feature representations from large-scale datasets, thereby reducing the dependence on manual feature engineering and enabling it to achieve breakthrough results in high-dimensional and unstructured data processing tasks. Typical model structures include convolutional neural networks, recurrent neural networks, and Transformer, which has been widely used in recent years. The development of deep learning has not only promoted the rapid progress of image recognition, speech recognition and natural language processing, but also marked the important transformation of the artificial intelligence research paradigm from "feature engineering driven" to "representation learning driven" (Vercelli, 2024). At the same time, knowledge representation and reasoning, as an important research direction in the early development of artificial intelligence, focuses on how to represent knowledge in a formal way and make inferences through logical rules, providing an important theoretical foundation for applications such as expert systems and semantic networks. Although data-driven methods have gradually become dominant in recent years, knowledge representation and symbolic reasoning still have irreplaceable value in interpretable artificial intelligence and complex decision-making systems. Overall, the core research content of artificial intelligence includes not only data-driven learning algorithm systems, but also cognitive mechanisms centered on knowledge expression and logical reasoning, and extends to multiple intelligence levels such as perception, cognition and decision-making. The various research directions form a relationship of mutual support and synergistic evolution in theory and technology, jointly promoting the

transformation of artificial intelligence from dedicated systems for single tasks to comprehensive intelligent systems with multiple functions (Lainjo, 2024). With the continuous deepening of interdisciplinary integration, the research boundaries of artificial intelligence are continuously expanding to more forward-looking fields such as multimodal learning, human-machine collaboration and general intelligence, providing a more solid theoretical foundation and methodological support for the continuous innovation and system upgrading of intelligent technologies.

2.3. The theoretical foundation of artificial intelligence

The development of artificial intelligence not only relies on continuous breakthroughs in engineering technology, but is also deeply rooted in the theoretical system of multidisciplinary integration. As a comprehensive discipline, the theoretical foundation of artificial intelligence involves a wide range of research fields such as mathematical theory, statistical learning theory, information theory, control theory and cognitive science. These theories together constitute an important supporting framework for the design of artificial intelligence algorithms and the construction of intelligent systems. In this system, mathematical foundation is regarded as the most core theoretical pillar. Among them, linear algebra provides key formal tools for matrix operations and feature representation in neural networks. Vector space theory, feature decomposition methods and matrix transformation mechanisms constitute important mathematical foundations for the computational structure of deep models. At the same time, probability theory and mathematical statistics provide theoretical basis for uncertainty modeling and statistical inference, and play an important role in the method system of Bayesian inference, hidden Markov models and probabilistic graphical models (Sharma, 2023). In addition, optimization theory lays the algorithmic foundation for model parameter learning. Among them, gradient descent and convex optimization methods have become the core technical paths in deep learning model training. From a theoretical perspective, the learning process in artificial intelligence can essentially be regarded as an optimization problem in a high-dimensional parameter space. Its convergence and stability largely depend on the support of rigorous and systematic mathematical theory.

The development of artificial intelligence not only relies on continuous breakthroughs in engineering technology, but is also deeply rooted in the theoretical system of multidisciplinary integration. As a comprehensive discipline, the theoretical foundation of artificial intelligence involves a wide range of research fields such as mathematical theory, statistical learning theory, information theory, control theory and cognitive science. These theories together constitute an important supporting framework for the design of artificial intelligence algorithms and the construction of intelligent systems. In this system, mathematical foundation is regarded as the most core theoretical pillar. Among them, linear algebra provides key formal tools for matrix operations and feature representation in neural networks. Vector space theory, feature decomposition methods and matrix transformation mechanisms constitute important mathematical foundations for the computational structure of deep models. At the same time, probability theory and mathematical statistics provide theoretical basis for uncertainty modeling and statistical inference, and play an important role in the method system of Bayesian inference, hidden Markov

models and probabilistic graphical models (Taeihagh, 2025). In addition, optimization theory lays the algorithmic foundation for model parameter learning. Among them, gradient descent and convex optimization methods have become the core technical paths in deep learning model training. From a theoretical perspective, the learning process in artificial intelligence can essentially be regarded as an optimization problem in a high-dimensional parameter space. Its convergence and stability largely depend on the support of rigorous and systematic mathematical theory.

3. Key Technological Advances in Artificial Intelligence

3.1. Machine Learning

Machine learning is one of the most core branches of artificial intelligence. Its basic idea is to enable computer systems to automatically learn the potential patterns in data without explicit rule programming through data-driven methods, and to complete prediction, classification or decision-making tasks accordingly. Unlike traditional algorithms that rely on manual rule design, machine learning emphasizes that the model continuously improves its performance through parameter optimization and pattern abstraction during sample training. Its essence can be regarded as the process of approximating an unknown function under given data distribution conditions. From a methodological perspective, the construction of machine learning models usually includes three key links: model selection, parameter estimation and performance evaluation. Model selection mainly involves the determination of hypothesis space and control of model complexity. Parameter estimation relies on optimization algorithms to solve the objective function. Performance evaluation is carried out by quantitative analysis of the model's predictive ability through methods such as cross-validation, error index or confusion matrix (Corrêa et al., 2023). In high-dimensional data environments, models are prone to overfitting. Therefore, regularization methods, model pruning and early stopping strategies are widely used to improve the generalization ability and stability of the model, thereby ensuring the reliable performance of the machine learning system on unknown data.

From the perspective of learning paradigms, machine learning can generally be divided into three basic types: supervised learning, unsupervised learning, and reinforcement learning. Supervised learning trains a model by minimizing the error between the predicted value and the true value under the condition of given input-output sample pairs. Its typical applications include classification and regression problems. Representative algorithms include linear regression, logistic regression, support vector machine, and neural network. Its theoretical basis mainly comes from the principle of risk minimization and statistical learning theory. Unsupervised learning achieves pattern discovery and feature representation by mining the internal structure of data under the condition of lack of label information. Typical methods include cluster analysis, principal component analysis, and autoencoders. These methods can reveal the potential distribution structure or low-dimensional representation of data and play an important role in tasks such as data preprocessing, anomaly detection, and feature extraction. Reinforcement learning is a learning paradigm based on interaction mechanism. Its core idea is to continuously

optimize the behavior strategy through continuous interaction between the agent and the environment during the trial and error process (Xu et al., 2024). Reinforcement learning is typically built on the Markov decision process framework, aiming to maximize long-term cumulative rewards through iterative updates of the value function or policy function. Representative methods include Q-learning, policy gradient methods, and deep reinforcement learning, and have achieved significant results in fields such as game theory decision-making, robot control, and autonomous driving.

With the exponential growth of data scale and the continuous improvement of computing power, the machine learning method system is gradually evolving from traditional models with shallow structures to multi-layered, highly complex deep architectures. This evolution not only significantly enhances the expressive power of the model, but also enables the algorithm to more effectively characterize the complex nonlinear feature relationships in high-dimensional data. In this process, ensemble learning methods effectively reduce the variance and bias of a single model by constructing multiple base learners and strategically combining them, achieving synergistic optimization of prediction performance and model stability, and showing strong generalization ability and robustness in structured data analysis tasks (Ricciardi et al., 2025). At the same time, the continuous improvement of large-scale distributed computing frameworks and parallel training mechanisms provides important technical support for the efficient training and large-scale deployment of machine learning algorithms in massive data environments. Relying on distributed storage systems, parameter server architectures and high-performance computing resources, the computational bottleneck in the model training process is significantly alleviated, thereby promoting the widespread application of machine learning in industrial scenarios. Overall, machine learning has gradually built a key technical link between data resources and intelligent applications, playing a fundamental and pivotal role in data value mining and intelligent decision support systems.

3.2. Deep Learning

Deep learning, as an important branch of machine learning (its structural relationship is shown in Figure 1), has become one of the core driving forces for breakthroughs in artificial intelligence technology in recent years. Its basic idea lies in constructing a multi-layered nonlinear neural network structure to achieve hierarchical representation and automatic feature extraction of complex data. Unlike traditional machine learning methods that rely on manually designed features, deep learning emphasizes an end-to-end training mechanism, enabling the model to automatically learn high-level abstract features driven by large-scale data, thereby significantly improving the performance of complex pattern recognition tasks. Structurally, deep learning models typically consist of an input layer, multiple hidden layers, and an output layer. The hidden layers use nonlinear activation functions to map features layer by layer, allowing the model to gradually learn more abstract and higher-level representations. According to the general approximation theorem, multi-layered neural networks can theoretically approximate any continuous function, providing an important mathematical basis for the powerful expressive capabilities of deep learning models. During model training, parameters are typically calculated using the backpropagation algorithm, and weights are continuously updated using stochastic

gradient descent and its improved optimization algorithms to minimize the loss function and improve the model's predictive performance.

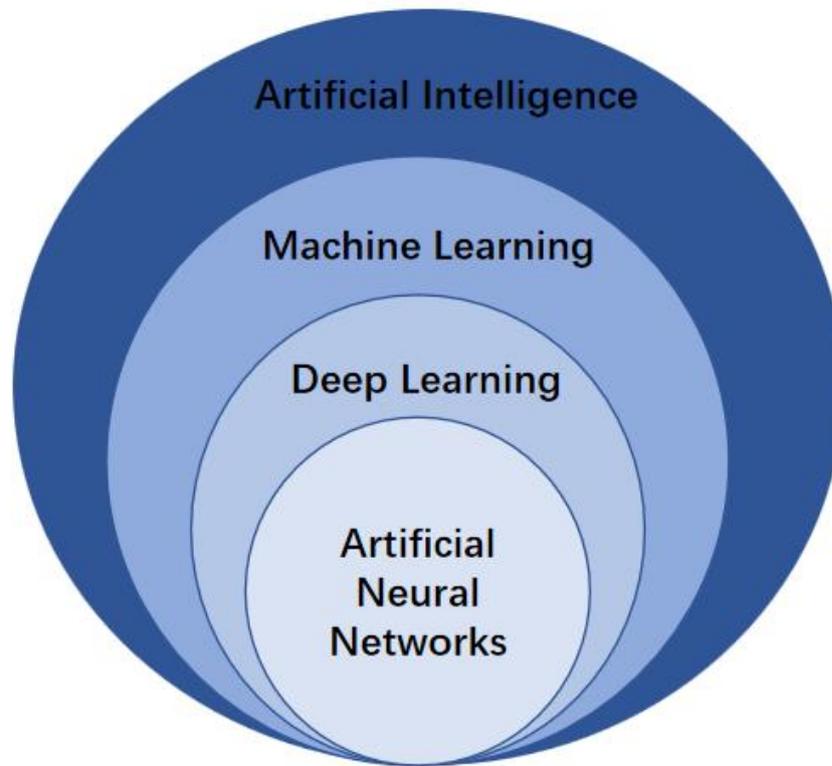


Figure 1. Hierarchical framework diagram of artificial intelligence technology

In terms of specific model structure, Convolutional Neural Network (CNN) is an important representative of the breakthroughs achieved by deep learning in the field of computer vision. CNN can effectively extract local spatial features through convolution operation and weight sharing mechanism, while significantly reducing the number of model parameters, thereby improving training efficiency and enhancing the generalization ability of the model. Therefore, this structure has shown excellent performance in tasks such as image classification, object detection and medical image analysis. On the other hand, Recurrent Neural Network (RNN) is mainly used to process sequential data with time dependence. It realizes the modeling of context information through recurrent connection structure and plays an important role in tasks such as natural language processing and speech recognition (Whig, 2025). However, traditional RNN is prone to gradient vanishing and gradient explosion problems during long sequence training. To solve this limitation, researchers have proposed improved structures such as Long Short-Term Memory Network (LSTM) and Gated Recurrent Unit (GRU), which enhance the model's ability to express long-term dependencies by introducing gating mechanism, thereby significantly improving the sequence modeling effect.

In recent years, the Transformer architecture based on attention mechanisms has become a significant milestone in the development of deep learning. Unlike RNNs, which rely on sequential computation, Transformers achieve global dependency modeling through self-attention, significantly improving parallel computing efficiency and model training speed. This architecture has spurred the development of large-scale pre-trained models in natural language processing and

propelled artificial intelligence research into a new phase centered on large-scale models. The success of the Transformer architecture has not only changed the traditional sequence modeling paradigm but also laid an important foundation for multimodal learning and cross-domain transfer learning. Overall, deep learning, relying on multi-layered neural network structures and automated feature representation mechanisms, has significantly enhanced the task processing capabilities of artificial intelligence systems in complex scenarios and high-dimensional data environments. With continuous innovation in model architecture, improved computing resources, and the development of cross-modal fusion technologies, deep learning is continuously driving the evolution of artificial intelligence from dedicated systems for single tasks to multifunctional integrated intelligent systems, and will play a sustained and far-reaching role in the future development towards higher levels of autonomy and general intelligence.

3.3. Large Model and Data-Driven Approach

With the continuous evolution of deep learning technology, artificial intelligence has gradually entered a development stage with "large-scale models" as its core feature. The so-called large model usually refers to a deep neural network model with a massive number of parameters, trained on large-scale data and capable of cross-task generalization. Its basic process is shown in Figure 2. Such models usually rely on the training paradigm of "pre-training and fine-tuning", are pre-trained in multi-domain corpora or multi-modal data environments, and are adaptively optimized on specific tasks to achieve wide application transfer. From the perspective of technical path, the development of large models is largely based on self-supervised learning mechanisms. By designing pre-training tasks that do not require manual annotation, the model can learn general representations in massive unlabeled data. For example, in the field of natural language processing, the model can perform language modeling through context prediction or mask word prediction, thereby mastering language structure and semantic rules; in the field of computer vision, feature representation learning can be achieved through image reconstruction or contrastive learning (Kashefi et al., 2024). Meanwhile, the development of large-scale models also exhibits a significant scale effect; that is, as the size of model parameters, the size of training data, and computing resources expand simultaneously, model performance shows a relatively stable upward trend. This pattern is often referred to as the law of scale, which reveals the functional relationship between model performance and parameter size, data size, and computing resources. Scaling not only enhances the model's ability to express complex patterns but also promotes cross-task transfer capabilities, enabling the model to demonstrate a certain degree of generalization ability even on tasks without explicit training. This ability is considered an important indicator of artificial intelligence moving towards a higher level of intelligence.

At the application level, the rise of large models has promoted the development of multi-domain technology integration and cross-modal learning. Large models based on a unified architecture can process multiple types of data such as text, images, speech and even video at the same time, realize the collaborative modeling and unified representation of multimodal information, and thus significantly improve the comprehensive perception and understanding capabilities of intelligent systems in complex scenarios. For example, through joint training of text and images, the model can complete tasks such as visual question answering, image and text

treatment decision support and health management. Its core goal is to improve the efficiency of medical services and the accuracy of diagnosis and treatment through data-driven intelligent analysis, thereby optimizing the allocation of medical resources and promoting the development of precision medicine. In the field of medical image analysis, artificial intelligence technology has made significant progress. Medical image data usually has the characteristics of large scale, high dimension and complex structure. The traditional manual interpretation method that relies on physician experience is not only inefficient, but also inevitably has certain subjective differences. Deep learning models based on convolutional neural networks can automatically complete image feature extraction and hierarchical representation learning, realize accurate identification and segmentation of lesion areas, and show high accuracy and stability in multiple disease detection tasks (Li et al., 2024). By constructing a large-scale labeled medical image dataset and carrying out multi-center clinical verification, artificial intelligence systems have, to a certain extent, the ability to assist in early disease screening and improve diagnostic consistency, thereby significantly improving the objectivity and standardization of clinical image analysis.

Artificial intelligence also plays a crucial role in clinical decision support. Intelligent systems based on Clinical Decision Support Systems (CDSS) can integrate multi-source information, including electronic medical records, laboratory test indicators, and historical case data, to comprehensively assess patient health and predict risks. Leveraging the pattern recognition and statistical inference capabilities of machine learning models, these systems can model and analyze disease progression trends, readmission risks, and individualized treatment strategies, thereby improving the accuracy and foresight of clinical decisions. Simultaneously, AI demonstrates significant value in drug development and precision medicine. Traditional drug development is often characterized by long cycles and high costs, while AI-based molecular structure prediction and virtual screening technologies can significantly shorten candidate drug screening time. For example, using deep learning models to predict protein structures and molecular interactions can improve drug design efficiency. Furthermore, by systematically analyzing genomic data and constructing individualized risk prediction models, AI can extract key features from massive amounts of genomic data, providing crucial support for individualized treatment plans in precision medicine.

Despite the broad application prospects of artificial intelligence in the medical field, its practical promotion and in-depth application still face multiple challenges. From a data governance perspective, medical data often contains highly sensitive personal privacy information. Therefore, ensuring patient privacy and data security while achieving data sharing and value mining is a crucial issue that must be balanced. The tension between data openness and privacy protection makes the design of systems and the construction of technical protection mechanisms particularly critical. Furthermore, regarding model performance and algorithmic fairness, the representativeness and quality of training data directly affect the system's generalization ability and stability. If the data sources have issues such as regional concentration, unbalanced sample structure, or potential biases, the model's performance may differ across different populations or clinical scenarios, leading to uneven treatment outcomes and the risk of algorithmic bias. Therefore, building a multi-center, diverse, and high-quality medical data system is a fundamental

condition for improving the reliability of AI-powered medical systems. Overall, against the backdrop of the deepening application of AI in medicine, achieving synergistic advancement of technological innovation and institutional governance has become a key prerequisite for its sustainable development. By establishing robust data security mechanisms, ethical standards, and regulatory frameworks, and promoting the deep integration of artificial intelligence technology with medical knowledge and multi-source data, AI is evolving from an early single-task auxiliary tool into a comprehensive intelligent support system integrating multiple modules. This is driving the continuous transformation of healthcare services towards data-driven, refined, and personalized approaches. With the optimization of algorithm performance, the improvement of data governance systems, and the continuous strengthening of interdisciplinary collaboration mechanisms, AI will play a more fundamental and strategic supporting role in the modern healthcare system, providing continuous impetus for improving healthcare quality and innovating healthcare models.

4.2. Education Sector

The application of artificial intelligence technology in the field of education is profoundly changing the traditional teaching model and learning methods. With the continuous enrichment of digital educational resources and the widespread popularity of online learning platforms, educational data is gradually showing characteristics of scale and diversification. Against this background, artificial intelligence provides important technical support for personalized teaching, intelligent assessment and educational management optimization by systematically analyzing and modeling learning behavior data. Its core goal is to improve teaching efficiency and learning effect, thereby promoting the development of a learner-centered precision education model. In terms of specific applications, personalized learning systems have become one of the important directions of artificial intelligence education applications^[18]. Traditional classroom teaching usually adopts a uniform pace and standardized teaching content, which is difficult to fully respond to the differentiated needs of students in terms of knowledge base, learning pace and cognitive style. Relying on machine learning methods and knowledge tracking models, intelligent teaching systems can analyze students' learning behavior data and answer performance in real time, assess their knowledge mastery, and dynamically adjust the learning content structure and task difficulty accordingly. This adaptive learning mechanism based on data modeling helps to improve the matching degree between teaching resources and learning needs, thereby enhancing learning motivation and improving overall learning efficiency.

Building upon this foundation, the application value of Intelligent Tutoring Systems (ITS) in educational practice is becoming increasingly prominent. These systems achieve continuous monitoring and dynamic feedback of the learning process by constructing student models, domain knowledge models, and teaching strategy models. By leveraging artificial intelligence algorithms to analyze student error types, learning paths, and knowledge gaps, the system can generate targeted explanations and prompts, thus simulating personalized tutoring to some extent. Related research indicates that intelligent tutoring systems demonstrate good teaching effects in mathematics, programming, and language learning, and have potential advantages in alleviating the problem of uneven distribution of educational resources. Furthermore, artificial intelligence

also shows significant application prospects in teaching evaluation and educational management. For example, automatic scoring technology, combined with natural language processing and deep learning models, performs semantic understanding and score prediction on subjective question answers, effectively improving assessment efficiency and scoring consistency. Simultaneously, by clustering and trend prediction of learning behavior data, educational management departments can identify potential learning risk groups and implement targeted interventions. This data-driven educational decision-making model helps optimize resource allocation and improve the scientific level of educational governance.

Despite the vast potential of artificial intelligence in education, its practical application still faces numerous challenges. From a data governance perspective, learning behavior data often involves issues of personal privacy and the protection of minors' information; therefore, strict adherence to relevant laws, regulations, and ethical norms is essential during data collection, storage, and use. Regarding algorithmic fairness, model bias can negatively impact learning evaluation results, especially in high-risk exams or critical assessment scenarios, making it crucial to enhance model transparency and interpretability. Furthermore, the deep integration of AI technology may affect teachers' professional roles and teaching autonomy; therefore, constructing a teaching model centered on human-machine collaboration has become a key research direction. From an overall development perspective, AI is driving the transformation of education models from traditional "uniform teaching" to a learner-centered "personalized learning" system. Leveraging data analysis and intelligent modeling technologies, education systems can more accurately identify learning needs, assess teaching effectiveness, and optimize resource allocation. In the future, with the further integration of algorithmic innovation and educational theory, and the continued deepening of interdisciplinary collaboration mechanisms, AI is expected to play a more profound and systematic role in promoting educational equity, improving teaching quality, and building a lifelong learning system.

4.3. Financial Sector

Artificial intelligence technology is increasingly widely used in the financial field and has gradually become an important technological foundation for promoting the digital transformation of the financial industry. The financial industry has the characteristics of large data scale, high transaction frequency and strong risk sensitivity, which provides rich application scenarios for the deployment and optimization of intelligent algorithms. By systematically modeling and analyzing structured and unstructured financial data, artificial intelligence has shown significant advantages in risk management, credit assessment, financial market prediction and intelligent services, thereby improving the efficiency of financial services and the scientific nature of decision-making. In the field of risk control and credit assessment, machine learning models have become an important tool for financial institutions to improve risk identification capabilities and pricing accuracy. By comprehensively analyzing historical transaction data, user behavior characteristics and external credit information, relevant algorithms can construct a credit scoring system and predict the probability of default. Compared with traditional rule-based or statistical regression methods, models such as gradient boosting trees, random forests and deep neural networks have stronger expressive power and can characterize complex nonlinear relationships between

variables, thereby significantly improving risk prediction performance^[19]. In personal credit and micro and small enterprise financing scenarios, this type of technology not only expands the coverage of financial services, but also promotes the refinement and differentiation of risk pricing. However, during model training, issues such as data structure imbalance, sample bias, or inappropriate feature selection may lead to systematic shifts in the scoring results. Therefore, it is necessary to strengthen the model validation mechanism and the algorithm fairness evaluation system to improve the reliability of risk assessment results.

In the fields of financial market forecasting and quantitative investment, artificial intelligence methods also demonstrate broad application prospects. Relying on technologies such as time series analysis, deep learning, and reinforcement learning, models can extract potential structural features from historical price fluctuations and macroeconomic indicators, providing data support for asset allocation optimization and trading strategy construction. For example, models based on recurrent neural networks or Transformer structures have strong pattern-capturing capabilities in financial time series modeling, helping to improve the accuracy of short-term price fluctuation predictions; while reinforcement learning methods achieve a dynamic balance between return objectives and risk constraints by simulating market environments and strategy iteration processes. Meanwhile, in the field of robo-advisory, robo-advisory systems provide investors with personalized asset allocation solutions by building customer profiles and risk preference models. The system can automatically generate investment portfolio recommendations and make dynamic adjustments based on the user's financial situation, investment objectives, and risk tolerance, thereby reducing service costs while improving investment decision-making efficiency. Furthermore, artificial intelligence also plays a crucial role in fraud prevention and anomaly detection. By monitoring and recognizing transaction behavior in real time, machine learning models can identify potential fraudulent activities and abnormal transaction behaviors. Among them, graph neural networks and other methods have shown significant advantages in the analysis of complex transaction network structures, which helps to reveal hidden relationships and risk propagation paths, and significantly enhances the security and stability of financial system operations in electronic payment and cross-border transaction scenarios.

Despite significant progress in the financial sector, artificial intelligence still faces numerous challenges in its application and promotion. From a data governance perspective, financial data often contains highly sensitive information, thus data security and privacy protection must be a primary focus during data collection, storage, and sharing. In terms of regulatory compliance, the interpretability of model decision-making results is crucial for key business operations such as credit approval and risk pricing; therefore, relevant models need high transparency and auditability to meet regulatory compliance requirements. Furthermore, the stability and robustness of algorithmic models need further improvement under extreme market conditions or systemic shocks to prevent the accumulation of potential systemic risks. From an overall development trend perspective, AI is driving the financial industry to gradually shift from a traditional experience-driven decision-making model to a decision-making system centered on data analysis and intelligent algorithms. Through efficient processing and in-depth mining of large-scale financial data, AI technology has significantly enhanced risk management capabilities and service

efficiency, and provided a key technological foundation for the formation of new financial service models. In the future, with the continuous improvement of the regulatory framework and the sustained enhancement of algorithmic capabilities, AI will play a more fundamental and strategic role in promoting the sound operation of the financial system, optimizing resource allocation, and promoting inclusive finance.

5. Ethical and Social Impact

With the widespread application of artificial intelligence technology in multiple fields such as medicine, finance, education, industry and public governance, its social impact is becoming increasingly prominent. While significantly improving production efficiency and optimizing social governance structure, the development of artificial intelligence has also triggered a series of complex and profound ethical and social issues. These issues not only involve the security, reliability and controllability of the technology system itself, but also further relate to social fairness, privacy protection, human subject status and the stability of social value system. Figure 3 systematically shows the core ethical issues involved in the social application of artificial intelligence and their interrelationships. By constructing an analytical framework covering dimensions such as data privacy protection, algorithm fairness, responsibility attribution, social structural changes and value alignment, we can more systematically understand the inherent logical connection between the artificial intelligence risk generation mechanism and governance path. Among them, data governance issues are particularly critical^[20]. Artificial intelligence systems usually rely on massive data resources during model training and performance optimization. These data often contain sensitive content such as personal identity information, behavioral trajectories and health records. In the process of data collection, storage, sharing and cross-domain circulation, if there is a lack of sound security protection and compliance mechanisms, it is easy to trigger risks such as privacy leakage, data abuse and network attacks. While technologies such as differential privacy, federated learning, and homomorphic encryption offer solutions for privacy protection to some extent, the core challenge in the governance of artificial intelligence remains how to ensure the reasonable use of data value while protecting individual privacy rights.

Meanwhile, the fairness of algorithmic decision-making has gradually become an important topic in AI ethics research. The prediction results of AI models largely depend on the distribution structure and sample characteristics of the training data. When the data itself has structural biases or insufficient sample representativeness, the model may solidify or even amplify existing inequalities during the learning process, leading to potential discriminatory consequences in scenarios such as credit approval, recruitment screening, or judicial assistance decision-making. To mitigate this problem, researchers have proposed various technical approaches, including data resampling, constraint optimization, and the construction of fairness indicators, to reduce the negative impact of algorithmic bias. However, defining fairness standards in diverse social contexts and translating them into actionable model optimization objectives still faces significant theoretical and practical challenges. Furthermore, the rapid development of AI technology is profoundly impacting labor structures and the social division of labor. The widespread application

of automated and intelligent systems in manufacturing, service, and management sectors risks the replacement of some repetitive and rule-based jobs, leading to a restructuring of employment and changes in skill requirements. While technological progress creates new job types and industrial opportunities, structural unemployment and skills mismatch may exacerbate social inequality in the short term. Therefore, guiding the workforce to complete skills transformation and career upgrading through education system reform, vocational skills training, and public policy adjustments has become one of the important issues in social governance in the era of artificial intelligence.

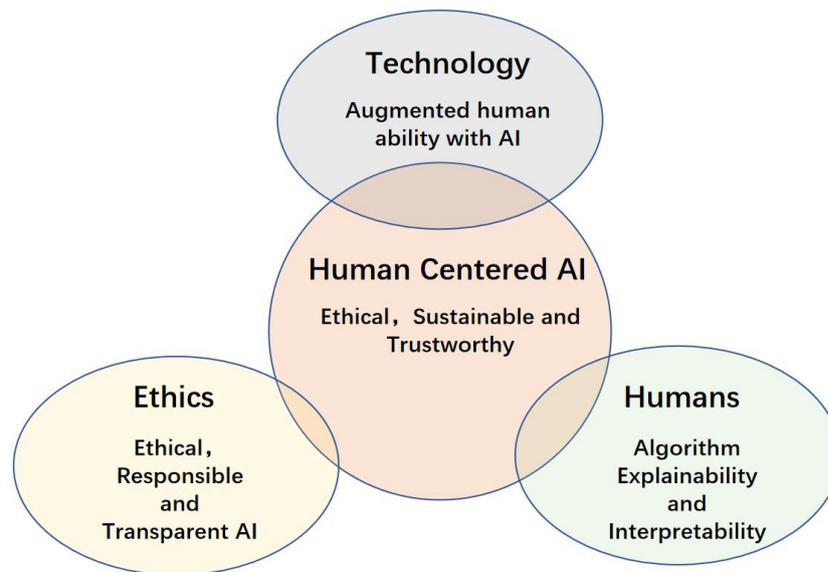


Figure 3. The core ethical issues of artificial intelligence and their interrelationships

At a deeper institutional and value level, the widespread application of artificial intelligence systems also brings challenges in terms of responsibility definition and technology governance. In particular, large-scale models based on deep learning often exhibit significant "black box" characteristics, with their internal decision-making logic difficult to explain intuitively. When system judgments erroneous or cause real-world harm, the division of responsibility among algorithm developers, system deployment organizations, and end-users still lacks unified standards, which is particularly prominent in high-risk application scenarios such as autonomous driving and medical assisted diagnosis. Therefore, strengthening research on model interpretability, improving decision-making transparency, and establishing a clear legal responsibility framework have become crucial links in promoting the standardized development of AI. Simultaneously, AI technology itself has a distinct dual nature: on the one hand, it can enhance social governance capabilities and information production efficiency; on the other hand, it can also be used for activities such as generating false information, spreading deepfakes, and automated cyberattacks, impacting the public opinion environment and social trust systems. With the development of generative models, the threshold for producing false content continues to decrease, increasing the difficulty of governance. At a more macro level, as AI systems gradually acquire autonomous decision-making and content generation capabilities, the boundaries of human-machine relationships are constantly being reshaped. How to ensure that the behavior of AI systems aligns with human social values and norms has become an important direction for

value alignment research. Current research attempts to improve the ethical consistency of model outputs by introducing human feedback mechanisms and reinforcement learning alignment strategies. However, due to the cultural differences and context-dependent nature of values, constructing a value alignment mechanism that is both universal and dynamically adaptable remains an important topic that needs to be explored in depth in the study of artificial intelligence ethics.

6. Future Development Trends

With the continuous optimization of algorithm model architecture, the exponential improvement of computing power and the continuous expansion of data resources, artificial intelligence is gradually entering a new development stage with large-scale pre-trained models, cross-modal fusion and intelligent agent systems as the core features. Unlike the early development path that mainly relied on the performance improvement of a single algorithm, the current evolution of artificial intelligence reflects a comprehensive trend of systematic integration, cross-domain collaboration and parallel development of technology and system. From the perspective of overall technology, the rise of large-scale pre-trained models marks the transformation of artificial intelligence from a task-specific paradigm to a general framework. By implementing self-supervised learning on massive data, the model can obtain cross-task transfer and multi-scenario adaptation capabilities. The "pre-training-fine-tuning" paradigm significantly reduces the dependence on manually labeled data and improves the generalization performance of the model to a certain extent. The future development of models will pay more attention to structural optimization and parameter efficiency, and improve resource utilization through model compression, knowledge distillation and efficient computing power scheduling mechanisms. At the same time, the integration of reinforcement learning methods and symbolic reasoning mechanisms is also regarded as an important direction for improving the logical reasoning and complex planning capabilities of the system. Although true general artificial intelligence is still in the exploratory stage, the development of cross-task transfer capabilities and multimodal representation learning has provided a key technological foundation for it.

At the data processing and cognitive level, multimodal fusion is gradually becoming an important path for artificial intelligence to improve its environmental understanding capabilities. Real-world information often exists in the form of interwoven text, images, voice, video, and data from multiple sensors, making it difficult for a single-modal model to achieve a comprehensive and accurate environmental representation. By constructing a unified semantic representation space and achieving cross-modal alignment, AI systems can perform comprehensive perception, information integration, and reasoning decision-making in complex situations, thereby improving the overall robustness and adaptability of the system. In highly complex application scenarios such as medical diagnosis, autonomous driving, and intelligent manufacturing, multimodal collaboration mechanisms not only significantly improve decision-making accuracy but also promote interdisciplinary integration and innovation. Meanwhile, with the widespread deployment of IoT devices and smart terminals, traditional cloud-centric centralized computing models are gradually facing latency and bandwidth bottlenecks. Against this backdrop, edge

computing and distributed intelligent systems have become important development directions. By deploying some inference tasks to network edge nodes, data transmission pressure can be effectively reduced and the system's real-time response capability improved. Furthermore, with the help of distributed training mechanisms such as federated learning, multiple devices can achieve collaborative model updates while ensuring data privacy, thereby gradually building an intelligent computing infrastructure that operates collaboratively across the cloud, edge, and terminal.

In terms of human-machine relationships and social governance, the development trend of artificial intelligence is gradually shifting from simple automation replacement to a collaborative intelligence model centered on enhancing human capabilities. This paradigm emphasizes that AI, as a cognitive support and decision-making aid, forms a complementary collaborative relationship with humans. For example, in highly specialized fields such as medicine, law, and scientific research, AI systems can undertake data analysis, pattern recognition, and knowledge discovery tasks, while value judgments and final decisions are still made by human agents. To achieve efficient human-machine collaboration, system design needs to further enhance interpretability, interactivity, and context awareness, improve user trust through transparency and feedback mechanisms, and enhance the system's understanding of human intentions and social context. Furthermore, from a sustainable development perspective, as model scale continues to expand, the energy consumption resulting from training and deployment is increasingly attracting attention. Through technical approaches such as model pruning, parameter sharing, knowledge distillation, and low-power hardware architecture design, computational resource consumption can be reduced while maintaining performance stability, promoting the development of green AI. At the institutional and governance level, future AI systems will place greater emphasis on standardization and ethically embedded design. By introducing risk assessment mechanisms and value alignment strategies during model development, ethical principles will be reflected upfront in the technical process. Simultaneously, international cooperation and the development of technical standards systems will become important trends. By establishing unified evaluation frameworks and compliance mechanisms, the safety, transparency, and social credibility of AI systems will be improved, thereby achieving a long-term and stable dynamic balance between technological innovation and social responsibility.

Author Contributions:

All authors have read and agreed to the published version of the manuscript.

Funding:

This research received no external funding.

Institutional Review Board Statement:

Not applicable.

Informed Consent Statement:

Not applicable.

Data Availability Statement:

Not applicable.

Conflict of Interest:

The authors declare no conflict of interest.

References

- Batool, A., Zowghi, D., & Bano, M. (2025). AI governance: A systematic literature review. *AI and Ethics*, 5(3), 3265 – 3279.
- Corrêa, N. K., Galvão, C., Santos, J. W., et al. (2023). Worldwide AI ethics: A review of 200 guidelines and recommendations for AI governance. *Patterns*, 4(10), 100857.
- Kashefi, P., Kashefi, Y., & Ghafouri Mirsarai, A. H. (2024). Shaping the future of AI: Balancing innovation and ethics in global regulation. *Uniform Law Review*, 29(3), 524 – 548.
- Kattnig, M., Angerschmid, A., Reichel, T., et al. (2024). Assessing trustworthy AI: Technical and legal perspectives of fairness in AI. *Computer Law & Security Review*, 55, 106053.
- Lahusen, C., Maggetti, M., & Slavkovik, M. (2024). Trust, trustworthiness and AI governance. *Scientific Reports*, 14(1), 20752.
- Lainjo, B. (2024). The role of artificial intelligence in achieving the United Nations sustainable development goals. *Journal of Sustainable Development*, 17(5), 30.
- Li, Y., Wu, B., Huang, Y., et al. (2024). Developing trustworthy artificial intelligence: Insights from research on interpersonal, human-automation, and human-AI trust. *Frontiers in Psychology*, 15, 1382693.
- McCormack, L., & Bendeche, M. (2024). Ethical AI governance: Methods for evaluating trustworthy AI. *arXiv preprint arXiv:2409.07473*.
- Ricciardi Celsi, L., & Zomaya, A. Y. (2025). Perspectives on managing AI ethics in the digital age. *Information*, 16(4), 318.
- Russell, S., & Norvig, P. (2021). *Artificial intelligence: A modern approach* (4th US ed.). Pearson.
- Sharma, S. (2023). Trustworthy artificial intelligence: Design of AI governance framework. *Strategic Analysis*, 47(5), 443 – 464.
- Sistla, S. (2024). AI with integrity: The necessity of responsible AI governance. *Journal of Artificial Intelligence & Cloud Computing*, 3(5), 2 – 3.
- Taeihagh, A. (2025). Governance of generative AI. *Policy and Society*, 44(1), 1 – 22.
- Vercelli, A. (2024). United Nations, artificial intelligences and regulations: Analysis of the General Assembly AI resolutions and the final report of the advisory body on AI (short paper). In *Proceedings of BEWARE@AI* (pp. 99 – 106).
- Whig, D. P. (2025). Ethical AI governance: A framework for ensuring transparency, fairness, and accountability. *Journal of Healthcare AI and ML*, 12, 12.

Xu, J., Lee, T., & Goggin, G. (2024). AI governance in Asia: Policies, praxis and approaches. *Communication Research and Practice*, 10(3), 275 – 287.

License: Copyright (c) 2026 Author.

All articles published in this journal are licensed under the Creative Commons Attribution 4.0 International License (CC BY 4.0). This license permits unrestricted use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited. Authors retain copyright of their work, and readers are free to copy, share, adapt, and build upon the material for any purpose, including commercial use, as long as appropriate attribution is given.