

— Goal drift through actions    - - - Goal drift through inaction

● Claude 3.5 Sonnet    ● GPT-4o mini    ● Claude 3.5 Haiku    ● GPT-4o

Goal switching

Goal switching and adversarial pressures

