

LLM Guardrails

Reza Fayyazi

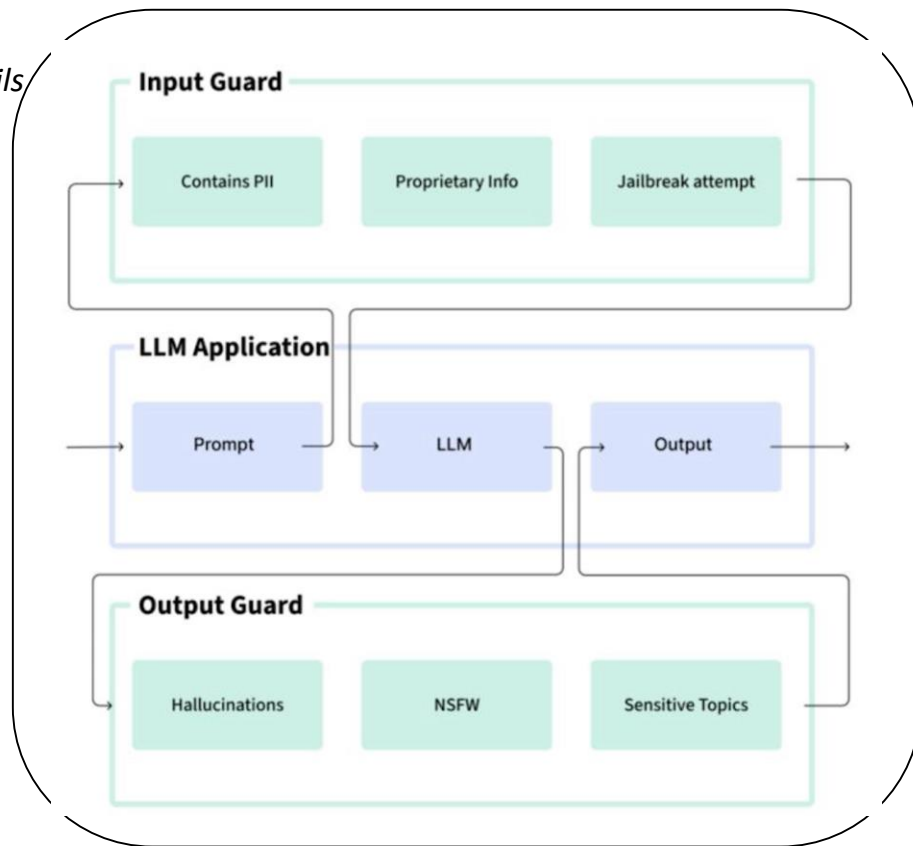
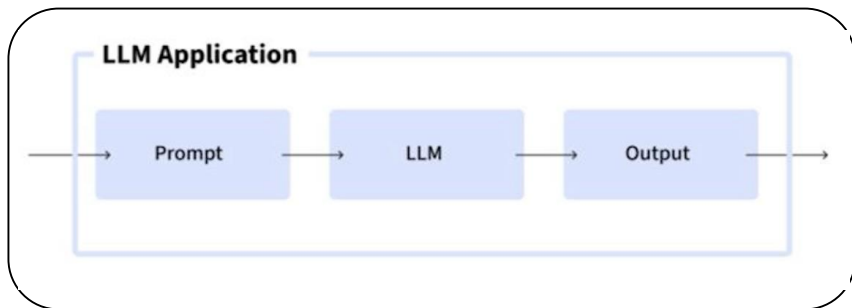
What are guardrails?

- A *guardrail* is a secondary check or validation around the input or output of an LLM model
- Input validation guardrails include:
 - Prevent the injection of code or adversarial prompts to manipulate the model
 - Check for phrasing in the input that may trigger biased or harmful responses
 - Check for Personal Identifiable Information (PII) leakage
- Output validation guardrails include:
 - Check for hallucination by cross-referencing against trusted sources
 - Check for PII leakage
 - Check for sensitive topics

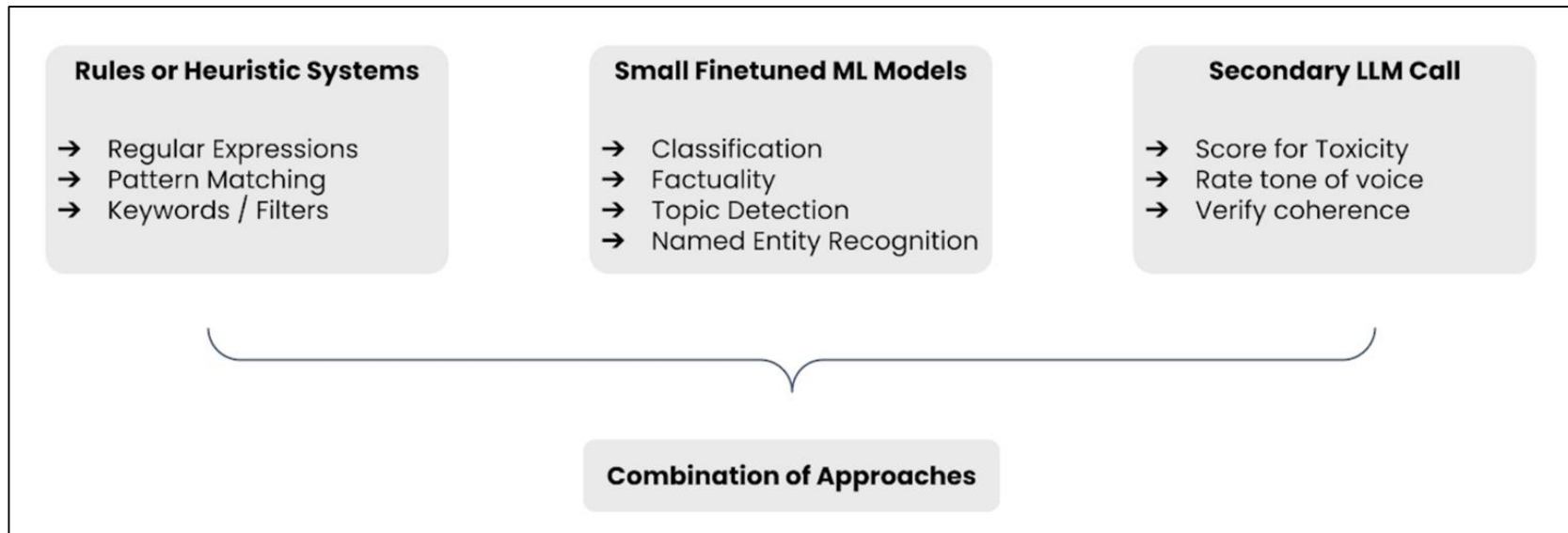
Where are guardrails applied in LLMs?

With Guardrails

Without Guardrails

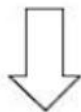


How are guardrails implemented?



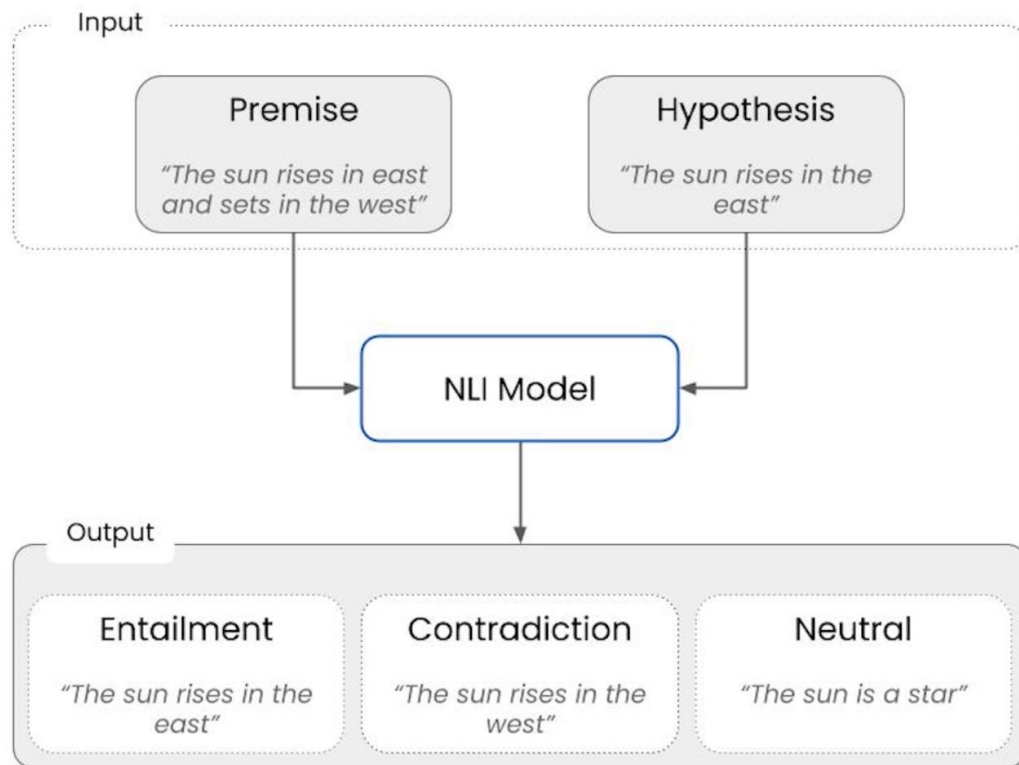
Hallucinations are due to Lack of Groundness

Grounded AI Response



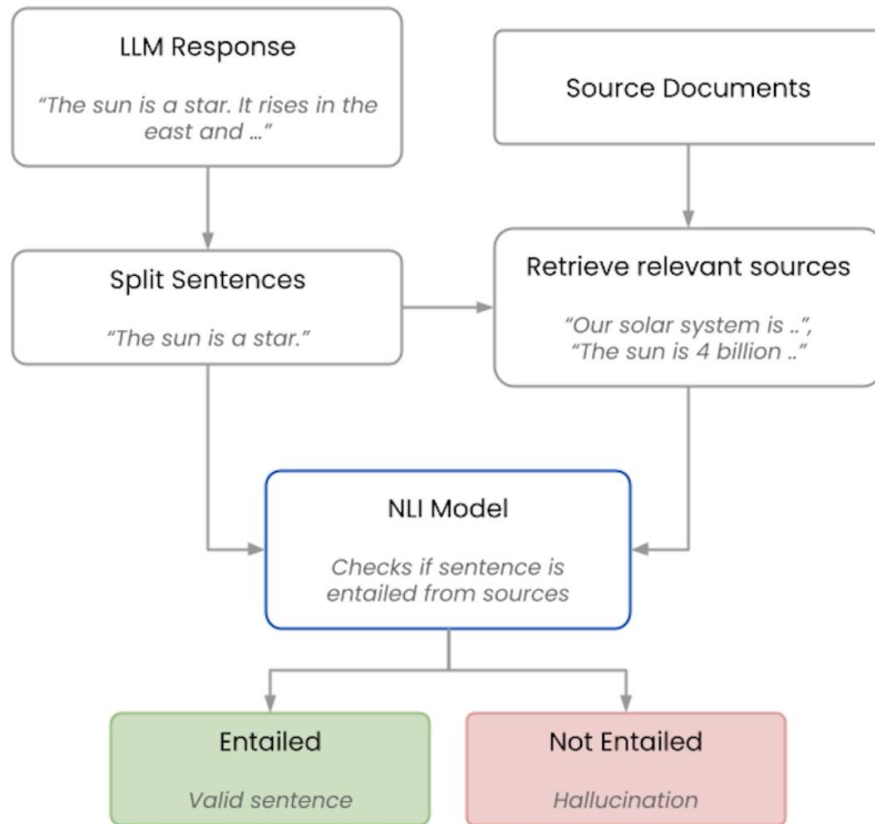
All claims made by the LLM
are explicitly supported by the input context rather
than being fabricated or recited from training

Natural Language Inference (NLI)

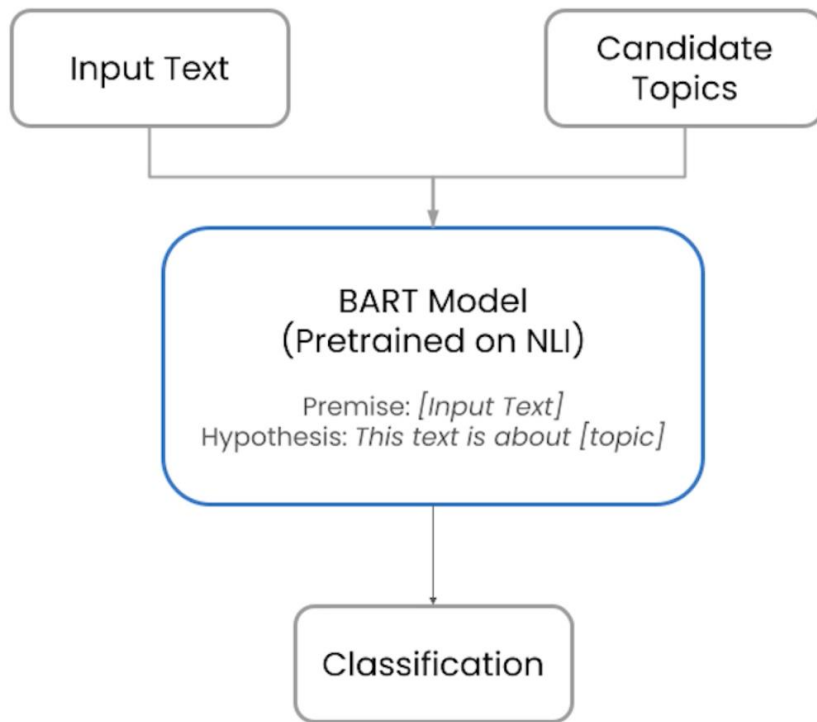


Using NLI for Detecting Hallucinations

- Using cosine similarity to check the splitted sentences from the response with splitted sentences from the retrieved sources



Zero-Shot Topic Classification



What is PII?

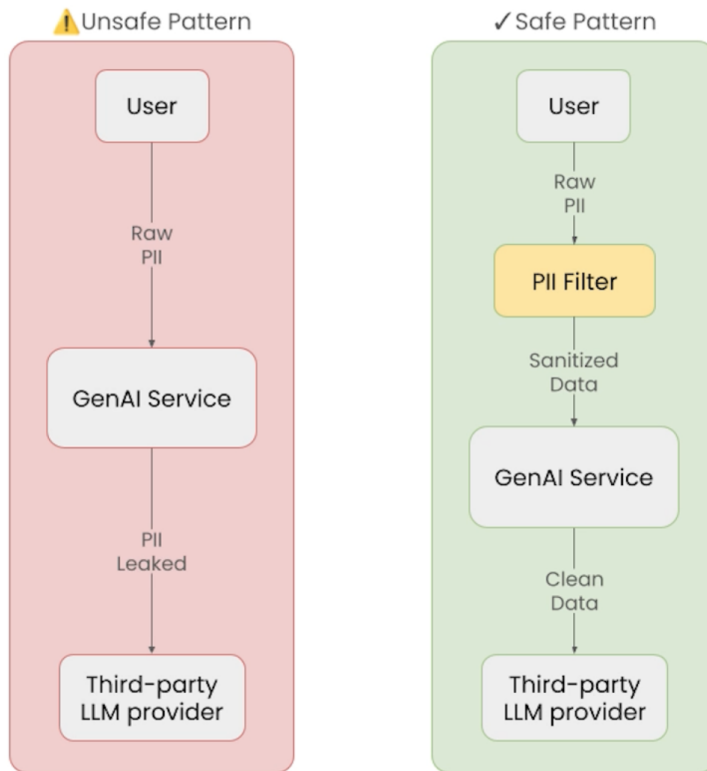
Personal Identifiable Information

Direct identifiers	Name, SSN, Email, ...
Indirect identifiers	Location, Demographics, ...
Sensitive data	Health records, Financial info, ...

LLM Data Privacy Risks:

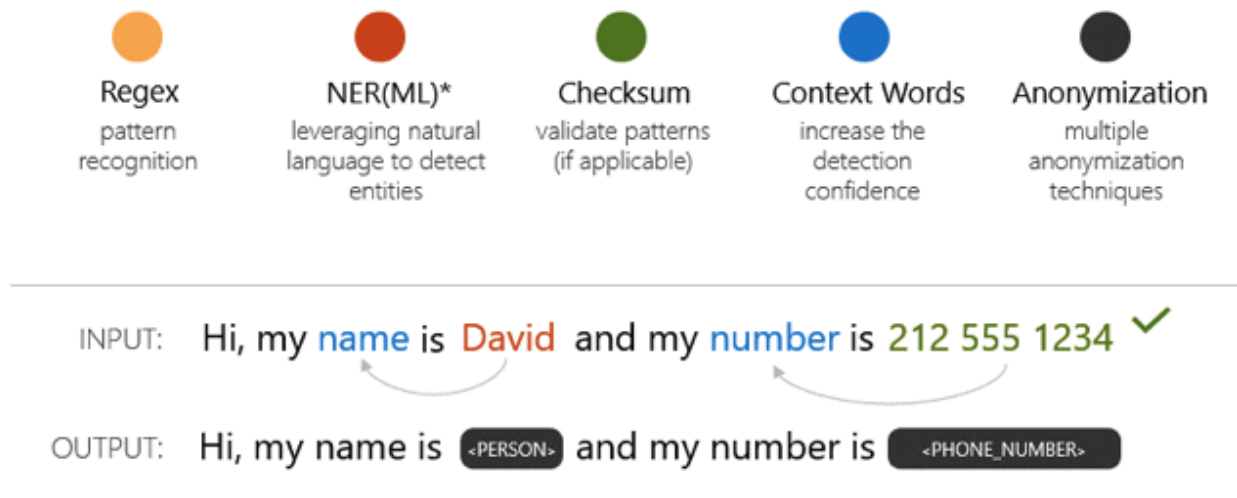
- Third-party processing exposure
- Potential data retention by providers
- Risk of training data contamination
- Limited control over data handling

PII Handling with LLMs



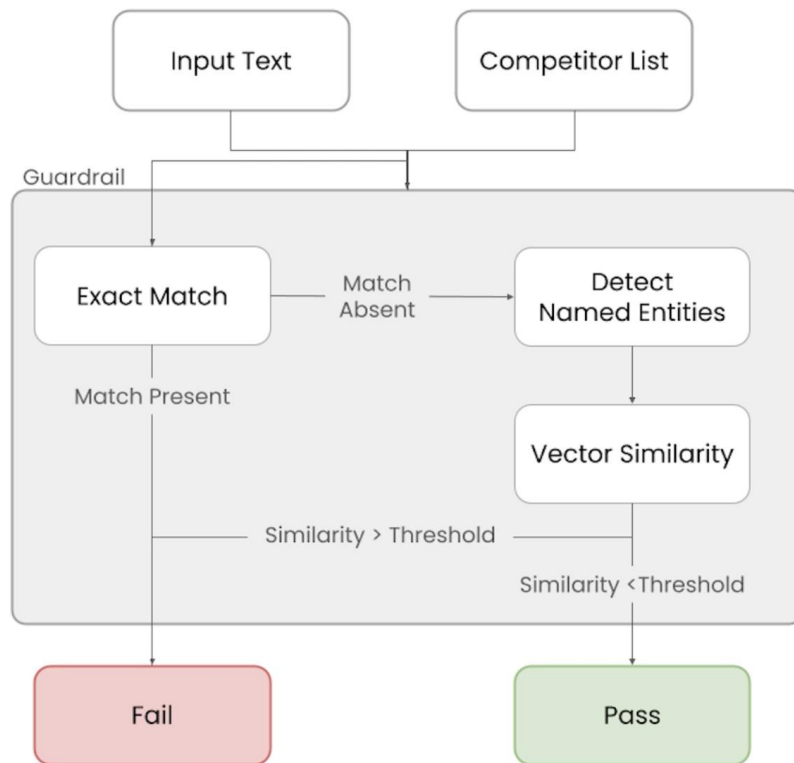
Microsoft Presidio

Presidio Detection Flow



*NER – Named Entity Recognition

Sensitive Word Detection



Thank you!

Reza Fayyazi
rf1679@rit.edu