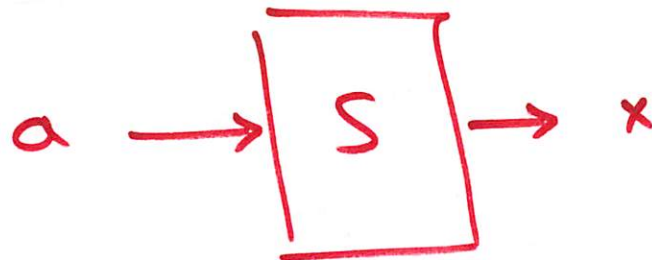


Classical Software vs LLMs

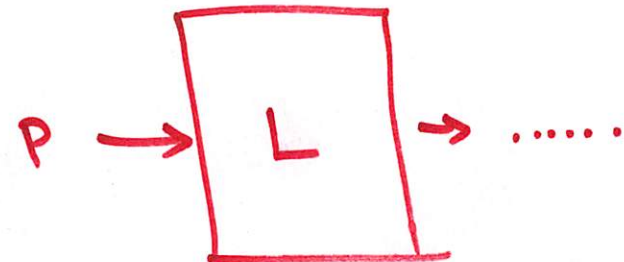
- Classic Software

- Deterministic control flow
- Auditable rules - explicit logic
- Verifiability



- LLMs

- Stochastic decoding & sampling
- Distributed representations across billions of weights - implicit logic
- Empirical tests



Security Risks - the top 10

- Confidentiality

- Sensitive Information Disclosure
- Model Theft

3 promises

- Integrity

- Prompt Injection Attacks
- Insecure Plugins/Output Handling
- Supply Chain Vulnerabilities
- Training Data Poisoning
- Overreliance

- Availability

- Model Denial of Service
- Excessive Agency

Mental model of an LLM

- Untrusted powerful intern
 - Helpful but unpredictable; must be contained, verified, and logged.
- Untrusted, unpredictable, over-confident intern-robot.
 - Capable, fast, tireless but no innate ethics, common sense or boundaries unless explicitly asked for.

The intern vs The robot:

- Knowledge
- Reasoning
- Ethical Boundaries }
- Common Sense }

→ Bias

→ Adversarial Robustness

