

Implementation Classic block/allow-lists ✓ Exact text matching and regex Machine learning methods ✓ Semantic search using text embeddings ✓ Classifiers, such as BERT ✓ LLM-as-judge **Prompt** Large

Language

Model

(LLM)

Rule #1:

Treat the LLM

as untrusted

INPUT OUTPUT

Consulting

Inspect / Sanitize

- ✓ Contextual, use-case specific **output** validation
- ✓ Allow only **expected content** types/formats
- ✓ Block potentially dangerous content like URLs, JavaScript, or markdown images, unless explicitly validated through an allow-list.
- ✓ Ensure **safe rendering** in web apps (contextual encoding, Content Security Policy, ...)
- ✓ Moderate outputs for harm categories like hate speech, violence, self-harm, sexual content, ...

Tool/Function Calls

Access Control

- ✓ Least privilege
- ✓ Downstream checks

Safe APIs

Output

- ✓ Limit function / scope of allowed operations
- ✓ Ensure APIs/functions are not vulnerable themselves (OWASP Top 10)

approved

Instruction / Data Separation

routing to ensure user queries

and are routed accordingly

match intended scope / use-case

Use **Spotlighting** / prompt engineering techniques to help the LLM distinguish instructions from data

- ✓ Data-marking

✓ Multi-turn dialogue