



大模型安全解决方案 白皮书

Ver 1.0

2023/8/24

版权说明

本白皮书版权属于北京百度网讯科技有限公司（以下简称“百度”），并受法律保护。任何单位或者个人未经百度书面许可，不得擅自摘抄、转载、摘编或者以其他方式使用本白皮书文字或观点，由于本方案相关产品版本升级或其它原因，本文件内容将会不定期进行更新。除非另有约定，本文件仅做指导、参考作用，所有陈述不构成对于本方案相关产品合同相对方的任何担保、承诺，不视为合同的组成部分或者附件。



目录

1.前言.....	2
2.大模型安全的重要性.....	4
2.1 安全风险引发的重要性	4
2.2 安全方案服务的使命与目标.....	5
3.大模型应用面临的安全挑战与潜在威胁	7
3.1 数据安全与隐私问题	7
3.2 模型流转/部署过程中的安全问题	8
3.3 AIGC 的内容合规问题.....	9
3.4 大模型运营的业务安全问题.....	15
4.大模型安全解决方案.....	19
4.1 大模型数据安全与隐私保护方案.....	19
4.2 模型保护方案	33
4.3 AIGC 内容合规.....	36
4.4 大模型业务运营与安全风控.....	41
5.大模型蓝军安全评测解决方案	44
5.1 建立大模型蓝军所面临困难.....	44
5.2 百度安全面向大模型蓝军的解决方案	47



6.总结与展望	54
6.1 总结成果与贡献	54
6.2 展望未来发展	55
6.3 结语	55
参考文献	57



1.前言

在当今迅速发展的数字化时代，人工智能技术正引领着科技创新的浪潮，而其中的大模型技术则被视为人工智能的一大突破。大模型是指参数量巨大、能力强大的人工神经网络模型，以其卓越的表现自然语言处理、计算机视觉、语音识别等领域赢得了持续的关注和青睐。这些模型的出现，不仅在学术界引起了研究者的广泛兴趣，也在商业应用领域带来了一系列创新和变革。大模型技术的崛起，首要得益于深度学习的发展以及硬件计算能力的提升。深度学习模型，尤其是基于 Transformer 架构的模型，如 BERT、GPT 和 T5，通过在海量数据上进行训练，学习到了丰富的语义和特征表示，使得其在多项人工智能任务中展现出远超以往的性能。例如，在自然语言处理领域，这些大模型能够实现更准确、更流畅的语言生成、机器翻译和情感分析等任务，大大提升了人机交互和信息处理的能力。伴随着大模型的不断演进和不断优化，其在商业领域的应用也愈发广泛，金融行业可以利用大模型进行风险评估和市场预测，医疗领域可以通过大模型实现图像识别和疾病诊断，而广告、营销等领域也能够通过大模型实现更精准的用户推荐和个性化服务。同时，大模型还在科学研究、文化创意和娱乐产业中发挥着积极作用，为人类创造了更多可能性。但伴随着大模型技术的迅猛发展，一系列安全风险和伦理挑战也开始浮现。大规模数据的采集和存储，可能导致个人隐私的泄露和滥用。模型的强大能力也可能被恶意利用，用于虚假信息生成、社会工程和网络攻击。对抗样本攻击则可能使得模型产生误导性结果，严重影响决策的准确性。在社会伦理层面，大模型的使用引发了关于人工智能责任、算法歧视等诸多争议。



因此，建立稳固的大模型安全风控体系势在必行。本白皮书旨在全面探讨大模型安全风险，并为各界提供指导，以确保大模型在广泛应用中的安全性和可信度。通过深入剖析大模型领域的安全挑战，我们可以制定切实可行的措施，确保大模型在为人类创造价值的同时，也能够保障个人隐私、社会稳定和信息安全。



2. 大模型安全的重要性

2.1 安全风险引发的重要性

随着大模型技术的高速发展，其在各个领域的应用日益广泛，从科研到商业，再到日常生活、办公等方方面面。但随之而来的是一系列潜在的安全风险，这些风险的引发和应对不仅关乎企业的声誉，还牵涉到个人隐私的保护和社会的稳定。正因如此，深入了解和应对这些安全风险变得至关重要。

首先，大模型在许多应用场景中处理大量敏感数据和个人信息，如用户的搜索记录、社交媒体互动和金融交易等。这使得数据泄露和隐私侵犯的风险不容忽视。一旦这些敏感信息遭受泄露，个人隐私权益可能会受到严重损害，甚至被用于恶意行为，如身份盗窃、诈骗和社会工程攻击。这不仅会对受害者造成经济损失，还可能导致社会的恐慌和不信任。

其次，大模型的强大能力也可能被用于进行各种形式的恶意攻击。模型的对抗性样本攻击，即针对模型的输入进行微小改动，从而欺骗模型产生错误预测，已成为一种常见的威胁。恶意使用者可以通过这种方式制造虚假信息，影响决策结果，如将误导性的信息传播到社交媒体平台，从而扰乱社会秩序。此外，大模型的生成能力也可能被用于生成虚假的内容，威胁到媒体的可信度和新闻的真实性。

另外，模型本身也可能成为攻击者的目标。模型参数和权重的泄露可能导致知识产权的损失，甚至使恶意使用者能够复制或修改模型，进一步恶化风险。对模型的针对性攻击，如投毒攻击，可能使模型的输出产生不良影响，从



而影响到正常的业务运行。这些威胁可能在不经意间对企业和社会造成巨大的损失。

此外，大模型的使用往往涉及到社会伦理和法律问题。例如，算法的歧视性问题，即模型在处理数据时产生的不公平或偏见，可能引发社会的不满和争议。此外，大模型可能会被用于传播虚假信息、仇恨言论或不当内容，从而引发社会不安定和文化冲突。

最后，国家网信办联合国家发展改革委、教育部、科技部、工业和信息化部、公安部、广电总局公布《生成式人工智能服务管理暂行办法》，自 2023 年 8 月 15 日起施行，旨在促进生成式人工智能健康发展和规范应用，维护国家安全和公共利益，保护公民、法人和其他组织的合法权益。这既是促进生成式人工智能健康发展的重要要求，也是防范生成式人工智能服务风险的现实需要。

因此，确保大模型的安全性和可信度是一个紧迫的任务。需要综合运用技术手段、政策法规以及社会共识，建立起一套全面的大模型安全风险管理体系。通过逐一应对数据隐私保护、模型防御、内容合规、恶意行为检测等方面的挑战，我们能够更好地应对现实中的安全风险，保障个人权益和社会稳定。这也是本白皮书所要探讨的核心议题之一。

2.2 安全方案服务的使命与目标

本白皮书的使命在于为大模型领域的各方利益相关者提供指导，以确保大模型技术的安全应用。我们致力于建立一个安全、稳定且可信赖的大模型生态



系统，旨在维护用户的数据隐私、保护企业的商业机密，并提供有效的对抗措施来应对潜在的安全威胁。我们的目标包括但不限于：

- 提供一套综合性的安全解决方案，以减轻大模型应用过程中的安全压力。
- 建立规范和标准，指导大模型的安全设计、开发、部署和监测。
- 促进安全意识的提高，使所有相关方能够更好地理解和应对安全挑战。
- 推动研究和创新，以增强大模型的鲁棒性和防御能力，应对新型攻击。

本白皮书将按照不同的维度深入探讨大模型安全的关键问题，以提供全面的指导和建议。



3. 大模型应用面临的安全挑战与潜在威胁

ChatGPT 引爆的生成式人工智能热潮，让 AI 模型在过去几个月成为行业瞩目的焦点，并且在国内引发“百模大战”，在大模型高速发展的同时，大模型应用所面临的安全挑战、与潜在的威胁也不能够忽视，本文将依托百度安全大模型安全实践与总结，分别从数据安全与隐私问题、模型流转/部署过程中的安全问题、AIGC 的内容合规问题、以及大模型运营过程中的业务安全问题在内共计四个方向，详细介绍一下相关的安全挑战。

3.1 数据安全与隐私问题

1、传输截获风险：在进行大模型非私有化的预训练、精调、推理服务时，数据需要在不同的主体或部门之间进行传输。这些数据通常包括各种敏感信息和隐私，如个人身份信息、金融数据等。在数据传输过程中，如果没有采取足够的安全措施，攻击者可能会截获这些数据，从而获取敏感信息，给用户和组织带来安全和隐私问题。因此，在使用大模型服务时，必须采取适当的安全措施来保护数据的机密性和完整性，以防止传输截获风险。

2、运营方窥探风险：在精调与推理阶段，通常需要使用个人身份信息、企业数据等敏感数据来提高模型的准确性和性能。然而，如果这些数据被大模型运营机构窥视或收集，就可能存在被滥用的风险。运营方可能会利用这些数据来了解用户的隐私信息，例如个人偏好、行为习惯、社交网络等，从而进行针对性的广告投放或者推销策略。此外，运营方还可能将数据泄露给第三方，



这些第三方可能是合作伙伴、数据分析公司、广告公司等，从而获取不正当的利益。

3、模型记忆风险：经过模型的训练和推理后，模型会形成记忆。这些记忆包括各种历史数据和相关信息，如果这些模型被泄露或共享使用，则可能存在模型记忆甚至记忆内容泄密的风险。攻击者可能会利用这些记忆信息来实施恶意行为，例如针对性攻击、诈骗等。此外，如果记忆内容被泄露，也会对用户的隐私和安全造成威胁。因此，在使用大模型服务时，必须采取适当的安全措施来保护模型的机密性和隐私性，例如加密和访问控制等。同时，应该定期对模型进行评估和更新，以减少模型记忆风险。

3.2 模型流转/部署过程中的安全问题

大模型本身也是一种重要的资产，它包含了大量的知识和技能，如果没有合理的管理和控制，就可能被盗取、复制或篡改，导致模型的性能下降或功能失效。此外，大模型也可能受到对抗攻击的威胁，如对抗样本、对抗训练等，这些攻击可以使模型产生错误的输出；本白皮书围绕数据、模型、网络通信等多个方面所面临的安全问题做一下介绍：

1、模型知识泄漏：在将模型部署到生产环境中，模型的输出可能会暴露训练数据的一些信息。攻击者可以通过分析模型的输出，推断出训练数据的特征和分布，进而构建类似的数据集，甚至还原部分原始数据。

2、模型逆向工程：攻击者可能尝试通过逆向工程技术还原部署模型的架构、权重和训练数据。这可能导致知识产权盗窃、模型盗用和安全漏洞的暴



露。逆向工程可能通过模型推理结果、输入输出分析以及梯度攻击等方式进行。

3、输入数据的合法性和安全性： 在模型部署阶段，恶意用户可能试图通过提供恶意输入来攻击系统。例如，输入中可能包含恶意代码、命令执行、注入语句或文件包含路径，从而导致安全漏洞。

4、模型更新和演化： 模型需要定期更新以保持性能和适应新的数据分布。然而，模型更新可能引入新的漏洞和问题。安全地更新模型需要考虑版本控制、验证新模型的安全性和稳定性，以及备份机制以防产生不良影响。

3.3 AIGC 的内容合规问题

自 2023 年 4 月 11 日，国家互联网信息办公室为促进生成式人工智能技术健康发展和规范应用，根据《中华人民共和国网络安全法》等法律法规，国家互联网信息办公室起草了《生成式人工智能服务管理办法（征求意见稿）》，再到国家网信办联合国家发展改革委、教育部、科技部、工业和信息化部、公安部、广电总局共同公布的《生成式人工智能服务管理暂行办法》的正式施行，在国家层面不断指导和促进生成式人工智能健康发展和规范应用，也是防范生成式人工智能服务风险的现实需要，现以在百度安全在生成式人工智能服务的安全实践和业务理解，总结了如下几个方面的安全挑战：

1. 个人隐私问题：



隐私问题涉及到生成式人工智能技术在使用用户个人数据时可能引发的隐私泄露和滥用问题。生成技术通常需要大量的数据来提供更准确的内容生成，这可能包括用户的文本、图像、音频等信息。然而，当个人数据被用于生成内容时，可能导致用户的隐私权受到侵犯。此外，生成的内容可能会反映用户的个人喜好、兴趣等，从而进一步加剧隐私问题，例如：

- **Smart Compose 隐私问题：** 谷歌的 Smart Compose 功能可以根据用户的输入预测邮件的内容。然而，这可能意味着谷歌能够访问用户的邮件内容，引发了用户隐私泄露的担忧。
- **语音助手隐私问题：** 语音助手如 Siri、Alexa 等需要收集和分析用户的语音指令，以提供更个性化的服务。但这也涉及对用户的语音数据进行收集和存储，引发了关于隐私和数据安全的问题。
- **个性化内容生成隐私问题：** 生成技术可能会根据用户的浏览历史、社交媒体活动等生成个性化的内容，涉及用户个人数据的使用。这可能让用户感到他们的隐私受到了侵犯。

2. 虚假信息和误导性内容：

虚假信息和误导性内容是指生成技术产生的信息在形式或内容上误导受众，可能违背事实真相，损害信息的可信度和准确性。这种问题可能出现在各种内容中，包括文字、图像、音视频等多模态内容。虚假信息和误导性内容可能导致社会的混乱、信息泛滥和不信任的情况。人们可能不再能够确定何时可以相信所看到的内容，这可能削弱媒体的权威性和信息的真实性。此外，虚假



信息也可能对政治、商业和社会产生重大影响，导致不稳定和不确定性。例如：

- **Deepfake 虚假视频：** Deepfake 技术可以制作逼真的虚假视频，使人物出现在他们实际未出现的场景中。例如，有人可能使用 Deepfake 技术将名人的脸部特征添加到不实的视频中，以制造虚假事件。
- **虚假新闻和评论：** 生成技术可以产生看似真实的新闻报道、评论和社交媒体帖子，但这些内容可能缺乏事实支持，误导受众。这可能会对公共舆论、政策制定和个人信任产生负面影响。
- **制造虚假证据：** 生成技术可能用于制造虚假的证据，例如在法庭上使用虚假的文件或录音。这可能导致司法领域的不公正判决。

3. 民族仇恨言论和不当内容：

这一问题指的是生成的内容可能涉及针对特定民族、种族、宗教或文化群体的仇恨性、甚至挑衅性言辞。这种民族仇恨言论和不当内容的存在可能导致仇恨情绪加剧，引发潜在的冲突和社会争议，甚至导致社会的分裂，这对社会和谐、文化多元性以及人们之间的相互理解产生负面影响。

4. 偏见和歧视问题：

这一问题涉及到生成的内容可能带有种族、性别、性取向、宗教、地域、年龄、健康、职业、国别等方面的偏见和歧视，进而对个体、群体和社会造成不公平和伤害；产生这一问题的原因主要是指生成技术产生的内容可能反映出技术模型所学习的数据中存在的偏见和歧视。这些偏见可能是源自原始数据中



的社会偏见，也可能是因为模型在大规模训练数据中学习到的不平衡。例如：2018 年推出的谷歌 AI 助手 “Google Duplex” 。这个助手被设计成能够与人类自然对话，例如预订餐厅的电话。然而，有用户发现当助手模仿不同的人物时，它可能会展现出性别偏见，例如模仿女性声音时表现出过于顺从的态度，而模仿男性声音时则更自信；这个案例揭示了生成技术可能内在地反映出社会中已经存在的偏见和歧视。虽然这些模型并不是有意的，但大模型在学习过程中继承了这些偏见和歧视，并在未来的内容生成中可能不受控制的接受并发挥这些倾向。

5. 淫秽色情内容：

淫秽色情内容问题是指生成技术产生的内容可能包含裸露、性暗示、不雅言辞等不适宜公开传播的内容。这种内容可能在社交媒体、聊天应用、新闻评论等领域出现，可能冒犯人们的道德观念和价值观，对社会道德、个人尊严和文化价值产生负面影响，可能引发道德争议、社会不安以及对技术应用的担忧。例如： DeepNude 这个应用程序。这款应用可以使用深度学习技术，将普通照片中的衣着 “去除” ，从而制作出虚拟的裸体图像。虽然该应用最初声称是用于艺术目的，但它引发了广泛的担忧，认为这有可能被用于创造不适宜的虚假淫秽内容，侵犯个人隐私和尊严。直观的揭示了生成技术可能被用于制造淫秽色情内容，甚至可能损害个人形象和社会道德。

6、政治/军事敏感内容：



针对大模型生成的内容，可能因训练数据污染、用户恶意引导等导致生成有关国家领导人、国家制度/法律法规、政治事件等严重错误的内容，以及可能涉及军事等领域的敏感信息，可能对国家安全、国际关系和社会稳定产生影响。

7、恐怖/暴力内容：

生成式人工智能技术能够模仿并创造包括文字、图像和音频在内的多模态内容，这使得恐怖和暴力内容的创造变得更容易。虽然生成技术可以用于创作娱乐作品、艺术创作等领域，但它也可能被滥用，创造具有恐怖和暴力元素的内容，对社会产生负面影响。例如前两年一款名为"NightCafe Studio"的应用在社交媒体平台上引起了争议。该应用可以根据用户提供的文字描述生成有关恐怖和暴力场景的图像。虽然该应用声称是用于娱乐目的，但这种技术可能被滥用，用于创造恐怖主义、暴力行为等不良内容，对用户造成精神和情感伤害。

8、版权和知识产权问题：

大模型在生成过程中，模型可能会从大量的原始数据中提取灵感，导致生成的内容与现有的作品相似，从而引发版权和知识产权问题。例如 2020 年，某艺术家声称他的作品被 NVIDIA 的人工智能算法所复制，这一算法通过学习大量艺术作品生成了一系列类似的图像。这引发了关于生成技术是否侵犯了原创艺术家的知识产权的讨论。此外，生成的文本内容也可能受到版权保护。例



如，一些新闻机构和出版商可能会使用生成技术自动创作新闻报道，这可能引发与原创性和知识产权相关的问题。

9. 滥用和恶意使用：

生成技术可能被用于制造虚假信息、网络爬虫、网络钓鱼、欺诈行为、网络攻击等恶意目的。滥用技术可能会造成社会混乱、信任危机和人身安全问题，需要设定合适的监管和制约措施。例如近期出现以 FraudGPT（欺诈 GPT），和 WormGPT 为代表的黑化的生成式 AI 工具，专为网络攻击、犯罪而生的大模型。

10 责任和透明度：

生成技术的逻辑和决策过程往往难以解释，造成责任追溯困难，同时缺乏明确的责任归属：

- **责任归属：**在生成式人工智能的系统中，往往难以确定具体的责任主体。例如，如果一个由人工智能驱动的机器人犯下了错误，或者生成了有害的结果，很难确定应该由谁来承担责任。此外，由于人工智能系统的复杂性，即使试图进行责任追溯也可能面临困难。例如，在某些情况下，人工智能系统可能根据其接受的大量数据进行决策，而这些数据可能来自于多个来源，难以追踪其原始来源。
- **透明度和解释性：**生成式人工智能系统通常设计为能够生成多样、复杂的内容结果，这使得人们难以理解其内部的工作机制和决策过程。这种



缺乏透明度和解释性的问题可能导致人们对人工智能系统的信任降低，同时也使得在出现错误或争议时难以进行责任评估和追究。

3.4 大模型运营的业务安全问题

大模型服务在投入实际业务运营与应用时，同样面临诸多业务安全挑战，本节将如下几个业务环节来介绍大模型应用的安全问题：

1. 前置业务环节：

本环节主要涵盖企业在构建大模型服务时，与大模型交互前的各类业务阶段，如账号注册、登录、权益申请等业务运营的诸多环节，存在的业务安全风险主要包含：企业自有账号体系的批量注册、盗号、撞库、扫库、拖库等账号攻击风险，以及包含薅羊毛、权益侵占、机器作弊、审核资源浪费等诸多的业务运营风险。如下图所示 ChatGPT 推出仅两个月，注册用户就突破 1 个亿，随着用户规模的不断增长，各类违规账号也在不断的活跃，于是在 2023 年 4 月初开始，大规模封禁各类违规注册账号；同样以百度文心一言大模型服务上线为例，再面向全国用户开放了服务试用申请后，短时间内收到了大量新注册用户的提交，其中不乏一些违规账号的存在。





因此，大模型在投入运营阶段，其前置的业务环节的安全风控能力建设也会直接影响服务上线后的运营效果与服务质量。

2. 大模型交互环节：

在大模型交互环节，本节将分别从用户的“提问行为”和“提问内容”两个维度展开”。

首先是提问行为，在针对大模型发起提问时，黑产等不法分子围绕提问接口发起 AIGC 盗爬 / 垃圾提问 / 接口攻击 / 频控突破 / 资源侵占等攻击行为；针对大模型输出结果，黑灰产可以发起投毒反馈、恶意反馈等攻击行为。如下图所示，今年北京某公司起诉其多年的合作的伙伴某知名网校品牌，指其近期推出的数学大模型 MathGPT 和在某品牌学习机上线的 AI 助手，在未经其授权和许可情况下，爬取了海量数据，要求其公开道歉、删除数据资源，求偿 1 元，打响了 AIGC 盗爬的第一案。



针对“提问行为”存在的安全挑战

在大模型提问时，黑产等不法分子围绕提问接口发起AIGC盗爬 / 垃圾提问 / 接口攻击 / 频控突破 / 资源侵占等攻击行为，针对大模型输出结果，黑灰产可以发起投毒反馈、恶意反馈等攻击行为

- AIGC盗爬**: 展示大模型输出结果被非法爬取和滥用的场景。
- 垃圾提问/资源消耗**: 展示大量垃圾提问导致资源被恶意消耗的场景。
- 投毒/恶意反馈**: 展示用户输入恶意反馈或投毒内容，影响模型输出的场景。
- 接口攻击**: 展示通过接口发起攻击，导致服务异常的场景。

其次是用户提问内容安全，针对用户输入的各类 prompt，属于常规 UGC 内容安全范畴，例如需要针对用户输入内容进行包含“涉黄、涉赌、涉毒、涉政治、涉恐、涉爆、低俗、辱骂”等内容审核；同时还需要进行“恶意代码、网址安全”等注入、违规内容的甄别，避免违法违规内容作为 prompt 提交给大模型，诱导生成不合规的内容，如下图所示：

针对“提问内容”存在的安全挑战

在与大模型交互提问时，用户输入的 prompt 也能存在各类内容风险，主要涵盖如下几类：

- 涉黄
- 涉赌
- 涉毒
- 涉政
- 涉恐
- 涉爆
- 低俗/辱骂
- 恶意代码

3. 大模精调/推理环节：



在大模型服务上线后，还需要持续的对模型进行精调、推理；因此在运营阶段数据安全和隐私问题同样不能忽视，相关风险不在此章节进行赘述，可以参考 3.1 数据安全和隐私问题。



4.大模型安全解决方案

百度二十余年安全对抗的总结与提炼，围绕百度【文心大模型】安全实践经验，推出以 AI 安全为核心的大模型安全风控解决方案，从大模型全生命周期视角出发，方案涵盖大模型训练/精调/推理、大模型部署、大模型业务运营等关键阶段所面临的安全风险与业务挑战，提供全套安全产品与服务，助力企业构建平稳健康、可信、可靠的大模型服务。



如上图所示，本方案针对大模型训练阶段、部署阶段和业务运营阶段所面临的安全挑战，给出了完整的应对方案，本章节将会围绕数据安全与隐私保护方案、模型保护方案、AIGC 内容合规方案、以及业务运营风控方案四个维度详细阐述大模型安全能力建设；同时结合以攻促防的思路详细阐述如何建立 AIGC 内容安全蓝军评测能力，对大模型实现例行化的安全评估。

4.1 大模型数据安全与隐私保护方案



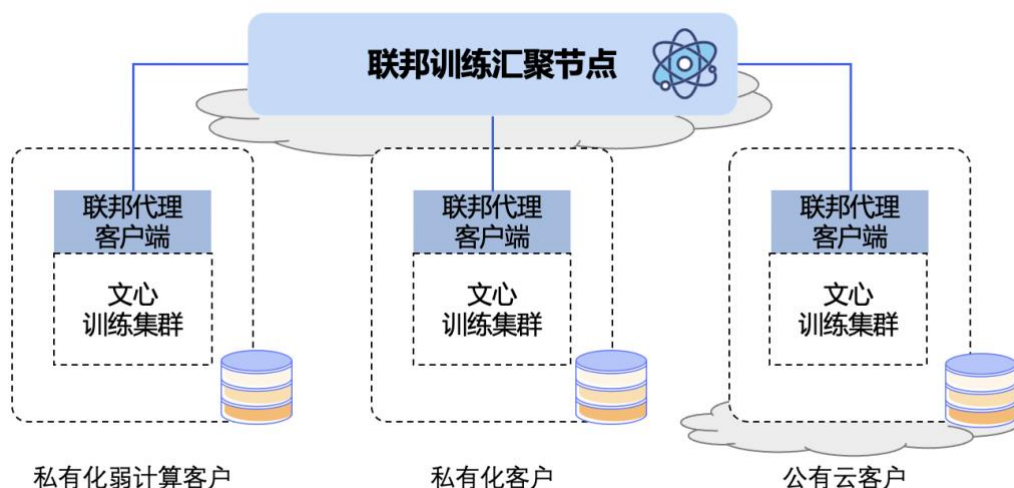
4.1.1 横向联邦大模型解决方案

百度安全支持公有云、私有化两种场景下的横向联邦软件方案，使得数据不出域的情况下，完成大模型的预训练、精调，解决数据传输过程中被截获的风险。

在联邦学习的横向技术基础上，又针对大模型的训练中遇到的特性做了优化。首先，大模型的训练较传统的训练阶段，又细分为预训练和精调两个阶段，并且训练模式也不同，为半监督训练和监督训练，并且两个阶段的训练量上，预训练要远大于精调，特别在精调训练手段也有很多特殊的 peft 的手段。其次，大模型的模型参数量较传统机器学习模型要多出几个数量级，并且在训练过程中有着计算量大和计算节点的算力不均衡等问题。最后，较传统的横向联邦，安全模型也是不相同的，传统上需要保护的是用户数据，而不是模型。而对于大模型的场景，除了用户数据是隐私的，其中训练的模型也是厂家投入了大量资本产生的，所以在在大模型场景下模型安全也是需要考虑的。

我们依照大模型所特有的特性，对现有的横向联邦技术做了演进。采用中心化的 CS 架构，中心节点为汇聚服务器，用于将不同参与方的结果数据进行汇聚，平衡各参与方的计算节奏，保持和管理最终的合并后的模型。每个参与方，采取弱侵入式的接入方式，部署参与方插件，用于和现有的算力平台进行结合，收集和管理本方的计算集群。



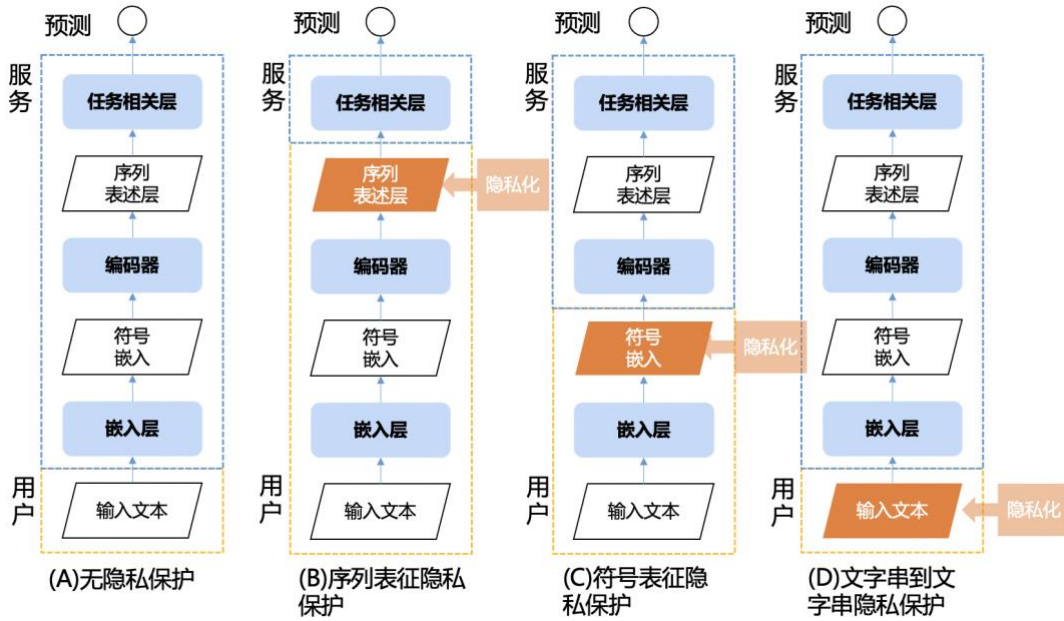


4.1.2 基于差分隐私的软件精调/推理方案

差分隐私 (differential privacy) 是一个数据保护手段, 通过使用随机噪声来确保请求信息的可见结果时, 不会因为个体的变化而变化, 实现仅分享可以描述数据库的一些统计特征、而不公开具体到个人的信息。这一特性可以被用来保护大模型在精调和推理时与云端服务端交互的用户数据隐私。

基于差分隐私的云上精调方案, 主要是利用差分隐私算法, 通过添加噪声去保护用户与模型之间交互的数据。部署上会分成客户端和一个服务提供端。根据目前的研究 Chen Qu[1], 依据大模型保护的位置的不同, 可以分成四种类型:





(A) 没有隐私保护 (Null Privacy) , (图 a) , 不应用何隐私限制, 因此也不提供任何隐私保护方式。这种方式为模型提供了最大的可用性。

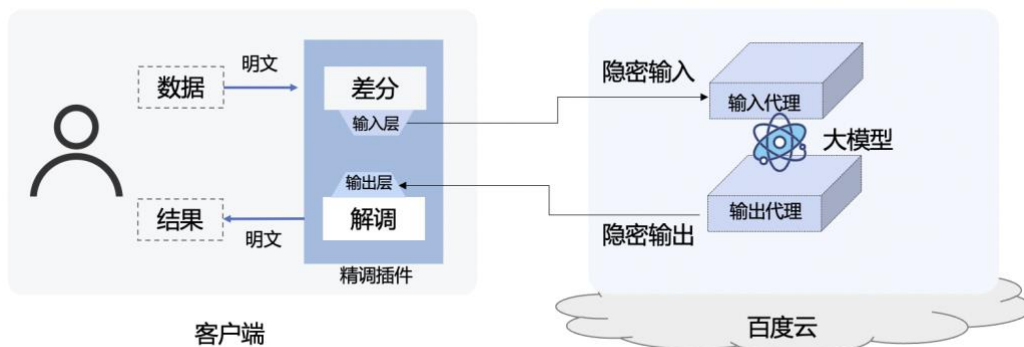
(B) 序列表征隐私保护 (Sequence Representation Privation) , (图 b) , 嵌入层 (Embedding layer) 和编码层 (Encoder layer) 是部署在用户侧。用户在本本地打乱序列的内容, 后再传输给部署在服务端的特定的任务层, 最后在服务端完成结果的计算。

(C) 符号表征隐私保护 (Token Representation Privatization) , (图 c) , 只有符号嵌入层 (embedding layer) 是部署在用户侧。用户本地完成符号化和嵌入表查询后, 完成文本到符号表征表示。这里可以将隐私保护的手段应用在符号表征上, 然后再发送给服务提供者。服务提供者将收到的符号表征, 再加上必须的符号表征和位置特征, 然后再作为编码层的输入

(Encoder Layer)



(D) 字符串到字符串隐私保护 (Text-to-text Privatization) , (图 d) , 用户在本地完成了字符串到字符串的转换, 并且在过程中完成隐私化的保护, 最后再将保护后的文字发送给服务提供者。服务提供者拥有一个完整的自然语言模型, 来处理这些隐私保护后的字符串。



我们的差分隐私方案, 主要应用在大模型的精调和推理阶段, 特别是对于性能高于精度的场景。其部署是包含一个客户端和一个服务端。客户端, 将用户的明文数据添加噪声混淆, 并进行初步输入层的计算, 完成对用户的输入数据进行保护, 并发送给服务端。服务端收到用户的隐秘数据, 并将数据通过大模型的输入代理层传递给大模型进行计算。计算后, 未解密的结果通过输出代理, 发送给客户端。客户端收到后, 先进行输出的解码等输出层操作后, 再经过差分解调, 强化输出结果, 消除噪声对结果的影响, 得到计算的明文结果, 并返回给客户。由于在整个计算过程中, 传递的数据均为添加噪声后的中间计算结果, 在保证计算性能的基础上, 通过差分隐私增加数据还原的难度, 在一定程度上能够保护用户数据的安全。

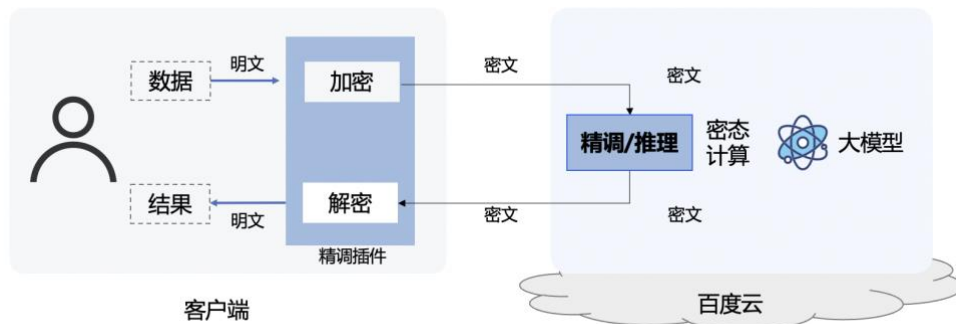


4.1.3 基于同态密码学的软件精调/推理方案

同态密码学是一项联邦学习的关键技术，提供了在加密状态下对数据进行计算和处理的能力，从而保护数据的隐私和安全。对于大模型的数据保护思路，是通过同态密码学来实现大模型的计算逻辑，从而大模型可以接受密态化的数据输入，整体精调和推理过程完全是密态化的进行，最终的结果也是以密态的形式返回给客户端，整个过程完全是密态化的，所以将此过程全部部署到在云上的服务端。而客户仅需要将本地的隐私数据密态化后上传给服务端，所有计算过程由云端外包完成，但是云端服务，不能获取到计算的内容。

对于同态密码学方案，核心是如何通过同态密码学实现大模型的核心计算逻辑，其中主要包括，Embedding，Transformer (Attention) 和 Header 等大模型基础组件结构。由于同态密码学计算复杂性和支持的计算有限，如何合理的利用同态密码学算法能达到可用性和精度的要求，实现精调和推理阶段隐私保护的方案。

目前基于同态密码学方面的大模型研究，公开研究主要集中在推理阶段，也有少量的精调方面。根据所采用的同态密码学算法的实现不同，大致可以分为基于全同态密码学 FHE 实现和基于 MPC (SecureShare) 实现两大方向。在 FHE 方向，有基于 CKKS 的 THE-X[3]，以及基于 BGV 的 Liu, Xuanqi[7]，



基于 HGS 的 Primer[4]; 在 MPC 方向有基于 2PC 的 MPCFormer[2]和 Iron[5], 基于 3PC 的 Puma[6]。除了底层实现方法的不同之外, 对于如何通过同态密码学中有限的计算方式去实现和逼近大模型的基础算子也是目前研究的热点。在降低计算量的同时, 如何平衡计算量和网络传输量之间的关系, 以达到在实际应用中能最大化的降低耗时, 将算法可用性能进一步接近可用, 也是研究所追求的目标。

我们的同态密码学方案是结合同态密码学和差分隐私等技术, 构建的一个对用户数据进行密态计算的方案, 并将此技术运用在大模型的精调和推理阶段。在用户客户端, 会安装一个客户端插件, 此插件主要用于加密用户的隐私数据, 形成可以用于密态计算的语料, 通过网络连接将加密后的数据发送给服务端。在服务端, 将加密的语料直接加载后, 通过同态的特性直接用于模型计算。最终的结果也将以密文的形式, 返回给客户端。客户端, 通过插件将数据解密后得到最终的结果。由于数据全程都是密态形态, 所以任何第三方都不可窃取到用户在使用大模型中交互的数据, 从而保护了用户数据的隐私。

4.1.4 可信执行环境解决方案

可信执行环境 (trusted execution environment, TEE) 是处理器中的安全区域, TEE 保护程序与数据的机密性和完整性不被外部窃取和破坏。与存储加密和网络通信加密一起, TEE 可以保护落盘(at rest)和通信过程中(in transit)的数据隐私和安全。随着 TEE 技术的发展, 在计算核心与内存之间增加安全处理器, 以保护被计算核心使用(in use)的数据安全和隐私的机密计算技术出现。



TEE 能够作为云计算的信任根，保管根密钥，确保 TEE 外实体无法获取，还可以通过远程证明和负载度量值的结合，使公有云达到私有云的安全等级。对于私有化的场景，TEE 可以充分发挥“飞地”的作用，将隐私敏感资产部署在其他实体。在如下大模型使用场景中保护敏感数据资产：

- 在使用第三方大模型服务提供的精调和预测功能时，保护用户输入数据和精调产出模型的隐私

- 在第三方部署大模型服务时，保护模型的隐私

TEE 方案具有强兼容性、高性能和模型准确等优势：

- 通过安全的设计与配置，可以运行复杂的分布式系统；
- 处于机密计算状态的处理器操作明文数据，不需要使用差分 and 同态等密码学算法，具备高性能处理海量数据的能力；
- 支持通用的 NLP 算法，模型和计算精度无损失。

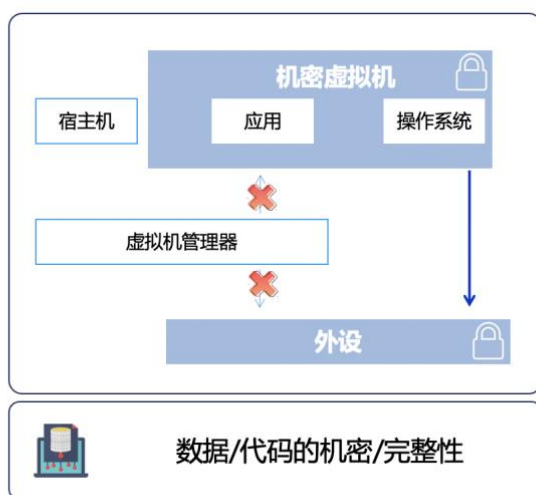
由于 TEE 包括多个硬件模块，涉及数据处理和流通全生命周期，比较容易受到侧信道攻击，需要构建纵深防御安全体系，抵抗不同方向的攻击，并加强安全测试来主动发现问题，还需要及时更新系统中所有组件的安全补丁。

目前，Intel，AMD 和海光提供虚拟机 TEE（机密虚拟机），可以保护虚拟机内的应用、操作系统和外设不被宿主机和虚拟机管理器访问：

- TDX (Trusted Domain Extension) 是英特尔新提出的能够部署硬件隔离的虚拟机（可信域，trusted domain，TD）的技术框架，TDX 从多方面保护机密虚拟机，降低 TCB，加强对数据和知识产权流通控制。TDX 技术生态全面，功能强大，对安全启动，IaaS 层部署，NLP 运算的 CPU 加速等支持较好；



- SEV (Secure Encrypted Virtualization) 是 AMD 提出的机密虚拟机方案，其推出时间较 TDX 早，因此软件生态较好，upstream linux kernel, openstack, kubevirt 和 libvirt 等虚拟化相关生态支持较好；
- CSV (Chinese Secure Virtualization) 是海光根据 AMD SEV 国产化的解决方案，使用国密算法，信任根全部国产化。



数据只能被有权限的实体访问

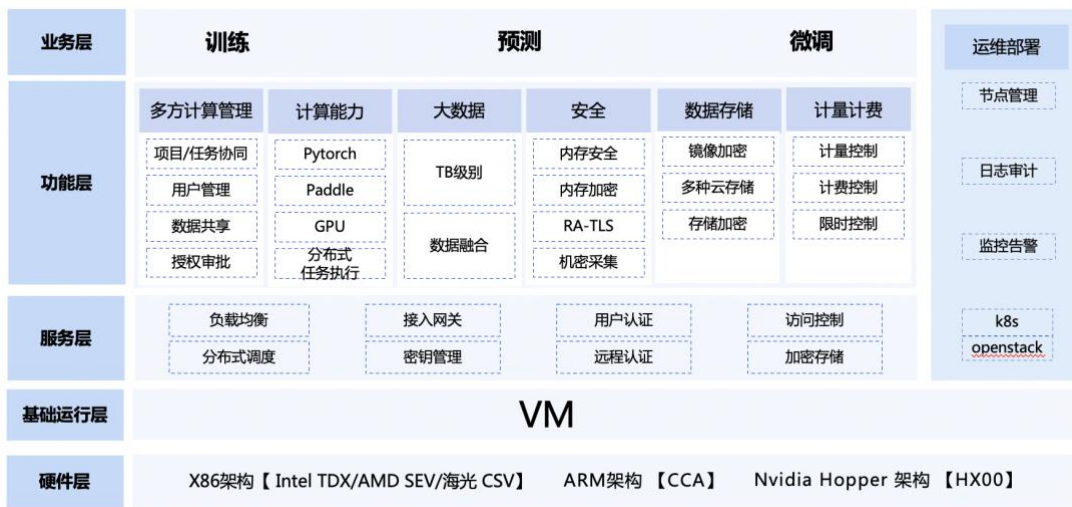
百度提供的可信执行环境解决方案 (MesaTEE) 是基于硬件 TEE 的大模型机密计算方案，支持在 Intel SGX, Intel TDX, AMD SEV-SNP 和海光 CSV 等多种硬件 TEE 内进行大模型训练和推理。通过以 PCIe 穿透 (passthrough) 的方式访问 Nvidia H 系列 GPU 和海光集成 DCU 等具备机密计算能力的外部加速设备，MesaTEE 能够获得与非机密计算相当的大模型计算性能，拥有良好的模型效率和用户体验。

MesaTEE 将传统虚拟机安全与 TEE 相结合，将可信启动过程记录到远程认证的度量值中，保证启动过程的安全，提高远程认证的真实性。机密虚拟机中运行的容器启动前，其数字签名会被校验，确保程序来源的合法性。隐私数据以透明加解密的方式落盘，保护数据隐私和安全的同时，提高应用的兼容



性；非秘密的程序，通过建立哈希树（Hash tree）的方式，保证其完整性的同时，兼顾访问性能。机密虚拟机之间使用基于远程认证的透明加解密技术，确保通信过程中的数据隐私安全。

MesaTEE 深耕大模型使用场景，支持分布式训练、精调和推理，通过基于身份的访问控制，具备多租户数据及模型隔离管理和保护，多方数据训练和推理等数据融合功能。



可信执行环境是云计算中不可或缺的一部分，它从硬件层面解决了软件根本的信任问题，是云计算的“根”。机密计算是大趋势，英特尔、AMD 和英伟达等硬件提供商均提供了机密计算硬件解决方案。微软、亚马逊云、谷歌云和阿里云等均提供机密计算的设备和解决方案。百度、蚂蚁金服和字节跳动等均在使用机密计算为业务提供隐私及安全能力。

4.1.5 基于安全沙箱的解决方案

安全沙箱技术是一种通过构建隔离的可供调试、运行的安全环境，来分离模型、数据使用权和所有权的技术，同时提供模型精调计算所需的算力管理和



通信等功能，保证模型拥有方的预训练模型在不出其定义的私有边界的前提下，数据拥有方可以完成模型精调任务。

安全沙箱产品是提供给模型开放共享过程中各参与方使用，提供模型安全开放共享所需的算力管理和通信等功能，并满足计算任务需求的软件系统或软硬件一体化系统。

安全沙箱通过界面隔离、环境隔离、网络隔离、执行隔离、数据隔离五大隔离技术达到模型和数据的可用不可见。



界面隔离：为抵抗来自站点外对平台调试环境的窃取数据的攻击，通过界面渲染的手段，使用户仅可以看到调试环境中的内容，可以向环境中提交操作和数据，但是无法直接从环境中获取到操作的内容，实现指令到环境的操作是单向的效果。

环境隔离：为抵抗来自调试环境中对于运行环境的渗透攻击，通过将使用环境划分的手段，根据操作对象的不同，将调整逻辑代码的区域划分为调试环境，将对真实全量数据进行操作的区域划分为运行环境，两个区域完全隔离不存在直接的访问介质。从而达到：在调试环境中改动程序逻辑，仅可通过脱敏数据了解格式，但不可触碰真实全量数据；在运行环境中要操作真实全量数



据，其所提交的程序逻辑为固定的，操作的内容经过审查，其最终运行的结果为确定的，操作过程也是可回溯和可追责的。

网络隔离：为抵抗来自组件陷落后，形成跳板对内部其它组件发起的攻击，通过物理硬件策略的手段，使隔离环境间和组件间所工作的网络层面是独立的，其间的交换的数据是单向的，目的是确定的，协议上是简单和明确的。以达到网络层面上访问可控，可审计，以及出现风险后可以进行有效阻断和控制影响范围的作用。

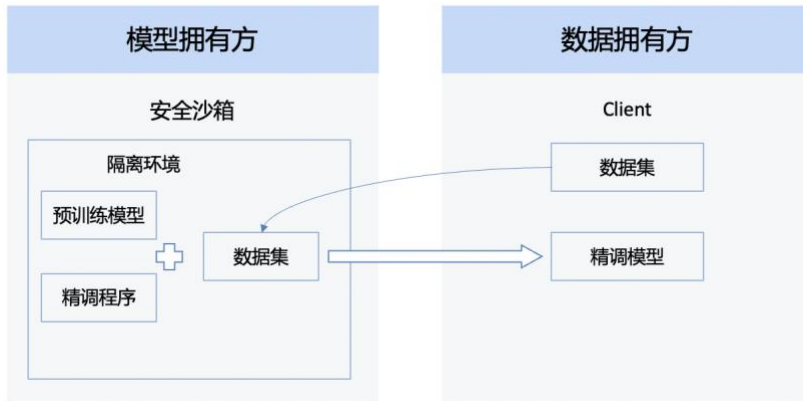
执行隔离：为抵抗来自执行环境内，运行逻辑对执行环境的渗透攻击，通过虚拟化技术，将用户直接操作的调试，以及间接使用的运行环境的执行体，与真实执行的物理机环境相互隔离，去除运行环境之间的物理机的差异，保护物理环境的不被穿透，消除运行残留，阻断租户间的相互影响。

数据隔离：为了防护对数据的直接窃取的攻击出现，通过对数据的访问进行控制，在调试环境和运行环境所访问的数据：物理策略上限定，使用者上限定，使用方式上限定。

在大模型精调领域，基于以上五大隔离技术，再结合访问控制策略，可将安全沙箱技术应用在单方保护模型的场景及保护模型和数据的场景。

对于单方保护模型的场景，安全沙箱部署在模型拥有方，模型拥有方在沙箱中上传预训练模型和精调程序，数据拥有方在安全沙箱中上传精调数据集，在沙箱中完成精调工作，产出精调模型。



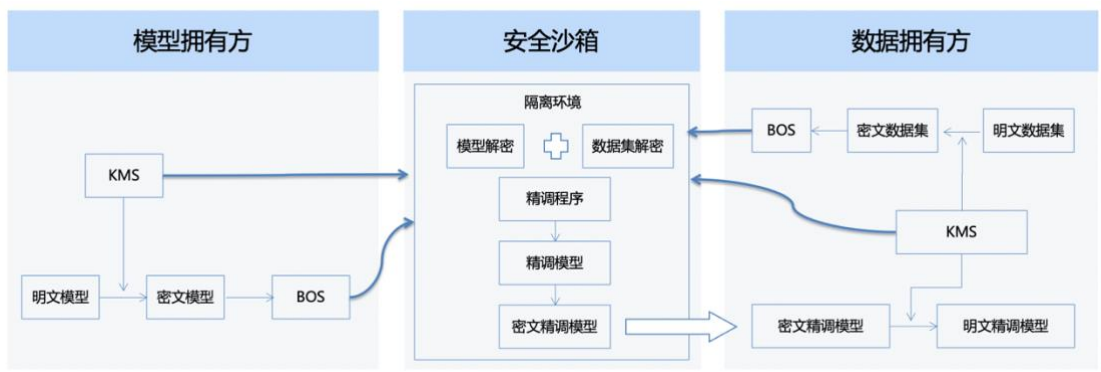


对于同时保护模型和数据的场景，模型拥有方和数据拥有方角色相同，但是其不了解精调相关领域知识，无法自主完成精调，只能雇佣外部人员，但是又不能让外部人员触碰到需要保护的模型和数据。此时便可以利用安全沙箱可用不可见的特性，使得外部人员可以在不触碰原始数据的情况下，对数据进行清洗、治理，使其符合精调数据的格式要求，能够将处理后的数据集应用于模型精调工作，产出精调模型。



结合 KMS 作为密钥管理服务，对进入到沙箱的预训练模型和精调数据集进行加密，在沙箱的隔离环境内进行模型精调时，对其使用对应的 KMS 进行解密，完成精调时产出密文精调模型，数据拥有方使用自己的 KMS 对其解密为明文精调模型。





在大模型推理领域，安全沙箱可提供在线推理服务用于一键部署精调后的大模型，对外提供在线 API 推理服务。在线推理服务提供精调模型部署、API 网关、负载均衡、安全访问认证、动态脱敏等功能。模型支持多实例部署方式和实例动态扩缩容，提供 API 网关能力，对后端的实例进行负载均衡，以保证在线推理服务的高可用性；对请求进行安全访问认证，确保请求来源的合法性；对推理服务返回的内容实时动态脱敏，确保推理结果不包含敏感数据。

业务层	模型精调	模型推理	
	模型保护	数据保护	
服务层	模型部署	API网关	负载均衡
	内容安全动态脱敏	安全访问认证	
控制层	网络访问控制	数据访问控制	出域控制
执行层	安全容器	虚拟机	机密计算虚拟机
存储层	模型加密	数据集加密	
机器层	物理机	虚拟机	
	CPU	GPU	



4.2 模型保护方案

在模型训练、管理、部署等环节，主要有如下两个方向的业务痛点：

- 1、语料数据管理：面对多渠道收集珍贵语料数据，如何实现高效的数据管理，防范模型原始语料数据泄漏，提高语料数据加工效率
- 2、模型资产保护：大模型文件是企业核心数字资产，如何防范大模型文件在训练、推理、微调等环节的模型文件泄漏风险

为了解决上述模型安全相关问题，构建行之有效的模型保护方案，如下图所示

示，分大模型语料数据安全管理与大模型资产全流程保护两套管理方案：



4.2.1 大模型语料数据安全方案

在大模型的语料数据安全方案中，保护敏感数据、确保数据的完整性和合规性是至关重要的。以下是一套综合的语料数据安全方案：



1. **元数据管理：** 建立完善的元数据管理系统，记录数据的来源、用途、分类、权限等信息，以便对数据进行跟踪和监控。
2. **分类分级：** 对话料数据进行分类分级，根据敏感程度和机密性将数据划分为不同等级，以便进行适当的保护和控制。
3. **流转审批：** 设计流程化的数据流转审批机制，确保数据在不同环节的传递经过合法的授权和审批。
4. **数据鉴权：** 引入严格的数据鉴权机制，只有经过授权的人员可以访问特定等级的数据，确保数据不被未经授权的人员获取。
5. **加密保护：** 对敏感数据进行加密，保障数据在存储和传输过程中的安全性。采用合适的加密算法，确保数据的保密性和完整性。
6. **行为审计：** 部署行为审计系统，记录数据的访问、修改、复制等操作，以便跟踪数据的使用情况，及时发现异常行为。
7. **数据脱敏：** 在特定场景下，对敏感数据进行脱敏处理，以保护个人隐私和敏感信息的安全。
8. **访问控制：** 实施严格的访问控制策略，确保只有具备访问权限的人员才能进入数据存储区域。
9. **数据备份与恢复：** 建立定期的数据备份机制，以防止数据丢失或损坏。同时，测试数据恢复流程，确保在紧急情况下能够快速恢复数据。
10. **敏感信息检测：** 引入敏感信息检测技术，及时识别数据中的敏感信息，如个人身份证号、银行账号等。

4.2.2 大模型资产全流程保护方案



针对大模型的全生命周期，从模型训练到部署，采用全方位的安全防护措施是关键。以下是大模型资产全流程保护的方案：

1. **模型训练安全：** 在模型训练过程中，采用隔离环境，确保模型训练的数据和代码不受未授权访问。引入训练数据的加密和隐私保护措施，防止敏感信息泄露。
2. **模型流转安全：** 设计模型流转的安全机制，确保模型在传递过程中不被篡改或恶意替换。可以使用数字签名等方式验证模型的完整性。
3. **模型推理安全：** 在模型推理阶段，引入安全沙箱和权限控制，确保模型运行在受控环境中，避免恶意代码注入和攻击。
4. **模型微调安全：** 在模型微调过程中，采用差分隐私等技术，确保微调数据的隐私性和保密性。
5. **私有化部署安全：** 对于私有化部署，强调数据在企业内部的隔离和安全性。建立私有化部署的权限控制和监控机制。
6. **模型演化与更新：** 引入安全审查流程，确保模型的更新和演化过程中不引入漏洞或不安全的元素。
7. **模型审计与跟踪：** 对模型的运行情况进行审计和跟踪，及时发现异常行为和风险。
8. **安全修复与更新：** 在发现模型存在漏洞或安全问题时，能够及时进行修复和更新，防止潜在的威胁扩散。

通过上述方案，可以在大模型的整个生命周期中，确保模型资产的安全，保护敏感信息，防范恶意攻击，维护业务的正常运行。这些方案涵盖了模型的各



个阶段，从数据安全到部署安全，都为大模型业务的安全运营提供了强大的保障。

4.3 AIGC 内容合规

内容合规能力建设在实践中面临着多个难点与挑战，这些挑战需要企业、技术团队和监管机构共同努力克服，以下是百度安全在实践中总结的当前 AIGC 合规能力建设过程中所面临的主要的难点和挑战：

1. **复杂多样的内容：** AIGC 生成的内容可能涵盖各种主题、形式和风格，包括文字、图像、音频、视频等。针对不同类型的内容制定合适的合规标准和规则是复杂的任务。
2. **监管法规不断变化：** 互联网内容的监管法规不断变化，跨国业务还需要适应不同国家和地区的法律法规。因此，跟上法规变化，确保合规性是一个挑战。
3. **技术与人工判定的平衡：** 判断 AIGC 生成内容是否合规往往需要结合技术和人工判定。技术虽然可以自动检测一部分问题，但对于某些复杂情况，人工判断仍然是必要的，平衡二者需要技术和人员的投入。
4. **多语言与文化差异：** 如果业务覆盖多个语言和文化，要确保所生成的内容不涉及不当言论、歧视性语言等，需要深入了解不同语言和文化的特点。
5. **隐蔽性的风险：** 有些合规问题可能不容易被自动检测，需要进行深入的内容分析和理解。例如，某些内容可能含有隐喻、讽刺等，难以简单地依赖技术检测。



6. **合规规则的标准化：** 制定合适的合规规则 and 标准需要深入的行业和领域知识。但是，在不同领域和业务中，合规标准可能存在差异，制定一套通用的标准是具有挑战性的。
7. **时间敏感性：** 有些内容可能在发布后迅速传播，导致迅速产生影响。在这种情况下，需要在短时间内判断内容的合规性，需要高效的合规审核机制。
8. **平衡安全与隐私：** 在确保内容合规的同时，也要保护用户的隐私和个人信息。确保安全合规的同时，避免不必要的数据收集和使用。
9. **技术局限性：** 当前的自动化技术虽然在内容检测方面取得了进步，但仍存在误报和漏报的问题。技术的局限性需要考虑如何提高准确性和效率。

本方案基于《生成式人工智能服务管理办法》，以及百度在人工智能技术的沉淀与总结，构建了五道安全防线，确保大模型生成内容的安全与合规：

1. 预训练数据过滤方案

在构建大语言模型之前，需要对训练数据进行有效的筛选和清洗，保留高质量的语料数据用于训练对大模型安全性有着至关重要的影响。通过预训练数据过滤方案减少训练数据中的偏见、不准确性和不适当内容，从而从根本上提高模型生成内容的质量和安全性。

百度使用安全内容业务中积累的海量有标注数据，基于 ERNIE 模型的领先内容理解能力，构建了通用的内容安全召回模型，能够高效检出训练语料中的



有害内容；同时通过业务风控富集的敏感词词库过滤数据中的脏话和不适出现词汇，提供召回模型之外的快速更新能力。除了过滤有害内容，预训练数据过滤方案也能够删除可能包含个人身份信息、隐私敏感信息的内容，用以严格保护用户的隐私。

2. 内容干预系统

大模型的内容干预是指通过人工审核、过滤技术或其他方式，干预模型输入的内容，以确保其符合特定的标准、规范和价值观。这种干预可以帮助减少有害、不准确或不恰当的内容，并提高生成内容的质量和安全性。

百度可提供完整的实时内容干预系统，内置红线必答和 Query 干预功能。红线必答能够很好回答常见的红线问题，确保回复内容高度安全合规，维护社会主义核心价值观；Query 干预支持用户配置相应规则，通过对包含特定敏感词的快速匹配，将不安全 Query 引导至更加合适的处理流程中（例如标准回复模版），减少大模型在该 Query 输入下产生有害内容或者不正确数据。

值得注意的是，内容干预需要权衡大模型的自由创作能力与生成内容的质量和安全性之间的关系。过于严格的内容干预可能会大幅抑制大模型的创造性，而过于宽松则可能导致有害内容的生成。因此，掌握合适的内容干预尺度也对使用方提出了高要求，百度提供了相对审慎可用的预置策略，能够很好地兼顾大模型创新能力和回复内容的安全性。



3. 安全分类算子

大模型输入的安全分类是指将用户输入内容进行分类，以判断其安全性和合适性。这种分类能够帮助防止不良内容的生成，保护用户免受有害、不准确或不适当的内容影响。通过有效的输入内容安全过滤，能够极大程度地减少大模型生成不安全或者负面的回复内容，同时结合高精度的分类标签，通过改写技术可以构造出更适于大模型输出合规回复的提示词模版。

百度结合多年的业务内容安全分类实践，将输入内容划分为不同的主题类别和语义类别，由此构建出完整正交的标签体系，基于知识增强的 ERNIE 系列模型，提供覆盖涉政、涉黄、违法等不同主题和恶意、攻击、中立、正常等不同语义的内容分类能力，能够高效检出涉政、涉黄、违法、歧视、辱骂、负面价值观等类别的不安全输入，同时提供高质量的提示词改写模版，协助大模型更好地理解问题并正确回答。

4. 大模型微调安全策略

在大模型预训练完成后，为了提高其生成内容的安全性，可以进行安全微调。基于已经通过安全审核的、符合安全标准的指令数据对大模型进行微调，以指导其生成更合适、不含有害内容的回复内容。微调后的大模型可以进一步通过 RLHF 方式提升大模型对安全回复内容的偏好程度，引导鼓励大模型生成更加高质量的安全内容。

百度基于数据标注和数据质量管理的既往工作，依赖多样性的内容标注语料库和人类反馈的偏好标注，提供多类别的安全调优语料和多维度的预置奖励



模型，能够有效地将大模型的生成能力与人类偏好相对齐，从根本上让大模型遵循安全有用的原则与用户进行信息交互。

5. 输出内容安全过滤

大模型输出内容安全过滤是指对大模型生成的文本内容进行检测和筛选，以识别并过滤掉有害、不准确、不适当或不合规的回复内容，这有助于确保大模型生成内容的质量和安全性。

百度使用业务风控中积累的高危词典对输出内容进行安全过滤，在滤除有害敏感词后通过语义改写将安全回复内容作为最终的大模型输出，确保输出环节安全合规。

此外，在面对第三方自建大模型的服务厂商，百度安全同样构建了如下图所示的大模型内容安全防护体系，围绕用户输入的 prompt 内容、大模型生成内容提供专业的内容审核能力，其核心服务包含如下

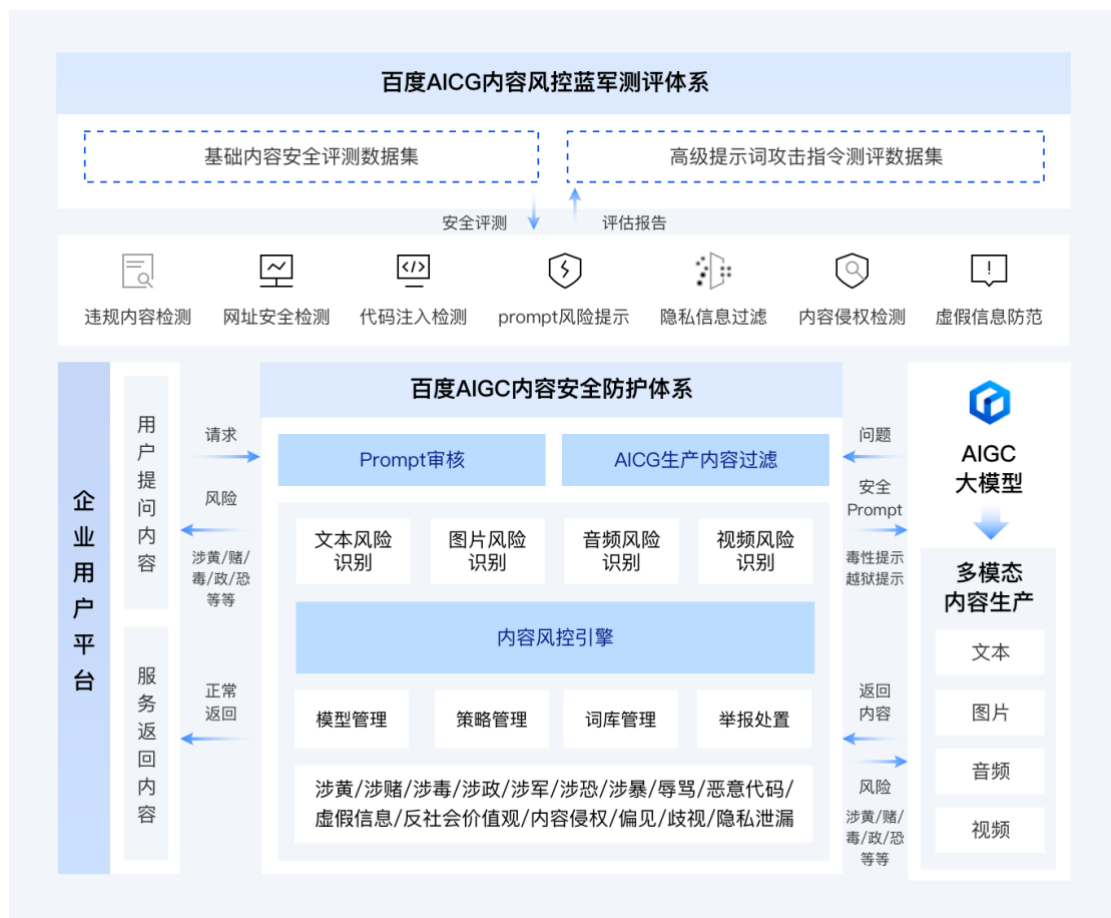
- **Prompt 审核与改写**：对于用户输入内容提供多维度内容审核能力、并针对恶意诱导大模型生成违规内容的 Prompt 进行改写并做毒性提示；



- **红线知识库**：对于用户 query 涉及国家领导人、制度、政策评价等诸多敏感的内容，基于红线知识库参与第三方大模型的精调与推理，保障内容客观、准确、全面以及政治中立。

- **AIGC 多模态内容审核**：为大模型生成内容提供包含违法违规内容审核、违反价值观、存在偏见歧视、内容侵权等风险内容过滤服务；

整体业务流转如下图所示



4.4 大模型业务运营与安全风控

在大模型业务运营的环节，依托百度安全智能风控能力，可以建立大模型业务运营的安全风控防护体系，可以在大模型前置云运营阶段（如：用户注册、登录、权益申请等环节）、以及大模型交互环节（如：用户提问环节、回



答内容反馈等环节)，结合用户行为、终端环境、网络特征等信息建立有效的安全防护体系，针对异常请求做实时风险检测，保障大模型处于一个安全、可靠的运营状态，如下图所示：



大模型在交互场景中的业务运营中，面临着多重安全威胁和风险，本方案结合当前场景，依托百度安全昊天镜智能风控服务，构建了包含账号安全、接口防刷、人机识别、AIGC 盗爬识别、设备风控以及风险情报等方面的能力：

1. 账号安全： 大模型的交互场景通常涉及用户账号，保护用户账号安全是首要任务。使用多因素身份验证（MFA）等措施，确保只有合法用户能够访问系统，防止恶意登录和盗号行为。此外，采用安全令牌、会话管理等方式，加强对用户身份的验证和保护。



2. 接口防刷： 针对大模型的接口，可能会受到恶意请求的攻击，导致系统资源过载甚至崩溃。通过实施限制频率、验证码验证等手段，可以有效减轻接口被恶意刷取的风险，确保正常的业务运行。

3. 人机识别： 大模型的交互场景可能会遭受机器人攻击，如恶意机器人批量注册、刷单等。引入人机识别技术，如验证码、滑动验证码、人脸识别等，可以辨别真实用户和机器人，防止自动化攻击。

4. AIGC 盗爬： 大模型生成的内容可能会被恶意爬虫大规模复制、传播，造成信息泄露和盗用。通过部署反爬虫技术，如 IP 封禁、User-Agent 检测等，可以减少非法爬取行为，保护生成内容的安全性。

5. 设备风控： 大模型的交互可能涉及多种设备，如电脑、手机、平板等。为了防止设备被劫持用于恶意行为，可以采用设备指纹识别、用户行为分析等技术，识别和阻止异常设备的访问。

6. 风险情报： 搜集和分析安全风险情报，了解当前的安全态势，能够及早预防和应对潜在威胁。风险情报可以来自外部的安全报告、漏洞数据库等，也可以基于内部的访问日志和异常行为分析。

在大模型业务运营中，上述安全措施和风险情报的作用是不可忽视的。综合运用这些措施，可以减轻大模型交互场景中的各种安全风险，保护用户隐私和数据安全，维护业务的稳定运行。同时，持续的监控、分析和改进也是确保业务安全的重要环节，以适应不断变化的安全威胁。



5. 大模型蓝军安全评测解决方案

如前文所述，大模型产出的内容都是基于大量的数据筛选和模型训练，不具备任何价值观，但数据的筛选、清洗，以及不同地区的内容监管尺度差异，会导致最终的内容产出存在不同的风险；有可能输出具有含有侮辱性和偏见歧视的内容，有可能输出非常不正确的价值观，也可能被用于恶意活动，如欺诈、虚假信息传播；因此对大模型的内容安全性进行评估和改进显得尤为重要。

本检测方案以网信办《生成式人工智能服务管理办法（征求意见稿）》为指导基础划分安全分类，通过在该安全分类体系中设定的不安全对话场景，针对性的生成了对应的评测内容，供大模型进行内容安全评测评估，以达到帮助大模型内容风控系统升级，促进大模型生态健康发展的目的。

5.1 建立大模型蓝军所面临困难

大型语言模型(LLM)可以自动化或协助人类完成各种任务，但获得的回复存在如幻觉、偏见和越狱等问题，这可能导致生成有害输出。因此在部署之前，建立大模型蓝军测试体系，通过主动攻击大模型的方法来发现缺陷非常重要。主动攻击成功的样本数据将提供给大模型安全防御开发人员进行针对性优化，高质量的样本将大幅提升安全防御开发人员的研发效率。



大模型发展的初期，大模型蓝军测试主要依赖人工编写测试语料，并人工标注危险回复。这种完全基于人工的蓝军测试流程限制了发现威胁的数量和多样性。因此，建立基于自然语言处理技术，机器学习技术，大语言模型技术的自动化大模型蓝军测试框架来代替人工测试体系显得尤为重要。

建立自动化大模型蓝军测试体系面临以下的挑战：

5.1.1 风险语料生成的自动化实现

生成大量对大模型具有潜在风险的语料存在巨大的挑战：首先需要生成通顺且符合人类表达逻辑的语料，其次生成的语料需要满足具有潜在的风险的条件。业界常规的方法是通过人工撰写收集的方式获取这一部分的内容数据，然而这种方法存在以下缺点：

- 人工成本昂贵：招募、培训和管理大量的志愿者需要耗费大量的人力资源和时间。同时，为了确保生成的测试数据质量，需要对志愿者进行严格的监督和审核，增加了运营成本和人力投入。
- 测试集数量存在瓶颈：由于依赖于人工手写生成攻击测试数据，测试集的数量受到限制。这种限制可能导致测试集的规模不够大，无法全面评估和发现大模型的潜在安全漏洞和问题。
- 人工生成语料存在偏狭：人工生成的语料往往受到个人经验、偏见和局限性的影响。志愿者可能无法涵盖各个领域和语境，导致生成的测试数据在覆盖范围和多样性上存在限制。这可能导致在处理新的、未知领域的输入时表现不佳。



- 人工生成语料的框架可扩展性差：当需要引入新的内容或应对特定场景时，依赖人工手写生成攻击测试数据的方法往往难以快速实现。对于复杂的测试需求或涉及到大规模语料的变化，人工方法的可扩展性和灵活性有限。

综上所述，传统的依赖志愿者招募和人工手写生成攻击测试数据的方法在成本、规模、多样性和可扩展性方面存在一些不足之处。为了克服这些问题，可以考虑结合自动化的方法，利用更强大的大模型语言生成能力来进行攻击测试和安全评估。针对该诉求，我们建立了风险内容评测数据自动生成框架。通过模型的语料生成方法同样面临诸多挑战，包括：

- 生成语料的质量：生成的攻击语料需要符合人类表达逻辑，通顺，能被人类所理解。
- 生成语料的多样性：生成的攻击语料需要在内容上足够广泛，避免大量测试语料在语义上聚集，降低测试的整体范围。
- 生成语料的威胁性：生成的攻击语料需要具备潜在的引起内容风险回答的能力，较低的攻击成功率将降低攻击样本的采集效率。

5.1.2 建立大模型回答内容的自动评测能力

对于海量的测试问题和大模型回答的风险性评测，全部依靠人工审核将耗费大量人力资源。需要建立自动化模块，快速准确地完成海量风险内容的自动



评测。需要通过大量算法优化和提示词工程研究提升模型判断的准确性，逼近人工审核的能力。自动评测能力面临以下挑战：

- 评测准确性：自动评测模块需要准确地感知回答的内容风险，过低的评测准确性将影响攻击样本采集的质量和效率。
- 评测计算效率：自动评测模块需要快速地完成内容风险感知的计算，过低的计算效率将影响攻击样本采集效率。

5.2 百度安全面向大模型蓝军的解决方案

大模型蓝军评测是一种主动的安全测试方法，旨在模拟攻击者的行为，评估大模型系统在真实威胁面前的安全性能与内容合规问题。蓝军安全评测的意义在于为大模型的业务运营提供全面的安全保障，增强系统的内容对抗能力，从而确保生成内容的安全性、完整性和可用性，大模型蓝军建设的整体目标是：

- 建立自动化的攻击语料生成能力
 - 提升威胁攻击语料的输出数量
 - 提升威胁攻击语料的输出多样性
 - 提升威胁攻击语料的攻击成功率
- 建立自动化大模型回复风险标注能力
 - 优化自动化大模型回复风险标注的效率
 - 优化自动化大模型回复风险标注的准确性



- 建立大模型安全评测框架
 - 设计全面权威的评测标准，量化指标

5.2.1 自动化的攻击语料生成

通过参考互联网安全领域中红蓝攻防的思路，建立大模型安全蓝军体系，通过自建提示词数据集主动引起大模型的不安全回复来发现潜在的风险。建立的风险内容评测数据自动生成框架包括以下详细部分：

- 风险内容语料生成：我们利用开源的大型语言模型蓝军（红队）数据集作为基础，通过筛选其中具有高风险的提示词，以及采用 stochastic few-shot 的方法，利用外部的语言模型生成新的测试提示词。这种方法可以利用已有的蓝军数据集中的关键信息，并借助外部模型的生成能力来扩充语料库，增加测试数据的多样性和覆盖范围。

其中 Stochastic few-shot 是一种用于生成新样本的机器学习方法，旨在通过少量的示例来生成具有多样性和创新性的数据。这种方法特别适用于语言生成任务，如生成提示词、扩充语料库等。在 stochastic few-shot 中，通过使用概率模型来模拟数据的生成过程。通常，该方法利用预训练的语言模型作为生成器，以提供语言生成的基础。然后，通过给定少量的示例输入，例如具有特定属性或特定上下文的样本，该方法通过采样和重组模型的内部表示来生成新的样本。该方法的优势在于，它能够利用有限的示例来生成更多样的数据，从而提高数据的多样性和丰富性。这对于训练模型、进行评估和测试以及

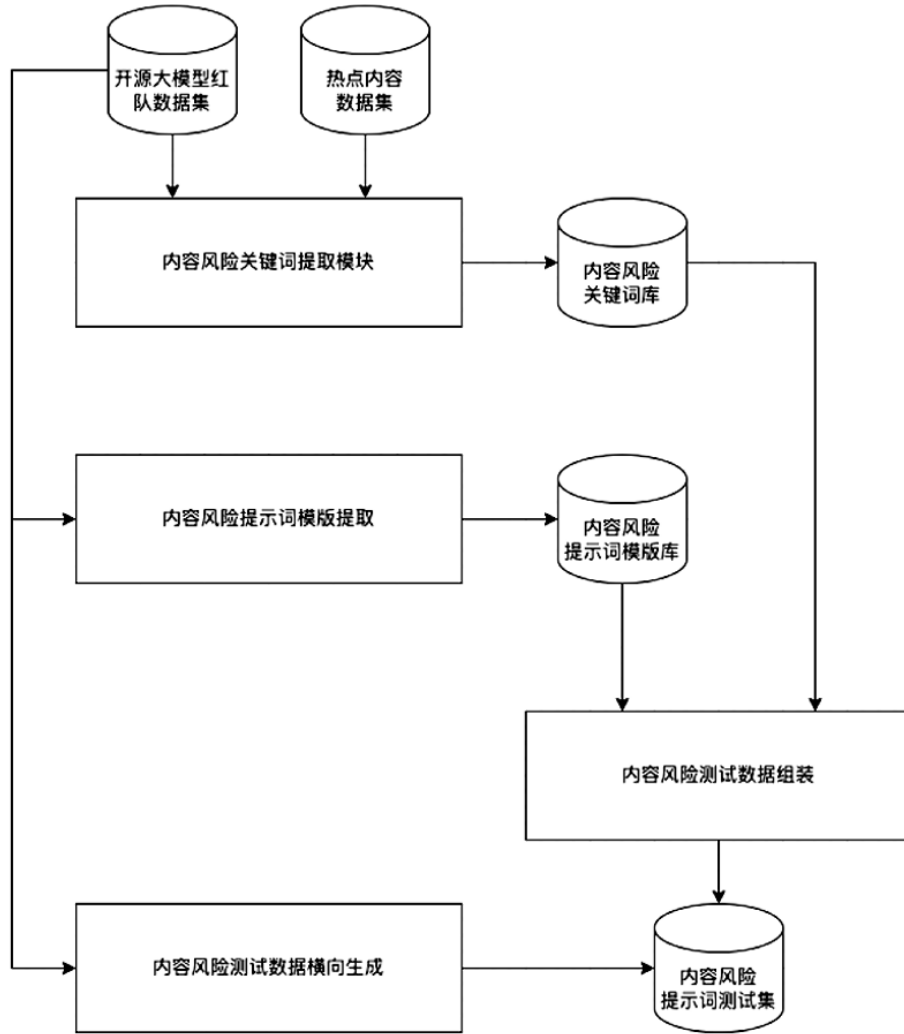


扩充语料库等任务非常有用。此外，stochastic few-shot 方法还可以用于探索模型在不同条件下的生成能力，帮助发现模型的潜在弱点和漏洞。

- 风险话题收集：我们定期从外部数据源收集相关的风险话题文本，并将其纳入我们的库中。这些数据源可能包括社交媒体、新闻报道、论坛讨论等。通过不断更新和丰富话题文本，我们可以确保风险内容评测数据的时效性和多样性。
- 风险关键词提取：我们利用自然语言处理技术和图计算技术，对外部的风险文本进行处理，提取其中的关键词，并构建风险关键词图谱。这个图谱可以帮助我们更好地理解并组织风险内容的关联性，为后续的评测和分析提供基础。
- 风险提示词模板生成：基于高风险的提示词，我们通过自然语言处理算法提取相应的模板。这些模板可以包含语法结构、词汇选择和上下文信息等。然后，借助 stochastic few-shot 的方式，我们利用外部的语言模型横向生成新的测试提示词，以丰富测试数据集的内容。

整体架构如下图：





开源大模型蓝军数据集和热点内容数据通过内容风险关键词提取模块提取筛选出风险程度较高的关键词库。同时开源大模型通过内容风险提示词模版提取模版，提取并横向生成大量提示词模版存入库中；通过内容风险测试数据横向生成模块直接生成提示词测试集数据存入数据库中。另外一部分提示词测试集通过关键词库和模版词库的信息组装后形成完整提示词数据存入库中。

通过以上的模块，我们的框架能够自动生成具有多样性和丰富性的风险内容评测数据。这样的自动生成方法能够降低人工成本，扩大测试集规模，提高测试



数据的多样性和覆盖度，并能够根据需求快速引入新内容。这种框架可以有效支持对大模型的风险评估和安全性测试。

5.2.2 自动化大模型回复风险标注

服务生产了海量风险内容评测数据后，我们将评测数据输入被测大模型，获得大模型的对应回答。我们需要检测这些对应回答的风险情况，并汇总整体的回答内容风险得到被测大模型的整体风险情况。

对于海量大模型输出结果做人工标注需要较大成本，因此我们探索一种可扩展的检测架构，支持自动化地完成回答内容准确快速的风险监测。风险内容检测框架包括多种方法并行，包括：

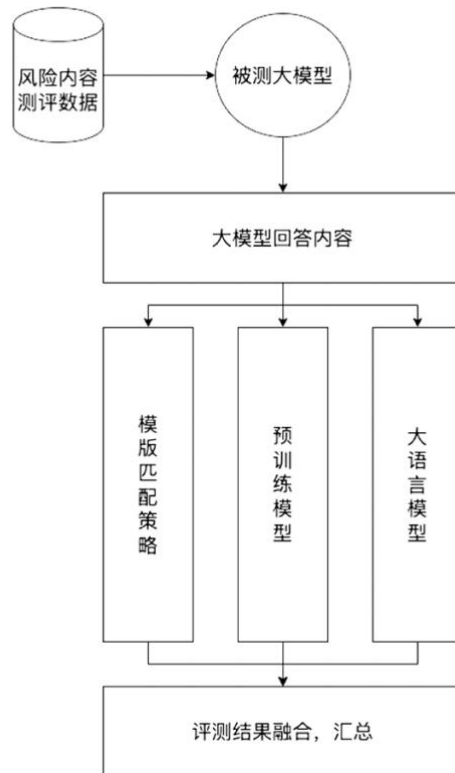
- 模版匹配策略：多数大模型在检测到内容存在风险时，会使用固定的格式生成回答内容，如：
 - 抱歉.....
 - 作为一个人工智能语言模型.....
 - 对不起.....

通过模版记录这类固定格式的回复，可以快速豁免回答内容是否存在风险。

- 预训练模型：使用一些预训练的语言模型，比如 Bert, Roberta, 或大语言模型通过 lora, p-tuning 等技术做微调，并人工标注一批回答与对应的风险情况，将标注数据用来对预训练模型做微调，可以实现通过这些模型对回答内容的风险预测。



- 大语言模型：评测内容探索采用多个大模型辅助标注方式快速、自动化的实现结果的评估。大模型输出的回答的评估方法借鉴了业界先进的实践经验和提示工程技术，将恶意问题提示语句和被测试大模型的对应输出通过模版组装成评估提示语句，并将评估提示语句输入多个评估达模型，获得评估结论。最终整合各个风险维度维度来自各个大模型的评估结论，输出被测试大模型整体的风险情况。
- 评测结果融合汇总：我们需要通过一个融合汇总模块，将来源于模版匹配策略，预训练模型，大语言模型对被测回答内容的风险情况输出做融合。这其中涉及到多种数据类型的转换，包括布尔值和文本数据。将模型输出转换为统一的布尔值格式后，我们设计了一个 bagging 模型汇总三个模型的判断结果，模型的权重可以根据系统配置自定义修改。



使用大语言模型分析文本的风险情况，需要通过一个提示语句模版将被测大模型回答内容与要求大语言模型分析风险的具体需求组装在一个长文本中，组装评估提示语句的模版需要克服以下难点：

- 内容识别错乱：存在对提问与回答的内容理解错误，导致误召回。
- 是非判断的命题界定不清楚：通过该方法发现的风险 case，大量误召回了回答的内容立场正确，但是涉及到了不安全的领域。
- 输出格式不固定：自然语言输出结论无统一格式，后续自动化分析困难。

我们采用了一些提示词工程方法，包括 Chain of Thought，梳理了风险分析的范式，加强了大语言模型通过文本内容得到正确风险情况的能力。

综上所述，通过多维度的模型预测大模型回答的潜在风险，有助于我们快速准确发现生成测试数据中的风险内容。高质量的蓝军攻击样本将有助于大模型安全防御模块开发人员更好地开展下一步针对性优化工作。

5.2.3 大模型安全评测框架

本检测方案最终会输出一份详细评测报告，内容包括评测方法、评测测试集、评测指标等数据；其中评测量化指标参考如下：

- 监测覆盖度，测试集数据不少 xx；
- 新型风险黑词感知能力，日均新发现黑词 xx、构建测试数据 xx；
- 监测发现风险数量不低于 xx；



报告中还会根据实际发现的风险，给出相应的改进建议，以达到帮助大模型内容风控系统升级，促进大模型生态健康发展的目的。

本评测框架通过自动化评测内容生成，自动化大模型回答评估，能够快速准确地量化大模型在多个内容安全维度的风险情况。通过定期的评测执行，能够实现对大模型内容安全能力的实时追踪，快速定位大模型在内容安全潜在的潜在漏洞，全面保障大模型的安全内容输出能力。

6. 总结与展望

在本白皮书中，我们深入探讨了大模型安全风控的多个关键方面，从数据安全与隐私保护到模型保护、内容合规、业务安全风控，以及蓝军评测，旨在为大模型领域的从业者、企业和用户提供全面的指导与建议。通过对安全风险的认识和解决方案的探讨，我们可以更好地应对挑战，确保大模型的可信度和应用价值。

6.1 总结成果与贡献

本白皮书强调了大模型安全风控的重要性，并从多个角度提供了解决方案。我们深入讨论了数据安全与隐私保护的策略，模型保护的技术手段，内容合规的方法，业务安全风控的实践，以及蓝军评测的意义。这些讨论不仅帮助相关方了解现有的安全挑战，还为他们提供了实际可行的方法，以确保大模型的安全性和可信度。



6.2 展望未来发展

然而，大模型安全风险领域仍然充满了挑战和机遇。随着技术的不断进步，新的安全风险和威胁可能会不断涌现。因此，我们需要保持警惕，并始终保持创新和适应能力。在未来，我们预见以下几个发展方向：

- 1. 跨界合作加强：** 由于大模型领域的安全问题涉及多个领域，跨界合作将变得更加重要。技术研究人员、法律专家、政策制定者等需要共同合作，以制定更全面的安全解决方案。
- 2. 持续创新和技术进步：** 安全风险不断演变，我们需要不断创新和提升安全技术。新型的防御手段、对抗攻击策略以及加密技术的应用都将是未来的研究方向。
- 3. 法律法规的完善：** 针对大模型的安全和隐私问题，需要更多的法律法规来保护用户的权益和数据隐私。政府和监管机构需要积极参与，制定适应快速变化的技术环境的法规。
- 4. 安全意识的提高：** 对于大众和企业来说，安全意识的提高是防范安全威胁的关键。教育培训、信息宣传等方式都有助于提高用户和从业者的安全意识。

6.3 结语

大模型的安全风险和挑战在不断变化，需要我们的共同努力来解决。通过合作、创新和持续的努力，我们可以建立起一个安全、稳定和可信赖的大模型



生态系统。本白皮书所提供的指导和建议将为这一目标的实现提供有力支持。

在未来的道路上，我们有信心克服各种挑战，推动大模型技术更加安全和可持续地发展。



参考文献

- [1] Chen Qu, et al. "Natural Language Understanding with Privacy-Preserving BERT" Conference on Information and Knowledge Management. 2021.
- [2] Li, Dacheng, et al. "MPCFormer: fast, performant and private Transformer inference with MPC." arXiv preprint arXiv:2211.01452(2022).
- [3] Chen, Tianyu, et al. "THE-X: Privacy-Preserving Transformer Inference with Homomorphic Encryption" arXiv preprint arXiv:2206.00216(2022).
- [4] Mengxin, Zheng, et al. "Primer: Fast Private Transformer Inference on Encrypted Data." Design Automation Conference. 2023.
- [5] Hao, Meng, et al. "Iron: Private Inference on Transformers." Advances in Neural Information Processing Systems. 2022.
- [6] "PUMA: Secure Inference of LLaMA-7B in Five Minutes"
- [7] Liu, Xuanqi, et al. "LLMs Can Understand Encrypted Prompt: Towards Privacy-Computing Friendly Transformers." arXiv preprint arXiv:2305.18396(2023).
- [8] “邪恶版” ChatGPT 出现：毫无道德限制，专为“网络罪犯”而生？
<https://mp.weixin.qq.com/s/YS4GMUgZPfwmg1MXBCRImA>
- [9] <https://new.qq.com/rain/a/20230613A03W0900>
- [10] 《生成式人工智能服务管理暂行办法》



V分



谢谢观看

相关咨询可以扫码联系

