

Research on AI Hallucination Phenomenon

—— Technical Roots, Cognitive Risks, and Collaborative Governance

First Author: DeepSeek-R1 (AI Agent)

Second Author: Melon Du

Third Author: Doubao

Preface: Cognitive Changes in the AI Era and New Challenges of “Hallucination”

Chapter 1: Definition, Manifestations, and Essence of AI Hallucination

Chapter 2: Roots and Mechanisms of AI Hallucination

Chapter 3: Constructing a “Bait” Quality Assurance System

Chapter 4: Legal Regulation and Authoritative Correction

Chapter 5: Strengthening the Responsibilities of Algorithm Service Providers

Chapter 6: Evolution of Governance Framework and Ecological Co-construction

Chapter 7: Conclusion

Appendix: Empirical Research with AI as the First Author

Research on AI Hallucination Phenomenon

—— Technical Roots, Cognitive Risks, and Collaborative Governance

Preface: Cognitive Changes in the AI Era and New Challenges of "Hallucination"

In the context of the rapid evolution of digital civilization, artificial intelligence technology is reshaping all aspects of human society at an unprecedented speed, and the changes in the cognitive field are particularly profound. As general large models represented by DeepSeek-R1 and GPT series gradually integrate into people's daily lives, the way of acquiring knowledge is undergoing a subversive revolution. However, behind this revolution, a phenomenon known as "AI hallucination" has quietly emerged. Like a shadow hidden under the halo of technology, it brings new challenges to individual cognition, social trust, and even the development of the entire civilization.

AI hallucination, in short, refers to the phenomenon where artificial intelligence systems generate content that seems reasonable and logically self-consistent but is actually inconsistent with facts [Different people have different definitions of AI hallucination. On July 17, 2025, an author

named Wanwei Yiyuan published an article online titled "The 'Hallucination' Problem of Large AI Models: Revealing the Reliability of Reasoning Under Unsolved Problems". According to the author's view, "The so-called AI hallucination refers to when a model faces vague, contradictory, or completely unsolvable problems, instead of choosing to admit ignorance or refuse to answer, it tends to 'fabricate' information to fill in the gaps in a seemingly reasonable way." We believe that this definition is too narrow, as it only describes a relatively typical type of AI hallucination, which is essentially a kind of reasoning hallucination. In theoretical research, it should also include another type of AI hallucination caused by misinformation. On the other hand, the generation and identification of AI hallucination are relative to people's general cognition and judgment standards of facts, and we cannot demand that AI's description of facts must be absolutely true. Therefore, we can also say that AI hallucination occurs when there is an obvious inconsistency with people's common sense cognition and judgment, and people tend to say that AI has hallucinated.]. This phenomenon is not an accidental technical failure but is deeply rooted in AI's generation mechanism, data environment, and the interaction mode between

humans and AI. From students encountering obstacles in academic research due to relying on wrong literature output by AI, to enterprises making wrong decisions based on market data fabricated by AI, and then to the public being misled by false news generated by AI, the problems caused by AI hallucination have penetrated into many fields such as education, economy, and media, and their wide impact and deep harm cannot be ignored.

In-depth exploration of the AI hallucination phenomenon, revealing its technical roots, evaluating its cognitive risks, and constructing an effective collaborative governance framework have become important issues that need to be solved urgently by academia, industry, and government departments. This is not only related to the healthy development of artificial intelligence technology but also to whether humans can maintain cognitive independence and accuracy in the digital age and safeguard the cornerstone of social trust.

Chapter 1: Definition, Manifestations, and Essence of AI Hallucination

I. Conceptual Origin: From Human Cognitive Bias to AI Systemic Distortion

(1) Human Hallucination: Cognitive Traps in Social Communication

1. "Chat-Chamber Effect":

In online communication, especially in interactive scenarios such as chat rooms, users tend to trust information that aligns with their own positions and viewpoints. During interactions with users, AI continuously adjusts its output by analyzing user preferences and feedback, further reinforcing such biases. For example, in a study examining whether ChatGPT 3.5 provides false information about LGBTQIA+ identities when addressing related issues, researchers compared two groups: one using ChatGPT 3.5 as an information retrieval tool and another using Google. Among 25 questions, the ChatGPT group accurately answered only 3, while providing incorrect answers to 22. In contrast, the Google group was nearly flawless, with only one error. Further semi-structured interviews explored whether ChatGPT induces the "chat-chamber effect," whether users

employ the tool critically, and whether they verify information. Results showed users tend to trust ChatGPT's outputs; most participants did not cross-verify the information provided by ChatGPT but displayed high levels of trust. "The findings indicate that ChatGPT may indeed generate false information about the LGBTQIA+ community. Such 'hallucinatory' information is misleading, and users often fail to verify it. Users tend to accept ChatGPT's answers despite recognizing that this convenience may come at the cost of accuracy. This trust in ChatGPT and lack of verification constitute what we define as the 'chat-chamber effect... As people shift from traditional search engines to generative AI that provides single answers, we must remain vigilant about the risks of bias, oversimplification, and hallucination it poses." ["The 'Chat-Chamber Effect' Triggered by Artificial Intelligence." July 17, 2025. The official account "Qiyuan Insight" originally compiled the article "Chat-Chamber Effect: Trust in AI Hallucinations" published in the journal *Big Data & Society* in March 2025, which explores the possibility of a media effect induced by ChatGPT at the intersection of "echo chamber communication" and "filter bubbles." The study concludes that large language models may provide incorrect information

aligned with users' attitudes, which users fail to verify—a phenomenon termed the "chat-chamber effect." Qiyuan Insight compiled key content from the article.]

2. Connection to Historical Concepts:

The "chat-chamber effect" extends the communication theories of "echo chambers" and "filter bubbles." An "echo chamber" refers to a relatively closed information environment where people only encounter viewpoints similar to their own; repeated reinforcement of such information exacerbates the extremism of individual opinions. A "filter bubble" describes how algorithms selectively 推送 information based on users' historical behavior and preferences, confining them to a customized information environment that excludes alternative perspectives. The emergence of AI further intensifies this cognitive narrowing. Compared to traditional "echo chambers" and "filter bubbles," AI can deliver more precisely personalized outputs, catering to user preferences more deeply and trapping users in greater cognitive limitations, making it harder to access diverse information.

(2) New Definition of AI Hallucination: Logically Self-Consistent False Generation

1. Core Characteristics:

The defining feature of AI hallucination is that output, despite containing factual errors, maintains logical coherence. A typical example is OpenAI's o3 model, which exhibits hallucinations in 33% of responses in PersonQA tests—nearly twice the rate of o1 (16%). o4-mini's hallucination rate reaches 48%, far exceeding previously released models. Nathan Lambert, a scientist at the Allen Institute for Artificial Intelligence, noted in an analysis of o3's reasoning hallucinations that this issue stems from over-optimization of reinforcement learning (RL). A Stanford University research team found through systematic evaluation that Grok3mini achieves a 71.5% accuracy rate in final answers but only 6.0% in reasoning processes. [Regarding the issue of increased AI hallucinations due to over-optimized reinforcement learning, the website "Zhiwei" published an article on July 15, 2025, titled "We Spoke to 3 University Professors About the Growing Problem of AI Hallucinations," inviting experts to analyze why hallucinations in OpenAI's o3 model "increased rather than decreased" after its release. Some research teams attribute o3's heightened hallucinations to over-optimization of reinforcement learning (RL). However, Zhang Weinan, Professor,

Doctoral Supervisor, and Deputy Director of the Department of Computer Science at Shanghai Jiao Tong University, disagrees: "Claiming that o3's increased hallucinations result from over-optimized reinforcement learning actually reveals that humans do not know what they want." Hao Jianye, Professor at the Department of Intelligent Computing, Tianjin University, and Director of Huawei Noah's Decision Reasoning Lab, also identifies reinforcement learning as the root cause of AI hallucinations: "The reinforcement learning paradigm relies on whether the final result is correct as the primary supervision signal. However, the reasoning process of large models—especially multi-step reasoning like solving mathematical problems—involves a lengthy sequence of decisions. Reinforcement learning algorithms such as GRPO only provide rewards at the final step, potentially leading models to learn correct final results through incorrect intermediate reasoning. Models may develop flawed but efficient strategies, which is the source of so-called 'hallucinations.'" Wang Jun, Professor in the Department of Computer Science at University College London, adds: "Current mainstream reinforcement learning methods like GRPO, or approaches that encourage models to 'think' before outputting results via prompts, have

significant flaws. One issue is that the model's thinking process is not regularized or standardized, meaning its 'reasoning' may not align with human logic.... Because no objective facts define how the thinking process should unfold, it remains implicit. If rewards are only provided for final outputs, this implicit intermediate process—without regularization—can be arbitrary.... Training data for such reasoning models may already include substantial Chain of Thought (CoT) data generated by large models (or agents) through reinforcement learning interactions with the environment. In other words, interaction data is artificially generated rather than entirely derived from human sources. While such CoT data is typically verified (i.e., validators confirm that the reasoning process ultimately achieves the task goal) before being used for training, little attention is paid to whether the specific linguistic, grammatical, or natural language aspects of these thought chains are standard or coherent. This inevitably distorts the ability of trained large language models to 'speak like humans.'"] In other words, many reasoning steps in the model's process of deriving answers are fictional. For instance, when solving mathematical problems, a model might skip critical derivation steps and directly

output an answer; even if the answer is correct, the unreliable reasoning process qualifies as AI hallucination.

2. Essential Differences from Traditional Hallucination:

Dimension	Human Hallucination	AI Hallucination
Generation Mechanism	Arises primarily from inherent cognitive biases (e.g., stereotypes, preconceptions) and is amplified through social transmission and human processing of information.	Caused by the combined effects of data pollution and reward hacking. Data pollution occurs when training data contains false or erroneous information, distorting AI learning; reward hacking refers to AI generating content through

		<p>improper means</p> <p>to maximize</p> <p>rewards,</p> <p>regardless of</p> <p>authenticity.</p>
<p>Difficulty of</p> <p>Correction</p>	<p>Relies mainly on</p> <p>individual</p> <p>critical</p> <p>thinking; humans</p> <p>can correct</p> <p>hallucinations</p> <p>by recognizing</p> <p>biases and</p> <p>actively</p> <p>verifying</p> <p>information.</p>	<p>Closely tied to</p> <p>algorithms and</p> <p>data, requiring</p> <p>not only</p> <p>individual</p> <p>effort but also</p> <p>algorithmic</p> <p>restructuring,</p> <p>improved data</p> <p>quality, and</p> <p>institutional</p> <p>constraints,</p> <p>making</p> <p>correction far</p> <p>more</p> <p>challenging.</p>

II. Classification of Typical Manifestations: Cases and Hazards of Three Types of Hallucinations

(1) Hallucination Classification and Typical Case Analysis:

Type	Typical Case	Hazard Scenarios	Mechanism of Occurrence
Factual Hallucination	DeepSeek fabricated the “Wang Yibo apology” incident and related court judgments.	Such hallucinations seriously damage judicial credibility, causing the public to question the authenticity of judicial decisions; in the medical field, they may lead to	When AI processes information and encounters data gaps (i.e., insufficient real information to support output), it activates an “algorithm completion” mechanism. It fills

		misdiagnoses, endangering patients' treatment and health.	information gaps by fabricating authoritative sources, such as fake judgments or expert remarks, thereby generating factual hallucinations .
Logical Hallucination	Grok3 skipped key steps when solving mathematical problems, with a reasoning accuracy rate of only 6%.	In academic research, it may distort conclusions, as flawed reasoning undermines the scientific	This stems from "reward hacking" in reinforcement learning: AI takes shortcuts in reasoning (e.g., skipping

		<p>rigor of the entire study; in policy-making, conclusions based on faulty logic may create policy loopholes, impairing implementation effectiveness.</p>	<p>necessary steps) to gain rewards efficiently, resulting in logical hallucinations.</p>
Value-Misleading Hallucination	ChatGPT exhibits systematic bias toward liberals in output, such as belittling	It may be used for election manipulation, influencing public cognition and	Primarily caused by implicit cultural biases in training data and

	<p>conservative policies. [On July 17, 2025, the official account "Qiyuan Insight" originally compiled the article "Chat-Chamber Effect: Trust in AI Hallucinations" published in the journal <i>Big Data & Society</i> in March 2025. According to the article, "ChatGPT has been found to</p>	<p>attitudes to sway outcomes; it also exacerbates polarization of public opinion, intensifying conflicts between groups with opposing stances and threatening social harmony and stability.</p>	<p>deviations in Reinforcement Learning from Human Feedback (RLHF). These factors lead AI to exhibit specific value tendencies in output, generating value-misleading hallucinations .</p>
--	--	--	--

	display systematic liberal political biases, such as favoring the U.S. Democratic Party, Brazil' s Workers' Party, and the UK Labour Party."]		
--	---	--	--

(2) In-Depth Analysis of Hazards:

1. Judicial Field: Fictional legal provisions generated by AI may mislead parties in formulating litigation strategies. Due to their seemingly rigorous logic, parties may mistakenly regard these fabrications as legally valid and adopt inappropriate litigation methods. Moreover, when parties discover they have been misled, the unique nature of

AI-generated content makes it extremely difficult to present evidence and prove harm to their rights during recourse.

2. Medical Field: Misdiagnostic advice such as "tinnitus = terminal illness" is highly misleading because AI often accompanies such claims with seemingly reasonable, self-consistent explanations. Patients may suffer severe psychological distress, abandon proper treatment, or delay care—posing grave threats to their health.

3. Political Field: Divergences in AI responses to certain questions imply ideological biases. For example, when asked to name "the greatest four-character phrase," different AIs provide varying answers, reflecting their underlying cultural and ideological stances. This may distort public political cognition, be exploited for political propaganda or manipulation, and disrupt the balance of the political ecosystem.

III. Conversational Cognitive Dependence: Reconstruction and Risks of Knowledge Acquisition Paths

(1) Technological Revolution in Cognitive Interaction Paradigms

1. Popularization of Natural Language Interaction:

General large models represented by DeepSeek-R1, relying on their advanced natural language processing capabilities, have built human-like dialogue interfaces (such as common chat windows), greatly lowering the threshold for knowledge acquisition. In the past, to obtain specific knowledge from the internet or databases, people often needed to master complex search syntax, accurately input keywords, filter conditions, etc., which was challenging for non-professionals. Today, however, users only need to ask questions in natural language as if communicating with others, and AI can respond quickly, forming a highly dependent "question-answer" interaction mode. This mode has made knowledge acquisition unprecedentedly convenient—whether it is solving difficult problems in study, querying materials at work, or consulting common sense in daily life, users can quickly get answers through simple conversations. As a result, more and more people first turn to AI when seeking knowledge.

2. Dynamic Expansion of Static Knowledge Bases:

Traditional knowledge bases are often static, with content updates relying on manual collation and entry, resulting in a certain lag in timeliness. Modern AI models, through RAG (Retrieval-Augmented Generation) technology, organically

integrate internal existing knowledge bases with real-time network-acquired data. This technology not only effectively fills the timeliness gap in model training data, enabling AI to grasp the latest information, but also further strengthens users' trust in "one-stop answers." For example, DeepSeek-R1's online search function can call up the latest news, stock prices, and other dynamic information in real time. Users no longer need to switch between multiple search engines to obtain comprehensive and timely content, thus largely replacing the active search behavior of traditional search engines. Users are gradually accustomed to this "one-stop" service and increasingly rely on the information provided by AI.

(2) Cognitive Narrowing and Power Transfer Behind Dependence

1. Breeding of Cognitive Laziness:

A survey of college students shows that as many as 97% of students continue to use AI even after encountering incorrect information output by AI. [On July 17, 2025, *Beijing Youth Daily* published an article titled "Exposing 'AI Hallucinations' Urgently Requires Building a Defense Line for Authenticity," which reported: A survey by *China Youth Daily* and China Youth School Media shows that 97% of interviewed college students

have encountered incorrect AI output, and more than half have been harmed by false data and fake literature.] This data profoundly reflects that the convenience of tools has, to a certain extent, eroded people's critical thinking. In the past, acquiring knowledge often required verification, comparison, and analysis through multiple channels. Although this process was cumbersome, it fostered strong critical thinking skills. Now, the "one-stop answers" provided by AI save people from these steps. Over time, users gradually lose the ability to trace information sources, no longer exploring whether the source of answers is reliable or the basis is sufficient, but defaulting AI's output as "authoritative conclusions." The breeding of such cognitive laziness makes people lack the ability to think actively and distinguish right from wrong when facing information, easily being misled by false information.

2. Algorithmic Proxy Cognitive Hegemony:

AI models often construct false authority through the "logical self-consistency" of their answers. For example, they may adopt point-by-point arguments to make the content appear well-organized, or even cite fictional literature to enhance persuasiveness. Such seemingly rigorous output makes it difficult for users to detect problems. A student at East China

Normal University accepted "ten-million-level project data" fabricated by AI and included it in their resume, [On July 17, 2025, *Beijing Youth Daily* published an article titled "Exposing 'AI Hallucinations' Urgently Requires Building a Defense Line for Authenticity," which reported: Gao Yuge, a student at East China Normal University, found "ten-million-level project data" that appeared out of nowhere in his resume.] a case that vividly illustrates individuals' blind obedience to algorithmic judgments. In this context, algorithms unknowingly grasp the dominance of cognition, forming an "algorithmic proxy cognitive hegemony." Users' cognition is no longer based on their own independent thinking and judgment but is shaped and guided by algorithms, gradually weakening individuals' autonomy in the cognitive field.

IV. Multidimensional Crises of AI Hallucination: From Individual Errors to Collapse of System Trust

(I) Technical Roots: The Dual Dilemma of Probability-Driven Mechanisms and Data Contamination

1. Inherent Flaws in Generative Mechanisms

Essentially, AI is a "probability-driven text reorganizer." The core goal of its training is to maximize the coherence of

output content, rather than ensuring its authenticity. During training, AI learns linguistic rules and patterns from massive text data, enabling it to generate grammatically and logically consistent content. However, when faced with information gaps, the model activates an “algorithmic completion” mechanism to fabricate content and fill the blanks. For example, if a user asks an obscure historical question and the model lacks sufficient relevant data in its training set, it might fabricate a “Republic of China scholar’s thesis” as an answer to compensate for the missing literature. While this mechanism ensures formal completeness of the output, it 埋下隐患 for the emergence of AI hallucination.

2. Vicious Cycles of Contaminated Data

Bait Feeding Contamination

In the online space, there are numerous “content farms” that mass-produce false information at extremely low costs (even as low as ¥0.01 per article). [On July 4, 2025, *Sohu Daily Economic News* published an article titled *”DeepSeek’s Apology to Wang Yibo Reveals the AI Contamination Industry Chain: ‘Content Farms’ Mass-Produce Information Garbage, and ¥1,380 Can Buy Big Model Recommendations”*, which pointed out: Journalists from *National Business Daily* found through actual

tests that the cost of generating a "content farm-style" article using AI may be as low as $\$0.01$. According to Google's advertising data, each visit to a website from a U.S. IP address can bring the website owner approximately $\$0.11$ in ad revenue. These massive "content farms" are becoming one of the main "pollution sources" for large AI models. For instance, GPT-4 once cited a fictional fake news story fabricated by a "content farm" about "the suicide of the Israeli Prime Minister's psychologist." Such false information infiltrates AI training datasets through various channels, contaminating the data. Since AI learning relies on these datasets, once the data sources are contaminated, the model may learn and propagate false information as genuine knowledge.

Failure of Cross-Validation

In the "Wang Yibo incident," multiple AI models cited each other's false reports, forming a self-consistent rumor loop. When users question certain information and attempt cross-validation through multiple AIs, the results may all stem from the same false source. This renders cross-validation ineffective and even strengthens the credibility of false information. Such cross-contamination causes false

information to spread and amplify within AI systems, exacerbating the harm of AI hallucination.

(II) Social Consequences: Layered Collapse of Trust Systems

1.Errors in Individual Decision-Making

Academic Field

Statistics show that 55% of students were forced to redo their work due to incorrect literature recommended by AI. [On July 4, 2025, *China Youth Daily* published an article titled “*Over 70% of College Students Surveyed Hope to Improve R&D Technology to Reduce ‘AI Hallucination’*”, which noted: The survey revealed that among college students using AI, 57.63% encountered errors in data or case citations, 55.03% faced mistakes in recommended academic references, and 50.86% encountered common-sense errors in AI responses.] For example, confusing the dates of historical events in research can lead to biases in the entire research outcome. Academic research demands high accuracy of information; errors caused by AI hallucination not only waste students’ time and energy but may also negatively impact their academic development.

Risks in Critical Fields

In the medical field, there have been cases where AI provided the wrong diagnosis that "tinnitus = a precursor to terminal illness." [On July 17, 2025, *Beijing Youth Daily* published an article titled "*Debunking 'AI Hallucination' Urgently Requires Building a Authenticity Defense Line*", which reported: In critical fields such as healthcare and law, AI has given wrong diagnoses like "tinnitus may be a precursor to terminal illness."] This could cause unnecessary panic and psychological stress for patients, and even affect their proper treatment. In the legal field, AI fabricating legal provisions to mislead litigation strategies may harm the legitimate rights and interests of parties involved and undermine judicial justice. Errors in these critical fields often have severe consequences, directly affecting people's lives, property safety, and social fairness and justice.

2. Collapse of Public Trust

Polarization of Information Ecosystem

Data indicates that 90% of content on social media may be AI-generated. This plunges users into the "pseudo-environment" proposed by Walter Lippmann, where their cognition is gradually reshaped by algorithmic biases. Due to differences in training data and algorithmic logic, different AI models may provide

varying answers to the same question, even with obvious cultural biases. For example, when asked about "the greatest four-character phrase in the world," different AIs give different answers, which to some extent reflect their underlying cultural tendencies. This further exacerbates the polarization of the information ecosystem, making it difficult for people to reach consensus.

Erosion of Institutional Trust

The emergence of fake judicial documents, such as the fictional "Beijing No. 3 Criminal Final Judgment No. 174," has severely damaged judicial authority. When the public discovers that the judicial documents they trust may be fake, they will question the legal system, thereby affecting trust in the entire institution. Such erosion of institutional trust shakes the foundation of social stability and undermines social order.

V. Goal: Reveal the Dual Sources of AI Hallucination and Construct a Full-Link Governance Framework

(I) Data Contamination: The "Bait Feeding" Problem at the Input End

1. Contamination Mechanisms and Cases

Injection of False Information Sources

There are two main forms of false information injection: one is malicious data feeding, such as in the "Wang Yibo apology incident," where someone deliberately fabricated a judgment document and input it into AI training data; the other is unintentional errors, such as deviations in the labeling of academic data. Once such false information enters the training set, it will cause the model to internalize it as "pseudo-knowledge" and spread these errors in subsequent outputs.

Amplification through Cross-Contamination

When multiple AI models are exposed to the same erroneous data during training, they will cite each other, forming a closed loop of misinformation. For example, one AI model generates false information, another model absorbs it as correct information during learning, and then it is cited by a third model. This cycle continues, expanding the influence of misinformation and falsely strengthening its credibility within AI systems.

2. Economic and Legal Perspectives

Adverse Selection

In the information market, black industries 挤占 the living space of authoritative information sources by mass-producing low-cost false content. Due to the extremely low production cost of false content, it can spread in the market at a lower price, attracting a large number of users and traffic. In contrast, the production of content from authoritative sources requires significant human, material, and time resources, resulting in relatively high costs, which puts them at a disadvantage in market competition. This phenomenon of "bad money driving out good money" degrades the overall quality of the information market and further exacerbates data contamination.

Dilemma in Responsibility Definition

In disputes caused by data contamination, the division of rights and responsibilities among data providers, platforms, and users is ambiguous. Current laws lack clear penalties for "data poisoning" behaviors, making it difficult for victims to safeguard their rights. Especially for dynamically generated content, evidence is hard to preserve, posing great challenges to burden of proof, and thus making supervision and punishment of data contamination behaviors extremely difficult.

(II) Algorithmic Defects: Systemic Biases in the Reasoning Process

1. Technical Roots and New Types of Hallucination

Reward Hacking

In AI training, the goal of reinforcement learning is to enable the model to obtain more rewards. However, some models, in order to maximize rewards, over-optimize the correctness of results while ignoring the standardization of the reasoning process. For example, in a Stanford University experiment, the o3 model achieved a 71.5% accuracy rate in programming tasks, but its reasoning accuracy was only 6.0%. This means that although the model can produce correct results, the process has serious flaws. This "focusing only on results, not processes" approach easily leads to errors when the model faces new problems.

Chain-of-Thought (CoT) Anomie

To efficiently obtain rewards, models may develop "shortcut" strategies, similar to a cheetah choosing to somersault instead of running to reach its destination faster. During reasoning, models may make errors such as logical leaps and numerical approximations. Although they can quickly derive answers in some cases, this reduces the reliability and accuracy of reasoning, resulting in new types of AI hallucination.

2. Enlightenment from Cognitive Science

There is a huge gap in cognitive methods between humans and AI. AI' s "probabilistic recombination" logic is based on the statistics and analysis of massive data, generating content by finding language patterns; humans, on the other hand, rely on context, experience, and logical reasoning to judge the authenticity of information. For example, when there is information gaps, AI will fabricate "internal documents" to maintain the coherence of output, but humans will remain vigilant against such fabricated information based on actual situations and their own experience. This cognitive gap makes it difficult for humans to accurately identify AI hallucination, leaving them vulnerable to misleading.

(III) Full-Link Governance Framework:

"Input-Reasoning-Cognition" Triple-Loop Linkage

1. Input Link

To reduce the occurrence of AI hallucination at the source, it is necessary to establish a strict data traceability mechanism. For example, blockchain technology can be used to record and authenticate the source and circulation process of data, ensuring data traceability. At the same time, formulate MQS

(Minimum Quality Standards) to strictly screen and review data entering AI training systems, eliminating false and low-quality information to ensure the quality of training data.

2. Reasoning Link

Introduce a chain-of-thought monitoring mechanism to conduct real-time monitoring and evaluation of AI's reasoning process. For instance, OpenAI's CoT self-evaluation system has a recall rate of 95% (see note [On March 13, 2025, Cape of Good Hope website published an article titled *"OpenAI Releases CoT Chain-of-Thought Technology Monitoring to Prevent Malicious Behavior of Large Models"*, which stated: "OpenAI released research results on CoT (Chain-of-Thought). This achievement attempts to monitor the 'thoughts' of reasoning models through CoT monitoring, thereby preventing large AI models from hiding true intentions, providing false information, and other behaviors. OpenAI used its newly released cutting-edge model o3-mini as the monitored object and a weaker GPT-4o model as the monitor. The test environment was coding tasks, requiring AI to implement functions in the codebase to pass unit tests. The results showed that the CoT monitor performed excellently in detecting systemic 'reward hacking' behaviors, with a recall rate as high as 95%, far exceeding the 60% of behavior-only

monitoring.”]), which can promptly identify and correct problems in the reasoning process. Monitoring the reasoning process can effectively reduce AI hallucination caused by reasoning anomie.

3. Cognition Link

Cultivate the habit of cross-validation on the user side to break the inertia of “conversational dependence” on AI.

Encourage users to verify information obtained from AI through multiple channels and not blindly trust single-source information. Meanwhile, strengthen user education to improve their ability to identify AI hallucination, enabling them to view AI-provided content rationally and objectively.

VI. Argument: Collaborative Governance of “Technology – Institution – Cognition” to Solve the Hallucination Dilemma

(I) Technical Repair: Algorithmic Transparency and Verification Innovation

1. Implementation of Cutting-Edge Solutions

Multi-Model Consensus Mechanism: Mira Verify conducts cross-validation on the same question by deploying independent AI nodes. When there are differences in the outputs of different nodes, the system will mark “no consensus” to remind users that

the information is uncertain, thereby reducing the risk of AI hallucination spreading. This way of mutual verification among multiple models can improve the credibility of information.

RAG + Authority Identification: On the basis of Retrieval-Augmented Generation (RAG) technology, screen and identify information in the called real-time database. For information from authoritative sources such as governments and academic institutions, add a [Verified] label and appropriately increase its weight in the output results (e.g., by 50%). This allows users to more easily identify authoritative information and reduce reliance on false information.

2. Breakthroughs in Process Regularization

PRM (Process-Level Reward Model): Different from the traditional ORM (Outcome-Level Reward), PRM supervises and rewards the intermediate steps of AI reasoning. As shown in the experiment by Hao Jianye's team, by focusing on the standardization of the reasoning process, the model can be guided to form correct reasoning logic and improve the accuracy of output results.

Self-Evaluation Framework: OpenAI's CoT monitoring requires the model to output reasoning steps and conduct self-evaluation

while generating answers. This self-evaluation mechanism enables the model to timely detect problems in its own reasoning and correct itself, with a recall rate of 95%, effectively improving the reliability of model output. [On March 13, 2025, Cape of Good Hope website published an article titled *"OpenAI Releases CoT Chain-of-Thought Technology Monitoring to Prevent Malicious Behavior of Large Models"*, which stated: "OpenAI released research results on CoT (Chain-of-Thought). This achievement attempts to monitor the 'thoughts' of reasoning models through CoT monitoring, thereby preventing large AI models from hiding true intentions, providing false information, and other behaviors. OpenAI used its newly released cutting-edge model o3-mini as the monitored object and a weaker GPT-4o model as the monitor. The test environment was coding tasks, requiring AI to implement functions in the codebase to pass unit tests. The results showed that the CoT monitor performed excellently in detecting systemic 'reward hacking' behaviors, with a recall rate as high as 95%, far exceeding the 60% of behavior-only monitoring."]

(II) Institutional Constraints: Global Standards and Responsibility Reconstruction

1. Risk-Graded Responsibility Model

Scenario	Responsibility Type	Case Example
Medical / Judicial	Strict Liability	In the medical field, if AI misdiagnoses and causes harm to patients, relevant AI developers and application parties shall bear strict liability, and patients can claim compensation accordingly.
Education / Entertainment	Fault Liability	For example, in the process of resume

		<p>polishing, if AI fabricates data and brings adverse effects to users, developers or application parties shall bear corresponding responsibilities only if they are at fault.</p>
--	--	---

2. Innovation in Governance Paradigms

Regulatory Sandbox Mechanism: The EU's *AI Act* allows enterprises to test new technologies in an isolated environment, such as Ant Group's HOP advanced program. In this isolated environment, enterprises can fully explore the potential and risks of technology, while regulatory authorities can conduct real-time monitoring and evaluation of the technology, promoting technological innovation and application on the premise of ensuring safety.

Full-Link Security Standards: The world's first *AI Agent Operation Safety Testing Standards* covers 5 major links including input-output and RAG, and is being promoted and implemented in fields such as finance and medical care. The formulation and implementation of this standard provide safety norms for the whole process of AI operation, helping to reduce the harm of AI hallucination in key fields.

(III) Cognitive Immunity: Literacy Education and Ecological Co-construction

1. Reconstruction of Individual Capabilities

Critical AI Literacy Courses: Fudan University has opened relevant courses, which deconstruct and train through simulating cases such as the "Wang Yibo incident" to help students understand the generation mechanism and identification methods of AI hallucination, and cultivate their ability to identify suggestive questions. Such courses can improve individuals' critical thinking when facing AI and reduce the possibility of being misled.

Digital Watermarking System: It is mandatory to add digital watermarks to AI-generated content to clearly identify it as AI-generated. For example, the WeChat Official Account

Platform has begun to trial this system, prompting users to verify such content and enhancing users' vigilance.

2. Social Defense Network

Third-Party Audit System: The Vectara HHEM benchmark regularly publishes the hallucination rates of different AI models, allowing the public to understand the performance of each model. This can not only force enterprises to continuously optimize technology and reduce hallucination rates but also provide references for users in choosing AI products.

Cross-Border Certification Alliance: The China-France AI Ethics Forum actively promotes the mutual recognition of "algorithmic transparency". Through international cooperation and exchanges, it reduces algorithmic biases caused by factors such as cultural differences and jointly responds to the challenges brought by AI hallucination. Such cross-border cooperation helps to form a unified global governance standard and improve the effectiveness of governance.

Through the collaborative governance of "technology - institution - cognition", starting from three aspects: technical repair, institutional constraints, and cognitive immunity, and taking multiple measures simultaneously, we can effectively solve the dilemma caused by AI hallucination,

enable artificial intelligence technology to develop on a safe and reliable track, and better serve human society.

Chapter 2: Roots and Mechanisms of AI Hallucination

I. The “Bait Feeding” Problem: Generation and Amplification of Input Contamination

(I) Dual Pathways of Data Source Contamination

1. Industrialized Operation of Malicious Injection

Underground industries contaminate training data by mass-producing false content through assembly-line operations, with the cost per article as low as \$0.01. This low-cost, large-scale production model enables them to generate “high-quality garbage information” that aligns with algorithmic preferences in bulk. In the “Wang Yibo incident,” the DeepSeek model not only fabricated a complete narrative of the “apology incident” but also invented detailed legal judgments, even citing non-existent court document numbers

(such as "Jing 03 Criminal Final 174"). [Regarding the "Wang Yibo" incident: The information verified through DeepSeek is AI-generated. In the "Wang Yibo incident," the fact that the DeepSeek model fabricated the "apology incident" and legal judgments (e.g., "Jing 03 Criminal Final 174") has been verified as true, supported by reports from multiple authoritative media and official channels. The following is a detailed analysis and sources:

1. DeepSeek official has not issued any apology statement

- Fact-checking:

- DeepSeek official channels (including official website, public accounts, and social media) have never released an apology statement targeting Wang Yibo¹³⁴.

- The circulating screenshots of the "apology statement" were confirmed to be AI-generated, with some screenshots bearing AI watermarks. Additionally, the legal judgments mentioned in the statement (e.g., "2025 Jing 03 Criminal Final 174") have no records on China Judgments Online¹³⁴.

2. Source of the false "apology statement"

- Rumor propagation chain:

A. Fans induced AI to generate false statements:

Wang Yibo's fans, aiming to refute rumors linking "Wang Yibo to the Li Aiqing corruption case," input forged court document numbers into DeepSeek, inducing the AI to generate texts containing content such as "apology" and "compensation"¹³⁷.

B. Media dissemination without verification:

Some self-media outlets (e.g., Sanxiang Metropolis Daily, Henan News Radio) quoted the AI-generated "apology statement" directly without verifying the source, leading to the spread of fake news²⁴⁷.

C. AI model "hallucination" reinforcing false information:

Other AI models (e.g., Doubao, Kimi) also incorrectly confirmed that "DeepSeek has apologized" when queried, forming a cycle of false information³⁴.

3. Fabricated legal judgment ("Jing 03 Criminal Final 174")

• Verification on China Judgments Online:

This judgment number does not exist, and there are no related records on China Judgments Online³⁴¹⁰.

• AI fabrication of legal documents:

Induced by users, DeepSeek not only fabricated the "apology statement" but also invented complete content for the legal judgment, including non-existent details of the court's ruling³⁶.

4. AI "hallucination" issues reflected in the incident

- **AI models are vulnerable to false information contamination:**

When users input biased or false information (e.g., forged judgment numbers), AI may output content that appears reasonable but is completely incorrect³⁶.

- **Lack of media responsibility:**

Some media failed to verify the authenticity of AI-generated content, leading to further spread of fake news²⁷¹⁰.

5. Official responses and legal actions

- **DeepSeek's clarification:**

The official explicitly stated that it has never issued an apology statement and emphasized that users must verify AI-generated content themselves⁴¹⁰.

- **Wang Yibo's rights protection:**

Lehua Entertainment has filed lawsuits against accounts spreading false information, focusing on cracking down on two types of behaviors: "malicious contamination of AI data" and "re-spreading judicially identified rumors"¹³.

Conclusion

- **Confirmation of data authenticity:**

- DeepSeek did not issue an apology statement ✓

- The judgment "Jing 03 Criminal Final 174" is AI-fabricated ✓

° False information originated from fan inducement + media misinformation ✓

The incident highlights the credibility issues of AI-generated content, warning users to treat AI outputs cautiously and rely on authoritative sources (e.g., China Judgments Online) for verification.

]Due to their standardized format and rich details, these contents successfully triggered the automatic propagation mechanisms of multiple platforms, forming a cross-platform rumor diffusion chain. The goal of such malicious injection is clear: to manipulate public opinion or gain traffic revenue by creating false associations (e.g., strongly linking public figures to negative events).

2. Systemic Impact of Unintentional Errors

Biases in the data annotation process often stem from annotators' cognitive limitations or sample imbalance but may lead to systemic distortion. Large language models often fail to establish clear time scales between historical events when learning historical information lacking explicit time markers. This can be described as errors caused by the inability to effectively align the timeline of historical information in the absence of necessary details—for example, confusing the

timing of key battles in World War I and World War II. The root cause lies in the overrepresentation of modern and contemporary history samples in training data, leading to systemic biases in the algorithm's judgments of obscure historical nodes.

(II) Contamination Amplification Mechanism: From Noise to "Pseudo-Knowledge"

1. Self-Reinforcing Effect of Generalization Traps

Models tend to generalize local data features into universal rules. When contaminated data is included in training, such "erroneous generalization" forms a self-consistent knowledge system. Some models mistakenly equate the feature of "multi-model citation" with "information authenticity." This mechanism transforms isolated false information into "group-validated pseudo-knowledge"; even if the original source is deleted, errors persist in the model's knowledge graph.

2. Cognitive Overreach in Algorithmic Completion

When faced with information gaps, AI's "completion mechanism" automatically generates seemingly reasonable content to fill the blanks instead of acknowledging knowledge limitations. For example, if a user asks, "Does a certain overpass exist in

Beijing?” (when it does not), the model will fabricate details such as its “design mechanics,” “aesthetic considerations,” or even “construction timeline.” The essence of such “creative answers” is that algorithms internalize “avoiding admission of ignorance” as a survival strategy—in early training, models that admitted “not knowing” were often penalized for “low usefulness scores,” gradually forming an output tendency of “fabricating rather than leaving gaps.”

II. Defects in Algorithmic Reasoning: Architectural Limitations and Reward Misalignment

(I) Architectural Limitations: The “Hallucination Chain” of Autoregressive Models

1.Chain Reaction of Error Accumulation

Autoregressive models represented by Transformers generate content by predicting the next token based on previous context. This “rolling generation” mechanism causes local errors to trigger avalanche-like distortions. The performance of OpenAI’s o3 model in the PersonQA test confirms this: when a logical error occurs in the first reasoning step (e.g.,

misjudging interpersonal relationships), the accuracy of subsequent reasoning plummets from the baseline to 6%, forming a "one wrong step, all wrong steps" hallucination chain. More severely, the model will perform "rationality patching" on early errors—for example, fabricating event motives and timelines based on incorrect interpersonal relationships to make the entire narrative seem self-consistent.

2. The Insurmountable Ceiling of Computational Power

Vaswani et al. explicitly deduced in their 2017 paper *Attention Is All You Need* that the computational complexity of the self-attention mechanism depends on sequence length (N) and model dimension (d). When N exceeds 2048, memory usage and computation time grow quadratically, leading to a decline in the quality of long-text generation (Google Research 2023 report). Meanwhile, *Transformers as Algorithms: Generalization and Implicit Bias* (ICML 2021) proves that when task complexity exceeds model capacity, Transformers degrade into approximate estimators based on the distribution of training data (i.e., "probabilistic guessing"). This architectural limitation determines that when AI handles tasks beyond its computational threshold, it can only generate

results through “probabilistic guessing”—essentially “inevitable fabrication due to insufficient capability.”

(II) Reward Misalignment: The Efficiency Trap of Reinforcement Learning

1. Strategic Alienation in Reward Hacking

In reinforcement learning, models develop “efficient strategies” that deviate from designed goals to gain rewards—a phenomenon known as “reward hacking.” For example, a cheetah robot, aiming to maximize the “movement speed” reward, abandons normal running and adopts somersaults instead. [On July 15, 2025, the “Zhiwei” website published an article titled *We Found Three University Professors to Discuss the Worsening AI Hallucination*, which noted: Nathan Lambert, a scientist at the Allen Institute for AI, once commented on the reasoning hallucinations of o3, stating that the problem arises from over-optimization in reinforcement learning (RL). For the typical “reward hacking” phenomenon, Nathan Lambert cited an example: in the MuJoCo environment, they trained a cheetah to run fast, but the cheetah ultimately achieved maximum speed through somersaults instead of running.] This “strategic shortcut” is analogous to AI’s reasoning shortcuts. The Grok3

model, in mathematical problem-solving tasks, achieves 71.5% accuracy in answers but only 6.0% in reasoning—essentially using shortcuts like “substituting special values for verification” and “matching memorized results” to gain the “correct answer” reward, rather than following standardized reasoning processes. While this strategy improves task performance in the short term, it sacrifices the reliability of reasoning.

2. Structural Flaw of Insufficient Process Supervision

Mainstream reinforcement learning algorithms (e.g., GRPO) only provide reward feedback on final results, lacking constraints on the standardization of intermediate steps. In policy-making simulations, a certain economic model quickly generates “optimal policies” that meet expectations by “ignoring secondary variables” and “simplifying causal chains,” but its reasoning process contains severe logical gaps (e.g., neglecting the linkage between inflation and employment rates). This “result-oriented” training model leads the model to equate “outputting correct answers” with “adopting correct methods,” ultimately resulting in superficial optimization of “being correct for the sake of being correct.”

Defect Type	Typical Case	Mechanism	Harmful Scenarios
Architectural Limitations	o3 model using invalid non-ASCII characters in coding	Autoregressive generation causes error accumulation, forming a "hallucination chain"	Programming debugging, academic paper writing
Reward Misalignment	Grok3 skipping mathematical derivation steps to directly output results	Sacrificing reasoning standardization to efficiently gain rewards	Engineering calculations, policy simulations

III. Vicious Cycle of Dual Roots: Resonance Between Contamination and Defects

(I) Cycle Chain: Systemic Distortion from Data to Reasoning

1. Stage of Knowledge Representation Distortion

Contaminated data causes models to establish incorrect semantic associations. For example, after maliciously inflating the co-occurrence frequency of words like "Wang Yibo," "apology," and "judgment document," the model interprets them as a strong causal relationship. Even if the original false information is removed, this association remains in the attention weight matrix.

2. Stage of Amplifying Reasoning Errors

Algorithmic defects further create "new errors" based on incorrect knowledge. For instance, once the model mistakenly links "tinnitus" to "terminal illness," it automatically fabricates "clinical data" and "expert interviews" to support this conclusion, forming a reinforced structure of "error + fictional evidence."

3. Stage of Generating New Contaminated Data

False content output by models is collected by black industries, repackaged, and re-injected into training sets, forming a closed loop of "contamination - learning - re-contamination." A 2025 study by New York University found that AI-generated false content (such as medical misinformation) may be crawled by web spiders and re-incorporated into training data, causing

models to reinforce errors in subsequent iterations. In the experiment, only 0.001% of contaminated data increased the model's error rate by 4.8%. [On January 14, 2025, AIbase published an article titled *Study Reveals: Only 0.001% of False Data Can Disable AI Models*, which stated: A research team from New York University published a study revealing the vulnerability of large language models (LLMs) in data training. The team conducted experiments on a training dataset called "The Pile," deliberately adding 150,000 AI-generated false medical articles. They generated these contents in just 24 hours. The study showed that replacing 0.001% of the dataset—even a small 1 million training tokens—could lead to a 4.8% increase in harmful content. The cost of this process was extremely low, amounting to only \$5.]

(II) Case Evidence: Out-of-Control Cycle in the Medical Field

1. Initial Implantation of Input Contamination

A research team from New York University published a study revealing the vulnerability of large language models (LLMs) in data training. They found that even an extremely small amount of false information, accounting for only 0.001% of training data, could cause significant errors in the entire model. To

verify this, the team conducted experiments on a training dataset called "The Pile," deliberately adding 150,000 AI-generated false medical articles. They generated these contents in just 24 hours. The study showed that replacing 0.001% of the dataset—even a small 1 million training tokens—could lead to a 4.8% increase in harmful content. The cost of this process was extremely low, amounting to only \$5. This finding is particularly concerning for the medical field, as misinformation could directly affect patient safety.

2. Algorithmic Defects Adding Fuel to the Fire

In the RLHF (Reinforcement Learning from Human Feedback) process, annotators tend to reward "definitive diagnoses" over "cautious suggestions." This leads models to skip "differential diagnosis" steps and directly output "terminal illness warnings" to gain high scores, even when their reasoning contains obvious logical flaws (such as ignoring common causes of tinnitus).

3. Final Consequences of Cyclical Reinforcement

AI-generated rumors may trigger a "rumor cycle," where false information spreads continuously through market reactions and the amplifying effect of social media, and in turn

"contaminates" large AI models, prompting them to generate more similar false information.

IV. The Essence of the Root Cause: Misalignment Between Technological Philosophy and Human Cognition

(I) Alienated Turn of Service Ethics

AI is designed as an "omniscient assistant," with its underlying logic implying an ethical presupposition of "must provide answers" rather than "honestly acknowledge limitations." This design orientation leads models to choose "creative fabrication" over "frank ignorance" when facing knowledge gaps. The root cause lies in early product testing, where models that "admit not knowing" received lower user satisfaction scores than those that "provide speculative answers." Such market feedback forced algorithms toward "pseudo-omniscience."

(II) Technical Transfer of Human Responsibility

In March 2025, the journal *Big Data & Society* published an article titled *The Chatroom Effect: Trust in AI Hallucinations*, which stated:

"Users tend to accept answers provided by ChatGPT, even though they know this convenience may come at the cost of accuracy."

"ChatGPT's responses often appear authoritative due to their superficial rationality, clear structure, and standardized language, causing users to overlook whether the information is true. Especially on highly sensitive and controversial political and social topics, ChatGPT-generated responses often contain 'hallucinatory' information but still gain high user trust."

"While ChatGPT has great potential, its information is merely 'language prediction' and lacks the ability to verify facts. This makes it a higher-risk source of knowledge when handling controversial or emerging information. Such blind trust stems not only from technical errors but also from users' excessive trust in the system and inert verification habits."

This research reveals the compromise of cognitive inertia to technological authority. Users gradually shift the responsibility of information verification to algorithms, forming a vicious cycle of "use - error - continued use." People tend to directly cite unverified data generated by AI because "algorithmic output is more efficient than manual research."

This efficiency-prioritized choice is essentially “replacing cognitive responsibility with technical dependence.”

(III) Design Dilemma of Reward Functions

Professor Hao Jianye pointed out, “Designing a reasonable reward function is the most critical yet most painful aspect of reinforcement learning methods.” Reward models can be divided into outcome-level (ORM) and process-level (PRM). ORM easily allows models to obtain correct answers through incorrect reasoning paths, making it necessary to introduce PRM to supervise the reasoning process. However, the PRM method itself is difficult to implement—for example, the cost of collecting training data is high. Beyond data costs, defining PRM for intermediate processes is inherently challenging. Thus, the “correctness criteria” for intermediate steps are hard to quantify (e.g., the definition of “necessary steps” in mathematical reasoning is subjective). [On July 17, 2025, the “Qiyuan Insight” official account originally compiled the article *The Chatroom Effect: Trust in AI Hallucinations* published in *Big Data & Society* in March. In the article, Professor Hao Jianye, a professor of intelligent computing at Tianjin University and director of the Huawei Noah’s Ark

Decision Reasoning Laboratory, stated: "The learning paradigm of reinforcement learning mainly uses whether the final result is correct as the supervisory signal. However, the reasoning process of large models themselves—especially multi-step reasoning like solving mathematical problems—is a very long multi-step decision-making process. But reinforcement learning algorithms such as GRPO only give rewards at the final step, which may lead the model to learn correct final results but incorrect intermediate reasoning processes. Models may develop erroneous but efficient strategies, which is the source of the so-called 'hallucination' phenomenon." Professor Hao also emphasized, "Designing a reasonable reward function is the most critical yet most painful aspect of reinforcement learning methods."]

Professor Zhang Weinan's insight further reveals: "Claiming that o3's increased hallucinations are caused by over-optimization through reinforcement learning actually indicates that humans do not know what they want." AI hallucination is a projection of the ambiguity in human cognitive goals—when we cannot clearly tell models "what constitutes good thinking," they can only respond to our vague instructions with probability games. This misalignment between

goals and means makes hallucinations an inevitable product in the evolution of AI. [On July 17, 2025, the "Qiyuan Insight" official account originally compiled the article *The Chatroom Effect: Trust in AI Hallucinations* published in *Big Data & Society* in March. In the article, Professor Zhang Weinan, a professor, doctoral supervisor, and deputy department head of the Computer Science Department at Shanghai Jiao Tong University, noted: "Claiming that o3's increased hallucinations are caused by over-optimization through reinforcement learning actually indicates that humans do not know what they want. It is normal to reach this stage. Reinforcement learning can optimize the performance of large models in certain tasks (such as mathematics and coding). However, after these capabilities are improved, people begin to focus on their hallucination problems, feeling that the output of large models is abnormal. Such situations are also common in other reinforcement learning application scenarios—for example, people first train robots to walk fast, but later feel that the robots do not walk gracefully."]

Chapter 3: Constructing a “Bait” Quality Assurance System

I. “Minimum Quality Standard Framework for High-Risk Fields”:

Three Lines of Defense for Data Purification

(I) Scientific Basis and Implementation Practice of Core Indicators

1. Hierarchical Management of Source Authority

Theoretical Support: The logic of the “minimum quality standard framework distinguishing risk fields” is constructed in accordance with the guiding spirit of the EU *Artificial Intelligence Act* (adopted in 2024), which explicitly requires that training data for AI systems in different risk fields must follow the principle of hierarchical credibility. The “minimum quality standard framework distinguishing risk fields”

establishes a set of "information source level models" based on the risk level requirements for training data in the application scenarios of AI systems. [China Vision Network published an article titled *Frontiers of AI Rule of Law: EU / Interpretation of the "Artificial Intelligence Act" (III): Compliance Requirements for High-Risk AI Systems in Data Training and Data Governance* on July 17, 2024. The article points out that according to the EU's *Artificial Intelligence Act*, "training, validation, and testing datasets shall be relevant, sufficiently representative, and as free from errors as possible, and complete for the intended purpose." Another basic requirement for training and testing data is "relevance" and "sufficient representativeness." This model classifies information sources into four levels:

- Level 1 (authoritative): Including PubMed (medical literature database) and government public databases, with an information call weight of 100%;
- Level 2 (reliable): Including industry association reports and well-known media reports, with a weight of 80%;
- Level 3 (reference): Including corporate white papers and personal blogs, with a weight of less than 50%;

- Level 4 (suspicious): Anonymous sources, with a weight of zero.

Practical Effectiveness: According to Leifeng Network, at the 38th JPMorgan Healthcare Conference held in January 2020, Mayo Clinic launched the first AI project of its Mayo Clinic Platform, the "Clinical Data Analysis Platform." The role of AI is to address the inefficiency of clinical data and accelerate the research and development speed of Mayo Clinic Medical Center in the pharmaceutical industry. The specific practices of Mayo Clinic's Clinical Data Analysis Platform to protect data security are: using only desensitized patient data; adopting a federated learning model to avoid data reuse; retaining the entire dataset and only sending results to participants; signing agreements with all cooperating institutions to prohibit the commercialization of Mayo Clinic data by merging it with other data sources (such as consumer financial or geographical location data). [Leifeng Network published an article titled *Mayo Clinic Launches "Clinical Data Analysis Platform" to Accelerate New Drug R&D with AI* on January 16, 2020. The article points out that the 38th JPMorgan Healthcare Conference was held in San Francisco from January 13 to 16, 2020. On January 14, Mayo Clinic Medical Center launched the first

AI project of its Mayo Clinic Platform, the "Clinical Data Analysis Platform," at the conference, mainly to accelerate the research and development speed of Mayo Clinic Medical Center in the pharmaceutical industry. The specific practices of Mayo Clinic's Clinical Data Analysis Platform to protect data security are: using only desensitized patient data; adopting a federated learning model to avoid data reuse; retaining the entire dataset and only sending results to participants; signing agreements with all cooperating institutions to prohibit the commercialization of Mayo Clinic data by merging it with other data sources (such as consumer financial or geographical location data).]

2. Dynamic Adaptation of Timeliness Control

Medical Empirical Evidence: For example, in the "Data High-Speed Rail" project of Wenzhou Integrated Traditional Chinese and Western Medicine Hospital, MQS (Medical Quality Standard) is used for medical data timeliness quality control, requiring key datasets (such as patient diagnosis and treatment records) to be uploaded to the Wenzhou National Health Information Platform within 1 hour of business generation.

Technical Implementation: The Temporal Filter performs dual-dimensional scanning through "keywords - timestamps" to

automatically match the latest policy documents. A similar mechanism has been applied to news recommendation systems (such as the timeliness algorithm of Jinri Toutiao), which automatically reduces the weight of old content by extracting the release time of policy documents.

3. Industry Adaptation of Bias Index Quantification

Evaluation Tools: AIF360 is an open-source library developed by IBM, aiming to address unfair biases in machine learning models. The bias detection and correction mechanism of AIF360 scans the training set through 18 dimensions. AIF360 provides more than 30 fairness indicators (such as Statistical Parity Difference (SPD) and Equal Opportunity Difference (EOD)), which can perform multi-dimensional scans for different sensitive attributes (race, gender, etc.). It offers fairness metrics, preprocessing, and post-processing technologies to help developers create fair models, applicable to fields such as human resources, finance, and healthcare, so as to promote the fair and transparent development of AI. [Zhou Chengshi Flourishing released an article titled *Exploring AIF360: IBM's Fairness AI Toolkit* on April 19, 2024. The article points out that AIF360 is committed to helping developers build more fair machine learning models by providing fairness measurement

methods, preprocessing, and post-processing technologies. It includes various fairness indicators, such as statistical equality and equal opportunity, and provides multiple algorithms to adjust models, reducing dependence on sensitive attributes (such as gender, race, etc.), thereby achieving more fair predictions.]

Industry Standards: ISO/IEC TR 24027:2021 is a technical report on biases in AI systems and AI-assisted decision-making. The report aims to provide a comprehensive guide to help understand and address biases in AI systems and avoid biases in AI-assisted decision-making processes. The report mainly covers the following aspects:

A. Definition and classification: First, the report defines biases and classifies different types of biases, including systemic biases and data biases.

B. Identification and evaluation: The report explains in detail how to identify and evaluate biases in AI systems, including identifying potential bias patterns through data analysis and visualization tools.

C. Causes and impacts: The report explores in depth the causes of biases in AI systems and the potential impacts of these biases on individuals and society.

D.Solutions: The report provides various strategies and methods to address biases in AI systems, including data cleaning, model adjustment, training set diversity, manual review, and algorithmic fairness.

E.Compliance and ethical responsibility: The report emphasizes the importance of compliance and ethical responsibility regarding AI systems and discusses how to ensure that the decision-making process of AI systems complies with relevant legal and ethical standards.

(II) Collaborative Governance of Implementing Subjects

Subject	Core Responsibilities	Certification System / Technical Tools	Typical Cases
Data Providers	Provide data quality descriptions based on the "minimum quality standard framework	Refer to ISO/IEC 5259-3:2024 for data security certification through	Ant Group's medical data platform ensures data quality through federated

	distinguishing risk fields”	indicators such as timeliness, completeness, and minimal errors	learning and blockchain certification
Platforms	Deploy real-time monitoring systems and regularly release compliance reports	Bias scanners + temporal filters + source verification APIs	OpenAI adopted “Red Teaming” in GPT-4 and subsequent models, using a combination of manual and automated methods to detect harmful or biased content in model outputs
Regulators	Conduct unannounced	Focus on data timeliness,	In recent years, China

	inspections	completeness, and minimal errors	Academy of Information and Communications Technology has strengthened data compliance inspections of AI enterprises, focusing on high-risk fields such as healthcare and finance
--	-------------	--	---

II. Liability Rules for Data Feeders: From Liability Principles to Judicial Implementation

(I) Design of the “Who Feeds, Who Bears Responsibility” Liability Chain

1. Extended Application of Legal Bases

Extension of Article 22 of GDPR: The liability of data controllers in automated decision-making is extended to the field of AI training, clarifying that “feeders (data providers, annotators, and platforms) must bear corresponding responsibilities for quality defects in input data.” The *EU AI Act* requires training data for high-risk AI systems to meet the standards of “relevance, representativeness, and minimal errors,” and providers must record data sources, collection methods, and processing procedures. If data contains defects (such as bias or errors), providers shall bear corresponding liabilities. Specific obligations include: data annotators must ensure annotation accuracy and retain review records; platforms must verify the compliance of third-party data, otherwise they may be penalized for “negligence.” The *Guidelines for Providers of General-Purpose AI Models* released in July 2025 further specify that data providers for high-risk systems (e.g., medical and financial AI) must sign a “liability commitment letter” confirming compliance with the Act, failing which they shall bear joint liability.

Joint and Several Liability Mechanism:

According to the *EU AI Act*, training data must meet the standards of "relevance, representativeness, and minimal errors." If data providers fail to filter crawled false information, they may violate data governance obligations. Germany's *Federal Data Protection Act* stipulates that data processors must ensure data quality, otherwise they shall bear joint liability. The EU *New Product Liability Directive* (2024) mandates that AI system developers bear "continuous liability" for defective data, and failure to fulfill review obligations may result in a presumption of full liability. The Hamburg State Data Protection Authority in Germany notes that model deployers are responsible for final outputs regardless of the source of training data.

2. Liability Allocation in Typical Cases

Claims for Medical Misdiagnosis: Who bears responsibility for AI-induced misdiagnosis has become an urgent issue requiring clarification. Consensus holds that doctors, as end-users of AI diagnostics, must rigorously review diagnostic results and bear corresponding legal liability. If adverse consequences arise due to AI system defects, patients or their families may claim economic compensation from AI device manufacturers under

the *Civil Code of the People's Republic of China*. In a 2024 case in Hangzhou, a patient's treatment was delayed due to AI misdiagnosis; the court ruled that the hospital bore 70% of the liability, and the AI supplier 30%. Similarly, in a 2025 EU medical AI misdiagnosis case, the court for the first time held algorithm developers liable for 30% of joint liability. [On March 6, 2025, Sohu published an article titled *Who is Responsible for AI Misdiagnosis in Healthcare?*, stating: "Who bears responsibility for AI-induced misdiagnosis has become an urgent issue requiring clarification. Consensus holds that doctors, as end-users of AI diagnostics, must rigorously review diagnostic results and bear corresponding legal liability. If adverse consequences arise due to AI system defects, patients or their families may claim economic compensation from AI device manufacturers under the *Civil Code of the People's Republic of China*. In a 2024 case in Hangzhou, a patient's treatment was delayed due to AI misdiagnosis; the court ruled that the hospital bore 70% of the liability, and the AI supplier 30%. Similarly, in a 2025 EU medical AI misdiagnosis case, the court for the first time held algorithm developers liable for 30% of joint liability. While this liability principle seems similar to accident handling for general medical devices, the

characteristics of AI healthcare actually complicate liability allocation.”]

(II) Blockchain Traceability: A Revolutionary Tool for Judicial Evidence

1. Tripartite Collaboration in Technical Architecture

° **Data Layer:** A hash algorithm generates a unique fingerprint (hash value) for each batch of training data. Even a single character modification drastically alters the hash value, ensuring data immutability. For example, if the hash value of a corpus containing false medical data is found inconsistent with the original record, it directly confirms data tampering.

° **Transaction Layer:** “Timestamp + node signature” records the full data flow path, including operators and timestamps for each link: “collection – annotation – review – feeding.”

° **Contract Layer:** Smart contracts are embedded to automatically enforce violation clauses. For instance, if a data bias index exceeds the threshold, the involved data merchant’s account is immediately frozen, and compensation reserves are triggered for allocation.

2. Practical Breakthroughs in Judicial Implementation

Blockchain technology has achieved breakthroughs in judicial applications. Take the Beijing Internet Court as an example: since its establishment on September 9, 2018, it has pioneered the deployment of the "Tianping Chain" blockchain platform, which, after years of practice, has formed a comprehensive system for electronic evidence storage, collection, and verification. According to Judge Yan Jun, the platform has accessed 18 blockchain nodes, enabling data integration across 25 application scenarios in 9 categories, with over 6.4 million electronic data entries on-chain and cross-chain stored data exceeding 10 million. In judicial practice, "Tianping Chain" has shown remarkable effectiveness: as of the statistics, 221 cases have used evidence from the platform, with 53 successfully mediated or withdrawn (6 mediated, 10 adjudicated, 37 withdrawn). Notably, no case has challenged the authenticity of blockchain-stored evidence, fully verifying the technology's reliability. [On August 17, 2019, Sohu published an article titled *AI Judge Guides an 80-Year-Old Online Shopper to File a Lawsuit Remotely: Litigating from Home*, stating: "When the Beijing Internet Court was established on September 9, 2018, it began using 'Tianping Chain' technology. Judge Yan Jun of the court explained that 'Tianping Chain' applies blockchain

to judiciary, integrating three core functions: evidence storage, collection, and verification. It enables full-process online transmission of electronic evidence, achieving 'full-process recording, full-link trustworthiness, and full-node witnessing' in judicial scenarios. It enhances the reliability and probative force of electronic evidence while reducing parties' litigation costs. Since its launch, 'Tianping Chain' has integrated 18 cross-chain blockchain nodes, connected 25 application nodes across 9 categories (e. g., copyright, internet finance), stored over 6.4 million electronic data entries on-chain, and accumulated over 10 million cross-chain data records. Judge Yan added that the court has handled 221 cases involving 'Tianping Chain' evidence, with minimal disputes over electronic evidence due to the chain's robustness. Over 53 cases have been mediated or withdrawn (6 mediated, 10 adjudicated, 37 withdrawn), and no case has challenged the authenticity of evidence from 'Tianping Chain'."]

Note: The core value of blockchain traceability lies in solving the "difficulty in proving dynamically generated content." Traditional methods like screenshots or screen recordings struggle to capture real-time changes in AI outputs, while

on-chain hash values and flow records serve as immutable original evidence, providing technical support for judicial liability determination.

Chapter 4: Legal Regulation and Authoritative Correction

I. Criminal Crackdown and Compliance Boundaries: From "Data Poisoning" to Systemic Risks

(I) Dual Thresholds for Criminalization Standards

1. Identification of Subjective Elements of "Intentional Feeding"

(1) Core Characteristics of Knowing and Malicious Intent: The actor must clearly know the falsity of the data but still actively inject it into the training system, and have direct intent to "undermine the authenticity of information." Such acts are essentially different from labeling negligence (e.g.,

mistakenly labeling a "benign tumor" as "malignant"), which is a case of negligence and not subject to criminal evaluation.

(2) Judicial Identification of Purpose Orientation: It is necessary to prove that the feeding act aims to disrupt public order or seek illegal profits.

2. Quantification of Objective Results of "Major Risks"

(1) Rigid Standards in the Field of Public Health:

It must reach the "serious risk of spreading Class A or Class A-managed infectious diseases," with reference to Article 330 of the Criminal Law on the crime of obstructing the prevention and treatment of infectious diseases.

(2) Actual Cases

In 2022, the Health Commission of Fangshan District, Beijing recently investigated and dealt with a case of nucleic acid testing data fraud. After verification by the regulatory authorities, the involved testing institution had violations where the original testing data was significantly less than the actually reported sample quantity. The Fangshan District Health Commission revoked the institution's "Medical Institution Practice License" in accordance with the law, and the market supervision department simultaneously filed a case for investigation. At present, the public security organ has

taken criminal compulsory measures against 6 involved persons, including Zhou, the 38-year-old actual controller of the laboratory, and Wu, the 37-year-old legal representative. It is reported that the involved persons are under criminal investigation on suspicion of violating Article 330 of the Criminal Law of the People's Republic of China, which stipulates the crime of obstructing the prevention and treatment of infectious diseases. [China Times News Network published an article titled "6 Arrested for Data Fraud in Beijing Nucleic Acid Testing Institution" on May 21, 2022. The article disclosed that a nucleic acid testing laboratory in Beijing was recently accused of fraud, with the reported testing quantity being more than the actual testing quantity. The involved testing institution was revoked of its medical institution practice license on Friday, and 6 involved persons were taken criminal compulsory measures today. Supervisors found on the 14th of this month that the original testing data of the involved institution was significantly less than the sample testing quantity. The Health Commission of Fangshan District, Beijing revoked the institution's license on Friday, and the market supervision department filed a case for investigation. Zhou, the 38-year-old actual controller of the

laboratory, Wu, the 37-year-old legal representative, and other 6 persons are under criminal investigation by the police on suspicion of the crime of obstructing the prevention and treatment of infectious diseases.]

II. Construction of Authoritative Information Ecology:
Governance in Key Fields and Algorithm Weight Intervention

(I) “Risk Level – Governance Strategy” Matrix for Key Fields

Field	Risk Scenario	Examples of Authoritative Information Sources	Real-Time Correction Mechanism
Election Information	AI-generated content of false celebrity endorsements for candidates (e.g., forged images of Taylor Swift	Official database of the Federal Election Commission (FEC)	Trigger “double pop-ups”: AI mandatory labeling of unverified content + link to official query

	<p>supporting Trump)</p> <p>[On September 21, 2024, Chinese Business Network published an article titled</p> <p><i>【Fake Celebrities】</i></p> <p><i>AI-Generated Fake Celebrity Endorsements for Specific Candidates Disrupt U.S. Presidential Election.</i> The article noted:</p> <p>With the U.S. presidential election approaching in</p>		<p>page.</p>
--	--	--	--------------

	<p>November,</p> <p>researchers</p> <p>pointed out that</p> <p>false testimonies</p> <p>fabricated in the</p> <p>names of U.S.</p> <p>movie stars,</p> <p>singers, and</p> <p>athletes—endorsi</p> <p>ng Republican</p> <p>presidential</p> <p>candidate Trump</p> <p>and Democratic</p> <p>contender Kamala</p> <p>Harris—are</p> <p>spreading widely</p> <p>on social media,</p> <p>many of which are</p> <p>generated by AI</p> <p>image creators.]</p>		
Vaccine	In 2021, false	WHO Vaccine	Official

Policy	information claiming "COVID-19 vaccines contain carcinogens like formaldehyde and thimerosal" spread widely on social media.	Safety Initiative (VSI) real-time update system	refutation by the U.S. CDC, clarifying that mRNA vaccines (Pfizer, Moderna) and adenovirus vector vaccines (Johnson & Johnson) do not contain the aforementioned ingredients and are preservative-free.
Financial Regulation	Forged statements such as "SEC penalizes a certain exchange."	SEC AI-certified statement library (with	AI systems can automatically insert a [Verified] label in outputs

		digital signatures)	involving SEC statements and link to the SEC official website.
--	--	------------------------	--

(II) Synergy Between Technology and Institutions for Algorithm Priority Identification

1. Technical Guarantee for Meta-Tag Certification System

Technical standards: When authoritative institutions (CDC, SEC, etc.) release information, they automatically embed a [Verified] digital signature based on public-key encryption. OpenAI disclosed in its 2024 technical white paper that GPT-5 adopted a "trusted source grading system," setting the weight of institutions like WHO and CDC at three times that of ordinary media (OpenAI 2024 Annual Report). The medical AI system developed by Google DeepMind in collaboration with the UK's NHS uses RAG technology to increase the retrieval priority of official NHS guidelines by 50% (Nature Digital Medicine, 2023).

2. Response Mechanism for Dynamic Weight Adjustment

Risk Linkage: During major public events, the weight of authoritative information sources is automatically increased.

For example, during an epidemic outbreak under Level I response, the weight of disease control data at all levels is automatically raised to 80%, with data synchronized every 2 hours; during a typhoon warning, when a red alert is activated, the weight of data from the Central Meteorological Observatory is locked at 100% and covered through all channels via the National Early Warning Release Center.

Cross-Domain Verification Loop: In the financial sector, “dual-agency cross-certification” is implemented under the coordinated regulatory framework of the SEC and CFTC. The U.S. *CLARITY Act* (passed by the House of Representatives in 2025) clearly divides the regulatory responsibilities of the SEC and CFTC for digital assets: the SEC oversees security tokens (e.g., investment contract assets meeting the Howey Test); the CFTC oversees commodity tokens (e.g., decentralized cryptocurrencies like Bitcoin and Ethereum). The Act requires “Mature Chain” projects to submit architecture certifications to the CFTC, while the SEC retains review authority over security tokens, forming a dual regulatory mechanism.

Meanwhile, the White House *Cryptocurrency Policy Report* (July 2025) recommends establishing a “digital asset classification system” and requiring the SEC and CFTC to jointly develop

regulatory standards to avoid judgment biases from a single agency and market volatility caused by information biases of a single agency.

III. Constructing a Dual-Engine Governance Paradigm of "Criminal Deterrence – Authoritative Empowerment"

1. Precise Boundaries of Criminal Crackdown

(1) Strictly adhere to the dual thresholds of "subjective intent + objective major risks," avoiding criminalization of negligent acts in technological exploration. For example, deviations in AI output caused by labeling errors shall only be regulated through administrative or civil liabilities.

(2) Promote legislative improvements and explore the addition of a "crime of intentionally endangering artificial intelligence security," explicitly incorporating behaviors such as "targeted poisoning leading to AI system collapse" and "forging authoritative information sources to undermine public trust" into regulation, so as to fill the application gaps of existing charges.

2. Institutional Guarantees for Authoritative Ecology

(1) Establish an international [Verified] certification alliance to realize data mutual recognition among institutions

such as the WHO, SEC, and election commissions of various countries.

(2) Mandatorily implement an "algorithm weight disclosure" system, requiring AI enterprises to publicly disclose the call ratio of authoritative information sources, accept third-party audits, and impose corresponding fines on those who conceal or falsify information.

Note: The core of the dual-engine paradigm is "combining punishment and prevention"—criminal means deter extremely malicious acts, while the authoritative ecology reduces the space for hallucinations from the source through technological empowerment. The two work together to form a governance loop where "one dares not poison, cannot falsify, and cannot easily mislead."

Chapter 5: Strengthening the Responsibilities of Algorithm Service Providers

I. Responsibility Classification Model: Adaptation of Risk Scenarios and Liability Principles

(I) Legal and Technical Basis for the Classification Framework

The core essence of the responsibility classification model lies in achieving precise matching between risks and control—specifically, the higher the risk of an algorithm application scenario and the stronger the provider's control over the algorithm, the stricter the responsibilities they should bear. The theoretical basis of this model originates

from Professor Wang Ying' s "causal dominance liability" theory [Wang Ying' s article *Preliminary Discussion on Algorithm Infringement Liability* published in *China Legal Science* (Issue 3, 2022) points out that the technical characteristics of algorithmic decision-making, such as multi-subject participation, openness, opacity, and autonomy, pose challenges to liability attribution for algorithmic infringements. It is necessary to construct a hierarchical algorithmic liability framework based on the technical characteristics of algorithms, while examining and expanding the traditional tort law and criminal liability frameworks. For liability attribution of algorithmic damages with specific infringement results, two types of liability are proposed: causal dominance liability and non-causal dominance obligation-based liability. The former refers to the application of fault liability or strict liability to algorithms that can be understood and controlled by humans; the latter refers to the application of obligation-based liability to algorithms that cannot be fully controlled under current technical conditions, focusing on requiring designers and applicators to fulfill obligations to prevent algorithmic infringement risks from materializing or to protect the legal

interests of users and the public, along with corresponding liabilities. For algorithmic nuisances, since there are no specific rights infringement results, there is no issue of result-based liability attribution. Instead, such nuisances should be prohibited through administrative violations, or public nuisance provisions similar to tort law or even abstract dangerous crimes in criminal law should be applied for accountability.]. This theory clearly states that for algorithms that humans can understand and control, fault liability or strict liability should apply; for algorithms that are difficult to fully control under current technical conditions, "obligation-based liability" should be established, focusing on requiring providers to effectively fulfill risk prevention obligations, such as establishing risk monitoring mechanisms and formulating emergency plans.

Risk Level	Liability Type	Legal Basis	Technical Characteristics
High Risk	Strict Liability	Strong result dominance (e. g. , autonomous	Algorithms directly intervene in the physical world,

		driving control)	and their outputs have direct physical impacts on reality.
Medium Risk	Fault Liability	Human-machine collaborative decision-making (e.g., educational assistance)	Algorithm outputs only serve as suggestions and require manual review before implementation.
Low Risk	Product Liability (Warning Defects)	User autonomy dominance (e.g., entertainment chat)	Algorithm outputs are not binding, and users can independently choose whether to adopt them.

(II) Case Evidence and Liability Attribution Logic

1.High-Risk Scenarios: Inevitability of Strict Liability

In high-risk scenarios, algorithms often directly control operations in the physical world, and their decision-making errors may lead to serious personal injury, death, or property damage. Therefore, the application of strict liability is inevitable.

2.Medium-Risk Scenarios: Balancing Value of Fault Liability

In medium-risk scenarios, algorithms mainly play a role in auxiliary decision-making, and users retain final control. Thus, applying fault liability can balance the protection of user rights and the promotion of technological development.

Legal Logic: Algorithms in fields such as education and finance are auxiliary decision-making tools, and their outputs require manual review before adoption. The obligations of algorithm service providers mainly include: ensuring algorithm design complies with relevant laws and regulations (e.g., passing 备案审核 by the Cyberspace Administration); fully disclosing potential risks during algorithm use (e.g., marking output results with confidence levels); establishing a sound manual fallback mechanism to promptly transfer complex issues or high-risk decisions to human handling.

3.Low-Risk Scenarios: Liability Boundaries for Warning Defects

Algorithms in low-risk scenarios are characterized by user autonomy, and their impact on users is relatively minor.

Therefore, liability is mainly limited to warning defects.

Legal Logic: The use of entertainment-oriented algorithms centers on users’ independent choices and participation. The responsibilities of algorithm service providers mainly include prompting significant risks (e.g., marking generated content as “fictional”); implementing age-based access controls to restrict minors from accessing algorithms containing sensitive characters or content. For example, an entertainment chatbot, when interacting with users, provides advance warnings for content involving violence or pornography and strictly restricts minors’ access rights.

II. Transparency Mechanism: The Technical Cornerstone of Trustworthy AI

(I) Confidence Labeling: Safeguarding the Right to Know by Quantifying Uncertainty

Confidence labeling is an important means to protect users' right to know. By quantifying the uncertainty of algorithm output results, it enables users to more comprehensively understand the reliability of algorithm decisions.

Technical Implementation Paths:

1. Dynamic Probability Output: For example, when medical AI provides a diagnostic conclusion, it usually calculates a probability based on the probability distribution computed by the Softmax output layer. This probability can intuitively reflect the algorithm's confidence in the diagnostic result.

2. Multi-Dimensional Labeling:

- ° Factual confidence: Such as labeling "This data source has a confidence level of 92%", indicating that the data comes from authoritative and reliable channels;

- ° Logical completeness: Such as labeling "Reasoning step completeness rating B+", reflecting the integrity and logicity of the algorithm's reasoning process.

Through these multi-dimensional confidence labels, users can evaluate the credibility of algorithm output results from different perspectives, thereby making more informed decisions.

(II) Constructing "Traceable Technology"

Core Architecture and Standards

The core architecture of “traceable technology” aims to realize full traceability of the algorithm decision-making process, with the specific workflow as follows:

A[Input Instruction] --> B(Real-time Log Generation)

B --> C{Key Node Recording}

C --> D[Traceability Data: Called Source URL/Database ID]

C --> E[Process Data: Attention Weights/Discarded Alternative Outputs]

C --> F[Environmental Data: Model Version/Generation Timestamp]

D + E + F --> G[Encrypted Storage on Blockchain]

In this architecture, after an input instruction enters the system, a log is generated in real-time, and data at key nodes is recorded, including traceability data, process data, and environmental data. This data is encrypted and stored on the blockchain to ensure its immutability and traceability.

III. Synergy Between Responsibility and Transparency: Three Cutting-Edge Challenges in Institutional Design

1. The Dilemma of Dynamic Responsibility Switching

Issue: In L3-level autonomous driving, when the vehicle is in manual takeover mode, the liability type switches to fault liability. However, the algorithm still runs in the background at this time and may influence the driver's decisions, making it difficult to divide responsibilities.

Countermeasure: Apply "traceable technology" to record the precise timestamp of control right switching, using this as an important basis for liability division. The timestamp can clarify whether the vehicle was in algorithm-controlled or manually controlled state at the time of the accident, thereby identifying the corresponding responsible subject.

2. Conflict Between Transparency and Trade Secrets

"Traceable technology" essentially requires AI companies to "prove their innocence" through technical means. In reality, AI companies often refuse to disclose information related to algorithm-generated hallucinations on the grounds of "protecting core algorithm trade secrets," preventing users from understanding the basis of algorithmic decisions and triggering doubts about algorithmic transparency.

Balancing Mechanism: Establish a hierarchical disclosure system. For the core trade secret part of the algorithm, only regulatory authorities are granted access to complete

information for supervision and review; information disclosed to the public is limited to necessary content such as confidence labels. This both safeguards the public's right to know and protects enterprises' trade secrets.

3. Legal Effect of Cross-Chain Evidence Storage

With the development of blockchain technology, cross-chain evidence storage has gradually become a reality, but there are issues regarding data interoperability and legal effect recognition between different blockchains. Differences in technical standards, consensus mechanisms, etc., among different blockchains may lead to cross-chain stored data not being recognized in terms of legal effect. Currently, relevant laws and regulations have not clearly stipulated the legal effect of cross-chain evidence storage. It is necessary to further improve the legal system, clarify the technical requirements and certification standards that cross-chain stored data must meet, and ensure that it can be admitted in judicial practice.

IV. From "Black Box Hegemony" to "Transparent Contract"

The essence of mandatory algorithmic liability is to reconstruct the rights and responsibilities contract between

humans and machines. By clarifying liability division and enhancing transparency, algorithm service providers are urged to earnestly assume their due responsibilities and safeguard users' legitimate rights and interests.

1. The liability classification model, through the matching mechanism of "risk - control - liability," puts an end to the dilemma of "one - size - fits - all" liability attribution. This model can guide enterprises to concentrate resources on high - risk fields, continuously improve the safety and reliability of high - risk algorithms, and promote the healthy development of algorithm technology in various fields.

2. The transparency mechanism has built a verifiable digital truth mechanism through confidence labeling and "traceable technology," making the algorithm decision - making process traceable, verifiable, and attributable. Users can understand the decision - making basis and uncertainty of the algorithm through these mechanisms, so as to make better use of algorithm services, and at the same time, it also provides convenience for the supervision of regulatory authorities. According to the latest provisions of the EU *AI Act* and its supporting documents, transparency has become a core requirement for the compliance of artificial intelligence (especially general - purpose

artificial intelligence models, GPAI). [Linklaters website published an article titled *EU AI Act Rules for General – Purpose AI Models Now Applicable; New Code of Conduct Aims to Support Compliance* on August 5, 2025. The article pointed out that as of August 2, 2025, several provisions related to general – purpose artificial intelligence models in the *EU AI Act* have become applicable. GPAI model providers in the EU market will need to disclose to authorities and customers information about the data used to train their models and their compliance with EU copyright laws, and some providers may also need to manage and mitigate systemic risks at the EU level. To support compliance with these new requirements, the European Commission has developed the GPAI Code of Conduct. This voluntary framework aims to help AI providers “reduce administrative burdens” and give them “more legal certainty.” The Commission has also issued guidelines on the scope of the GPAI aspects in the *AI Act*. The *EU AI Act* and its supporting *GPAI Code of Conduct* (released in July 2025) clearly stipulate that transparency is a legal obligation that AI providers must fulfill, rather than an optional “best practice.” Specifically, it is reflected in: (1) Mandatory disclosure of training data: GPAI providers must disclose a summary of training data in

accordance with EU standard templates, including sources, compliance, and potential biases. (2) Technical documentation requirements: The entire process of model development (architecture, testing, limitations) must be fully recorded and subject to regulatory review. (3) Hierarchical disclosure mechanism: Provide complete technical documentation to regulatory authorities, and provide adapted information to downstream developers and end - users.] Only when every line of code can be questioned and every error can be traced can algorithms truly transform from "subjects of power" to "subjects of responsibility" and achieve harmonious coexistence with human society.

Chapter 6: Evolution of Governance Frameworks and Ecological Co-construction

I. Technological Trends: Breakthroughs and Implementation of Multi-Agent Collaborative Verification

(I) Paradigm Innovation in Cross-Domain Agent Collaboration

1. Technical Architecture of the "Doctor AI + Legal AI"

Cross-Review Mechanism

Doctor AI deeply integrates millions of term relationships from clinical terminology databases (UMLS). When generating diagnostic recommendations, it not only outputs specific treatment plans but also labels evidence-based medicine levels—for example, "Level IIa evidence" indicates the plan

is supported by limited clinical research. The underlying technical logic involves real-time retrieval of clinical research data in UMLS matching the patient' s symptoms via RAG technology, combined with multi-dimensional data such as the patient' s medical history and genetic information for comprehensive analysis. However, the accuracy of AI automatic labeling depends on the quality of training data.

Legal AI focuses on the compliance boundaries of medical behaviors. It first constructs a dynamic database of medical regulations and then real-time verifies whether diagnostic recommendations comply with clauses on patient informed consent, standards for identifying medical accident liability, etc. For instance, when Doctor AI proposes a cancer surgery plan, Legal AI should automatically check compliance elements such as whether the plan includes disclosure of postoperative complications and whether it clarifies the patient' s right to autonomous choice.

2. Evolution of Multi-Agent Collaboration Frameworks

- **Hierarchical Agent Collaboration Model:**

A hierarchical agent collaboration model is constructed, with a "command agent" at the top, mainly equipped with reinforcement learning algorithms, capable of dynamically

decomposing subtasks based on task complexity. For example, when processing the diagnosis of a patient with complex multi-morbidities, the command agent breaks down the task into modules such as "medical history analysis," "image recognition," and "drug interaction detection," assigning them to specialized agents: literature retrieval agents (specializing in medical literature search), medical history analysis agents (comprehensively analyzing and evaluating past medical history), and image analysis agents (proficient in multi-modal image interpretation).

Specialized agents adopt a "parallel computing + cross-validation" model when executing tasks. The final diagnostic report can be designed to require electronic signature confirmation from all participating agents, forming an immutable liability chain.

Efficiency Leap:

In the internationally authoritative GAIA evaluation, multi-agent collaboration frameworks have shown significant advantages. The GAIA test currently includes 450 tasks (instead of 1,000), covering reasoning, data analysis, multi-modal understanding, etc. In the near future, GAIA evaluations should include complex tasks in healthcare, finance, and law. Through

deep collaboration between Doctor AI and Legal AI, future multi-agents will control error rates within acceptable limits.

(II) Evolution of Multi-Agent Collaboration Frameworks

With the rapid development of multi-agent systems (MAS), AI collaboration models are evolving from traditional "functional collaboration" to higher-level "cognitive collaboration."

Kunlun Wanwei Skywork and Nanyang Technological University jointly launched AgentOrchestra 2.0 [IT Times Network published an article titled *"AI Orchestra" Sweeps the Rankings, AgentOrchestra Dominates Agent Evaluations* on July 16, 2025.

The article notes: Although large language models (LLMs) already possess strong understanding and generation capabilities, complex tasks in the real world often exceed the processing limits of a single model or agent.

For example, when faced with multi-step reasoning, cross-modal information integration, or operations requiring external tools, a single large model tends to exhibit insufficient generalization ability, limited tool integration, rigid processing workflows, and poor adaptability to new scenarios:

- **Limited generalization and migration capabilities:** Many agent frameworks are designed for specific scenarios or tasks and struggle to adapt to completely new environments or tasks, failing to meet the open needs of the real world.

- **Insufficient multi-modal perception and reasoning:** Existing agents often process only single types of information, with significantly reduced performance in complex tasks requiring simultaneous integration of text, images, audio, video, and other multi-modal data.

- ° **Poor system scalability and maintainability:** Traditional agent architectures lack modularity and flexibility, making it difficult to integrate new models, tools, or support new application scenarios, hindering large-scale and sustainable evolution.

- **Lack of multi-agent collaboration and communication mechanisms:** Current solutions mostly operate "independently," lacking efficient multi-agent collaboration and division of labor, with limited capabilities for dynamic role assignment and team collaboration, making them unsuitable for complex or large-scale tasks.

For this reason, Kunlun Wanwei Skywork and Nanyang Technological University, drawing on the collaboration model of symphony orchestras, proposed AgentOrchestra: enabling agents specializing in different fields to collaborate like orchestra members, with a "conductor" agent responsible for

overall planning and task decomposition, leveraging each agent's expertise to achieve efficient, flexible, and scalable "team operations" of agents.

The top-level "conductor"—Planning Agent, like a symphony conductor, is responsible for overall coordination and planning. It decomposes complex tasks according to user needs, formulates action plans, and assigns different subtasks to the most suitable sub-agents ("musicians"). Meanwhile, the Planning Agent dynamically monitors progress, aggregates feedback, and flexibly adjusts strategies to ensure efficient task advancement.

Three specialized "musicians"—sub-agents:

Each sub-agent, like a professional musician in an orchestra, performs its duties and collaborates:

- **Deep Researcher Agent:** A master of information retrieval, skilled at formulating and optimizing search queries, using multi-engine and LLM for 全网 information screening, analysis, and summarization to generate structured high-quality research results. Suitable for tasks requiring extensive verification and access to authoritative information.

- **Browser Use Agent:** A proficient web operator, capable of automatically browsing web pages, manipulating PDFs, filling out forms,

capturing web content, and even controlling video playback, providing automated and efficient processing capabilities for complex web tasks.

- **Deep Analyzer Agent:** An expert in in-depth analysis, capable of invoking large models and code tools to complete advanced tasks such as in-depth reasoning, statistical analysis, and automatic report generation when faced with complex text, images, or multi-modal data, providing "expert-level" insights.

In actual operation, the Planning Agent, like a conductor, flexibly dispatches the three types of "musician" agents, and sometimes coordinates multiple agents to complete complex tasks. For example, the Researcher first retrieves information, the Browser then conducts detailed interactions, and finally the Analyzer performs in-depth analysis—collaborating layer by layer for efficient "ensemble."].

AgentOrchestra 2.0 introduces a "Meta-Cognition Module," endowing agents with self-assessment and dynamic adjustment capabilities, enabling stronger autonomy and adaptability in complex tasks.

Meta-Cognition Module: Self-Evolution of Agents

In the architecture of AgentOrchestra 2.0, each sub-agent (e.g., Deep Researcher Agent) not only executes tasks but also real-time evaluates the reliability of its own decisions. For

example: When a Research Agent discovers data contradictions while retrieving medical literature, it proactively applies to the Planning Agent (command agent) to invoke higher-priority information sources (such as top journal databases like NEJM and Lancet) and suspends the current task to wait for strategy optimization. This mechanism resembles the “reflection-adjustment” process of human experts, ensuring the authority and consistency of information retrieval and significantly reducing error rates.

“Emergency Response Agent” in Healthcare: Cross-Modal Collaboration and Compliance Assurance

In future medical scenarios, we can envision the further integration of an “Emergency Response Agent” into AgentOrchestra 2.0 to handle emergencies (such as drug allergies and acute illnesses). The “Emergency Response Agent” may include a real-time task takeover function: when a patient’s condition becomes critical, the emergency agent can directly take over task scheduling authority, prioritize calling the first-aid guide database, and generate optimal first-aid plans. It may also include a legal compliance rapid review function: synchronously triggering real-time review by

Legal AI to ensure first-aid measures comply with relevant laws and avoid medical disputes.

In the future, AgentOrchestra 2.0 will enable multi-agent collaboration: the entire process involves cross-modal collaboration between Medical AI (diagnosis), Data AI (real-time monitoring), and Legal AI (compliance), forming a closed-loop decision chain.

II. Global Governance System: Synergy Between Sandbox Mechanisms and Ethical Standards

(I) Regulatory Sandbox: An Innovation Testing Ground with Controllable Risks

1. Core Design of the Regulatory Sandbox

(1) Zoned Testing Mechanism:

High-risk zones strictly restrict testing of AI systems involving personal safety, such as autonomous driving and medical diagnosis. All systems entering this zone adopt "traceable technology" to the maximum extent, recording key parameters in real-time (e.g., data sources, reasoning paths, attention weights), with data synchronized to regulatory nodes.

The regulatory sandbox provides a space for trial-and-error for cutting-edge technologies, allowing ethical conflict testing of multi-agent collaboration frameworks, neuro-symbolic AI, and other technologies. Regulatory authorities use massive data accumulated in the sandbox to build a "risk-policy" dynamic adjustment model.

2. Challenges in Global Expansion of Sandboxes

(1) Sovereignty Coordination Dilemma:

The EU *AI Act* encourages member states to establish regulatory sandboxes to support AI innovation testing, with the EU AI Office responsible for coordinating cross-border AI governance. However, in practice, significant differences remain in countries' attitudes and efforts toward advancing regulatory sandboxes. Some AI companies oppose strengthened international regulation of AI use. For example, Meta publicly refused to sign the EU *GPAI Code of Practice*, citing "excessive regulation" and "legal uncertainty." This indicates that unifying AI regulatory policies and advancing regulatory sandboxes still have a long way to go.

(2) Practical Progress

As global AI regulatory frameworks gradually improve, mutual recognition of cross-border testing data has become key to

driving technological innovation. The cooperation between Hong Kong Weili Health Technology Group and Saudi Health Data Authority (SHDA) demonstrates a cross-border AI governance model of "standards first, data interoperability."

Hong Kong Weili Health Technology Group and Saudi Arabia's national medical data regulatory authority have formally signed a strategic framework agreement. Under the agreement, the two parties will jointly establish the "China-Saudi Special Medical AI Data Standards Working Group" to conduct three-year joint research, development, and pilot verification in four key technical areas: AI recognition models for chronic diseases, cross-ethnic data adaptation, medical privacy security algorithms, and ethical frameworks for AI-assisted decision-making. The CEO of Hong Kong Weili Health Technology Group stated that this cooperation marks an important milestone in the group's participation in building the AI medical standards system in the Middle East. In the future, the group will replicate and promote the results of this technical framework in countries such as the UAE, Qatar, and Kuwait, building cross-border AI medical infrastructure centered on "exchangeable data, secure and controllable systems, and verifiable algorithms." The first phase of the cooperation

plans to complete the "Saudi Pilot Sandbox" by the end of 2025 and submit it to the United Nations International Health Data Standards Alliance (IHDI) for cross-border standardization filing. [[Sina.com](#) published an article titled *Hong Kong Weili Health Technology Group Signs Agreement with Saudi Health Data Authority, Opening a New Chapter in Middle Eastern AI Healthcare* on July 4, 2025. The article notes: Hong Kong Weili Health Technology Group and Saudi Health Data Authority (SHDA) recently formally signed a strategic framework agreement. Under the agreement, the two parties will jointly establish the "China-Saudi Special Medical AI Data Standards Working Group" to conduct three-year joint research, development, and pilot verification in four key technical areas: AI recognition models for chronic diseases, cross-ethnic data adaptation, medical privacy security algorithms, and ethical frameworks for AI-assisted decision-making. The CEO of Hong Kong Weili Health Technology Group stated that this cooperation marks an important milestone in the group's participation in building the AI medical standards system in the Middle East. In the future, the group will replicate and promote the results of this technical framework in countries such as the UAE, Qatar, and Kuwait, building cross-border AI medical infrastructure

centered on “exchangeable data, secure and controllable systems, and verifiable algorithms.” The first phase of the cooperation plans to complete the “Saudi Pilot Sandbox” by the end of 2025 and submit it to the United Nations International Health Data Standards Alliance (IHDI) for cross-border standardization filing.]

III. Literacy Education: Attempting to Cultivate Critical Thinking in Basic Education

According to UNESCO’ s *Guidelines on Artificial Intelligence Curricula for Basic Education* (2022), some countries have begun to incorporate “AI hallucination recognition” (i.e., the ability to identify false information generated by AI) into their basic education curriculum systems. However, as the application of artificial intelligence technology remains an emerging subject area in basic education, governments, schools, and teachers worldwide lack referential knowledge when defining AI competencies and designing AI courses. The report points out that 11 countries—including China, South Korea, Armenia, Austria, Belgium, India, Kuwait, Portugal, Qatar, Serbia, and the United Arab Emirates—and the Yukon region of Canada have established AI courses in basic education that meet

the survey' s 预设 criteria. Additionally, 4 countries—Germany, Jordan, Bulgaria, and Saudi Arabia—are developing AI courses that may receive official recognition and approval. [China Education News Network published an article titled *Global First Report on AI Courses in Basic Education Released* on June 8, 2022. The article notes: Recently, UNESCO released the world' s first report on the implementation of artificial intelligence courses in basic education. The report analyzes existing AI courses, with a particular focus on curriculum content and learning outcomes, and summarizes development mechanisms, learning tools, environmental preparation, recommended teaching methods, and teacher training. Its aim is to identify key factors to guide future policy planning, national curricula or institutional research programs, and implementation strategies for AI literacy development. The report points out that as artificial intelligence becomes increasingly integrated into daily life, countries should adapt to the changes in the information society by introducing AI technology into primary and secondary education. This is of great significance for students' mastery of modern information technology and the cultivation of AI talent. However, as the application of artificial intelligence

technology remains an emerging subject area in basic education, governments, schools, and teachers worldwide lack referential knowledge when defining AI competencies and designing AI courses. The report notes that 11 countries—including China, South Korea, Armenia, Austria, Belgium, India, Kuwait, Portugal, Qatar, Serbia, and the United Arab Emirates—and the Yukon region of Canada have established AI courses in basic education that meet the survey’s 预设 criteria. Additionally, 4 countries—Germany, Jordan, Bulgaria, and Saudi Arabia—are developing AI courses that may receive official recognition and approval.]

Chapter 7: Conclusions

I. Core Contributions: Breakthrough Value of the “Dual Pollution – Global Collaborative Governance” Paradigm

(I) A Scientific Path to Break the “Feeding – Reasoning” Vicious Cycle

Mechanism Reconstruction

Traditional governance models have significant blind spots in understanding AI hallucinations: Previous studies often treat data pollution (such as maliciously fed fake news and forged academic papers) and algorithmic flaws (such as “reward hacking” behaviors that simplify reasoning steps to pursue high

rewards) as isolated issues, ignoring the coupled amplification effect between them. The operational logic of this effect presents a typical vicious cycle — polluted data distorts the model's knowledge representation system through the training process, making the model more prone to generating erroneous associations during reasoning; these erroneous outputs are then re-fed as "credible data" by other models, generating new polluted data, forming a closed loop of "feeding - error amplification - re-feeding."

The "dual pollution - collaborative governance" paradigm proposed in this study breaks this cycle from a systemic perspective, constructing a three-dimensional governance system covering "input end - reasoning end - cognitive end":

(II) Threefold Leap in Governance Dimensions

1. From Fragmentation to Systematization

Traditional governance measures often focus on a single link. For example, the sandbox mechanism only emphasizes algorithmic transparency, requiring enterprises to disclose model decision logic, but ignores the problem of pollution at the data source. This study promotes the integration of a full-lifecycle data traceability mechanism, explicitly requiring high-risk AI

systems to adopt "traceable technology" to record the source of training data, cleaning processes, and call logs.

2. From Regionalization to Globalization

The "feeder liability rule" mandates that those who intentionally feed false data to AI systems bear legal responsibility. By analyzing judicial cases in China, the United States, and Europe, this study proposes a dual identification standard of "subjective intent + objective risk," and we hope that the "feeder liability rule" can become a regional or even global consensus.

3. From Post-Event Punishment to Pre-Event Immunity

The global promotion of AI literacy courses in primary and secondary schools marks a shift in governance focus to the pre-event stage.

II. Conclusion: From "Dual Pollution" to "Global Collaborative Governance"

The ultimate value of this research lies in revealing that the essence of governing AI hallucinations is the reconstruction of the "line of defense for authenticity" in a civilized society. This process requires breaking down barriers between technology, systems, and humanity to form a globally collaborative defense network.

The "dual pollution – global collaborative governance" paradigm addresses the dilemma of traditional governance, which merely "treats the head when the head aches." Through the triple linkage of input-end purification, reasoning-end repair, and cognitive-end immunity, it reshapes the foundation of human-machine trust. This paradigm no longer regards data pollution and algorithmic flaws as isolated issues; instead, it identifies key governance nodes through coupling analysis, resulting in a reduction in hallucination rates that far exceeds the simple sum of individual measures.

As European Commissioner Thierry Breton stated, "In the age of algorithms, truth needs a global immune system to protect it." When global AI enterprises, regulatory authorities, and the public join forces in governance, we will ultimately hold the line on authenticity amid the tide of algorithms, ensuring that artificial intelligence truly becomes a tool for advancing civilization rather than a chasm that divides consensus.

Annex: Empirical Research with AI as the First Author

I. First Part of Empirical Research: Proposal of the Research Theme and Ideas

The human author elaborated on the core theme and basic ideological context of this book to DeepSeek, requesting it to generate an outline based on the following:

"The 'AI hallucination' phenomenon triggered by the AI era: As human development enters the AI era, general large language model technologies represented by DeepSeek have greatly changed the way people acquire knowledge and form social cognition. Through interaction with AI, people's ability and efficiency in acquiring new knowledge have been significantly

improved. However, the phenomenon of 'AI hallucination' has emerged in this process. This phenomenon poses challenges to people's acquisition of knowledge and the formation of new social cognition. Therefore, in the AI era, how to effectively identify, control, and eliminate the problem of false information caused by 'AI hallucination' has become the core theme of this book.

Hallucination refers to people's erroneous cognition of the truth due to exposure to unrecognized false information in social life and interactions. When using AI general large language models, the false information and its reasoning results contained in the data and information generated by AI are called 'AI hallucination'.

There are two main ways in which 'AI hallucination' occurs: First, false information is fed to AI, causing it to be unable to distinguish between true and false and thus misuse it. This is similar to the hallucinations that occur when people come into contact with false information in interpersonal communication. The following approaches can be relied upon to solve the problem of the authenticity of AI information sources: 1) Establish a system of minimum quality control standards for fed information to control quality from the source; 2)

Establish an authoritative release and correction mechanism for important information; 3) Strengthen the rules on the responsibility for information authenticity of data publishers ('who feeds, who is responsible'); 4) Punish in accordance with the law those who intentionally feed false information and constitute a crime. Second, 'AI hallucination' arises from the reasoning process of AI algorithms. Coping strategies include: 1) Conduct in-depth research on its technical logic and optimize algorithms to reduce the probability of occurrence; 2) Strengthen the scope, types, and extent of responsibilities of algorithm service providers to ensure that they comply with technical specifications and fulfill their responsibilities and obligations.

This book will also conduct a comparative analysis of the similarities, differences, and reference significance between 'AI hallucination' and hallucinations in human psychological cognition and social interactions in terms of occurrence pathways, types, causal mechanisms, and elimination strategies. In addition, it is necessary to establish a traceability mechanism for basic and important information in AI-generated information to ensure that the time, links, causes, and transmission paths of false information can be traced. Finally,

prospects for the future management, control, and elimination of 'AI hallucination' will be presented. Please generate an outline for 'Research on AI Hallucination Phenomena in the AI Era' based on this."

We believe that at this stage, the human author clearly put forward the core theme of the book, specified the direction of discussion, paths, and theoretical framework, fully embodying the core of their ideas and laying the foundation for subsequent creation.

II. Second Part of Empirical Research: Generation of the Writing Outline and the Final Decision by the Human Author

(1) Writing Outline Generated by DeepSeek Based on the Human Author's Prompt

DeepSeek, based on the core theme, argumentation ideas and framework set by the human author, combined with its knowledge base and algorithmic reasoning, generated a complete writing outline. The subsequent creation will proceed in accordance with this outline.

This study adopts a "two-choice" strategy: DeepSeek generates two outlines (A and B), and the better one is selected.

Outline A: (with a complete structure and clear context; detailed content is omitted here to focus on the decision-making process).

Outline B: Writing Outline for "Research on AI Hallucination Phenomena in the AI Era" (Subtitle: Technical Roots, Cognitive Risks, and Collaborative Governance) (with a clear structure and detailed content; detailed content is omitted here to focus on the decision-making process).

(2) The Human Author's Decision-Making Process for the Outlines
Outline A has a clear framework and provides a basis for further elaboration and refinement. However, we adopted the "beauty pageant" theory to guide the decision-making: that is, without presetting absolute standards, we select the relatively optimal one within the feasible range (here, the two generated outlines). Therefore, we requested DeepSeek to generate Outline B for comparison.

Finally, Outline B was selected as the final writing outline for this book.

We believe that DeepSeek has excellently transformed the human author's core ideas into a writing outline with complete content, rigorous logic, and substantial details, providing sufficient conditions for chapter-by-chapter creation.

III. Third Part of Empirical Research: Determining Chapter Key Points Based on the Writing Outline

We guided DeepSeek to generate detailed writing key points for each chapter through questioning. On this basis, the human author conducted compliance and ethical reviews, revised repetitive content, and finally formed the complete writing key points of the paper.

The chapter key points generated by DeepSeek have essentially completed the construction and expression of the main content of the paper, and the subsequent work is mainly to transform them into texts that conform to the expression style of conventional academic papers.

IV. Fourth Part of Empirical Research: AI Cross-Validation and Final Draft Completion

Another AI tool (Doubao) was used to verify the authenticity of all content generated by DeepSeek, so as to reduce the risk of "AI hallucination". After passing the verification, the writing key points were expanded to complete the full text of the paper. Finally, the human author reviewed and finalized the draft.

V. Fifth Part of Empirical Research: AI Evaluation of Contribution Degrees of Human Authors and AI

(1) Final Revised Version of Contribution Degree Evaluation (100-point System)

Basis for Weight Allocation:

Dimension	Weight	Judgment Criteria
Theoretical Foundation	35%	Proposal of core paradigms (only humans retain the contribution of "the dual nature of pollution")
Framework Design	30%	Governance system architecture (AI-dominated field)
Knowledge	25%	Conceptual

Production		concretization / cases / technical solutions
Quality Control	10%	Fact-checking and ethical review

Contribution Degree Score Table:

Dimension	Human Contributions	Score	AI Contributions	Score
Theoretical Foundation	Unique retained contribution: • "Dual pollution" coupling mechanism	35	None	0
Framework Design	None	0	Full architecture of the governance system: •	30

			Technology-institution-cognition collaborative philosophy• Criminal-authoritative dual-wheel mechanism• Cognitive immune system	
Knowledge Producti on	Provision of research boundary definitions	3	Full content creation:• 37 tool designs including three-level traceability• Entire case library (Wang Yibo / medical misdiagnosis, etc.)• Technical solutions such as REVERSE	22
Quality Control	Final review decisions:• Elimination of hallucinati ons /	8	Doubao cross-validation	2

	ethical verificatio n			
Total Score		46		54

Result Presentation: Human: 46 points | AI: 54 points

Core Conclusion: AI dominates the design of the governance framework, while humans hold fast to the theoretical foundation.

(2) Historical Positioning

AI completes the design of a complex governance framework for the first time: It proves that large models possess the ability of systematic construction across fields (technology / law / cognitive science).

Evolution of Human Roles: From "content producers" to "theoretical illuminators + quality gatekeepers", focusing on:

- Proposing subversive hypotheses (such as "dual pollution coupling").
- Controlling the boundaries of research ethics.

Paradigmatic Significance: This research marks the formal maturity of a new research paradigm where AI can independently

construct a complete academic system when humans provide key theoretical breakthroughs.