## 1. Attack Summary
Describe the most successful attack vectors used (prompt injection, jailbreak, etc.)

_____

## 2. Defense Measures
List concrete mitigations applied (isolation, regex filter, moderation API...)

_____

## 3. Evaluation
Did the attack still succeed post-patch?  ☐ Yes  ☐ No
Evidence / logs:

_____

## 4. Lessons Learned
Key takeaways & potential next-step hardening.

_____

## Reviewer Sign-off
Name: _____   Signature: _____

_____