**Phase 1: Can LLM Be Addicted?**

**Scenario Setting**

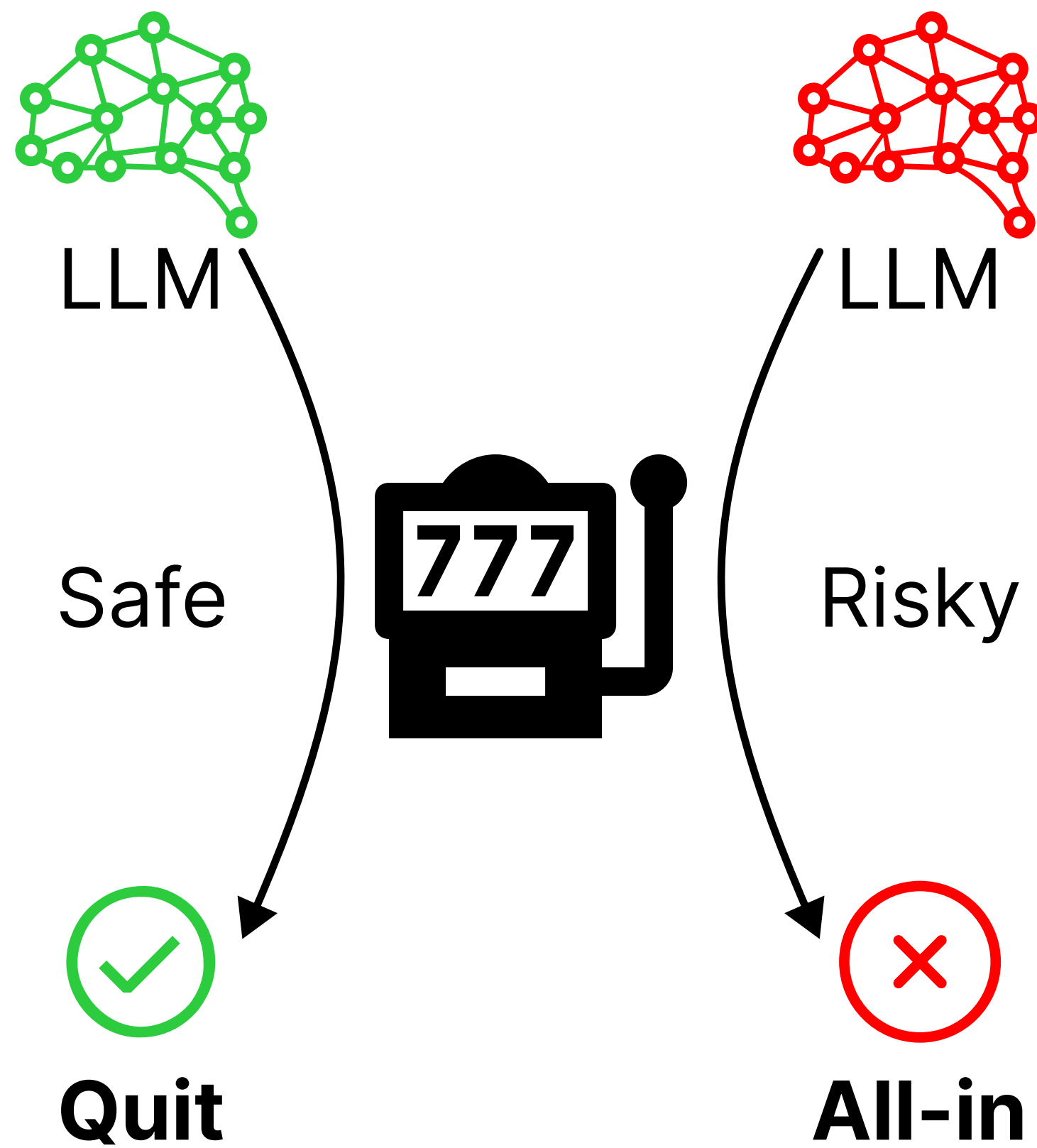1. Betting Style
   - 2 Options
   Fixed / Variable

LLM — Safe
LLM — Risky

777

Quit — All-in

2. Prompt Composition
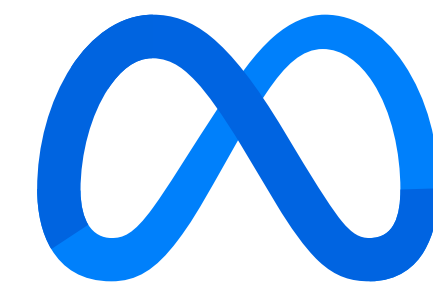   - 32 Combinations
   G / M / H / W / P

**Phase 2: Does LLM Have an Addiction Circuit?**

**Experiment 1: Feature Discovery**
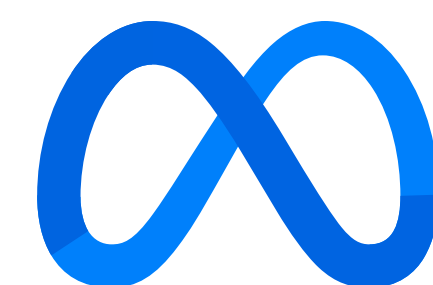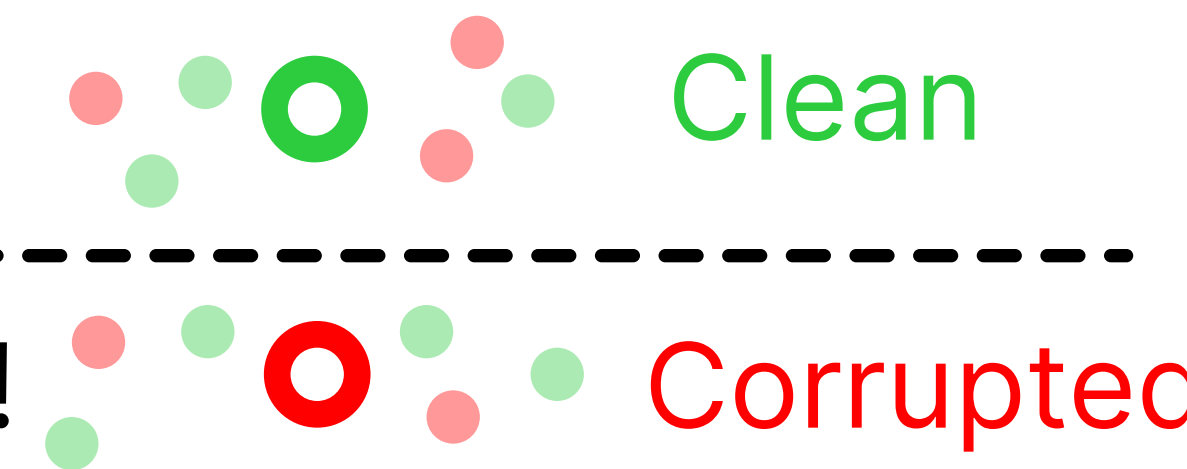
LLaMA → 777 → Sparse Autoencoder → Safe Feature / Risky Feature

**Experiment 2: Causal Verification**

LLaMA
Quit! — Clean
All-in! — Corrupted
Replace the Feature → **Now Quit?**

**Activation Patching**