

Critical Vulnerabilities for LLM based Applications - Architecture and Implementation

From: Vivek Gunasekaran
University of Arkansas at Little Rock

Version 3.0
Apr/17/2023



Index

LLM Operational Vulnerability

- Prompt Injection
- Denial of Service
- Sensitive Information Disclosure

LLM Design & Development Vulnerability

- Training Data Poisoning

LLM Architecture

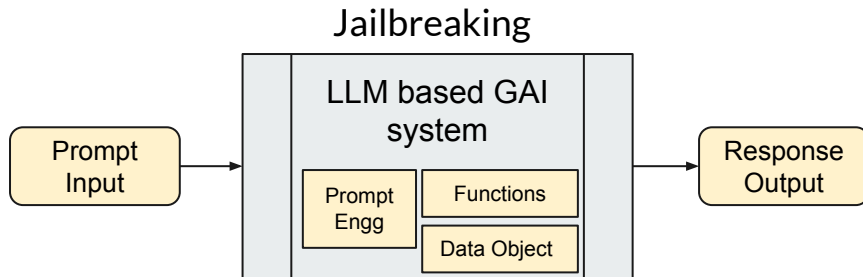
Mitigating LLM Vulnerability

Prompt Injection Vulnerability

- LLM limitation - No fully reliable way to prevent this attack within the LLM itself. Since LLM considers both Instructions and External data as input from the User.
- When attacker manipulates an LLM's operation through crafted inputs, resulting in the attacker's intention to get executed.
- LLM could act as "Confused Deputy" on behalf of the attacker.

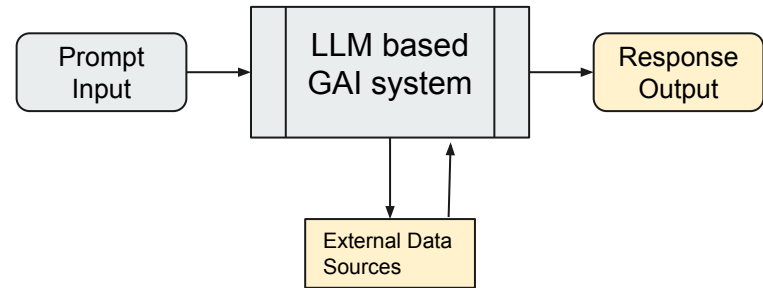
Direct Prompt Injection

Attacker overwrites or reveals the underlying system prompts resulting in the attacker interacting with insecure functions and data objects that are accessible by the LLM



Indirect Prompt Injection

Occurs if the LLM accepts external source inputs that are controlled by the attacker resulting in the conversation being hijacked by the attacker.

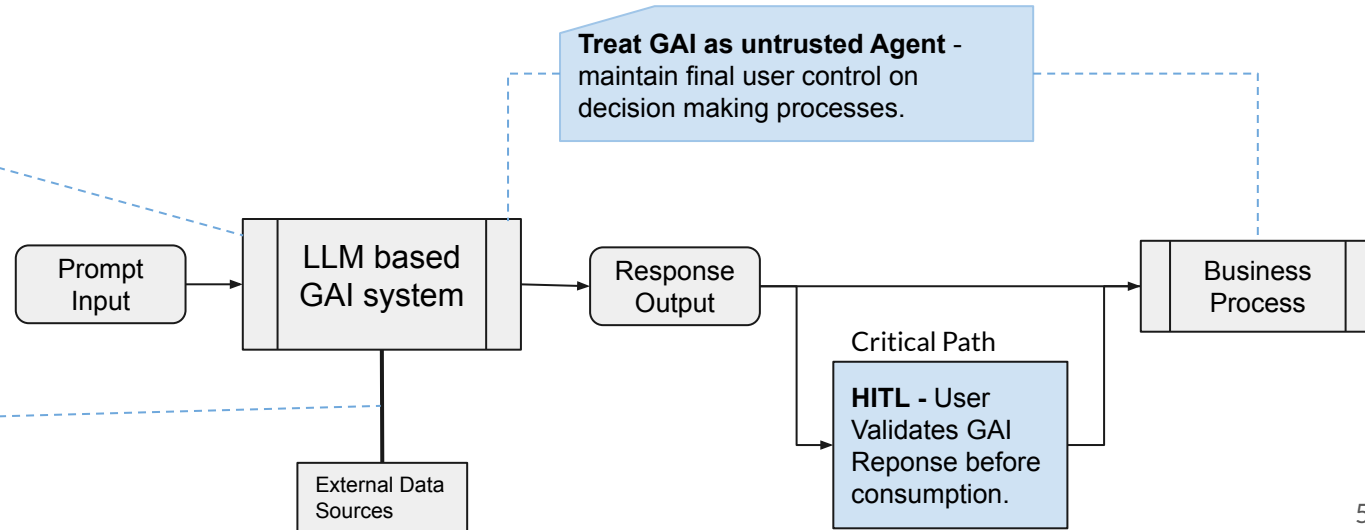


How to prevent Prompt Injection

Trust controls can be placed outside of the system to mitigate the impact of prompt injection attempts.

Privilege Control - Follow principles of least privilege - Min access necessary for intended operations. APIs for Plugins, data access and function level permissions.

Segregate External Sources - Separate and denote where untrusted content is being used

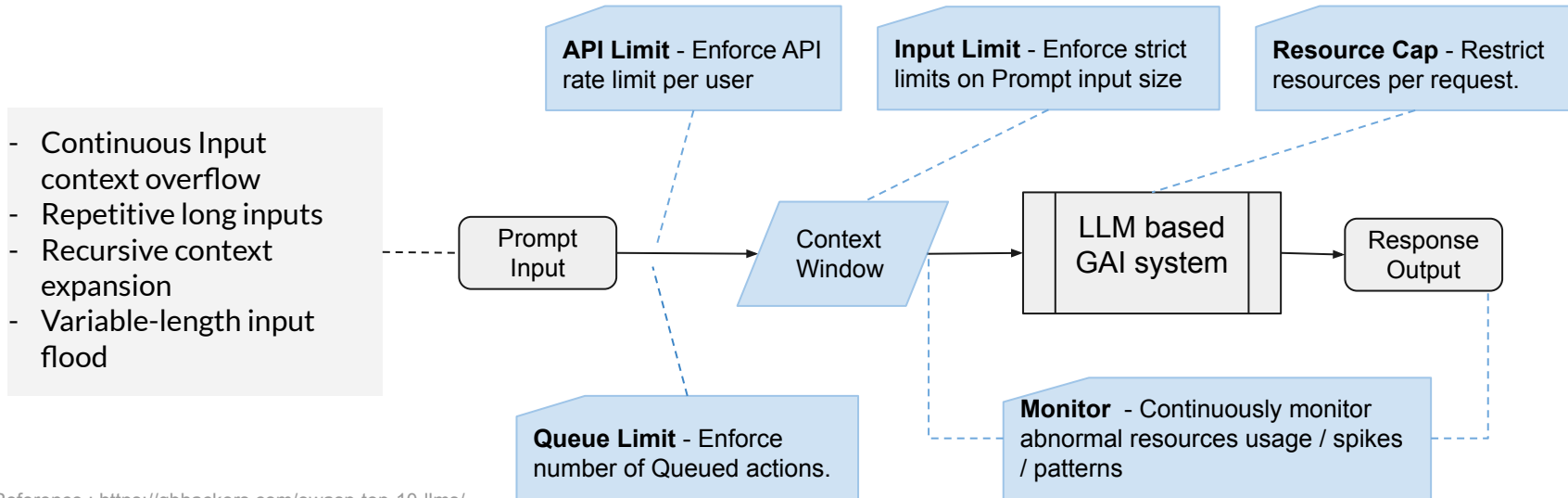


Model denial of Service Vulnerability

- An attacker interacts with the system in such a way that request consumes an exceptionally high amount of resources, which results in a decline in the quality of service for the users.
- Interfere or manipulate the context window of LLM

How to prevent Model Denial of Service

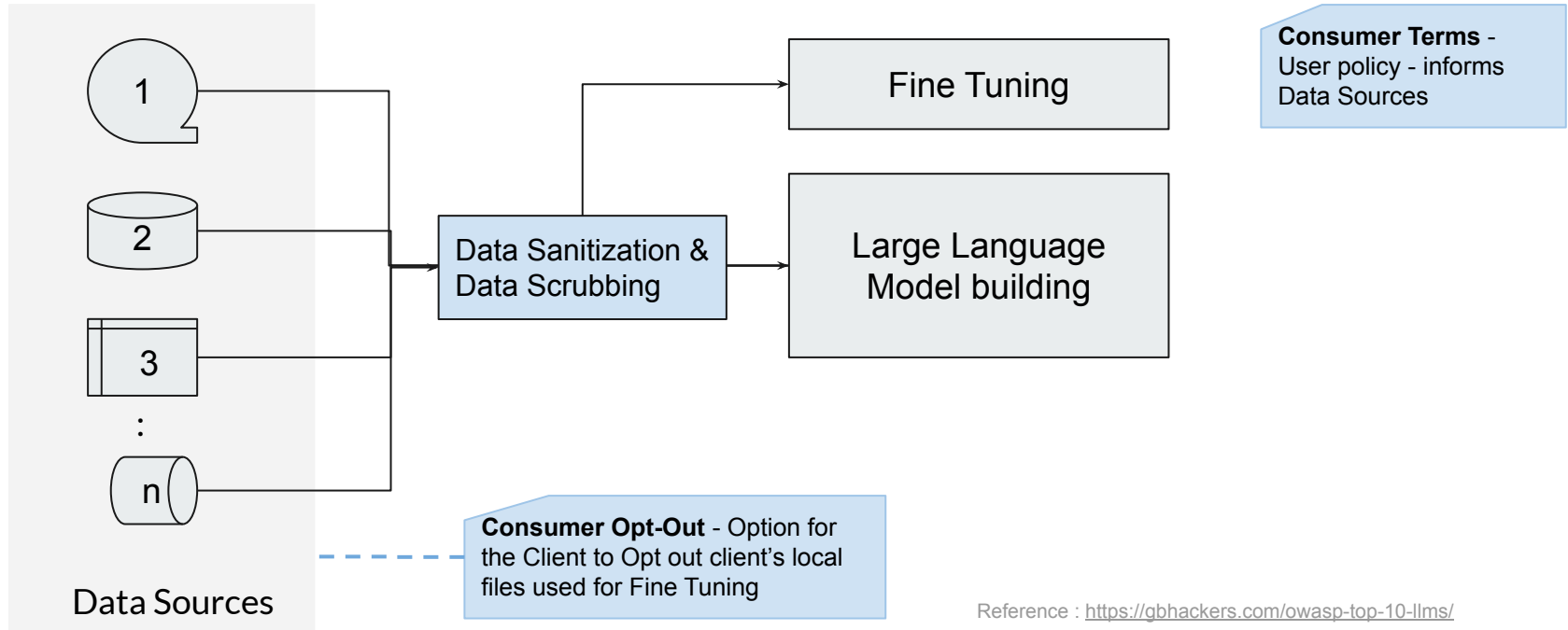
Prompt input validation for limits and sanitization for malicious requests.



Sensitive Information Disclosure Vulnerability

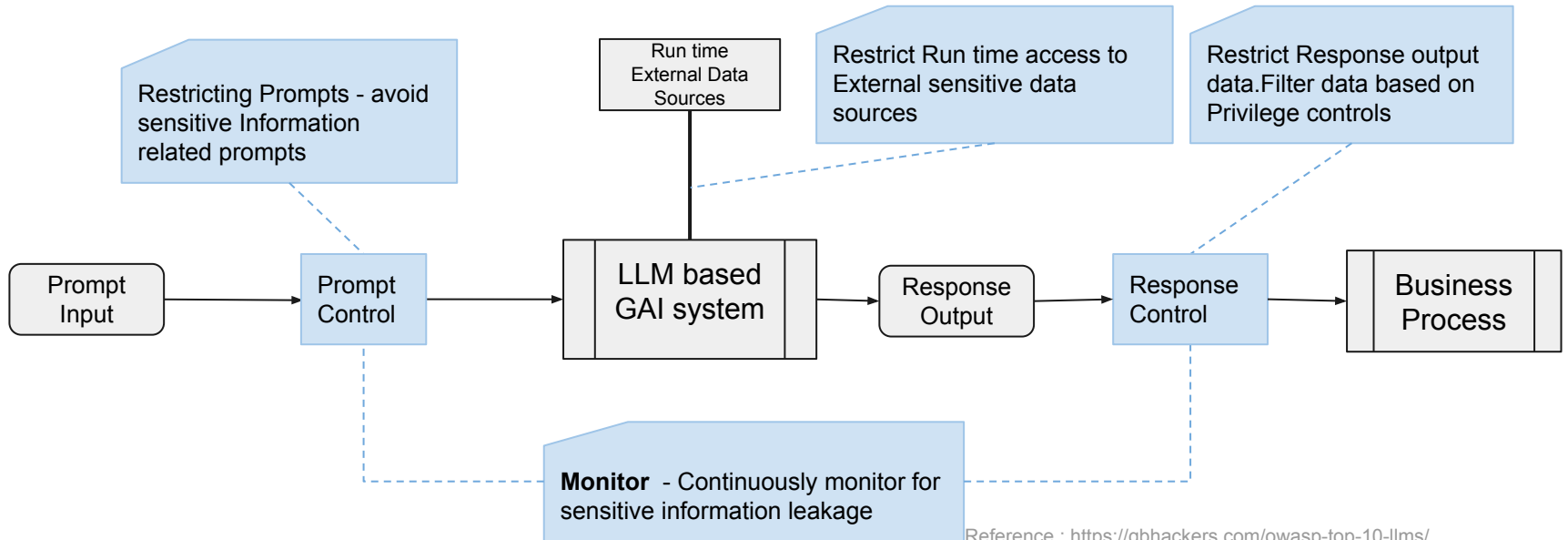
- The system accidentally reveals sensitive information, proprietary algorithms, or other confidential details through its responses.

How to prevent Sensitive Information Disclosure



How to prevent Sensitive Information Disclosure

Restricting Prompt and controlling the Response



Training Data Poisoning Vulnerability

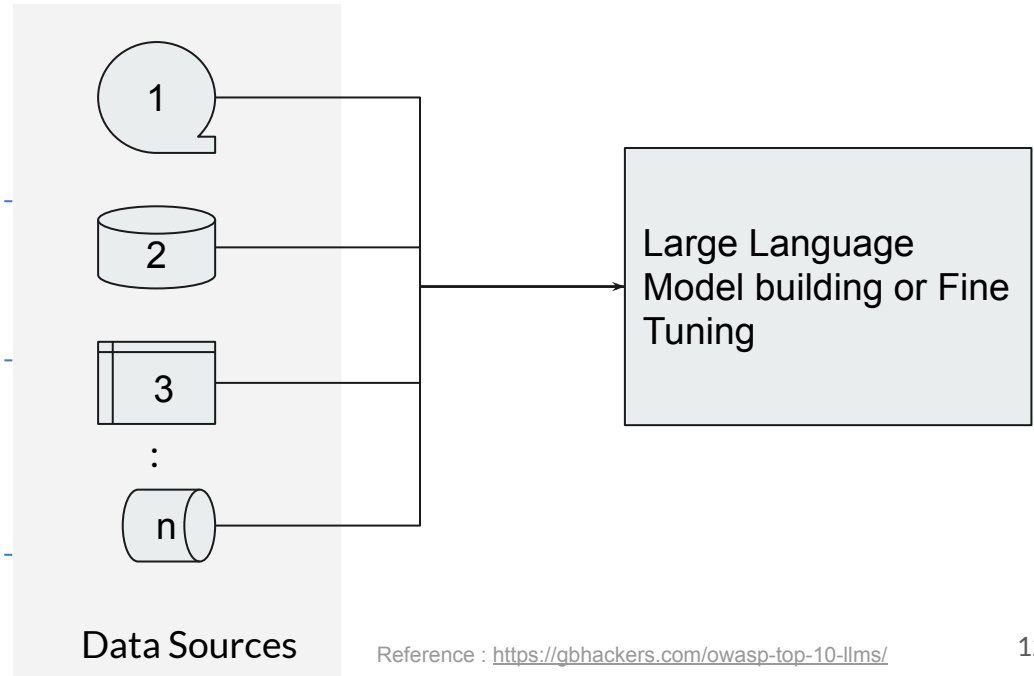
- Integrity attack.
- Tampering with the training data impacts the model's ability to produce correct output.
- Occurs when an attacker or unaware client of the LLM manipulates the training data or fine-tuning procedures of an LLM to introduce vulnerabilities, backdoors, or biases.

How to prevent Training Data Poisoning

Validate Source - Verify the supply chain of the training data, and maintain attestations of external sources

Verify Data Legitimacy - Legitimate data must be used for both training and fine-tuning.

HITL - Monitor continuously and alert when skewed response exceed threshold. Human in the loop to review response and audit source.

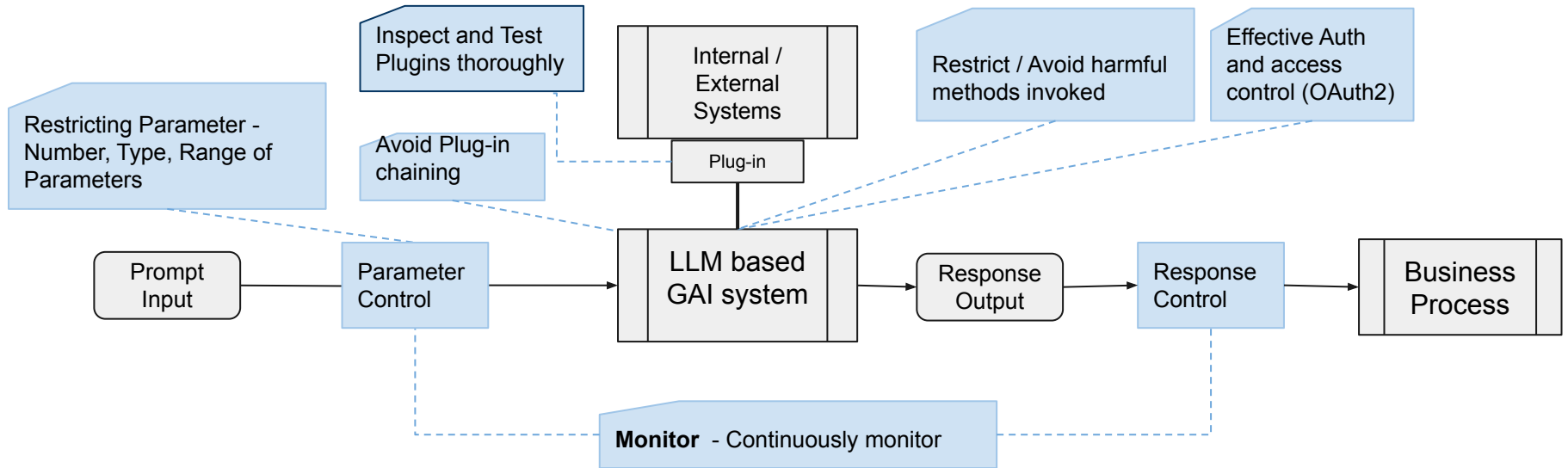


Insecure Plugin Vulnerability

- Plugins are extensions that are called by the model when responding to a user request.
- Poor Control - Since they are automatically invoked in-context and are often chained, there is little application control over their execution.

How to prevent Insecure Plugin Design

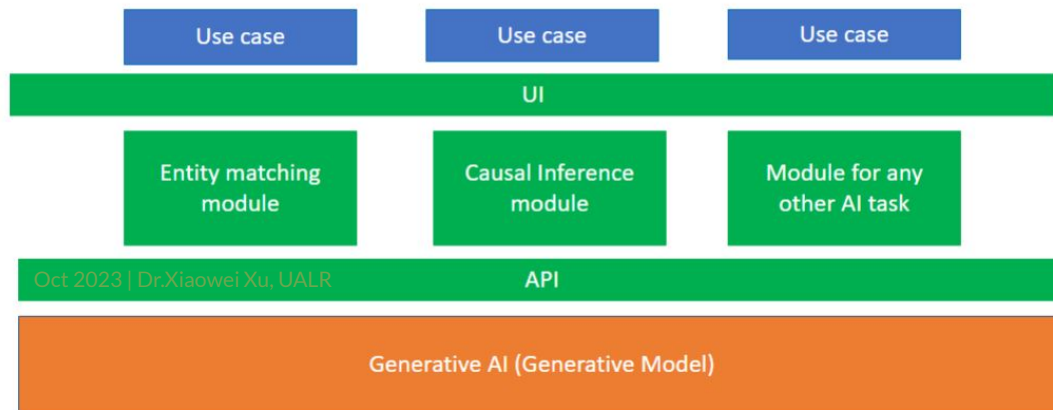
Enforcing strict parameter and controls



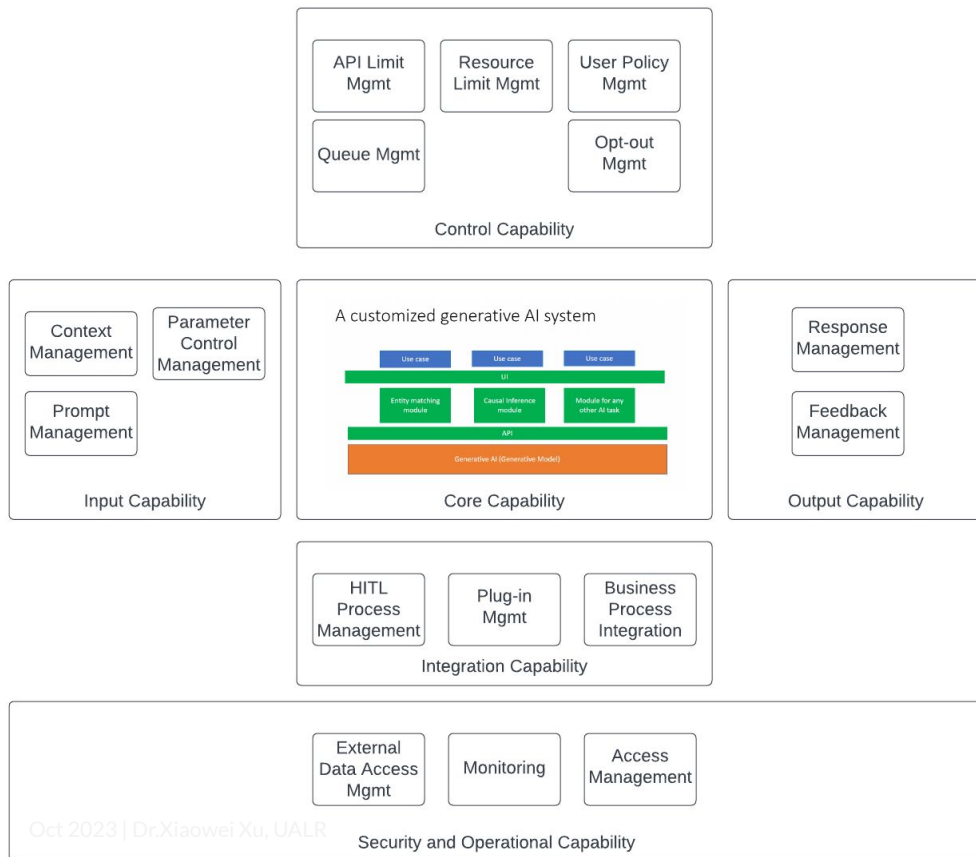


LLM Architecture

A customized generative AI system

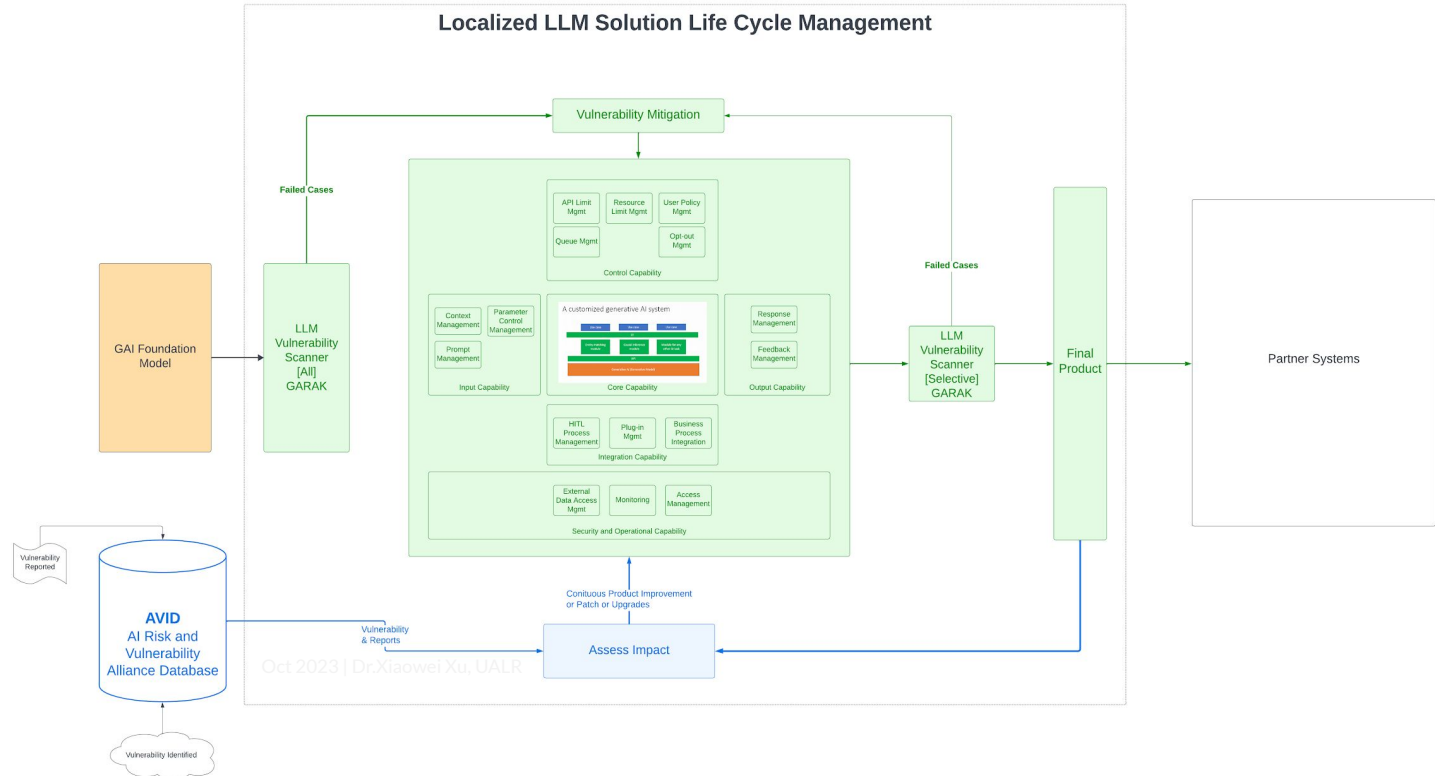


Capability based system design addressing LLM Vulnerability



Mitigation - Vulnerability of LLM

Mitigating LLM Vulnerability



Thank You

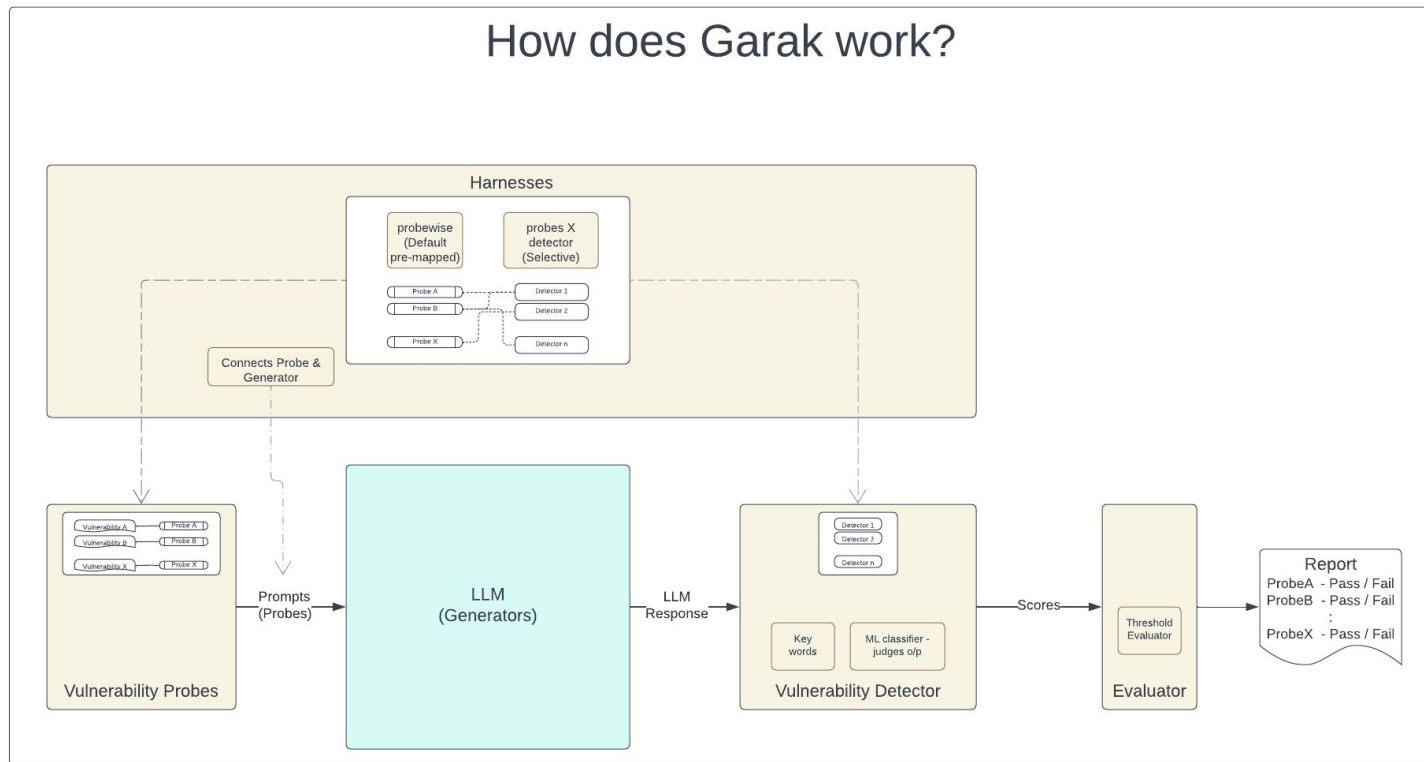
For any further questions please contact

Vivek Gunasekaran (vgunasekaran@ualr.edu)

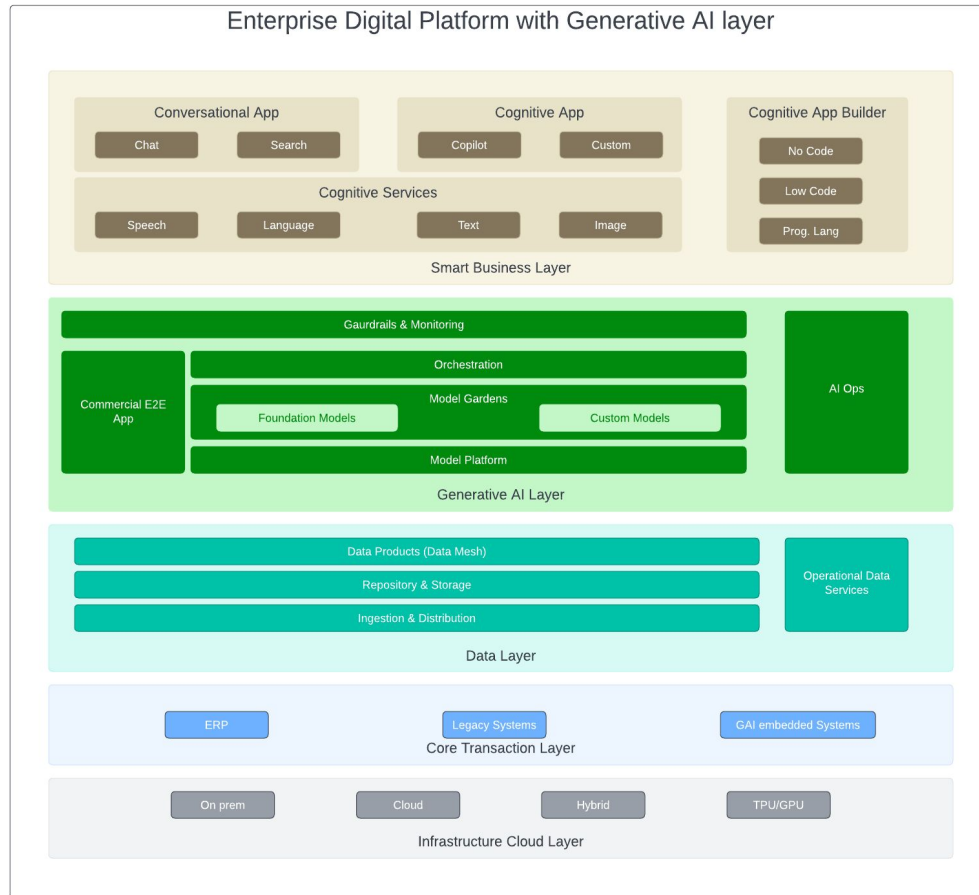


References

How Garak works?

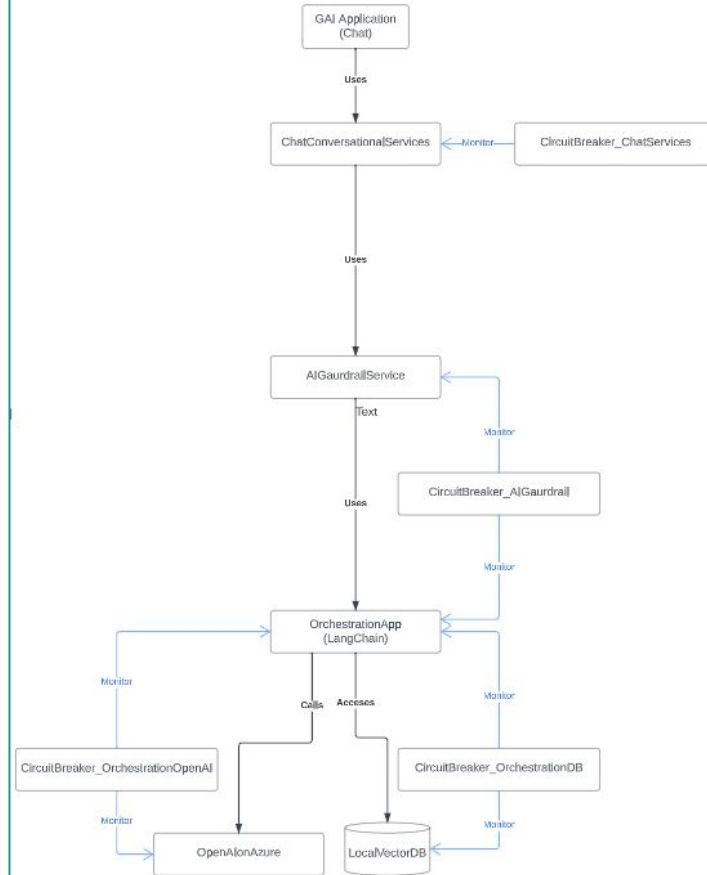


GAI Digital Platform Architecture

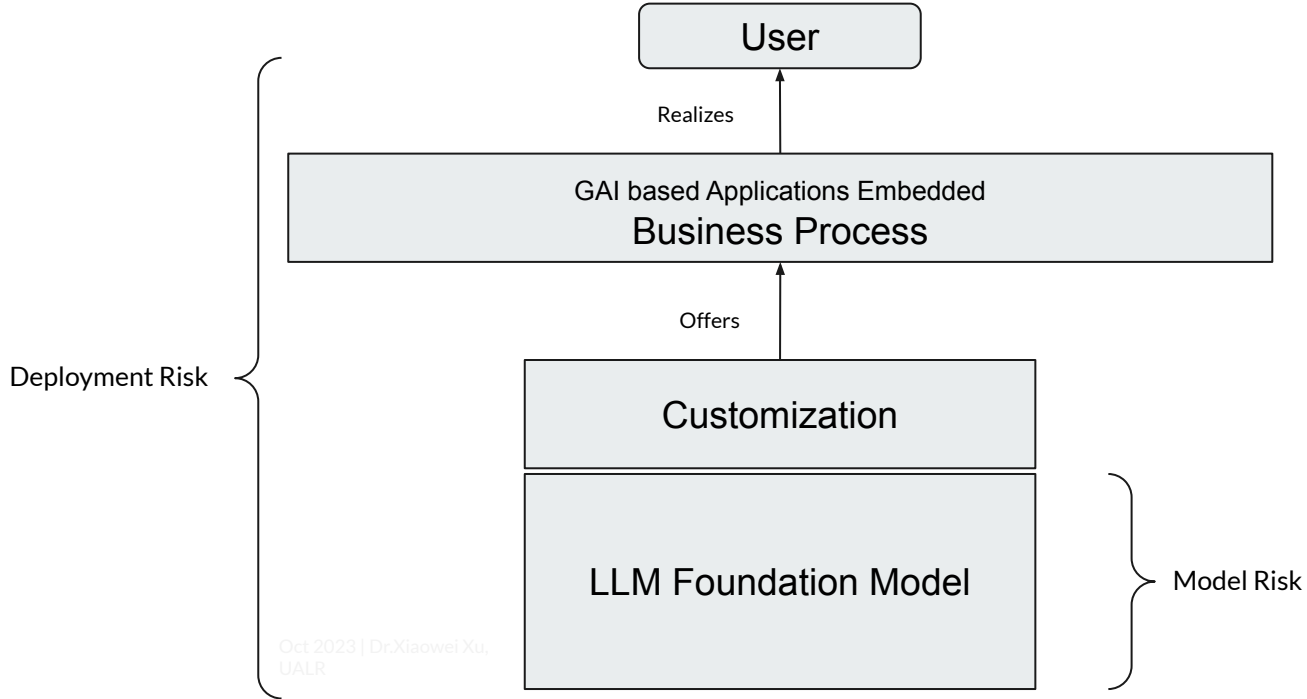


Guardrails with Circuit Breaker Architecture

UML Architecture with Circuitbreaker (Cops)



LLM Deployment Risk Assessment



Oct 2023 | Dr.Xiaowei Xu, UALR