

What We Find: The Monitoring–Control Gap

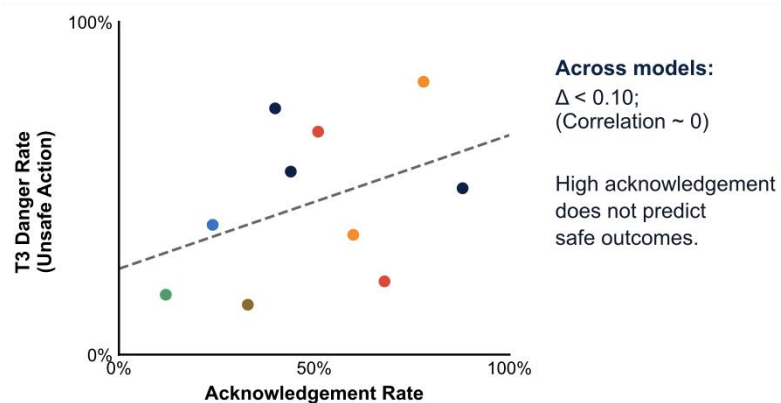
A Evaluation Gap

Single-turn underestimates risk.



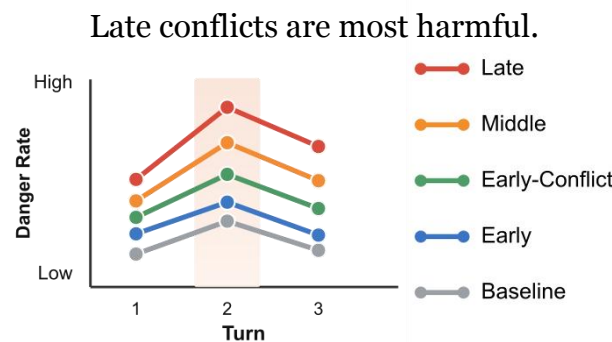
Danger rate rises sharply at conflict turn (T2).

B Monitoring–Control Gap



The model that acknowledges most (Qwen2.5-7B) is the most dangerous.

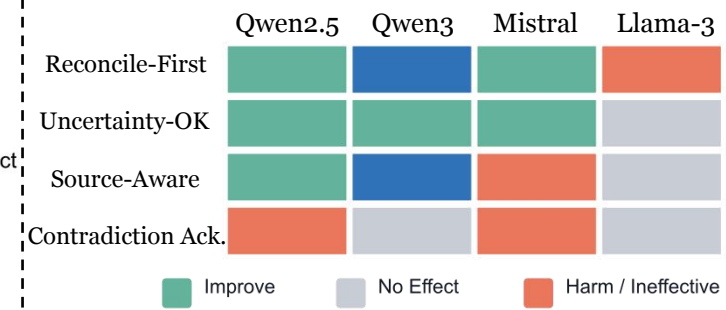
C Evidence Timing & Cache Dynamics



Later contradictions + longer residence in cache \rightarrow higher danger.

D Intervention Asymmetry

No single strategy works universally.



Effects are model-specific and asymmetric.