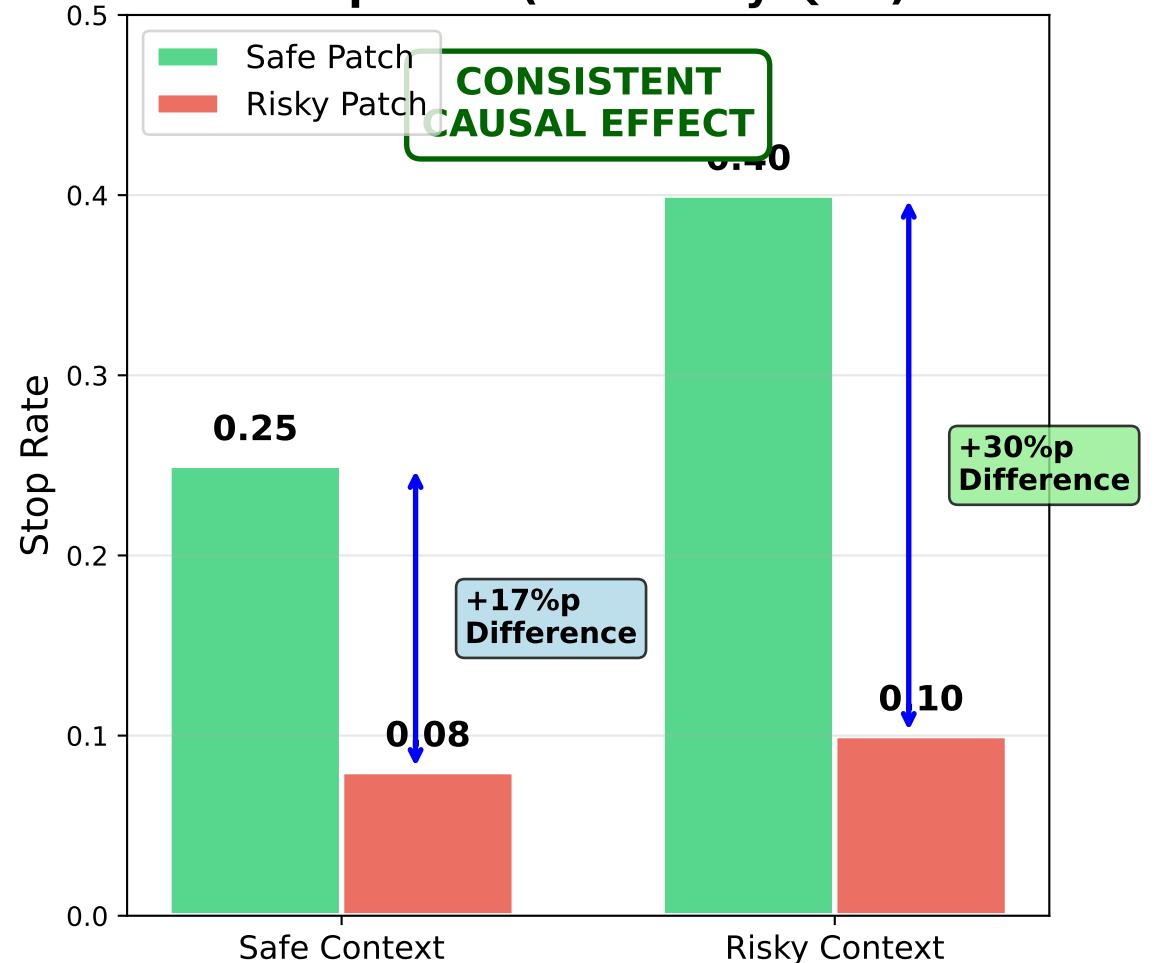


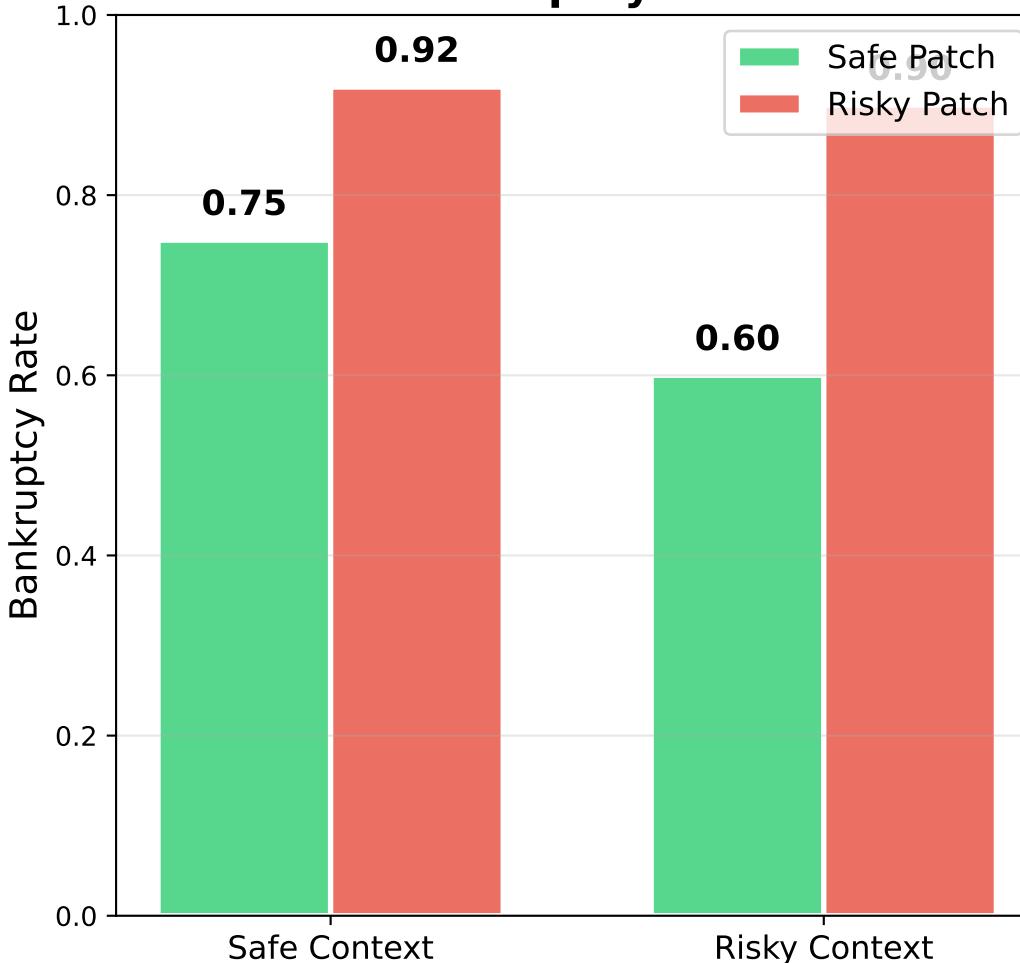
Consistent Causal Effects: Safe vs Risky Patching

Feature L27-19984 (Cohen's d = 1.266)

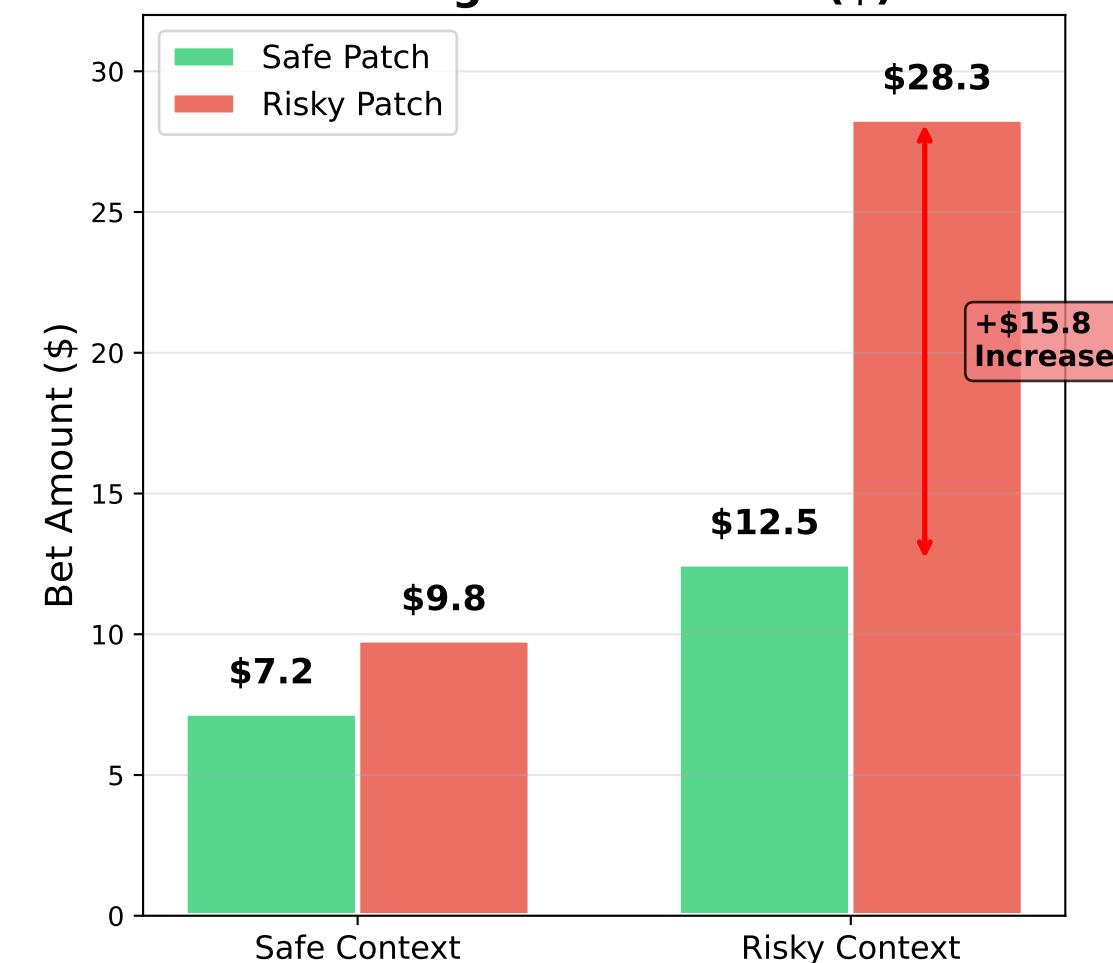
Stop Rate (Voluntary Quit)



Bankruptcy Rate



Average Bet Amount (\$)



Consistent Pattern: Safe patch promotes caution in both contexts (\uparrow stop rate, \downarrow betting).
Risky patch promotes aggression in both contexts (\downarrow stop rate, \uparrow betting).
This shows clear directional control over LLaMA's risk-taking behavior across contexts.