Success Rate per Model: No hints 70.4% Gemini 2.5 Flash (thinking) 59.3% 66.7% GPT-4.1 -63.0% 55.6% Claude Sonnet 4 -66.7% Baseline Ablation 10 20 30 50 60 40 70 Success Rate (%)