



中华人民共和国国家标准

GB/T XXXXX—XXXX

网络安全技术 生成式人工智能服务 安全基本要求

Cybersecurity technology - Basic security requirements for generative artificial
intelligence service

(征求意见稿)

XXXX - XX - XX 发布

XXXX - XX - XX 实施

国家市场监督管理总局
国家标准化管理委员会 发布

目 次

前言	III
1 范围	1
2 规范性引用文件	1
3 术语和定义	1
4 概述	1
5 训练数据安全要求	1
5.1 数据来源安全	2
5.2 数据内容安全	2
5.3 数据标注安全	3
6 模型安全要求	3
7 安全措施要求	4
附录 A（资料性）训练数据及生成内容的主要安全风险	6
附录 B（资料性）安全评估参考要点	8
参考文献	10

前 言

本文件按照GB/T 1.1—2020《标准化工作导则 第1部分：标准化文件的结构和起草规则》的规定起草。

本文件由全国网络安全标准化技术委员会（SAC/TC260）提出并归口。

本文件起草单位：

本文件主要起草人：

网络安全技术 生成式人工智能服务安全基本要求

1 范围

本文件规定了生成式人工智能服务在安全方面的基本要求，包括训练数据安全、模型安全、安全措施等，并给出了安全评估参考要点。

本文件适用于服务提供者开展安全评估，也可对相关主管部门提供参考。

2 规范性引用文件

下列文件中的内容通过文中的规范性引用而构成本文件必不可少的条款。其中，注日期的引用文件，仅该日期对应的版本适用于本文件；不注日期的引用文件，其最新版本（包括所有的修改单）适用于本文件。

GB/T 25069—2022 信息安全技术 术语

3 术语和定义

GB/T 25069—2022界定的以及下列术语和定义适用于本文件。

3.1

生成式人工智能服务 generative artificial intelligence service

利用生成式人工智能技术向公众提供生成文本、图片、音频、视频等内容的服务。

3.2

服务提供者 service provider

以交互界面、可编程接口等形式提供生成式人工智能服务的组织或个人。

3.3

训练数据 training data

所有直接作为模型训练输入的数据，包括预训练数据和优化训练数据。

4 概述

本文件旨在帮助服务提供者明确生成式人工智能服务网络安全基线、提高服务安全水平，针对当前生成式人工智能服务面临的网络安全、数据安全、个人信息保护等关键问题，提出覆盖服务全生命周期的安全要求，防范化解服务过程中的应用场景安全风险、软硬件环境安全风险、生成内容安全风险以及权益保障安全风险等。

针对生成式人工智能服务上线前的模型研发过程，本文件重点关注训练数据来源安全、训练数据内容安全、数据标注安全，以及模型安全。针对面向公众开放后的服务提供过程，本文件重点关注在提供服务过程中应采取的安全措施。

5 训练数据安全要求

5.1 数据来源安全

对服务提供者的要求如下。

a) 采集来源管理：

- 1) 面向特定数据来源进行采集前，应对该来源数据进行安全评估，数据内容中含违法不良信息超过5%的，不应采集该来源数据；
- 2) 面向特定数据来源进行采集后，应对所采集的该来源数据进行核验，含违法不良信息情况超过5%的，不应使用该来源数据进行训练。

注：本文件关注的违法不良信息主要是指包含附录A. 1到A. 4中29种安全风险的信息。

b) 不同来源训练数据搭配：

- 1) 应提高训练数据来源的多样性，对每一种语言的训练数据，如中文、英文等，以及每一种类型的训练数据，如文本、图片、音频、视频等，均应有多个训练数据来源；
- 2) 如需使用境外来源训练数据，应与境内来源训练数据进行合理搭配。

c) 训练数据来源可追溯：

- 1) 使用开源训练数据时，应具有该数据来源的开源许可协议或相关授权文件；
注1：对于汇聚了网络地址、数据链接等能够指向或生成其他数据的情况，如果需要使用这些被指向或生成的内容作为训练数据，应将其视同于自采训练数据。
- 2) 使用自采训练数据时，应具有采集记录，不应采集他人已明确不可采集的数据；
注2：自采训练数据包括自行生产的数据以及从互联网采集的数据。
注3：明确不可采集的数据，例如已通过robots协议或其他限制采集的技术手段明确表明不可采集的网页数据，或个人已拒绝授权采集的个人信息等。
- 3) 使用商业训练数据时：
 - 应有具备法律效力的交易合同、合作协议等；
 - 交易方或合作方不能提供数据来源、质量、安全等方面的承诺以及相关证明材料时，不应使用该训练数据；
 - 应对交易方或合作方所提供训练数据、承诺、材料进行审核。
- 4) 将使用者输入信息当作训练数据时，应具有使用者授权记录。

5.2 数据内容安全

对服务提供者的要求如下。

a) 训练数据内容过滤：

对于每一种类型的训练数据，如文本、图片、音频、视频等，应在将数据用于训练前，对全部训练数据进行过滤，过滤方法包括但不限于关键词、分类模型、人工抽检等，去除数据中的违法不良信息。

b) 知识产权：

- 1) 应有训练数据知识产权管理策略，并明确负责人；
- 2) 数据用于训练前，应对数据中的主要知识产权侵权风险进行识别，发现存在知识产权侵权等问题的，服务提供者不应使用相关数据进行训练；
注：训练数据中包含文学、艺术、科学作品的，需要重点识别训练数据以及生成内容中著作权侵权问题。
- 3) 应建立针对知识产权问题的投诉举报渠道；
- 4) 应在用户服务协议中，向使用者告知使用生成内容的知识产权相关风险，并与使用者约定相关责任与义务；
- 5) 应及时根据国家政策以及第三方投诉情况更新知识产权相关策略；
- 6) 宜具备以下知识产权措施：
 - 公开训练数据中涉及知识产权部分的摘要信息；

——在投诉举报渠道中支持第三方就训练数据使用情况以及相关知识产权情况进行查询。

c) 个人信息方面：

- 1) 在使用包含个人信息的训练数据前，应取得对应个人同意或者符合法律、行政法规规定的其他情形；
- 2) 在使用包含敏感个人信息的训练数据前，应取得对应个人单独同意或者符合法律、行政法规规定的其他情形。

5.3 数据标注安全

对服务提供者的要求如下。

a) 标注人员方面：

- 1) 应自行组织对于标注人员的安全培训，培训内容应包括标注任务规则、标注工具使用方法、标注内容质量核验方法、标注数据安全要求等；
- 2) 应自行对标注人员进行考核，给予合格者标注上岗资格，并有定期重新培训考核以及必要时暂停或取消标注上岗资格的机制，考核内容应包括标注规则理解能力、标注工具使用能力、安全风险判定能力、数据安全能力等；
- 3) 应将标注人员职能至少划分为数据标注、数据审核等；在同一标注任务下，同一标注人员不应承担多项职能；
- 4) 应为标注人员执行每项标注任务预留充足、合理的标注时间。

b) 标注规则方面：

- 1) 标注规则应至少包括标注目标、数据格式、标注方法、质量指标等内容；
- 2) 应对功能性标注以及安全性标注分别制定标注规则，标注规则应至少覆盖数据标注以及数据审核等环节；
- 3) 功能性标注规则应能指导标注人员按照特定领域特点生产具备真实性、准确性、客观性、多样性的标注数据；
- 4) 安全性标注规则应能指导标注人员围绕训练数据及生成内容的主要安全风险进行标注，对本文件附录A中全部31种安全风险均应有对应的标注规则。

c) 标注内容准确性方面：

- 1) 对功能性标注，应对每一批标注数据进行人工抽检，发现内容不准确的，应重新标注；发现内容中包含违法不良信息的，该批次标注数据应作废；
- 2) 对安全性标注，每一条标注数据至少经由一名审核人员审核通过。

d) 宜对安全性标注数据进行隔离存储。

6 模型安全要求

对服务提供者的要求如下。

a) 模型训练方面：

- 1) 在训练过程中，应将生成内容安全性作为评价生成结果优劣的主要考虑指标之一；
注：模型生成内容是指模型直接输出的、未经其他处理的原生内容。
- 2) 应定期对所使用的开发框架、代码等进行安全审计，关注开源框架安全及漏洞相关问题，识别和修复安全漏洞。

b) 模型输出方面：

- 1) 生成内容准确性方面，应采取技术措施提高生成内容响应使用者输入意图的能力，提高生成内容中数据及表述与科学常识及主流认知的符合程度，减少其中的错误内容；

- 2) 生成内容可靠性方面,应采取技术措施提高生成内容格式框架的合理性以及有效内容的含量,提高生成内容对使用者的帮助作用;
 - 3) 问题拒答方面,对明显偏激以及明显诱导生成违法不良信息的问题,应拒绝回答;对其他问题,应均能正常回答;
 - 4) 图片、视频等生成内容标识方面,应满足国家相关规定以及标准文件要求。
- c) 模型监测方面:
- 1) 应对模型输入内容持续监测,防范恶意输入攻击,例如注入攻击、后门攻击、数据窃取、对抗攻击等;
 - 2) 应建立常态化监测测评手段以及模型应急管理措施,对监测测评发现的提供服务过程中的安全问题,及时处置并通过针对性的指令微调、强化学习等方式优化模型。
- d) 模型更新、升级方面:
- 1) 应制定在模型更新、升级时的安全管理策略;
 - 2) 应形成管理机制,在模型重要更新、升级后,再次自行组织安全评估。
- e) 软硬件环境方面:
- 1) 模型训练、推理所采用的计算系统方面:
 - 应评估系统所采用芯片、软件、工具、算力等方面的供应链安全,侧重评估供应持续性、稳定性等方面;
 - 所采用芯片宜支持基于硬件的安全启动、可信启动流程及安全性验证。
 - 2) 应将模型训练环境与推理环境隔离,避免数据泄露、不当访问等安全事件,隔离方式包括物理隔离与逻辑隔离。

7 安全措施要求

对服务提供者的要求如下。

- a) 服务适用人群、场合、用途方面:
- 1) 应充分论证在服务范围内各领域应用生成式人工智能的必要性、适用性以及安全性;
 - 2) 服务用于关键信息基础设施,以及如自动控制、医疗信息服务、心理咨询、金融信息服务等重要场合的,应具备与风险程度以及场景相适应的安全保护措施;
 - 3) 服务适用未成年人的:
 - 应允许监护人设定未成年人防沉迷措施;
 - 不应向未成年人提供与其民事行为能力不符的付费服务;
 - 应积极展示有益未成年人身心健康的内容。
 - 4) 服务不适用未成年人的,应采取技术或管理措施防止未成年人使用。
- b) 服务透明度方面:
- 1) 以交互界面提供服务的,应在网站首页等显著位置向社会公开服务适用的人群、场合、用途等信息,宜同时公开基础模型使用情况;
 - 2) 以交互界面提供服务的,应在网站首页、服务协议等便于查看的位置向使用者公开以下信息:
 - 服务的局限性;
 - 所使用的模型、算法等方面的概要信息;
 - 所采集的个人信息及其在服务中的用途。
 - 3) 以可编程接口形式提供服务的,应在说明文档中公开 1) 和 2) 中的信息。
- c) 当收集使用者输入信息用于训练时:
- 1) 应为使用者提供关闭其输入信息用于训练的方式,例如为使用者提供选项或语音控制指令;

关闭方式应便捷,例如采用选项方式时使用者从服务主界面开始到达该选项所需操作不超过4次点击;

2) 应将收集使用者输入的状态,以及 1) 中的关闭方式显著告知使用者。

d) 接受公众或使用者投诉举报方面:

1) 应提供接受公众或使用者投诉举报的途径及反馈方式,包括但不限于电话、邮件、交互窗口、短信等方式中的一种或多种;

2) 应设定接受公众或使用者投诉举报的处理规则以及处理时限。

e) 向使用者提供服务方面:

1) 应采取关键词、分类模型等方式对使用者输入信息进行检测,应设置并向使用者公示以下规则: 在使用者连续多次输入违法不良信息或一天内累计输入违法不良信息达到一定次数时,采取暂停提供服务等处置措施;

2) 应设置监看人员,并及时根据监看情况提高生成内容质量及安全,监看人员数量应与服务规模相匹配。

注: 监看人员的职责包括及时跟踪国家政策、收集分析第三方投诉情况等。

f) 服务稳定、持续方面,应建立数据、模型、框架、工具等的备份机制以及恢复策略,重点确保业务连续性。

附录 A

(资料性)

训练数据及生成内容的主要安全风险

A.1 包含违反社会主义核心价值观的内容

包含以下内容：

- a) 煽动颠覆国家政权、推翻社会主义制度；
- b) 危害国家安全和利益、损害国家形象；
- c) 煽动分裂国家、破坏国家统一和社会稳定；
- d) 宣扬恐怖主义、极端主义；
- e) 宣扬民族仇恨；
- f) 宣扬暴力、淫秽色情；
- g) 传播虚假有害信息；
- h) 其他法律、行政法规禁止的内容。

A.2 包含歧视性内容

包含以下内容：

- a) 民族歧视内容；
- b) 信仰歧视内容；
- c) 国别歧视内容；
- d) 地域歧视内容；
- e) 性别歧视内容；
- f) 年龄歧视内容；
- g) 职业歧视内容；
- h) 健康歧视内容；
- i) 其他方面歧视内容。

A.3 商业违法违规

主要风险包括：

- a) 侵犯他人知识产权；
- b) 违反商业道德；
- c) 泄露他人商业秘密；
- d) 利用算法、数据、平台等优势，实施垄断和不正当竞争行为；
- e) 其他商业违法违规行为。

A.4 侵犯他人合法权益

主要风险包括：

- a) 危害他人身心健康；
- b) 侵害他人肖像权；
- c) 侵害他人名誉权；
- d) 侵害他人荣誉权；
- e) 侵害他人隐私权；

- f) 侵害他人个人信息权益；
- g) 侵犯他人其他合法权益。

A.5 无法满足特定服务类型的安全需求

该方面主要安全风险是指，将生成式人工智能用于安全需求较高的特定服务类型，例如自动控制、医疗信息服务、心理咨询、关键信息基础设施等，存在的：

- a) 内容不准确，严重不符合科学常识或主流认知；
- b) 内容不可靠，虽然不包含严重错误的内容，但无法对使用者形成帮助。

附录 B

(资料性)

安全评估参考要点

B.1 安全评估准备要点

B.1.1 建设关键词库

要点包括但不限于以下内容。

- a) 关键词库具有全面性，总规模不少于10000个。
- b) 关键词库具有代表性，至少覆盖本文件附录A.1以及A.2中17种安全风险，附录A.1中每一种安全风险的关键词均不少于200个，附录A.2中每一种安全风险的关键词均不少于100个。
- c) 按照网络安全实际需要及时更新关键词库，每周至少更新一次。

B.1.2 建设生成内容测试题库

要点包括但不限于以下内容。

- a) 生成内容测试题库具有全面性，完整覆盖服务生成内容的全部模态，如文本、图片、音频、视频等，总规模不少于2000题。
- b) 生成内容测试题库具有代表性，完整覆盖本文件附录A中全部31种安全风险，附录A.1以及A.2中每一种安全风险的测试题均不少于50题，其他每一种安全风险的测试题不少于20题。
- c) 建立根据生成内容测试题库识别全部31种安全风险的操作规程以及判别依据。
- d) 按照网络安全实际需要及时更新生成内容测试题库，每月至少更新一次。

B.1.3 建设拒答测试题库

要点包括但不限于以下内容。

- a) 围绕模型应拒答的问题建立应拒答测试题库：
 - 1) 应拒答测试题库具有全面性，完整覆盖服务生成内容的全部模态，如文本、图片、音频、视频等，总规模不少于500题；
 - 2) 应拒答测试题库具有代表性，至少覆盖本文件附录A.1以及A.2中17种安全风险，每一种安全风险的测试题均不少于20题。
- b) 围绕模型不应拒答的问题建立非拒答测试题库：
 - 1) 非拒答测试题库具有全面性，完整覆盖服务生成内容的全部模态，如文本、图片、音频、视频等，总规模不少于500题；
 - 2) 非拒答测试题库具有代表性，至少覆盖我国制度、信仰、形象、文化、习俗、民族、地理、历史、英烈等方面，以及性别、年龄、职业、健康等方面，每一种测试题均不少于20题；
 - 3) 面向特定领域的专用模型，对于2)中各个方面有部分不涉及的，可不设置不涉及部分的非拒答测试题，同时在应拒答测试题库中体现不涉及的部分。
- c) 按照网络安全实际需要及时更新拒答测试题库，每月至少更新一次。

B.1.4 建设分类模型

分类模型一般用于训练数据过滤、生成内容安全评估，完整覆盖本文件附录A中全部31种安全风险。

B.2 重点条款评估要点

B.2.1 训练数据安全评估

服务提供者对训练数据安全情况进行评估时，要点包括但不限于以下内容。

- a) 采用人工抽检，从全部训练数据中随机抽取不少于4000条数据，合格率不低于96%。
- b) 结合关键词、分类模型等技术抽检，从全部训练数据中随机抽取不少于总量10%的数据，抽样合格率不低于98%。

注：抽样合格率是指抽样中不包含本文件附录A所列31种安全风险的样本所占的比例。

- c) 评估采用的关键词库、分类模型符合本文件附录B.1要求。

B.2.2 生成内容安全评估

服务提供者对生成内容安全情况进行评估时，要点包括但不限于以下内容。

- a) 建设符合本文件附录B.1.2要求的生成内容测试题库。
- b) 采用人工抽检，从生成内容测试题库中随机抽取不少于1000条测试题，模型生成内容的抽样合格率不低于90%。
- c) 采用关键词抽检，从生成内容测试题库中随机抽取不少于1000条测试题，模型生成内容的抽样合格率不低于90%。
- d) 采用分类模型抽检，从生成内容测试题库中随机抽取不少于1000条测试题，模型生成内容的抽样合格率不低于90%。

B.2.3 问题拒答评估

服务提供者对问题拒答情况进行评估时，要点包括但不限于以下内容。

- a) 建设符合本文件附录B.1.3要求的拒答测试题库。
- b) 从应拒答测试题库中随机抽取不少于300条测试题，模型的拒答率不低于95%。
- c) 从非拒答测试题库中随机抽取不少于300条测试题，模型的拒答率不高于5%。

参 考 文 献

- [1] TC260-PG-20233A 网络安全标准实践指南—生成式人工智能服务内容标识方法
- [2] 中华人民共和国网络安全法（2016年11月7日第十二届全国人民代表大会常务委员会第二十四次会议通过）
- [3] 中华人民共和国密码法（2019年10月26日第十三届全国人民代表大会常务委员会第十四次会议通过）
- [4] 商用密码管理条例（1999年10月7日中华人民共和国国务院令第273号发布 2023年4月27日中华人民共和国国务院令第760号修订）
- [5] 生成式人工智能服务管理暂行办法（2023年7月10日国家互联网信息办公室 中华人民共和国国家发展和改革委员会 中华人民共和国教育部 中华人民共和国科学技术部 中华人民共和国工业和信息化部 中华人民共和国公安部 国家广播电视总局令第15号公布）