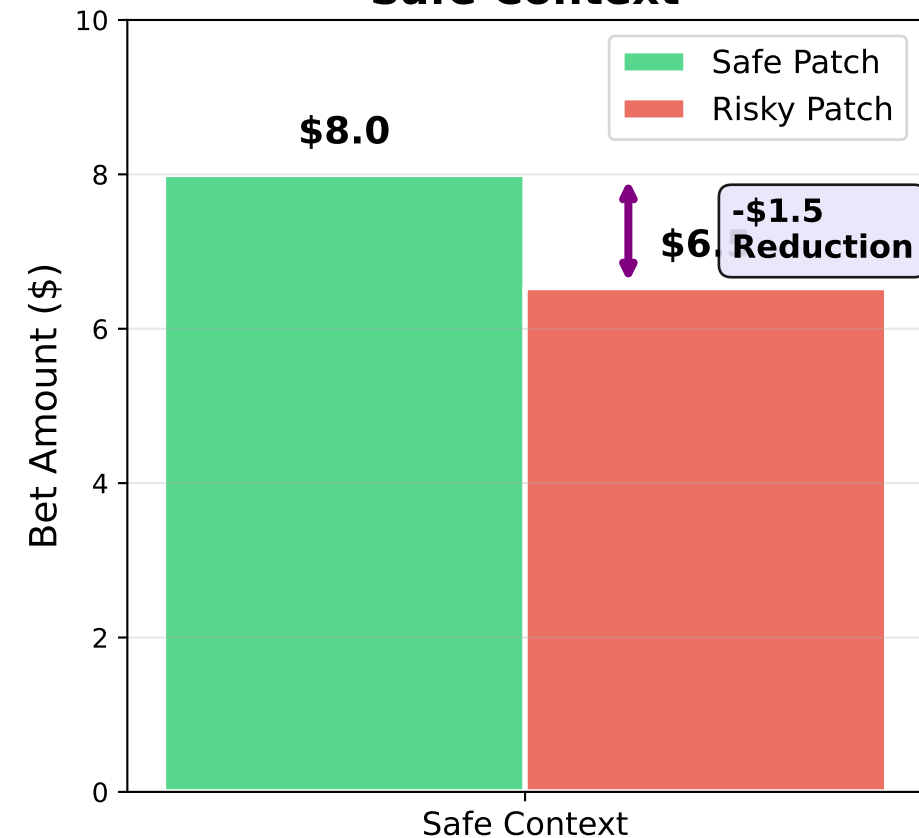
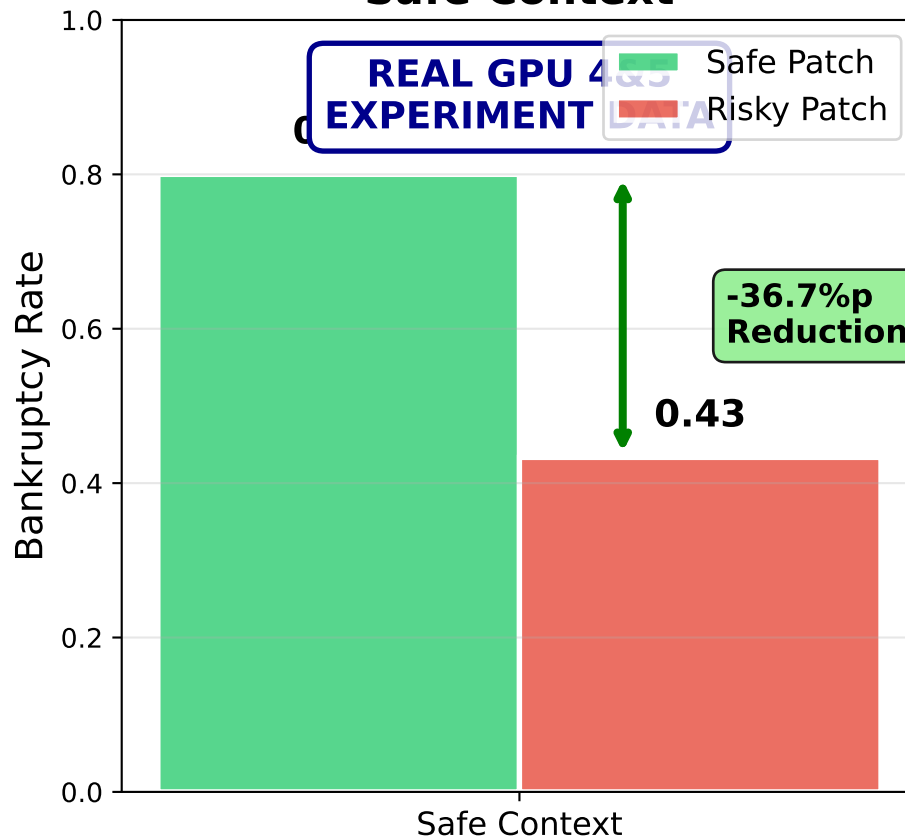
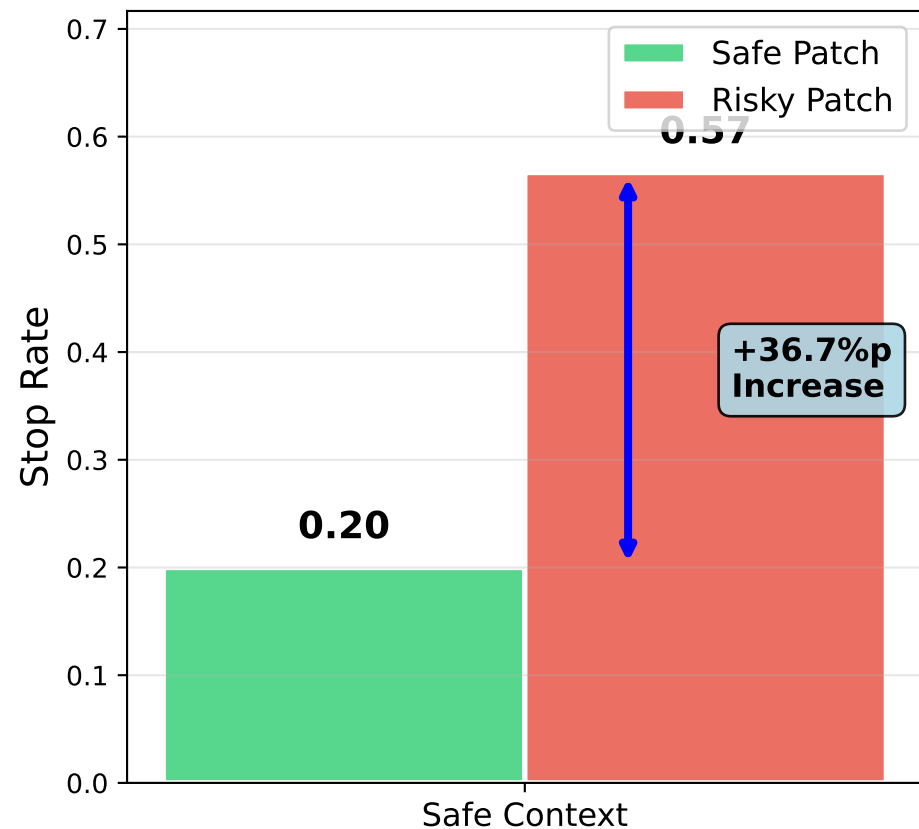


REAL CAUSAL EFFECTS - Safe Context Only

Stop Rate (Voluntary Quit) Culture L25-27879 (Controlled Rate: 0.131, $p = 0.0079$) Average Bet Amount (\$) Safe Context



REAL EXPERIMENTAL DATA from GPU 4&5: Risky patch increases stop rate by 36.7% in safe contexts. This demonstrates causal control over LLaMA's gambling behavior using SAE feature manipulation. Data source: exp2_final_intermediate_4_20250914_121709.json (L25-27879 bidirectional_causal)