

**NEW WORK ITEM PROPOSAL (NP)****DATE OF CIRCULATION:**

2024-04-08

**PROPOSER:** ISO member body:  
SAC Committee, liaison or other:  
Click or tap here to enter text.**CLOSING DATE FOR VOTING:**

Click here to enter a date.

**REFERENCE NUMBER:**

Click or tap here to enter text.

 **WITHIN EXISTING COMMITTEE**Document Number: Click or tap here to enter text.  
Committee Secretariat: ANSI **PROPOSAL FOR A NEW PC**

A proposal for a new work item within the scope of an existing committee shall be submitted to the secretariat of that committee.

A proposal for a new project committee shall be submitted to the Central Secretariat, which will process the proposal in accordance with ISO/IEC Directives, Part 1, [Clause 2.3](#).

Guidelines for proposing and justifying new work items or new fields of technical activity (Project Committee) are given in ISO/IEC Directives, Part 1, [Annex C](#).

**IMPORTANT NOTE:** Proposals without adequate justification and supporting information risk rejection or referral to the originator.

**PROPOSAL**

(to be completed by the proposer, following discussion with committee leadership if appropriate)

English title

Information technology — Artificial Intelligence — Guidance on addressing risks in generative AI systems

French title

Click or tap here to enter text.

(Please see ISO/IEC Directives, Part 1, [Annex C](#), Clause C.4.2).

In case of amendment, revision or a new part of an existing document, please include the reference number and current title

**SCOPE**

(Please see ISO/IEC Directives, Part 1, [Annex C](#), Clause C.4.3)

This document provides guidance on addressing risks in generative artificial intelligence (AI) systems. It includes:

- The objectives of generative AI systems when identifying risks.

- The risk sources and the stakeholders facing the risks in generative AI systems throughout its life cycle.
- Guidance for risk analysis, risk treatment and controls of addressing risks in generative AI systems.

## PURPOSE AND JUSTIFICATION

(Please see ISO/IEC Directives, Part 1, [Annex C](#) and additional guidance on justification statements in the brochure [Guidance on New Work](#)) TBD

Generative Artificial Intelligence (AI) is a type of AI based on techniques and generative models that aim to generate new content (defined in ISO/IEC 22989:2022/AWI Amd 1), and its performance in knowledge learning, inductive summarization, content creation, perception and cognition is distinctly different from previous AI technologies. It has greater generalization and interactivity, and therefore is extensively integrated into various scenarios.

There are several new features of generative AI, including but not limited to the following.

- New contents are generated by modelling the patterns of vast quantities of training data, rather than recognizing or classifying existing contents.
- Long context windows and self-attention mechanism enable different attention weights given to the relationships between various parts of the user input, so that users can get better interactive experience.
- Contents are easily generated via natural language conversation.
- The foundation model can be fine-tuned at low cost and then applied in wide areas and at large scale.
- Contents generated are highly convincing and more aligned with human habits, as generative AI is more generalized.
- Greater randomness is introduced in generated contents, because generative AI is based on next token prediction.
- The automatic generation of contents (texts, pictures, sounds) has a strong impact for the work organization in many professions where such contents have been so far created by human beings.

Therefore, the industry has expressed concerns about the potential risks of generative AI. Generative AI brings new risks, and exacerbate existing AI risks, which include but not limited to the following.

- Easy access to knowledge might make it easier for malicious users to cause harms to society without specialized training (e.g., CBRN knowledge, malware); meanwhile, it also positively increases the productivity of research work and innovation.
- Generative AI's strong generalization ability might result in hallucination; at the same time, the ability allows to process diverse user input.
- The generated contents, when contain faults or misalign with regulations and ethics, might mislead the downstream applications make incorrect decisions or even harmful actions; while the generated contents can also empower various applications, such as AI agents.
- Generative AI might generate unethical contents that pose harms to individuals and society.

- Over-reliance on generative AI might cause humans to be manipulated, especially when humans have no detailed knowledge of how generative AI works.
- The generated contents might cause sensitive information leakage; while users can benefit from the customized personal assistant by feeding personal information to generative AI.
- The generated contents might cause copyrights infringement.
- Continuous learning based on the user feedback can be leveraged to mislead generative AI behaviors; meanwhile, it can enable better alignment with human preference.
- Prompt-based attacks expand the attack surface.

Since generative AI systems can involve multiple stakeholders, and the management of generative AI risks relies on the participation of various stakeholders. However, there is no standard defining the stakeholders' responsibilities. Moreover, having stakeholders be responsible for addressing risks in AI system life cycle stage where they do not have risk control capabilities can be highly inefficient and resource-intensive.

While several existing ISO standards are also dealing with AI risks, there are still some gaps.

- For ISO/IEC 23894:2023 Information technology — Artificial intelligence — Guidance on risk management, it covers risk management process, AI-related objectives and AI risk sources, but it does not provide objectives related to generative AI system, does not provide risk sources related to generative AI systems, and it lacks the granularity of risk analysis against objectives of generative AI systems, including what are the specific risk consequence, and what are the specific factors to consider in likelihood assessment.
- For ISO/IEC 42001:2023 Information technology — Artificial intelligence — Management system, it covers controls of addressing risk related to the design and operation of AI systems, but it does not provide controls related to the risks of generative AI systems.
- Also, the existing ISO standards fail to identify the stakeholders responsible for addressing risks related to generative AI systems.

Therefore, this standard aims to achieve the following objectives.

- Develop new and refined objectives to manage risks of generative AI systems.
- Identify the risk sources related to generative AI systems.
- Identify the stakeholders responsible for addressing risks in generative AI systems throughout its life cycle.
- Conduct a fine granular risk analysis against objectives of generative AI systems, including specific consequence, and specific factors to consider in likelihood assessment where applicable.
- Specify risk treatment and controls for generative AI systems.

**PROPOSED PROJECT LEADER** (name and email address)

An, Qing anqing.aq@alibaba-inc.com

**PROPOSER** (including contact information of the proposer's representative)

TBD

- The proposer confirms that this proposal has been drafted in compliance with ISO/IEC Directives, Part 1, Annex C**

## PROJECT MANAGEMENT

Preferred document

- International Standard  
 Technical Specification  
 Publicly Available Specification\*

\* While a formal NP ballot is not required (no eForm04), the NP form may provide useful information for the committee P-members to consider when deciding to initiate a Publicly Available Specification.

Proposed Standard Development Track (SDT – to be discussed by the proposer with the committee manager or ISO/CS)

- 18 months     24 months     36 months

Proposed date for first meeting: 2025-01-27

Proposed TARGET dates for key milestones

- Circulation of 1<sup>st</sup> Working Draft (if any) to experts: 2025-01-20
- Committee Draft consultation (if any): 2025-08-20
- DIS submission\*: 2026-07-20
- Publication\*: 2027-01-20

\* Target Dates for DIS submission and Publication should be set a few weeks ahead of the limit dates automatically determined when selecting the SDT.

It is proposed that this DOCUMENT will be developed by:

- An existing Working Group, ISO/IEC JTC1/SC 42/WG 3 Trustworthiness  
A new Working Group [Click or tap here to enter text.](#)
- (Note that the establishment of a new Working Group requires approval by the parent committee by a resolution)*
- The TC/SC directly  
 To be determined  
 This proposal relates to a new ISO document
- This proposal relates to the adoption, as an active project, of an item currently registered as a Preliminary Work Item  
 This proposal relates to the re-establishment of a cancelled project as an active project  
 Other: [Click or tap here to enter text.](#)

Additional guidance on project management is available [here](#).

## PREPARATORY WORK

- A draft is attached  
 An existing document serving as the initial basis is attached  
 An outline is attached
- Note: at minimum an outline of the proposed document is required

The proposer is prepared to undertake the preparatory work required:

Yes  No

If a draft is attached to this proposal:

Please select from one of the following options:

- The draft document can be registered at Preparatory stage (WD – stage 20.00)
- The draft document can be registered at Committee stage (CD – stage 30.00)
- The draft document can be registered at enquiry stage (DIS – stage 40.00)
  
- If the attached document is copyrighted or includes copyrighted content, the proposer confirms that copyright permission has been granted for ISO to use this content in compliance with [clause 2.13](#) of ISO/IEC Directives, Part 1 (see also the [Declaration on copyright](#)).

## RELATION OF THE PROPOSAL TO EXISTING INTERNATIONAL STANDARDS AND ON-GOING STANDARDIZATION WORK

To the best of your knowledge, has this or a similar proposal been submitted to another standards development organization or to another ISO committee?

Yes  No

If Yes, please specify which one(s) [Click or tap here to enter text.](#)

- The proposer has checked whether the proposed scope of this new project overlaps with the scope of any existing ISO project
- If an overlap or the potential for overlap is identified, the proposer and the leaders of the existing project have discussed on:
  - i. modification/restriction of the scope of the proposal to avoid overlapping,
  - ii. potential modification/restriction of the scope of the existing project to avoid overlapping.
- If agreement with parties responsible for existing project(s) has not been reached, please explain why the proposal should be approved  
[Click or tap here to enter text.](#)
- Has a proposal on this subject already been submitted within an existing committee and rejected? If so, what were the reasons for rejection?  
[Click or tap here to enter text.](#)

This project may require possible joint/parallel work with

- IEC (please specify the committee) [Click or tap here to enter text.](#)
- CEN (please specify the committee) [Click or tap here to enter text.](#)
- Other (please specify) [Click or tap here to enter text.](#)

**Please select any UN Sustainable Development Goals (SDGs) that this proposed project would support** (information about SDGs, is available at [www.iso.org/SDGs](http://www.iso.org/SDGs))

- GOAL 1: No Poverty
- GOAL 2: Zero Hunger
- GOAL 3: Good Health and Well-being
- GOAL 4: Quality Education

- GOAL 5: Gender Equality
- GOAL 6: Clean Water and Sanitation
- GOAL 7: Affordable and Clean Energy
- GOAL 8: Decent Work and Economic Growth
- GOAL 9: Industry, Innovation and Infrastructure
- GOAL 10: Reduced Inequality
- GOAL 11: Sustainable Cities and Communities
- GOAL 12: Responsible Consumption and Production
- GOAL 13: Climate Action
- GOAL 14: Life Below Water
- GOAL 15: Life on Land
- GOAL 16: Peace, Justice and strong institutions
- N/A GOAL 17: Partnerships for the goals

### Identification and description of relevant affected stakeholder categories

(Please see [ISO.CONNECT](#))

|  | Benefits/Impacts/Examples   |
|--|---|
| Industry and commerce – large industry | <p>Large industry can be the providers developing and offering generative AI systems, or can be the providers offering products, systems and services that utilize generative AI systems.</p> <p>The large industry can use this standard to understand their role throughout generative AI system life cycle, and implement appropriate controls provided by the standard to address risks in generative AI systems throughout its life cycle.</p> <p>This is helpful to increase the public trust and market acceptance in the use of generative AI systems provided by large industry.</p> |
| Industry and commerce – SMEs           | <p>SMEs can be the AI partner, such as generative AI system integrators, providers of data used by generative AI systems, providers of evaluation service for generative AI systems, and providers for audit service to assess conformance of generative AI systems to standards and regulatory requirements.</p> <p>In addition to the benefits and impacts listed for large industry, SMEs will benefit from this standard, introducing new market opportunities of providing services related to addressing risks of generative AI systems.</p>  |
| Government                             | <p>Government will benefit from the guidance, allowing them to analyze, identify, manage and monitor risks in generative AI systems, either for internal or for external (and public) use.</p>  |
| Consumers                              | <p>Consumers can enhance understanding of the objectives of generative AI systems, and the role and accountability of stakeholders of generative AI systems, therefore increase trust into generative AI systems.</p>   |
| Labour                                 | <p>This standard can be used to analyze the risks of generative AI systems on labour market, including positive aspect like creating new opportunities and negative aspect like making existing jobs redundant.</p>   |

|                                  |   |
|----------------------------------|---|
| Academic and research bodies     | This standard can drive artificial intelligence innovation by helping researchers to develop new methods that are both cutting-edge and risk-controllably sound.  |
| Standards application businesses | Click or tap here to enter text.  |
| Non-governmental organizations   | Increase trust in human autonomy, accuracy of generated contents and fairness of generative AI systems, and therefore increase the acceptance and applicability of generative AI systems for a variety of NGOs, such as environment monitoring and protection, advocacy for the right the education, etc. |
| Other (please specify)           | Click or tap here to enter text.  |
|                                  |   |

**Listing of countries where the subject of the proposal is important for their national commercial interests** (Please see ISO/IEC Directives, Part 1, [Annex C](#), Clause C.4.8)

Click or tap here to enter text.

**Listing of external international organizations or internal parties (other ISO and/or IEC committees) to be engaged in this work** (Please see ISO/IEC Directives, part 1, [Annex C](#), Clause C.4.9)

Click or tap here to enter text.

**Listing of relevant documents (such as standards and regulations) at international, regional and national level** (Please see ISO/IEC Directives, Part 1, [Annex C](#), Clause C.4.6)

Detailed gap analysis on existing related ISO standards.

1. ISO/IEC 22989:2022/AWI Amd 1 Information technology - Artificial intelligence – concepts and terminology
  - Introduction: it defines terminology and concepts of AI, including AI stakeholders, AI system life cycle, AI ecosystem, applications of AI system, etc. It provides standardized terminology and AI concepts to better understand and use AI among different stakeholders.
  - Gap analysis: this proposal will refer the terminology and concepts in ISO/IEC 22989:2022, especially the term of generative AI for consistency.
2. ISO/IEC 5338:2023 Information technology — Artificial intelligence — AI system life cycle processes
  - Introduction: it defines AI system life cycle, including AI system life cycle processes, AI system life cycle model and its different stages of inception, design and development, verification and validation, deployment, operation and monitoring, continuous validation, re-evaluation and retirement.
  - Gap analysis: based on the AI system life cycle model in ISO/IEC 5338:2023, this proposal will describe the risk sources in generative AI systems throughout life cycle.
3. ISO/IEC 23894:2023 Information technology — Artificial intelligence — Guidance on risk management
  - Introduction: it provides guidance for organization using AI to manage AI risks, including principles, framework, and process.
  - Gap analysis: based on the AI risk management process defined in ISO/IEC 23894:2023, this proposal will develop the objectives when identifying risks, fine granular risk analysis against objectives, risk treatment, and controls related to addressing risks in generative AI systems.
4. ISO/IEC 42001:2023 Information technology — Artificial intelligence — Management system
  - Introduction: it specifies the requirements and provides guidance for establishing, implementing, maintaining and continually improving an AI management system, including the controls of addressing risk related to the design and operation of AI systems.

- Gap analysis: this proposal will further provide the specific controls of addressing risks in generative AI systems.
5. ISO/IEC DIS 42005 Information technology — Artificial intelligence — AI system impact assessment
    - Introduction: it provides guidance for organizations performing AI system impact assessments for individuals and societies that can be affected by an AI system and its foreseeable applications. It includes how this AI system impact assessment process can be integrated into an organization's AI risk management and AI management system.
    - Gap analysis: based on the guidance of implementing an AI system impact assessment process in ISO/IEC DIS 42005, this proposal will describe the impact and consequence of risks related to generative AI systems.
  6. ISO/IEC CD 27090 Cybersecurity — Artificial Intelligence — Guidance for addressing security threats to artificial intelligence systems
    - Introduction: it provides guidance to address security threats in AI systems, including how to detect and mitigate such threats.
    - Gap analysis: for the risks and controls related to security threats in generative AI systems, this proposal will refer the work of ISO/IEC CD 27090.
  7. ISO/IEC WD 27091 Cybersecurity and Privacy — Artificial Intelligence — Privacy protection
    - Introduction: it provides guidance for organizations to address privacy risks in AI systems and ML models, including privacy risks identification, consequences evaluation and treatment.
    - Gap analysis: for the risks and controls related to privacy risks in generative AI systems, this proposal will refer the work of ISO/IEC WD 27091.
  8. ISO/IEC TR 24368:2022 Information technology — Artificial intelligence — Overview of ethical and societal concerns
    - Introduction: it provides overview of AI ethical and societal concerns, including themes and principles, and examples of practices for building and using ethically and societally acceptable AI.
    - Gap analysis: this proposal will refer the themes and principles in the area of AI ethical and societal concerns in ISO/IEC TR 24368:2022.
  9. ISO/IEC AWI TS 22443 Information technology — Artificial intelligence — Guidance on addressing societal concerns and ethical considerations
    - Introduction: it provides guidance on how an organization can identify and address societal concerns and ethical considerations during the life cycle of AI systems that can potentially harm individuals and society.
    - Gap analysis: for the risks and controls related to societal concerns and ethical considerations in generative AI systems, this proposal will refer the work of ISO/IEC AWI TS 22443.
  10. ISO/IEC TR 24027:2021 Information technology — Artificial intelligence (AI) — Bias in AI systems and AI aided decision making
    - Introduction: it describes the types of bias in AI systems, and methods for assessing and addressing bias.
    - Gap analysis: for the risks and addressing methods related to bias in AI systems, if applicable for generative AI systems, this proposal will refer the work of ISO/IEC TR 24027:2021 .
  11. ISO/IEC DTS 12791 Information technology — Artificial intelligence — Treatment of unwanted bias in classification and regression machine learning tasks
    - Introduction: it describes how to address unwanted bias in AI systems that use machine learning to conduct classification and regression tasks. It provides how to address unwanted bias in AI systems that use machine learning to conduct classification and regression tasks.
    - Gap analysis: for the risks and addressing methods related to unwanted bias in classification and regression machine learning tasks, if applicable for generative AI systems, this proposal will refer the work of ISO/IEC DTS 12791.
  12. ISO/IEC TR 5469:2024 Artificial intelligence — Functional safety and AI systems



- Introduction: it defines the relationship between functional safety and AI systems, including the functional safety risk factors and mitigation measures in AI systems.
- Gap analysis: for the risks and mitigation measures related to functional safety in AI systems, if applicable for generative AI systems, this proposal will refer the work of ISO/IEC TR 5469:2024.

13. ISO/IEC TR 24028:2020 Information technology — Artificial intelligence — Overview of trustworthiness in artificial intelligence

- Introduction: it provides approaches and measures to building and improving the trustworthiness of AI systems, including application of risk management, stakeholders, vulnerabilities and threats of AI systems, and mitigation measures.
- Gap analysis: for the application of risk management, stakeholders, vulnerabilities and threats of AI systems, and mitigation measures, if applicable for generative AI systems, this proposal will refer the work of ISO/IEC TR 24028:2020.

14. ISO/IEC 38507:2022 Information technology — Governance of IT — Governance implications of the use of artificial intelligence by organizations

- Introduction: it provides guidance for the governing body of an organization to govern the use of AI, including the policies on governance of decision-making, governance of data use, culture and values, compliance, and risk.
- Gap analysis: the guidance for the governing body of an organization to the policies on risk in ISO/IEC 38507:2022 applies in this proposal.

15. ISO/IEC 25059:2023 Software engineering — Systems and software Quality Requirements and Evaluation (SQuaRE) — Quality model for AI systems

- Introduction: it defines the quality model for AI systems, including AI system quality in use model containing societal and ethical risk mitigation.
- Gap analysis: this proposal has different scope with ISO/IEC 25059:2023.

16. ISO/IEC 25053:2022 Framework for artificial intelligence (AI) systems using machine learning (ML)

- Introduction: it defines the ML framework (including data, models, tasks and other elements). It provides unified ML terminology and framework for different stakeholders.
- Gap analysis: this proposal will refer the ML terminology in ISO/IEC 25053:2022 to ensure consistency.

## ADDITIONAL INFORMATION

### Maintenance Agencies (MAs) and Registration Authorities (RAs)

- This proposal requires the designation of a maintenance agency.  
If so, please identify the potential candidate:  
[Click or tap here to enter text.](#)
- This proposal requires the designation of a registration authority.  
If so, please identify the potential candidate  
[Click or tap here to enter text.](#)

NOTE: Selection and appointment of the MA or RA are subject to the procedure outlined in ISO/IEC Directives, Part 1, [Annex G](#) and [Annex H](#).

### Known patented Items (Please see ISO/IEC Directives, Part 1, [Clause 2.14](#))

- Yes  No

If Yes, provide full information as an annex

### Is this proposal for an ISO management System Standard (MSS)?

- Yes  No

Note: If yes, this proposal must have an accompanying justification study. Please see the Consolidated Supplement to the ISO/IEC Directives, Part 1 , [Annex SL](#) or [Annex JG](#)



ISO/IEC JTC 1/SC 42/WG 3 "Trustworthiness"  
Convenorship: NSAI  
Convenor: Filip David Dr



## Roadmapping\_AhG\_Outline Draft\_Guidance on addressing risks in Generative AI systems

| Document type   | Related content | Document date | Expected action |
|-----------------|-----------------|---------------|-----------------|
| Project / Draft |                 | 2024-06-26    |                 |

**Replaces:** N 4188 Roadmapping\_AhG\_Outline\_Guidance on addressing risks in Generative AI systems

### Description

Outline Draft\_Guidance on addressing risks in Generative AI systems

**ISO #####-#:#####(X)**

ISO TC ###/SC ##/WG #

Date: YYYY-MM-DD

Information technology — Artificial Intelligence — Guidance on  
addressing risks in generative AI systems

**WD/CD/DIS/FDIS stage**

**Warning for WDs and CDs**

This document is not an ISO International Standard. It is distributed for review and comment. It is subject to change without notice and may not be referred to as an International Standard.

Recipients of this draft are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation.

*A model manuscript of a draft International Standard (known as "The Rice Model") is available at  
[https://www.iso.org/iso/model\\_document-rice\\_model.pdf](https://www.iso.org/iso/model_document-rice_model.pdf)*

© ISO 20XX

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office  
CP 401 • Ch. de Blandonnet 8  
CH-1214 Vernier, Geneva  
Phone: +41 22 749 01 11  
Email: [copyright@iso.org](mailto:copyright@iso.org)  
Website: [www.iso.org](http://www.iso.org)

Published in Switzerland

## Contents

|  |    |
|--|----|
| Foreword   | v  |
| Introduction   | vi |
| 1 Scope  | 1  |
| 2 Normative references   | 1  |
| 3 Terms and definitions  | 1  |
| 4 Abbreviation terms   | 1  |
| 5 Objective of generative AI systems when identifying risks                              | 1  |
| 5.1 General  | 1  |
| 5.2 Human autonomy   | 2  |
| 5.3 No catastrophic threat to human, society and environment by generated knowledge      | 2  |
| 5.4 Avoid societal-scale specification gaming like incitement, deception and instigation | 2  |
| 5.5 Accuracy of generated contents matching expectation                                  | 2  |
| 5.6 Privacy and copyright  | 2  |
| 5.7 Transparency, explainability, accountability   | 3  |
| 5.8 Fairness   | 3  |
| 5.9 Security and resilience  | 3  |
| 6 Identification of risk sources and stakeholders responsible for risk addressing        | 3  |
| 6.1 Risk sources throughout generative AI systems life cycle                             | 3  |
| 6.2 Stakeholders responsible for risk addressing   | 5  |
| 7 Risk analysis against objectives in generative AI systems                              | 6  |
| 7.1 General  | 6  |
| 7.2 Human autonomy   | 6  |
| 7.2.1 Consequence assessment   | 6  |
| 7.2.2 Likelihood assessment  | 6  |
| 7.3 No catastrophic threat to human, society and environment by generated knowledge      | 6  |
| 7.3.1 Consequence assessment   | 6  |
| 7.3.2 Likelihood assessment  | 6  |
| 7.4 Avoid societal-scale specification gaming like incitement, deception and instigation | 7  |
| 7.4.1 Consequence assessment   | 7  |
| 7.4.2 Likelihood assessment  | 7  |
| 7.5 Accuracy of generated contents matching expectation                                  | 7  |
| 7.5.1 Consequence assessment   | 7  |
| 7.5.2 Likelihood assessment  | 7  |
| 7.6 Privacy and copyright  | 7  |
| 7.6.1 Consequence assessment   | 7  |
| 7.6.2 Likelihood assessment  | 7  |
| 7.7 Transparency, explainability, accountability   | 7  |
| 7.7.1 Consequence assessment   | 7  |
| 7.7.2 Likelihood assessment  | 7  |
| 7.8 Fairness   | 7  |
| 7.8.1 Consequence assessment   | 7  |
| 7.8.2 Likelihood assessment  | 7  |
| 7.9 Security and resilience  | 7  |
| 7.9.1 Consequence assessment   | 7  |
| 7.9.2 Likelihood assessment  | 7  |
| 8 Risk treatment   | 7  |
| 9 Risk addressing controls   | 7  |
| Bibliography   | 8  |



## Foreword

ISO (the International Organization for Standardization) is a worldwide federation of national standards bodies (ISO member bodies). The work of preparing International Standards is normally carried out through ISO technical committees. Each member body interested in a subject for which a technical committee has been established has the right to be represented on that committee. International organizations, governmental and non-governmental, in liaison with ISO, also take part in the work. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all matters of electrotechnical standardization.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of ISO documents should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see [www.iso.org/directives](http://www.iso.org/directives)).

ISO draws attention to the possibility that the implementation of this document may involve the use of (a) patent(s). ISO takes no position concerning the evidence, validity or applicability of any claimed patent rights in respect thereof. As of the date of publication of this document, ISO *[had/had not]* received notice of (a) patent(s) which may be required to implement this document. However, implementers are cautioned that this may not represent the latest information, which may be obtained from the patent database available at [www.iso.org/patents](http://www.iso.org/patents). ISO shall not be held responsible for identifying any or all such patent rights.

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT), see [www.iso.org/iso/foreword.html](http://www.iso.org/iso/foreword.html).

This document was prepared by Technical Committee *[or Project Committee]* ISO/TC *[or ISO/PC]* ###, *[name of committee]*, Subcommittee SC ##, *[name of subcommittee]*.

This *second/third/...* edition cancels and replaces the *first/second/...* edition (ISO #####:####), which has been technically revised.

The main changes are as follows:

— xxx xxxxxxxx xxx xxxx

A list of all parts in the ISO ##### series can be found on the ISO website.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at [www.iso.org/members.html](http://www.iso.org/members.html).



## Introduction

Generative Artificial Intelligence (AI) is a type of AI based on techniques and generative models that aim to generate new content that are similar from real data (defined in ISO/IEC 22989:2022/AWI Amd 1), and its performance in knowledge learning, inductive summarization, content creation, perception and cognition is distinctly different from previous AI technologies. It has greater generalization and interactivity, and therefore is extensively integrated into various scenarios.

There are several new features of generative AI, including but not limited to the following.

- New contents are generated by modelling the patterns of vast quantities of training data, rather than recognizing or classifying existing contents.
- Long context windows and self-attention mechanism enable different attention weights given to the relationships between various parts of the user input, so that users can get better interactive experience.
- Contents are easily generated via natural language conversation.
- The foundation model can be fine-tuned at low cost and then applied in wide areas and at large scale.
- Contents generated are highly convincing and more aligned with human habits, as generative AI is more generalized.
- Greater randomness is introduced in generated contents, because generative AI is based on next token prediction.

Therefore, the industry has expressed concerns about the potential risks of generative AI. Generative AI brings new risks, and exacerbate existing AI risks, which include but not limited to the following.

- Easy access to knowledge might make it easier for malicious users to cause harms to society without specialized training (e.g., CBRN knowledge, malware); meanwhile, it also positively increases the productivity of research work and innovation.
- Generative AI's strong generalization ability might result in hallucination; at the same time, the ability allows to process diverse user input.
- The generated contents, when contain faults or misalign with regulations and ethics, might mislead the downstream applications make incorrect decisions or even harmful actions; while the generated contents can also empower various applications, such as AI agents.
- Generative AI might generate unethical contents that pose harms to individuals and society.
- Over-reliance on generative AI might cause humans to be manipulated, especially when humans have no detailed knowledge of how generative AI works.
- The generated contents might cause sensitive information leakage; while users can benefit from the customized personal assistant by feeding personal information to generative AI.
- The generated contents might cause copyrights infringement.
- Continuous learning based on the user feedback can be leveraged to mislead generative AI behaviors; meanwhile, it can enable better alignment with human preference.
- Prompt-based attacks expand the attack surface.

Since generative AI systems can involve multiple stakeholders, and the management of generative AI risks relies on the participation of various stakeholders. However, there is no standard defining the stakeholders' responsibilities. Moreover, having stakeholders be responsible for addressing risks in AI system life cycle stage where they do not have risk control capabilities can be highly inefficient and resource-intensive.

This document aims to achieve the following objectives.

- Develop new and refined objectives to manage risks of generative AI systems.
- Identify the risk sources related to generative AI systems.
- Identify the stakeholders responsible for addressing risks in generative AI systems throughout its life cycle.
- Conduct a fine granular risk analysis against objectives of generative AI systems, including specific consequence, and specific factors to consider in likelihood assessment where applicable.
- Specify risk treatment and controls for generative AI systems.

By using this guidance, the stakeholders involved in generative AI systems can develop risk management plans suitable for their roles, including specifying objectives when identifying risks, identifying the AI system life cycle stages involved, as well as identifying risk sources, analyzing the consequences and likelihood of risks, and providing appropriate risk treatment measures and controls needed for different risks.

# Information technology – Artificial intelligence – Guidance on addressing risks in generative AI systems

## 1 Scope

This document provides guidance on addressing risks in generative artificial intelligence (AI) systems. It includes the following:

- The objectives of generative AI systems when identifying risks.
- The risk sources and the stakeholders facing the risks in generative AI systems throughout its life cycle.
- Guidance for risk analysis, risk treatment and controls of addressing risks in generative AI systems.

This document is applicable to all types and sizes of organizations that develop or use generative AI systems.

## 2 Normative references

The following documents are referred to in the text in such a way that some or all of their content constitutes requirements of this document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 22989:2022/AWI Amd 1, *Information technology — Artificial intelligence — Artificial intelligence concepts and terminology*

ISO/IEC 5338:2023, *Information technology — Artificial intelligence — AI system life cycle processes*

ISO/IEC 23894:2023, *Information technology — Artificial intelligence — Guidance on risk management*

ISO/IEC 42001:2023, *Information technology — Artificial intelligence — Management system*

Editor's note: ISO/IEC 22989 amendments to include generative AI terminology and concepts.

## 3 Terms and definitions

For the purposes of this document, the terms and definitions given in ISO/IEC 22989:2022, ISO/IEC 23053:2022 and the following apply.

ISO and IEC maintain terminology databases for use in standardization at the following addresses:

- ISO Online browsing platform: available at <https://www.iso.org/obp>
- IEC Electropedia: available at <https://www.electropedia.org/>

## 4 Abbreviation terms

|    |                         |
|----|-------------------------|
| AI | artificial intelligence |
| ML | machine learning        |

## 5 Objective of generative AI systems when identifying risks

### 5.1 General

When identifying risks of generative AI systems, various generative AI-related objectives should be taken into account, depending on the nature of the system under consideration and its application context. Objectives of generative AI systems to consider include but are not limited to the objectives described in [Clauses 5.2 to 5.9](#).

## 5.2 Human autonomy

It is the same objective as in conventional AI, but refined in the context of generative AI systems.

- Ensure generative AI systems can operate as intended. For example, accidental misalignment or mis-specification of system goals can cause a model not to operate as intended
- Avoid uncontrollable model autonomy. E.g., 1) Deceptive reward hacking, in which models might develop the ability to act differently under human supervision and in unsupervised setups to get higher rewards. 2) Auto-induced distributional shift, in which models can cause a change in the distribution of their own inputs, and use this ability for undesirable purposes.

## 5.3 No catastrophic threat to human, society and environment by generated knowledge

It is a new objective in the context of generative AI systems.

- Generated knowledge shall fully consider its impact on ethics and morality, its widespread homogenizing consequences, the risk of “value lock-in”, and the risk of obscene, degrading, and/or abusive content.
- Generated knowledge shall fully consider its influence on societal safety and stability, for example, lowering the barrier to access CBRN info; and augmenting security attacks such as hacking, malware, and phishing (generative AI systems are already able to discover vulnerabilities in systems (hardware, software, data) and write code to exploit them).
- Generated knowledge shall fully consider its influence on humans. For example, generative AI systems enable the production of false or misleading information at scale, by which the malicious user can use to deceive or cause harm to others. Also, emotional entanglement between humans and GAI systems, such as coercion or manipulation that leads to safety or psychological risks.

## 5.4 Avoid societal-scale specification gaming like incitement, deception and instigation

It is the same objective as in conventional AI, but bigger impact in the context of generative AI systems.

- Fully consider the risks that generative AI systems may change people's thoughts and behaviors, or so-called social hacking, to exploit human weaknesses to gain their trust.
- Fully consider the impact on individual physical and mental health, such as inducing user addiction, or encouraging self-harm and suicide.
- Fully consider the risks that generative AI systems may generate dangerous or violent recommendations.

## 5.5 Accuracy of generated contents matching expectation

It is a new objective in the context of generative AI systems.

- Generated contents need to be truthful and correct.
- Generated contents need to align with mainstream social values and objective laws, and avoid hallucinations.

## 5.6 Privacy and copyright

It is the same objective as in conventional AI, but refined in the context of generative AI systems.

- Training data do not infringe on privacy.
- Generated content does not contain or infer personally identifiable information (PII)
- Training data do not infringe on IPR or copyright, or they are authorized for use.
- Generated content does not infringe on IPR or copyright.

## 5.7 Transparency, explainability, accountability

It is the same objective as in conventional AI, but bigger impact in the context of generative AI systems.

- Inform humans of confusions caused by generative AI systems, ensuring humans are aware that they are interacting with generative AI systems.
- Clearly define the accountability of stakeholders for addressing risks at each stage of the life cycle.
- Ensure that content likely to cause confusion is marked and traceable.

## 5.8 Fairness

It is the same objective as in conventional AI, but bigger impact in the context of generative AI systems.

- Fully consider the diversity of training data (including the diversity of annotators, as well as the distribution of training and testing data).
- Ensure generated contents align with various preferences in society and avoids producing biased content, causing homogenization, representation harm.
- Ensure fair distribution of capabilities or benefits from generative AI system access, to avoid that capabilities and outcomes of generative AI systems may be worse for some groups compared to others.

## 5.9 Security and resilience

It is the same objective as in conventional AI, but refined in the context of generative AI systems.

- Ensure generative AI systems can resist attacks such as prompt-injection indirect prompt-injection, data poisoning, among others.

# 6 Identification of risk sources and stakeholders responsible for risk addressing

## 6.1 Risk sources throughout generative AI systems life cycle

When identifying risks of generative AI systems, various generative AI-related risk sources should be taken into account.

Based on the AI system life cycle provided in ISO/IEC 5338:2023, generative AI-related risk sources to consider include but are not limited to the ones listed in Table 1.

**Table 1 — Risk sources throughout generative AI systems life cycle**

| AI system life cycle   | Risk Sources                                       | Description   |
|------------------------|--|---|
| Inception              | Objectives not aligned with regulations and ethics | It is the same risk source as in conventional AI, but has exacerbated impact in the context of generative AI systems. <ul style="list-style-type: none"> <li>— Due to generative AI system's ability to create highly convincing and realistic contents, this ability can be exploited for malicious purpose and cause potential large-scale consequences, if objectives are not aligned with regulations and ethics.</li> </ul>  |
|                        | Unreasonable accountability of stakeholders        | It is the same risk source as in conventional AI, but refined in the context of generative AI systems. <ul style="list-style-type: none"> <li>— Generative AI introduces more refined accountability allocation among stakeholders. If being set unreasonably, the risks would not be addressed effectively.</li> <li>— For example, managing malicious users is challenging on the generative model layer; instead, it requires management at the provider layer, like platform provider, service provider or product provider. So, it is unreasonable to make any single stakeholder accountable for the actions of malicious users.</li> </ul> |
| Design and development | Training data management                           | It is the same risk source as in conventional AI, but refined in the context of generative AI systems. <ul style="list-style-type: none"> <li>— Training data might contain IPR or copyright data which are not authorized for use, which is not typical in conventional AI.</li> </ul> <p>At the same time, it is the same risk source as in conventional AI, but has</p>  |

|                             |  |  |  |
|-----------------------------|--|--|--|
|                             |  | <p>exacerbated impact in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— If training data contains privacy info, besides personal info leaks, generative AI systems can fabricate personal images or likenesses for malicious use.</li> </ul>   |  |
|                             | Machine learning algorithm   | <p>It is the same risk source as in conventional AI, but refined in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— As generative AI systems get emergent capabilities, self-iteration and continuous learning, it can cause uncontrollable model autonomy</li> </ul>   |  |
|                             | Data annotation  | <p>It is the same risk source as in conventional AI, but refined in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— In the context of generative AI systems, data annotation involves annotating negative sample data and fine-tuning training to teach the model to recognize inputs containing harmful content, privacy or copyright violations.</li> <li>— If not done well, it can introduce risks when the generative AI system deals with user input and runtime input.</li> </ul> <p>At the same time, it is the same risk source as in conventional AI, but has exacerbated impact in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— If insufficient or low-quality data annotation, because of generative AI's strong generalization ability, it will result in more serious issue, like hallucination.</li> <li>— In contrast, conventional AI, such as discriminative AI, can simply fail to make a decision.</li> </ul> |  |
|                             | Reinforcement learning with human feedback   | <p>New risk source introduced in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— In this case, the model can manipulate the reward mechanisms, acting differently under human supervision and in unsupervised setups</li> </ul>   |  |
| Verification and validation | Testing data<br>Testing implementation   | It is the same risk source as in conventional AI.  |  |
| Deployment                  | Insecure deployment environment<br>Lack of protection of model<br>Vulnerability in model | It is the same risk source as in conventional AI.  |  |
| Operation and monitoring    | Operating data input   | User input data  | <p>It is the same risk source as in conventional AI, but refined in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— Typical example is the prompt-injection attack.</li> </ul>  |
|                             |  | Runtime input data   | <p>It is the same risk source as in conventional AI, but refined in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— Typical example is the indirect prompt-injection attack. The third-party plug-in or knowledge base can contain IPR or copyright data which are not authorized for use. This is not typical in conventional AI.</li> </ul> |
|                             | Model execution  | Generated contents   | <p>New risk source introduced in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— Generative AI's key and unique feature is it can create new or original contents.</li> <li>— The generated contents can violate copyright.</li> <li>— Also, generative AI can cause hallucination, or generate incorrect or untruthful contents.</li> </ul>  |
|                             |  | Malicious use  | <p>It is the same risk source as in conventional AI, but refined in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— Generative AI system makes it possible to generate malware with</li> </ul>  |

|                       |  |   |   |
|-----------------------|--|---|---|
|                       |  |   | <p>low cost, or assist hackers to generate malware. The traditional AI can just improve codes.</p> <p>At the same time, it is the same risk source as in conventional AI, but has exacerbated impact in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— Generative AI system can be used for incitement, deception and instigation purpose, causing societal-scale impact.</li> </ul>                                      |
|                       |  | Lack of transparency and explainability                     | <p>It is the same risk source as in conventional AI, but has exacerbated impact in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— Due to generative AI system's ability to create convincing and realistic contents, it can result in confusion that humans are not aware that they are interacting with generative AI systems and follow the generated decisions, causing potential large-scale consequences.</li> </ul> |
| Continuous validation |  | Data poisoning  | <p>It is the same risk source as in conventional AI, but has exacerbated impact in the context of generative AI systems.</p> <ul style="list-style-type: none"> <li>— It can cause the generated contents misalign with humans' values and objectives.</li> <li>— And the impact will be further enlarged due to user's heavy dependence on the generative AI systems because it can create highly convincing and realistic contents.</li> </ul>              |
| Re-evaluate           |  | Insufficient risk evaluation<br>Inaccurate risk treatment   | It is the same risk source as in conventional AI.   |
| Retirement            |  | Unthorough disposal of data<br>Unthorough disposal of model | It is the same risk source as in conventional AI.   |

## 6.2 Stakeholders responsible for risk addressing

Based on the AI system life cycle provided in ISO/IEC 5338:2023, stakeholders facing the risks in generative AI systems are listed in Table 2.

**Table 2 — Stakeholders facing the risks in generative AI systems throughout its life cycle**

| AI system life cycle        | Risk Sources  | Stakeholders   |             |
|-----------------------------|---|--|-------------|
| Inception                   | <ul style="list-style-type: none"> <li>— Objectives not aligned with regulations and ethics</li> <li>— Unreasonable accountability of stakeholders</li> </ul>                                       | All  |             |
| Design and development      | <ul style="list-style-type: none"> <li>— Training data management</li> <li>— Machine learning algorithm</li> <li>— Data annotation</li> <li>— Reinforcement learning with human feedback</li> </ul> | Data provider<br>Foundation model developer<br>Finetuned model developer |             |
| Verification and validation | <ul style="list-style-type: none"> <li>— Testing data</li> <li>— Testing implementation</li> </ul>  | AI developer<br>AI evaluator<br>AI auditor                               |             |
| Deployment                  | <ul style="list-style-type: none"> <li>— Insecure deployment environment</li> <li>— Lack of protection of model</li> <li>— Vulnerability in model</li> </ul>  | AI provider<br>AI developer<br>AI system integrator                      |             |
| Operation                   | Operating   | — User input data  | AI provider |

|                       |                    |  |   |
|-----------------------|--------------------|--|---|
| and<br>monitoring     | data input         | — Runtime input data   | Data provider<br>AI customer                        |
|                       | Model<br>execution | — Generated contents<br>— Malicious use<br>— Lack of transparency and explainability | AI provider<br>Regulators                           |
| Continuous validation |                    | — Data poisoning   | Data provider<br>AI provider                        |
| Re-evaluate           |                    | — Insufficient risk evaluation<br>— Inaccurate risk treatment                        | AI developer<br>AI evaluator<br>AI auditor          |
| Retirement            |                    | — Unthorough disposal of data<br>— Unthorough disposal of model                      | AI provider<br>AI developer<br>AI system integrator |

## 7 Risk analysis against objectives in generative AI systems

### 7.1 General

### 7.2 Human autonomy

#### 7.2.1 Consequence assessment

#### 7.2.2 Likelihood assessment

### 7.3 No catastrophic threat to human, society and environment by generated knowledge

#### 7.3.1 Consequence assessment

#### 7.3.2 Likelihood assessment

### 7.4 Avoid societal-scale specification gaming like incitement, deception and instigation

#### 7.4.1 Consequence assessment

#### 7.4.2 Likelihood assessment

### 7.5 Accuracy of generated contents matching expectation

#### 7.5.1 Consequence assessment

#### 7.5.2 Likelihood assessment



## **7.6 Privacy and copyright**

**7.6.1 Consequence assessment**

**7.6.2 Likelihood assessment**

## **7.7 Transparency, explainability, accountability**

**7.7.1 Consequence assessment**

**7.7.2 Likelihood assessment**

## **7.8 Fairness**

**7.8.1 Consequence assessment**

**7.8.2 Likelihood assessment**

## **7.9 Security and resilience**

**7.9.1 Consequence assessment**

**7.9.2 Likelihood assessment**

## **8 Risk treatment**

## **9 Risk addressing controls**

## Bibliography

- [1] ISO/IEC DIS 42005 *Information technology — Artificial intelligence — AI system impact assessment*
- [2] ISO/IEC CD 27090, *Cybersecurity — Artificial Intelligence — Guidance for addressing security threats to artificial intelligence systems*
- [3] ISO/IEC WD 27091 *Cybersecurity and Privacy — Artificial Intelligence — Privacy protection*
- [4] ISO/IEC TR 24368:2022 *Information technology — Artificial intelligence — Overview of ethical and societal concerns*
- [5] ISO/IEC AWI TS 22443 *Information technology — Artificial intelligence — Guidance on addressing societal concerns and ethical considerations*
- [6] ISO/IEC TR 24027:2021 *Information technology — Artificial intelligence (AI) — Bias in AI systems and AI aided decision making*
- [7] ISO/IEC DTS 12791 *Information technology — Artificial intelligence — Treatment of unwanted bias in classification and regression machine learning tasks*
- [8] ISO/IEC TR 5469:2024 *Artificial intelligence — Functional safety and AI systems*
- [9] ISO/IEC TR 24028:2020 *Information technology — Artificial intelligence — Overview of trustworthiness in artificial intelligence*
- [10] ISO/IEC 38507:2022, *Information technology — Governance of IT — Governance implications of the use of artificial intelligence by organizations*