**1. Why do Large Language Models (LLMs) require guardrails for safe and responsible use?**

LLMs can generate inaccurate, biased, or unsafe responses because they don't have true understanding or intent. Guardrails such as content filters, policies, and monitoring help ensure outputs remain safe, reliable, and aligned with ethical and organizational standards.

**2. What types of security risks are most common in LLMs?**

LLMs face three primary categories of risks:

- Confidentiality risks - Leakage of sensitive data.
- Integrity risks - Prompt injection or manipulated inputs altering behavior.
- Availability risks - Attacks or misuse that disrupt service.

**3. How do confidentiality, integrity, and availability concerns apply to LLMs?**

- Confidentiality: Models may expose sensitive training data or user information.
- Integrity: Attackers may use prompt injection to bypass rules or introduce false outputs.
- Availability: Adversarial inputs or spam can overload the system, making it unusable.

**4. What are some real-world examples of security risks in generative AI systems?**

- A chatbot accidentally revealed customer PII (confidentiality breach).
- Malicious prompts override system instructions to extract hidden data (integrity breach).
- Spam or adversarial queries cause downtime (availability risk).
- AI-generated harmful or offensive responses are damaging brand reputation.

**5. What layers of risk mitigation can organizations implement to secure LLMs?**

Organizations can implement multi-layer defenses, such as:

- Input filtering - Block malicious prompts.
- Output monitoring - Detect unsafe or biased responses.
- Human oversight - Add checks for sensitive use cases.
- Authentication and rate-limiting - Prevent misuse of APIs.
- Red-teaming and audits - Test models for vulnerabilities.

**6. What key considerations should be kept in mind when designing a secure generative AI solution?**

- Continuous monitoring: Update defenses as threats evolve.
- Responsible design: Balance innovation with safety, fairness, and transparency.