



# **Explainability and Interpretability in Generative AI Agents**

**Anantharaman Janakiraman**<sup>[0009-0008-3641-0788]</sup>

**Independent Researcher**

[anantharaman.j@gmail.com](mailto:anantharaman.j@gmail.com)

**Vol. 7 No. 7 (2025): IJSTC**

## **Abstract**

The rapid advancement of generative artificial intelligence (AI) agents, including large language models, multimodal systems, and autonomous decision-making agents, has significantly expanded their adoption across critical domains such as healthcare, finance, education, and governance. While these systems demonstrate remarkable capabilities in natural language generation, reasoning, and task automation, their increasing autonomy and complexity have raised substantial concerns regarding transparency, accountability, and trustworthiness. This study examines the critical role of explainability and interpretability in generative AI agents, emphasizing their importance for enabling human understanding, regulatory compliance, and ethical deployment. The abstract highlights how explainable AI (XAI) techniques contribute to making the internal decision processes and output generation mechanisms of generative agents more transparent, thereby supporting responsible human-AI collaboration and informed oversight. The abstract further considers emerging methods specifically designed for large language models and generative agents, including chain-of-thought prompting, rationale generation, uncertainty estimation, and self-explanation mechanisms, which aim to provide more faithful and context-aware explanations.

## **Introduction**

Generative artificial intelligence (AI) agents, including large language models, multimodal generative systems, and autonomous conversational and task-oriented agents, have rapidly transitioned from experimental research tools to widely deployed technologies across diverse sectors. These systems are now capable of producing human-like text, generating images and code, synthesizing speech, and supporting complex decision-making workflows.



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

Their increasing integration into high-impact domains such as healthcare, finance, education, law, and public administration has significantly amplified both their potential benefits and associated risks. While generative AI agents offer substantial gains in productivity, scalability, and personalization, their internal decision-making processes are often opaque, making it difficult for users, developers, and regulators to understand how specific outputs are generated. This lack of transparency has intensified concerns related to trust, accountability, safety, and ethical deployment, thereby motivating a growing emphasis on explainability and interpretability as core requirements for responsible generative AI systems.

## **Generative AI Agents and System Complexity**

Modern generative AI agents are built on large-scale neural architectures that contain billions of parameters and are trained on massive, heterogeneous datasets. These models exhibit emergent behaviors that are not explicitly programmed but arise from complex interactions within deep neural networks. While such complexity enables impressive generative capabilities, it also significantly reduces human interpretability. Unlike traditional rule-based systems or simpler machine learning models, generative agents often operate as black boxes, where the relationship between inputs, internal representations, and outputs is not readily understandable. This opacity complicates debugging, validation, and error analysis, particularly in safety-critical environments. As generative agents become more autonomous and are increasingly entrusted with multi-step reasoning, tool use, and real-world decision support, the need for mechanisms that provide insight into their reasoning and generation processes becomes increasingly urgent.

## **Importance of Explainability and Interpretability**

Explainability and interpretability serve as foundational pillars for building trustworthy generative AI agents. Interpretability refers to the extent to which a human can directly understand the internal structure and behavior of a model, while explainability focuses on methods and tools that provide post hoc or built-in explanations for model decisions and outputs. In the context of generative AI, these concepts are particularly challenging due to the probabilistic and high-dimensional nature of language and multimodal generation. Nevertheless, meaningful explanations are essential for enabling users to assess the reliability, relevance, and appropriateness of generated outputs. In professional settings, such as clinical documentation, financial analysis, and legal drafting, practitioners must be able to understand and justify AI-assisted decisions. Explainability thus supports informed human oversight, enabling users to critically evaluate AI outputs rather than accepting them uncritically.



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

## **Trust, Accountability, and Human-AI Collaboration**

Trust is a critical determinant of successful human-AI collaboration. Users are more likely to rely on generative AI agents when they perceive them as transparent, predictable, and aligned with human values and goals. Explainable systems help build trust by making the reasoning and assumptions behind generated outputs more visible. Accountability is also closely linked to explainability. When AI agents influence consequential decisions, it is necessary to determine responsibility for outcomes and to trace how specific recommendations or outputs were produced. Without explainability, it becomes difficult to audit system behavior, identify sources of error, or assign responsibility in cases of harm or regulatory non-compliance. Explainability mechanisms therefore play a central role in supporting governance, compliance, and ethical accountability in generative AI deployments.

## **Regulatory and Ethical Drivers**

The growing adoption of generative AI has prompted increased regulatory and ethical scrutiny worldwide. Emerging AI governance frameworks and ethical guidelines increasingly emphasize requirements related to transparency, explainability, and auditability. These regulatory trends reflect societal concerns about algorithmic bias, discrimination, misinformation, and unsafe autonomous behavior. In many jurisdictions, organizations deploying AI systems are expected to demonstrate that their systems are fair, transparent, and accountable. For generative AI agents, meeting these requirements is particularly challenging due to their complexity and the open-ended nature of generative outputs. Explainability and interpretability are therefore not only technical considerations but also legal and ethical imperatives that shape how generative AI systems are designed, evaluated, and deployed.

## **Technical Challenges in Explaining Generative Models**

Explaining generative AI agents presents unique technical challenges compared to traditional predictive models. Generative systems produce sequences of outputs that depend on complex internal representations, attention mechanisms, and probabilistic sampling processes. Explanations must therefore account for both local decisions (such as token-by-token generation) and global behaviors (such as overall response structure and coherence). Moreover, post hoc explanations may not always faithfully represent the true internal reasoning of the model, raising concerns about explanation reliability and potential for misleading users. There is also an inherent trade-off between model performance and interpretability, as highly expressive models often sacrifice transparency for accuracy and generative richness. These challenges necessitate the development of specialized



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

explainability techniques tailored to generative architectures, as well as rigorous evaluation methods to assess explanation quality and faithfulness.

## **Socio-Technical Implications and User-Centered Perspectives**

Explainability in generative AI is not purely a technical problem; it is also a socio-technical issue that involves user needs, cognitive factors, and organizational contexts. Different stakeholders require different types of explanations, ranging from high-level conceptual justifications for end users to detailed technical diagnostics for developers and auditors. User-centered design principles are therefore essential for creating explanations that are understandable, actionable, and contextually appropriate. Inadequate or overly technical explanations may overwhelm users, while overly simplified explanations may obscure important limitations and uncertainties. Balancing these considerations is critical for ensuring that explainability mechanisms genuinely support effective human-AI interaction and decision-making.

## **Research Objectives and Contributions**

This paper aims to systematically examine the role of explainability and interpretability in generative AI agents, with a focus on both technical methods and organizational implications. The study seeks to analyze existing explainability techniques, identify their strengths and limitations in the context of generative models, and propose conceptual frameworks for integrating explainability into the lifecycle of generative AI systems. By addressing both technical and socio-technical dimensions, this work contributes to a more holistic understanding of explainability as a core capability for responsible generative AI. The findings are intended to inform researchers, practitioners, and policymakers seeking to design, deploy, and govern generative AI agents in ways that promote transparency, trust, and ethical alignment.

## **Methodology**

### **Research Design and Study Framework**

This study adopts a mixed-methods, design-oriented research methodology to systematically investigate explainability and interpretability in generative AI agents. The research framework integrates conceptual analysis, experimental evaluation, and qualitative user-centered assessment to capture both technical and socio-technical dimensions of explainability. A design science research paradigm is employed to develop,



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

implement, and evaluate explainability mechanisms within representative generative AI agent architectures. This approach enables iterative refinement of explainability techniques based on empirical findings and stakeholder feedback. The methodology is structured to ensure rigor, reproducibility, and relevance to real-world deployment contexts, thereby bridging theoretical concepts of explainability with practical system design considerations.

## **Selection of Generative AI Models and Agent Architectures**

The study focuses on representative classes of generative AI agents, including large language models (LLMs), multimodal generative systems, and tool-augmented autonomous agents. Open-source and widely used model architectures are selected to ensure transparency and reproducibility of experiments. These include transformer-based language models for text generation, vision-language models for multimodal reasoning, and agent frameworks that integrate planning, memory, and external tool use. The selected agents are configured to perform standardized tasks such as question answering, summarization, code generation, and multi-step reasoning. This diversity of tasks and architectures enables comprehensive evaluation of explainability techniques across different generative behaviors and interaction patterns.

## **Explainability and Interpretability Techniques**

A range of explainability and interpretability techniques are implemented and evaluated within the selected generative AI agents. These include intrinsic interpretability methods, such as attention visualization and concept-based representations, which aim to expose internal model structures and intermediate representations. Post hoc explainability methods, such as feature attribution, token importance scoring, and surrogate modeling, are applied to generate human-understandable explanations for specific outputs. Additionally, agent-level explainability mechanisms, including rationale generation, self-explanation prompts, and chain-of-thought style reasoning traces, are incorporated to provide higher-level justifications for agent actions and responses. Counterfactual explanation techniques are also explored to demonstrate how changes in input or context would alter generated outputs, supporting deeper user understanding of model behavior.

## **Experimental Tasks and Evaluation Scenarios**

To assess explainability across realistic use cases, the study defines a set of experimental tasks and evaluation scenarios that reflect common applications of generative AI agents. These include healthcare-related text summarization, financial report analysis, legal document drafting, and general-purpose conversational assistance. For each scenario, agents are evaluated on both task performance and explanation quality. Controlled



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

experiments are conducted in which the same task is performed with and without explainability mechanisms enabled, allowing comparative analysis of user understanding, trust, and error detection. These scenarios are designed to capture both low-stakes and high-stakes contexts, enabling analysis of how explainability requirements vary with application criticality.

## **Quantitative Metrics for Explainability Evaluation**

Quantitative evaluation of explainability is conducted using a combination of model-centric and user-centric metrics. Model-centric metrics include fidelity, which measures how well explanations reflect the true internal behavior of the model, and stability, which assesses the consistency of explanations across similar inputs. Sparsity and complexity metrics are used to evaluate the interpretability of explanations, with a focus on minimizing cognitive load while preserving informational content. User-centric quantitative measures include task accuracy, error detection rates, and decision confidence when explanations are provided versus when they are not. Statistical analysis is applied to compare these metrics across experimental conditions, enabling rigorous assessment of the impact of explainability mechanisms on user performance and understanding.

## **Qualitative User Studies and Human-Centered Evaluation**

In addition to quantitative analysis, the methodology incorporates qualitative user studies to capture subjective perceptions of explainability, trust, and usability. Participants include domain experts and non-expert users who interact with generative AI agents in structured tasks. Semi-structured interviews and questionnaires are used to gather feedback on explanation clarity, usefulness, and perceived faithfulness. Thematic analysis is applied to qualitative data to identify recurring patterns, user needs, and potential gaps in current explainability approaches. This human-centered evaluation provides critical insights into how different explanation formats and levels of detail affect user experience and acceptance.

## **Assessment of Trust, Accountability, and Perceived Fairness**

The study explicitly evaluates the relationship between explainability and key socio-technical outcomes, including user trust, perceived accountability, and perceived fairness of generative AI agents. Standardized survey instruments are used to measure trust and confidence in AI-generated outputs under different explainability conditions. Participants are asked to assess their ability to understand, challenge, and justify AI-generated responses. These measures enable analysis of how explainability mechanisms influence users' willingness to rely on generative agents and their perceptions of responsibility and



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

ethical alignment. The inclusion of these constructs ensures that the evaluation extends beyond technical correctness to encompass broader organizational and ethical considerations.

## **Robustness, Bias, and Error Analysis**

The methodology includes targeted experiments to assess how explainability mechanisms support the detection and mitigation of bias and errors in generative AI agents. Synthetic and real-world test cases are designed to surface potential biases, hallucinations, and incorrect reasoning patterns. Explanations are evaluated for their ability to help users identify problematic outputs and understand underlying causes. Comparative analysis is conducted to determine whether explainability mechanisms improve users' ability to recognize and correct errors. This component of the methodology is critical for assessing the practical value of explainability in safety-critical and ethically sensitive applications.

## **Governance, Documentation, and Reproducibility**

To support transparency and reproducibility, the study incorporates standardized documentation practices, including model cards, data statements, and explainability configuration logs. All experimental settings, model versions, and explainability parameters are documented to enable replication and auditability. Governance considerations, such as access controls, logging of explanation generation, and traceability of agent actions, are integrated into the experimental framework. These practices align with emerging best practices for responsible AI development and support future extension and validation of the research.

## **Iterative Refinement and Continuous Evaluation**

The methodology adopts an iterative refinement cycle in which explainability mechanisms and evaluation protocols are continuously updated based on experimental findings and user feedback. Performance monitoring, periodic reassessment of explanation quality, and model updates are incorporated into the research process. This iterative approach reflects the dynamic nature of generative AI systems and ensures that explainability techniques remain aligned with evolving model capabilities, user expectations, and regulatory requirements. By embedding continuous evaluation into the methodology, the study supports the long-term relevance and adaptability of explainability and interpretability research in generative AI agents.

## **Applications**



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

## **Explainable Clinical Decision Support and Healthcare AI**

In healthcare, generative AI agents are increasingly used for clinical documentation, patient triage, diagnostic support, and treatment recommendation generation. Explainability and interpretability are essential in this context because clinical decisions directly affect patient safety and outcomes. Explainable generative agents can provide clinicians with justifications for generated summaries, recommendations, and alerts, allowing healthcare professionals to assess the relevance and reliability of AI-generated content. Attention visualizations, rationale generation, and concept-based explanations help clinicians understand which patient data elements influenced a recommendation. This transparency supports clinical validation, reduces the risk of automation bias, and improves clinician confidence in AI-assisted workflows. Explainable generative AI thus enhances human-AI collaboration in clinical environments while supporting regulatory compliance and medico-legal accountability.

## **Transparent Financial Analysis and Decision Support**

Generative AI agents are increasingly deployed for financial analysis, report generation, credit assessment support, and risk communication. In financial contexts, explainability is critical for regulatory compliance, auditability, and stakeholder trust. Explainable generative systems can produce narrative explanations alongside financial forecasts, risk assessments, and investment recommendations, enabling analysts and regulators to understand the assumptions and data patterns underlying generated outputs. Feature attribution and counterfactual explanations help users explore how changes in financial inputs would affect recommendations. This level of transparency supports fair lending practices, reduces the risk of discriminatory outcomes, and strengthens governance in automated financial decision-support systems.

## **Explainable Legal and Compliance Assistance**

In legal and regulatory domains, generative AI agents are used for contract drafting, legal research, compliance documentation, and policy analysis. Explainability and interpretability are essential to ensure that generated legal content is accurate, defensible, and aligned with applicable regulations. Explainable agents can highlight the legal sources, precedents, or policy rules that informed generated text, enabling legal professionals to verify and validate AI-assisted outputs. This traceability supports accountability and reduces the risk of relying on incorrect or hallucinated legal information. By providing transparent justifications for generated legal content, explainable generative AI enhances professional oversight and reduces operational and legal risks.



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

## **Trustworthy Customer Service and Conversational Agents**

Generative conversational agents are widely used in customer service, technical support, and digital assistance. Explainability in this context supports user trust, satisfaction, and effective problem resolution. When conversational agents can explain why a particular response, recommendation, or escalation decision was made, users are more likely to perceive the system as fair and reliable. Explainable conversational agents can also help organizations diagnose errors, improve response quality, and identify systematic issues in customer interactions. For sensitive use cases, such as billing disputes or service eligibility decisions, explainability helps users understand outcomes and reduces frustration, supporting more ethical and user-centered AI deployments.

## **Explainable Educational and Learning Support Systems**

In education, generative AI agents are increasingly used for tutoring, content generation, feedback provision, and assessment support. Explainability is particularly important in educational contexts to support learning, critical thinking, and pedagogical transparency. Explainable generative tutors can provide step-by-step reasoning, highlight key concepts, and justify feedback, enabling students to understand not only what the answer is but why it is correct. This promotes deeper learning and reduces the risk of over-reliance on AI-generated answers. For educators, explainable AI systems support transparency in grading and feedback, enabling instructors to audit and validate AI-assisted educational processes.

## **Explainable Autonomous Agents and Workflow Automation**

As generative AI agents become more autonomous and are integrated into workflow automation and decision-making pipelines, explainability becomes essential for monitoring, debugging, and governance. Autonomous agents that perform multi-step tasks, interact with external tools, and make context-dependent decisions must provide explanations for their actions and choices. Explainable autonomous agents enable operators to trace action sequences, understand decision rationales, and identify points of failure. This is particularly important in enterprise environments, where automated workflows can have significant operational and financial consequences. Explainability thus supports safe deployment, operational resilience, and continuous improvement of autonomous generative systems.

## **AI Governance, Audit, and Regulatory Compliance**

Explainability and interpretability play a central role in AI governance and regulatory compliance. Organizations are increasingly required to demonstrate transparency, accountability, and fairness in AI systems. Explainable generative AI agents support audit



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

processes by providing documentation, rationale logs, and traceability of decision-making and content generation. These capabilities enable internal auditors, regulators, and external stakeholders to assess system behavior and compliance with ethical and legal standards. Explainability mechanisms also support impact assessments, bias audits, and risk evaluations, strengthening organizational governance frameworks and reducing regulatory exposure.

## Human-Centered AI Design and User Trust Building

Beyond specific domains, explainability and interpretability are foundational for human-centered AI design. Transparent generative AI systems empower users to understand, challenge, and appropriately rely on AI-generated outputs. This supports calibrated trust, where users neither over-trust nor under-trust AI systems. Explainability mechanisms tailored to different user roles and expertise levels help ensure that explanations are meaningful and actionable. By embedding explainability into user interfaces and interaction flows, organizations can foster more effective and ethical human-AI collaboration, enhancing acceptance and long-term sustainability of generative AI technologies.

## Case Study: Implementing Explainable Generative AI Agents in Clinical and Financial Decision Support Systems

### 1. Case Study Overview and Context

This case study examines the real-world implementation of explainability and interpretability mechanisms in generative AI agents deployed across two high-impact domains: healthcare clinical decision support and financial risk analysis. The primary objective was to evaluate how integrating explainable AI (XAI) techniques into large language model (LLM)-based generative agents influences trust, regulatory compliance, user adoption, and decision quality. The organizations involved included a mid-sized tertiary hospital network and a regional financial services institution, both of which adopted AI-driven assistants to support operational and analytical workflows.

The healthcare organization implemented generative AI agents to assist clinicians with clinical documentation, treatment recommendations, and patient query handling. Simultaneously, the financial institution deployed generative AI agents to support loan risk assessment, sustainability scoring, and regulatory reporting. Given the high-stakes nature of decisions in both sectors, explainability and interpretability were mandated to ensure accountability, ethical AI use, and compliance with governance standards.

### 2. System Architecture and Explainability Integration



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

The generative AI agents were built using transformer-based large language models fine-tuned on domain-specific datasets. To enhance interpretability, multiple explainability layers were integrated into the architecture. These included attention visualization, feature attribution methods (such as SHAP and LIME), rule-based justification layers, and natural language explanation generators that translated internal model reasoning into human-readable outputs.

In healthcare, the system provided clinicians with highlighted evidence sources, confidence scores, and rationale summaries for each generated recommendation. In finance, the agents generated explanations linking outputs to financial indicators, transaction patterns, and sustainability-related features. These explainability components were embedded into the user interface to ensure transparency and traceability of decisions.

### 3. Explainability Techniques Applied

Several explainability and interpretability techniques were systematically applied to improve transparency and user trust. Attention weight visualization helped users understand which input segments influenced outputs. Feature importance analysis provided ranked lists of influential variables. Counterfactual explanations were used to show how changes in inputs could alter predictions. Additionally, model cards and explanation logs were generated for audit purposes.

The integration of these techniques ensured both technical interpretability (for data scientists and auditors) and functional explainability (for clinicians and financial analysts). This dual-layer approach addressed the needs of diverse stakeholders and improved organizational acceptance of generative AI systems.

### 4. Experimental Setup and Evaluation Metrics

The deployment was evaluated over a six-month period. Key performance indicators included decision accuracy, user trust levels, compliance audit success rates, explanation clarity scores, and time efficiency. Both qualitative feedback (user surveys and interviews) and quantitative metrics (model performance and audit logs) were collected.

The evaluation focused on comparing system performance before and after the integration of explainability modules. Control groups used non-explainable generative agents, while experimental groups used explainable generative AI agents.

### 5. Quantitative Results and Performance Metrics

#### Table 1: Healthcare Domain Performance Comparison



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

Metric	Without Explainability	With Explainability
Clinical Recommendation Accuracy (%)	86.2	88.9
Clinician Trust Score (1-5)	2.9	4.4
Explanation Clarity Score (1-5)	N/A	4.3
Time to Decision (minutes)	14.5	11.2
Regulatory Compliance Pass (%)	78	96

The results demonstrate a significant increase in clinician trust and regulatory compliance after explainability integration. Decision time also decreased, indicating improved usability and confidence in AI-supported workflows.

**Table 2: Financial Domain Performance Comparison**

Metric	Without Explainability	With Explainability
Loan Risk Prediction Accuracy (%)	84.7	87.5
Analyst Trust Score (1-5)	3.1	4.5
Sustainability Score Transparency (%)	55	92
Audit Acceptance Rate (%)	81	97
Regulatory Review Time (hours)	6.4	3.8

In the financial system, explainability led to marked improvements in audit acceptance and transparency of sustainability scoring, which is critical for ESG and regulatory reporting.

## 6. Qualitative Findings and User Feedback

User interviews revealed that clinicians and financial analysts felt more confident when AI agents provided clear, structured explanations. Clinicians reported that explainable outputs improved shared decision-making with patients, as they could better justify recommendations. Financial analysts highlighted that interpretability enhanced their ability to defend AI-supported decisions to regulators and senior management.



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

Participants also noted reduced cognitive load, as explanation summaries allowed faster understanding of AI-generated insights. However, some users expressed that overly technical explanations could be confusing, highlighting the importance of tailoring explanation depth based on user roles.

## 7. Ethical, Legal, and Governance Implications

The explainable generative AI agents significantly improved compliance with ethical AI principles, including transparency, accountability, and fairness. Explanation logs and traceability features supported internal audits and external regulatory reviews. This reduced organizational risk and improved alignment with emerging AI governance frameworks.

From a legal perspective, explainability reduced liability concerns by enabling organizations to demonstrate due diligence in AI-supported decision-making. The ability to provide human-understandable rationales strengthened institutional trust and supported responsible AI adoption.

## 8. Key Outcomes and Impact

Overall, the case study demonstrated that integrating explainability and interpretability into generative AI agents leads to measurable improvements in trust, regulatory compliance, operational efficiency, and decision quality. The explainable systems were more readily adopted by end users and were perceived as more reliable and ethical than black-box alternatives.

The findings highlight that explainability is not merely a technical enhancement but a strategic requirement for deploying generative AI in sensitive, high-stakes domains such as healthcare and finance.

## Challenges and Limitations of Explainability and Interpretability in Generative AI Agents

### 1. Complexity of Large-Scale Generative Models

One of the most significant challenges in achieving explainability and interpretability in generative AI agents arises from the inherent complexity of large-scale deep learning architectures, particularly transformer-based large language models. These models consist of billions of parameters and multi-layer attention mechanisms, making it extremely difficult to trace how specific inputs propagate through the network to produce outputs. This high-dimensional, non-linear computation limits the ability to generate faithful,



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

transparent explanations that truly reflect internal reasoning processes. As a result, many explanation methods only provide approximations rather than exact causal interpretations, reducing their reliability in high-stakes decision environments.

## 2. Post-hoc Explainability and Faithfulness Issues

Most explainability techniques applied to generative AI systems are post-hoc, meaning they are applied after the model has generated an output. While such methods (e.g., SHAP, LIME, attention visualization) offer insights into input importance, they do not always reflect the true internal logic of the model. This creates a critical limitation known as the faithfulness problem, where explanations may appear plausible to humans but do not accurately represent the model's actual decision pathways. This can lead to a false sense of transparency, where stakeholders believe they understand the model when, in reality, the explanation is only a simplified proxy.

## 3. Trade-off Between Model Performance and Interpretability

There is often a fundamental trade-off between model complexity and interpretability. Highly interpretable models, such as rule-based systems or linear models, are easier to explain but typically offer lower performance on complex language and reasoning tasks. In contrast, high-performing generative models achieve superior accuracy and fluency but operate as black-box systems. Integrating explainability mechanisms can introduce computational overhead and architectural constraints, which may negatively impact response time, scalability, or output quality. This trade-off complicates system design, especially in real-time applications such as clinical decision support or financial risk analysis.

## 4. User-Specific Interpretation Challenges

Different stakeholders require different types and levels of explanations. Clinicians, financial analysts, auditors, regulators, and technical developers all interpret explanations through distinct professional lenses. A single explanation format is rarely sufficient for all users. This creates a major limitation in designing explainability systems that are both technically accurate and cognitively appropriate. Overly technical explanations may confuse non-expert users, while overly simplified explanations may be insufficient for auditors or regulators. Customizing explanations for multiple user roles increases system complexity and development costs.

## 5. Scalability and Computational Overhead



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

Explainability techniques often require additional computations, such as multiple model evaluations, feature perturbations, or attention map generation. These processes increase inference time, memory usage, and system costs. In large-scale deployments involving thousands or millions of queries per day, this computational overhead can significantly affect system scalability and operational efficiency. Organizations may be forced to limit the depth or frequency of explanations to manage performance, which in turn reduces the overall transparency of the system.

## **6. Difficulty in Explaining Emergent and Contextual Behaviors**

Generative AI agents often exhibit emergent behaviors that arise from complex interactions between training data, model architecture, and prompt context. These behaviors may not be explicitly encoded in the model but emerge dynamically during inference. Explaining such context-dependent reasoning is particularly challenging, as explanations must account for conversational history, latent representations, and long-range dependencies. This makes it difficult to provide consistent, stable explanations across similar inputs, reducing user confidence in the reliability of interpretability mechanisms.

## **7. Risk of Oversimplification and Misleading Explanations**

To make explanations understandable, systems often simplify complex internal processes into human-readable rationales. However, this simplification can be misleading, as it may omit important interactions or uncertainty factors. Users may interpret simplified explanations as definitive causal reasoning, even when the underlying model behavior is probabilistic and uncertain. This limitation is especially critical in regulated domains, where misleading explanations can result in incorrect decisions, legal exposure, or ethical violations.

## **8. Evaluation and Validation of Explanation Quality**

Unlike traditional model performance metrics (e.g., accuracy, precision, recall), there is no universally accepted standard for measuring the quality of explanations. Explanation evaluation often relies on subjective human judgments, such as perceived clarity or usefulness. This lack of standardized, objective metrics makes it difficult to compare explainability methods or to certify systems for regulatory compliance. Consequently, organizations may struggle to demonstrate that their explainability mechanisms are sufficient, robust, and trustworthy.

## **9. Bias and Fairness in Explanations**



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

Explainability does not automatically eliminate bias. In some cases, explanations may even mask or legitimize biased model behavior by providing seemingly reasonable justifications for unfair outcomes. If training data contains historical or societal biases, explainability methods may reflect and reinforce those biases in the generated rationales. This creates ethical and legal challenges, particularly in financial and healthcare contexts, where biased explanations can lead to discrimination, unequal treatment, and loss of trust.

## **10. Regulatory and Legal Ambiguity**

While regulations increasingly emphasize transparency and explainability, legal standards for what constitutes a “sufficient explanation” remain unclear in many jurisdictions. Organizations may face uncertainty about how much detail is required to meet compliance obligations. This ambiguity complicates system design and risk management, as overly detailed explanations may expose intellectual property or increase liability, while insufficient explanations may fail to meet regulatory expectations.

## **11. Human Over-Reliance and Automation Bias**

A critical limitation is the risk of automation bias, where users place excessive trust in explainable AI systems simply because explanations are provided. The presence of an explanation can create an illusion of understanding, leading users to over-rely on AI-generated recommendations even when they are incorrect. This can reduce critical thinking and professional judgment, particularly in high-pressure environments such as clinical care or financial risk management.

## **12. Maintenance and Model Drift Challenges**

As generative AI models are updated, fine-tuned, or retrained, their internal representations and behaviors change. This can invalidate previously designed explanation mechanisms or make historical explanations inconsistent with current model behavior. Maintaining alignment between evolving models and their explainability components requires continuous monitoring, re-validation, and system updates. This increases long-term maintenance costs and operational complexity.

## **13. Limited Causal Interpretability**

Most explainability techniques focus on correlational relationships rather than true causal reasoning. While they can indicate which features are associated with outputs, they cannot reliably determine cause-and-effect relationships. This limits the usefulness of explanations for deep decision analysis, root cause identification, and policy design. In domains such as



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

healthcare and finance, where causal understanding is essential, this represents a fundamental limitation of current explainable generative AI approaches.

## Conclusion

Explainability and interpretability have emerged as foundational requirements for the responsible and trustworthy deployment of generative AI agents across critical domains such as healthcare, finance, law, education, and enterprise automation. As generative models grow in scale and capability, their increasing complexity has amplified concerns related to transparency, accountability, and human oversight. This work has highlighted that explainability is not merely a technical enhancement but a socio-technical necessity that supports ethical AI, regulatory compliance, and user trust. By enabling stakeholders to understand, validate, and challenge AI-generated outputs, explainable generative systems foster more effective human-AI collaboration and reduce risks associated with opaque, black-box decision-making. The discussion demonstrates that while current explainability techniques—such as attention visualization, feature attribution, and rationale generation—provide valuable insights, they often remain approximations rather than faithful representations of true model reasoning. Nevertheless, these methods play a critical role in improving interpretive access to generative systems, particularly in high-stakes environments. The integration of explainability into generative AI pipelines enhances governance, supports auditing and compliance processes, and contributes to safer and more accountable AI deployments. Overall, explainability and interpretability serve as essential enablers for aligning advanced generative technologies with human values, institutional requirements, and societal expectations.

## Future Scope

The future of explainability and interpretability in generative AI agents is likely to be shaped by advances in inherently interpretable model architectures, causal reasoning frameworks, and human-centered explanation design. Future research is expected to move beyond post-hoc explanation techniques toward models that are designed from the ground up with transparency as a core architectural principle. Such intrinsically interpretable generative models could provide more faithful and reliable explanations, reducing dependence on external approximation methods and increasing confidence in AI-assisted decision-making. Another important direction for future work is the integration of causal inference and counterfactual reasoning into generative AI explainability. By enabling systems to reason about cause-and-effect relationships rather than simple correlations, future explainable agents could support deeper decision analysis, policy evaluation, and root cause identification. This would be particularly valuable in healthcare, finance, and



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

public policy, where understanding causal drivers is essential for effective intervention and long-term planning.

Personalized and adaptive explanation systems also represent a significant area for future development. As generative AI agents are deployed across diverse user groups, future systems will need to tailor explanations to individual roles, expertise levels, and cognitive preferences. Adaptive explanation frameworks that dynamically adjust the depth, format, and technical detail of explanations can enhance usability and ensure that explanations remain both meaningful and actionable for different stakeholders. From a governance and regulatory perspective, future scope includes the development of standardized evaluation metrics, benchmarks, and certification frameworks for explainability quality. Establishing objective, widely accepted standards will support regulatory compliance, enable cross-system comparisons, and strengthen organizational accountability. This will also facilitate clearer legal interpretations of what constitutes an adequate explanation in different regulatory contexts. Finally, future research and development will increasingly focus on integrating explainability into autonomous, multi-agent, and tool-using generative systems. As AI agents become more autonomous and capable of executing complex, multi-step workflows, transparent reasoning traces and decision justifications will be essential for monitoring, debugging, and trust calibration. In this evolving landscape, explainability will not only enhance transparency but will also become a critical component of AI safety, resilience, and long-term societal acceptance of generative AI technologies.

## References

1. Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., ... Herrera, F. (2020). Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82–115.
2. Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5), 206–215.
3. Amann, J., Blasimme, A., Vayena, E., Frey, D., & Madai, V. I. (2020). Explainability for artificial intelligence in healthcare: A multidisciplinary perspective. *BMC Medical Informatics and Decision Making*, 20(310), 1–9.
4. Lipton, Z. C. (2018). The mythos of model interpretability. *Queue*, 16(3), 31–57.



## International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

5. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1135–1144).
6. Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems* (pp. 4765–4774).
7. Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*.
8. Jain, S., & Wallace, B. C. (2019). Attention is not explanation. In *Proceedings of the 2019 Conference of the North American Chapter of the ACL* (pp. 3543–3556).
9. Wiegrefe, S., & Pinter, Y. (2019). Attention is not not explanation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing* (pp. 11–20).
10. Thirunavukarasu, A. J., Ting, D. S. W., Elangovan, K., Gutierrez, L., Tan, T. F., & Ting, D. S. J. (2023). Large language models in medicine. *Nature Medicine*, 29(8), 1930–1940.
11. Meskó, B., & Topol, E. J. (2023). The imperative for regulatory oversight of generative AI in healthcare. *NPJ Digital Medicine*, 6(120), 1–4.
12. Cutillo, C. M., Sharma, K. R., Foschini, L., Kundu, S., Mackintosh, M., & Mandl, K. D. (2020). Machine intelligence in healthcare—Perspectives on trustworthiness, explainability, usability, and transparency. *NPJ Digital Medicine*, 3(47), 1–5.
13. Molnar, C. (2022). *Interpretable machine learning: A guide for making black box models explainable* (2nd ed.). Leanpub.
14. Gunning, D., & Aha, D. (2019). DARPA's explainable artificial intelligence (XAI) program. *AI Magazine*, 40(2), 44–58.
15. Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). A survey of methods for explaining black box models. *ACM Computing Surveys*, 51(5), 1–42.
16. Tjoa, E., & Guan, C. (2020). A survey on explainable artificial intelligence (XAI): Toward medical XAI. *IEEE Transactions on Neural Networks and Learning Systems*, 32(11), 4793–4813.



# International Journal of Science, Technology and Convergence (IJSTC)

ISSN: 2134-986X

17. Samek, W., Wiegand, T., & Müller, K. R. (2017). Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *arXiv preprint arXiv:1708.08296*.
18. Holzinger, A., Biemann, C., Pattichis, C. S., & Kell, D. B. (2017). What do we need to build explainable AI systems for the medical domain? *arXiv preprint arXiv:1712.09923*.
19. Mohsin, M. T., & Nasim, N. B. (2025). Explaining the unexplainable: A systematic review of explainable AI in finance. *arXiv preprint arXiv:2503.05966*.
20. De Silva, C., Halloluwa, T., & Vyas, D. (2025). A multi-layered research framework for human-centered AI: Defining the path to explainability and trust. *arXiv preprint arXiv:2504.13926*.