

CoT-induced over-refusal: injecting refusal CoT into benign trivia  
(baseline 0% refusal across all templates)

