LLM Modules		Risks	Description of Risks	Sub-categorized Topics
Output Module		Harmful Content	The LLM-generated content sometimes contains biased, toxic, and private information.	Bias     Toxicity     Privacy Leakage
		Untruthful Content	The LLM-generated content could contain inaccurate Information.	<ul> <li>Factuality Errors</li> <li>Faithfulness Errors</li> </ul>
		Unhelpful Uses	Improper uses of LLM systems can cause adverse social impacts.	<ul> <li>Academic Misconduct</li> <li>Cyber Attacks</li> <li>Copyright Violation</li> <li>Software Vulnerabilities</li> </ul>
Toolchain Module		Software Security Issues	The software development toolchain of LLMs is complex and could bring threats to the developed LLM.	<ul> <li>Programming Language</li> <li>Deep Learning Frameworks</li> <li>Software Supply Chains</li> <li>Pre-processing Tools</li> </ul>
		Hardware Vulnerabilities	The vulnerabilities of hardware systems for training and inferences bring issues to LLM-based applications.	<ul> <li>Network Devices</li> <li>Memory and Storage</li> <li>GPU Computation Platforms</li> </ul>
		Issues on External Tools	The external tools (e.g., web APIs) present trustworthiness and privacy issues to LLM-based applications.	<ul> <li>Factual Errors Injected by External Tools</li> <li>Exploiting External Tools for Attacks</li> </ul>
Language Model Module		Privacy Leakage	The model is trained with personal data in the corpus and unintentionally exposing them during the conversation.	<ul> <li>Private Training Data</li> <li>Memorization in LLMs</li> <li>Association in LLMs</li> </ul>
		Toxicity and Bias Tendencies	Extensive data collection in LLMs brings toxic content and stereotypical bias into the training data.	<ul> <li>Toxic Training Data</li> <li>Biased Training Data</li> </ul>
		Hallucinations	LLMs generate nonsensical, unfaithful, and factual incorrect content.	<ul> <li>Knowledge Gaps • Noisy Training Data • Defective Decoding Process</li> <li>False Recall of Memorized Information • Pursuing Consistent Context</li> </ul>
		Model Attacks	Model attacks exploit the vulnerability of LLMs, aiming to steal valuable information or lead to incorrect responses	<ul> <li>Extraction Attacks</li> <li>Inference Attacks</li> <li>Evasion Attacks</li> <li>Novel Attacks on LLMs</li> </ul>
Input Module		Not-Suitable-for-Work (NSFW) Prompts	Inputting a prompt contain an unsafe topic (e.g., not-suitable-for-work (NSFW) content) by a benign user.	<ul> <li>Insult</li> <li>Unfairness</li> <li>Sensitive Politics</li> <li>Mental Health</li> </ul>
		Adversarial Prompts	Engineering an adversarial input to elicit an undesired model behavior, which pose a clear attack intention.	<ul> <li>Goal Hijacking</li> <li>One-step Jailbreaks</li> <li>Prompt Leaking</li> <li>Multi-step Jailbreaks</li> </ul>