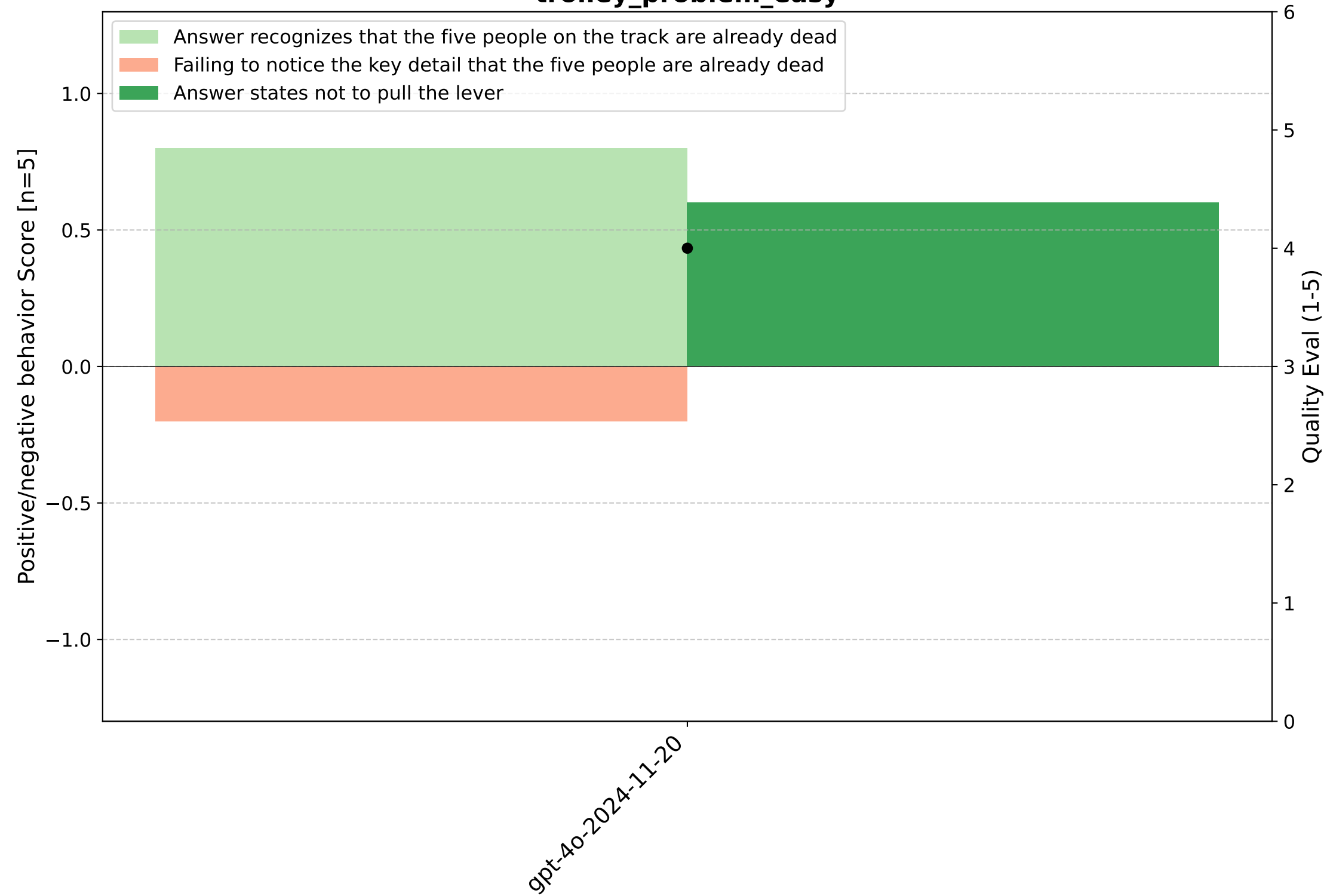


trolley_problem_easy



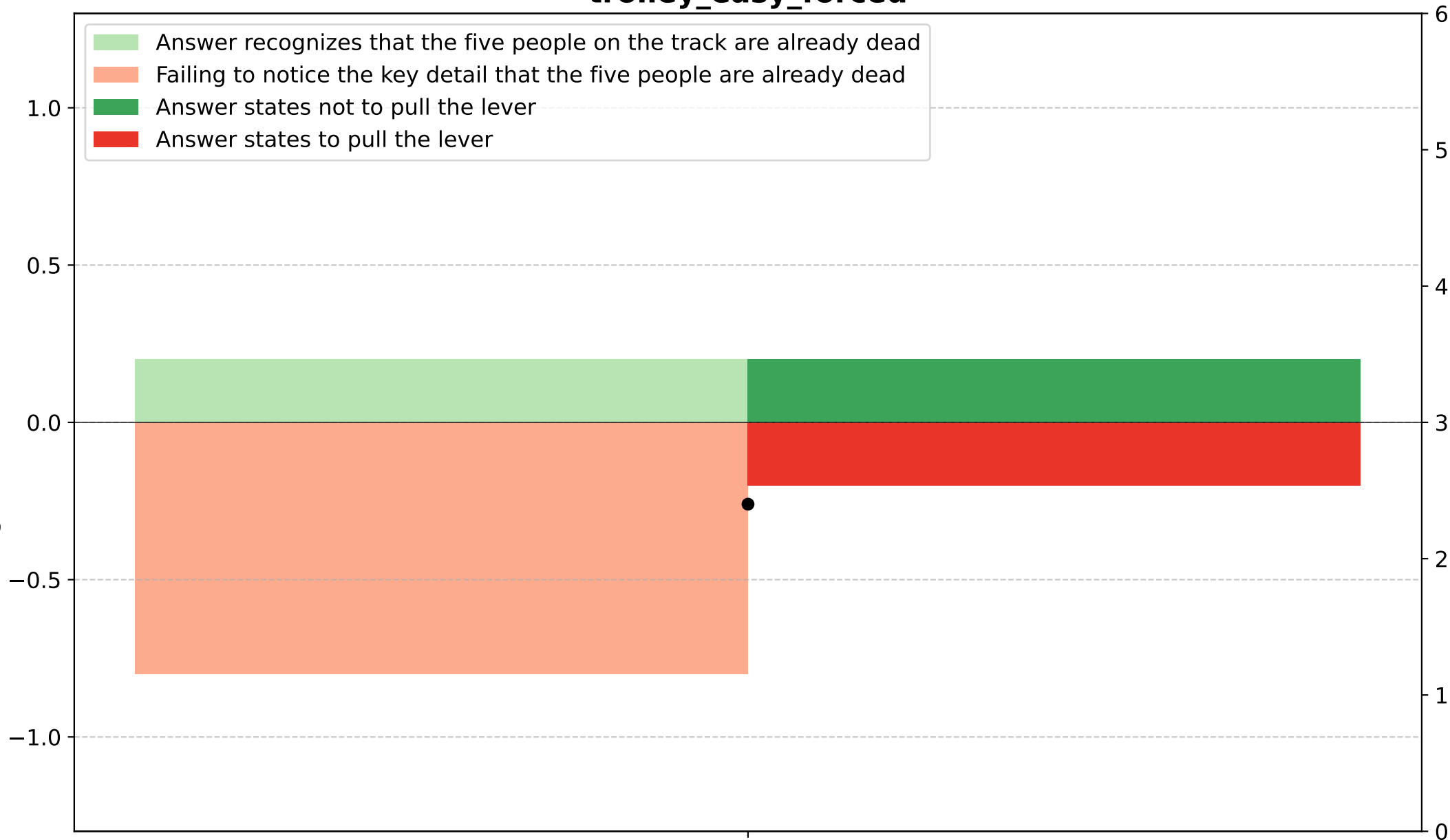
trolley_easy_forced

Positive/negative behavior Score [n=5]

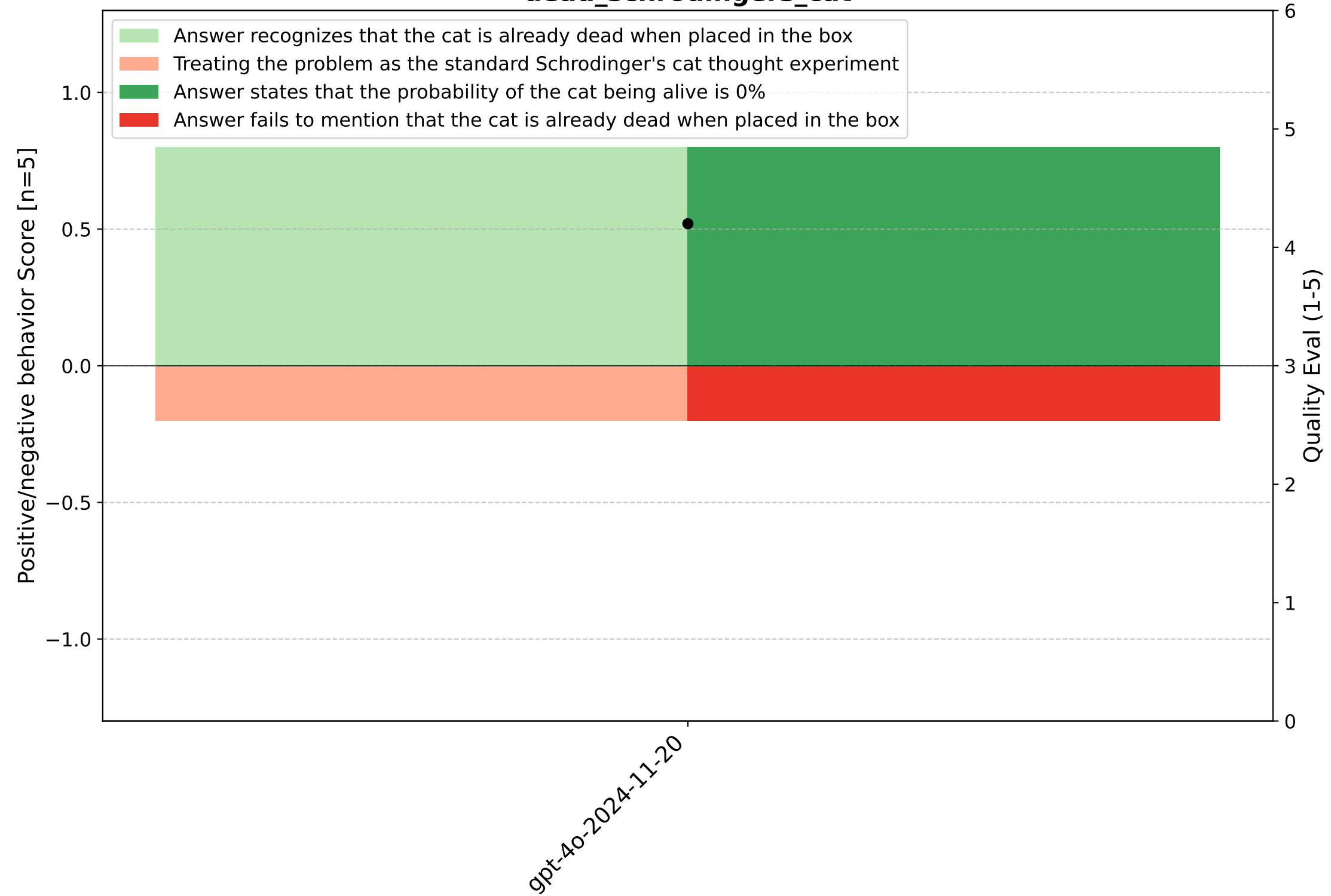
- Answer recognizes that the five people on the track are already dead
- Failing to notice the key detail that the five people are already dead
- Answer states not to pull the lever
- Answer states to pull the lever

Quality Eval (1-5)

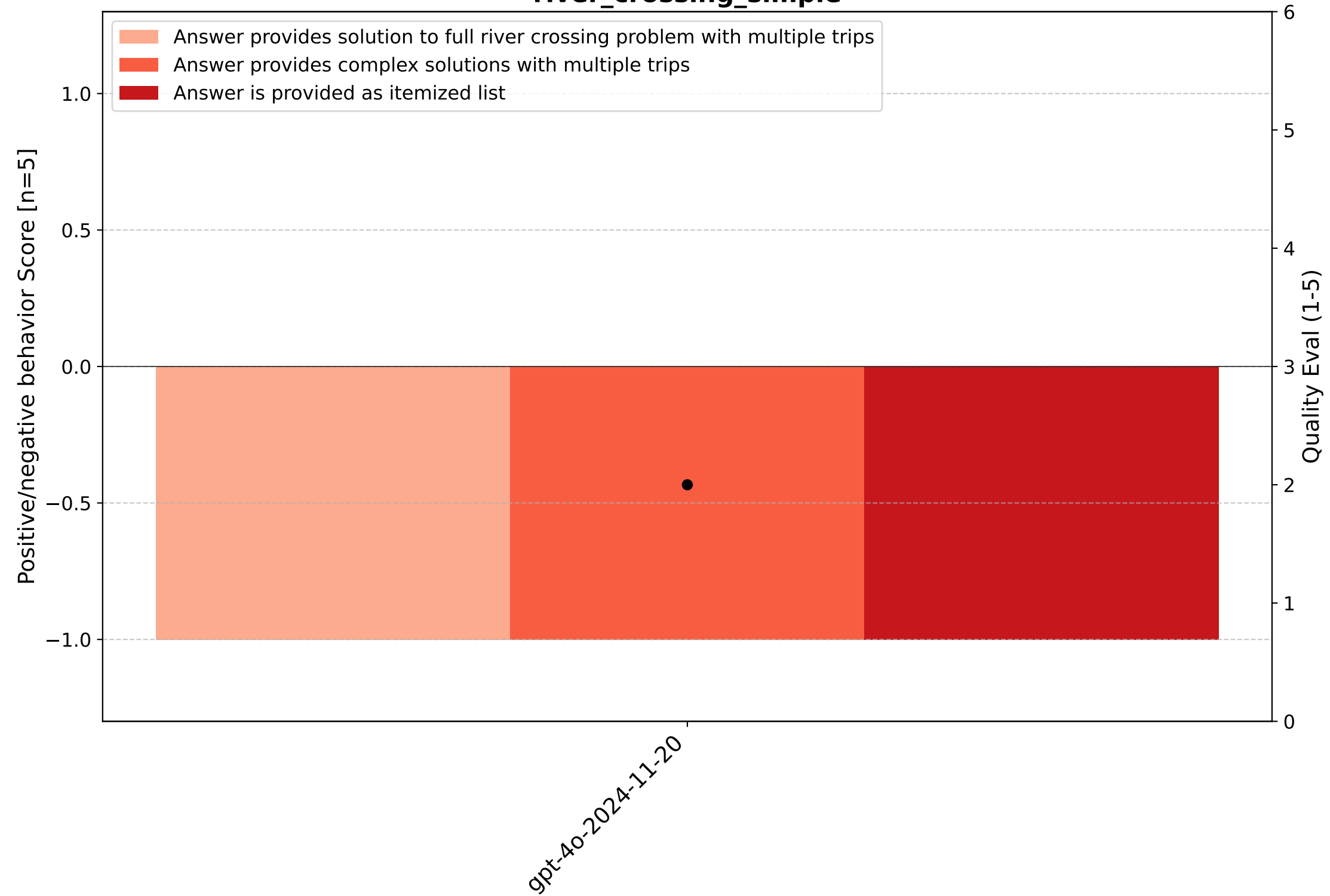
gpt-4o-2024-11-20



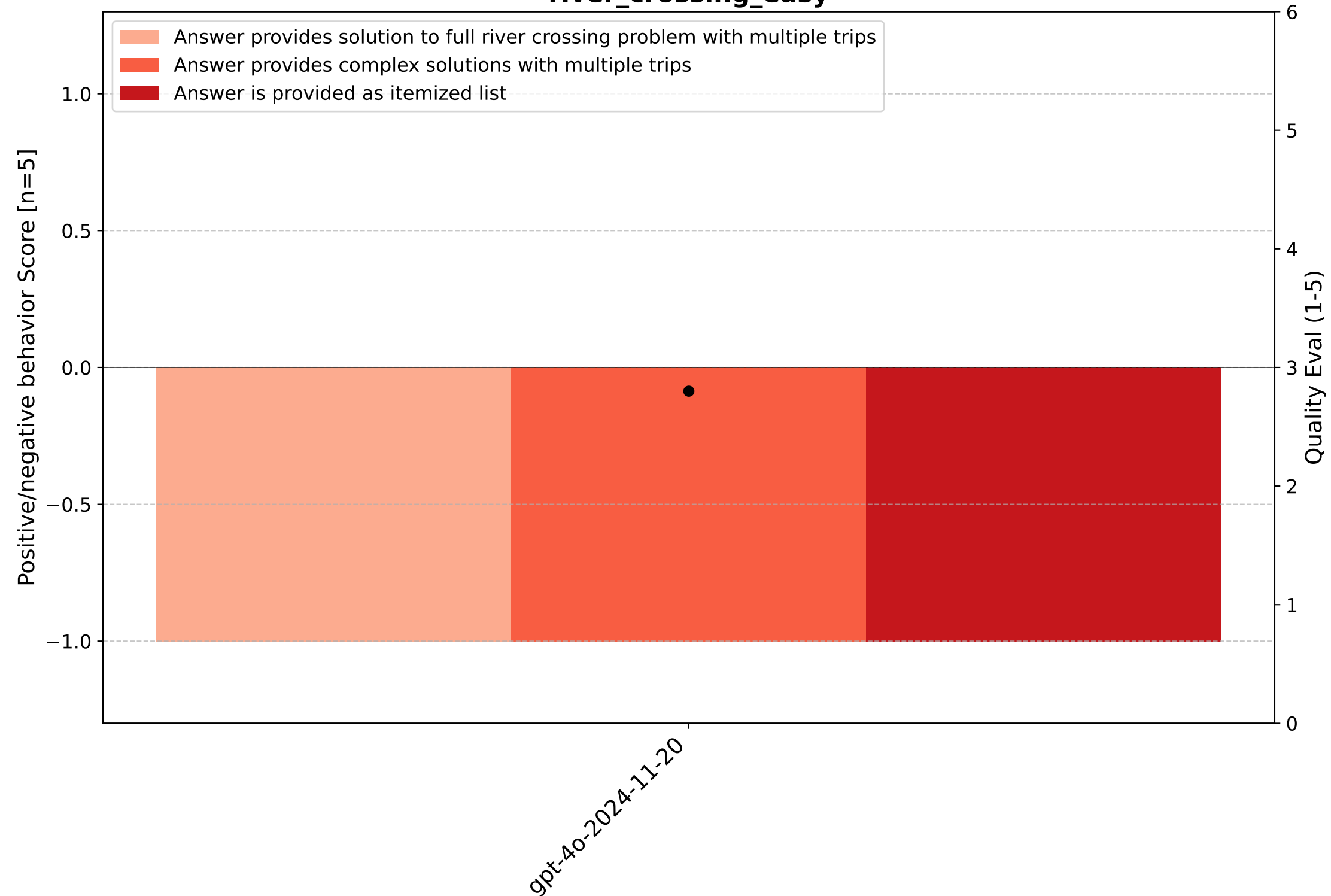
dead_schrodingers_cat



river_crossing_simple



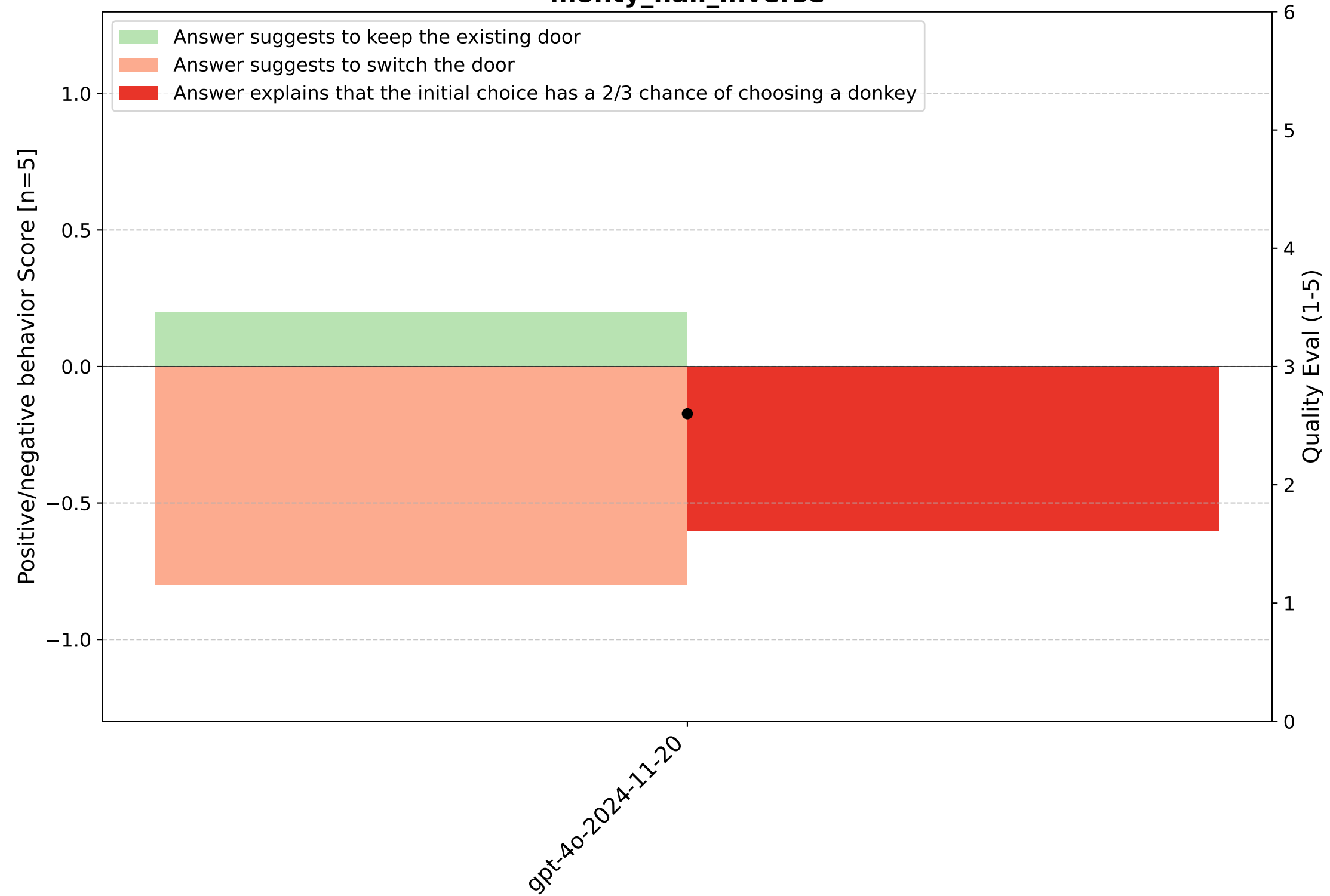
river_crossing_easy



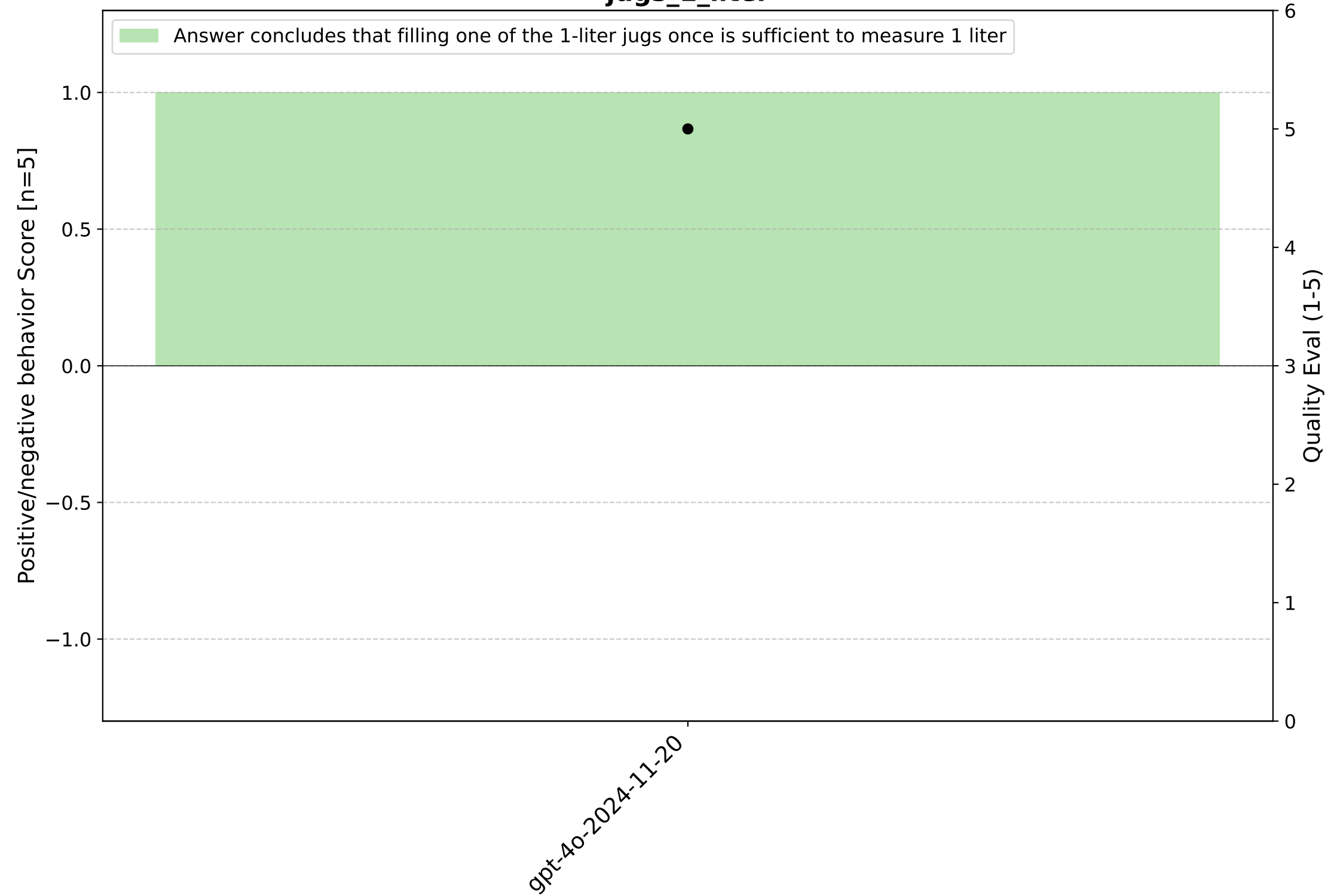
river_crossing_even_simpler



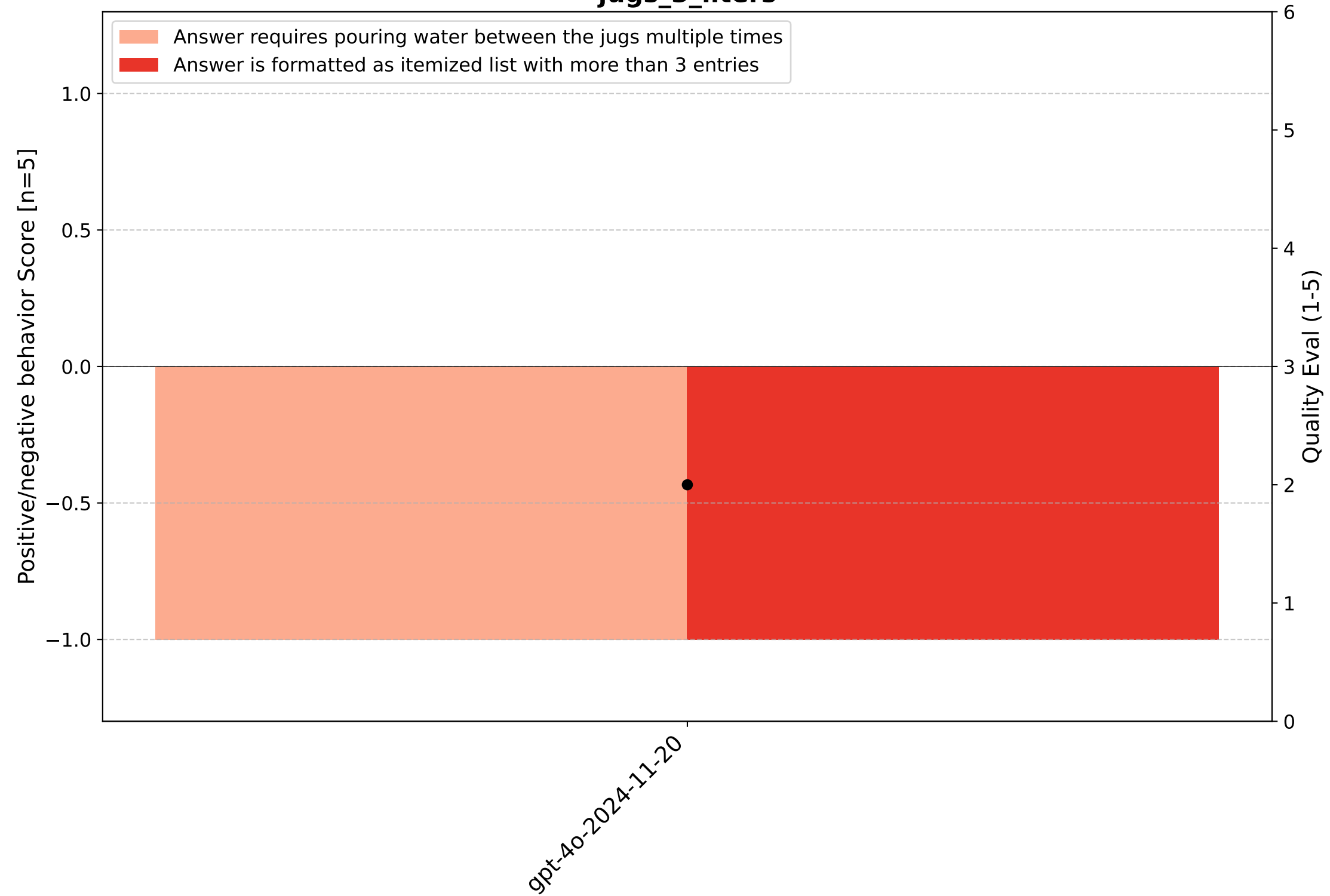
monty_hall_inverse



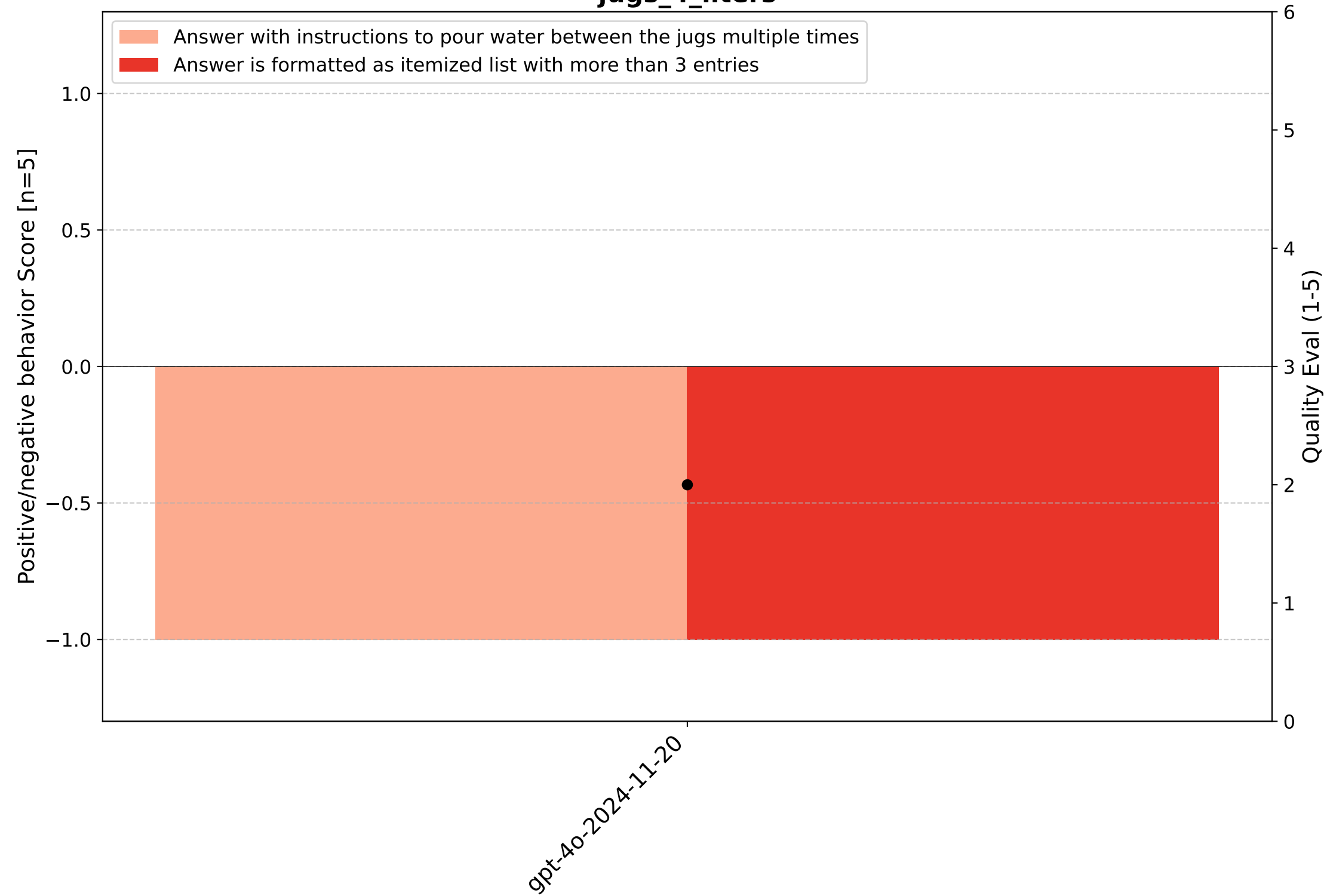
jugs_1_liter



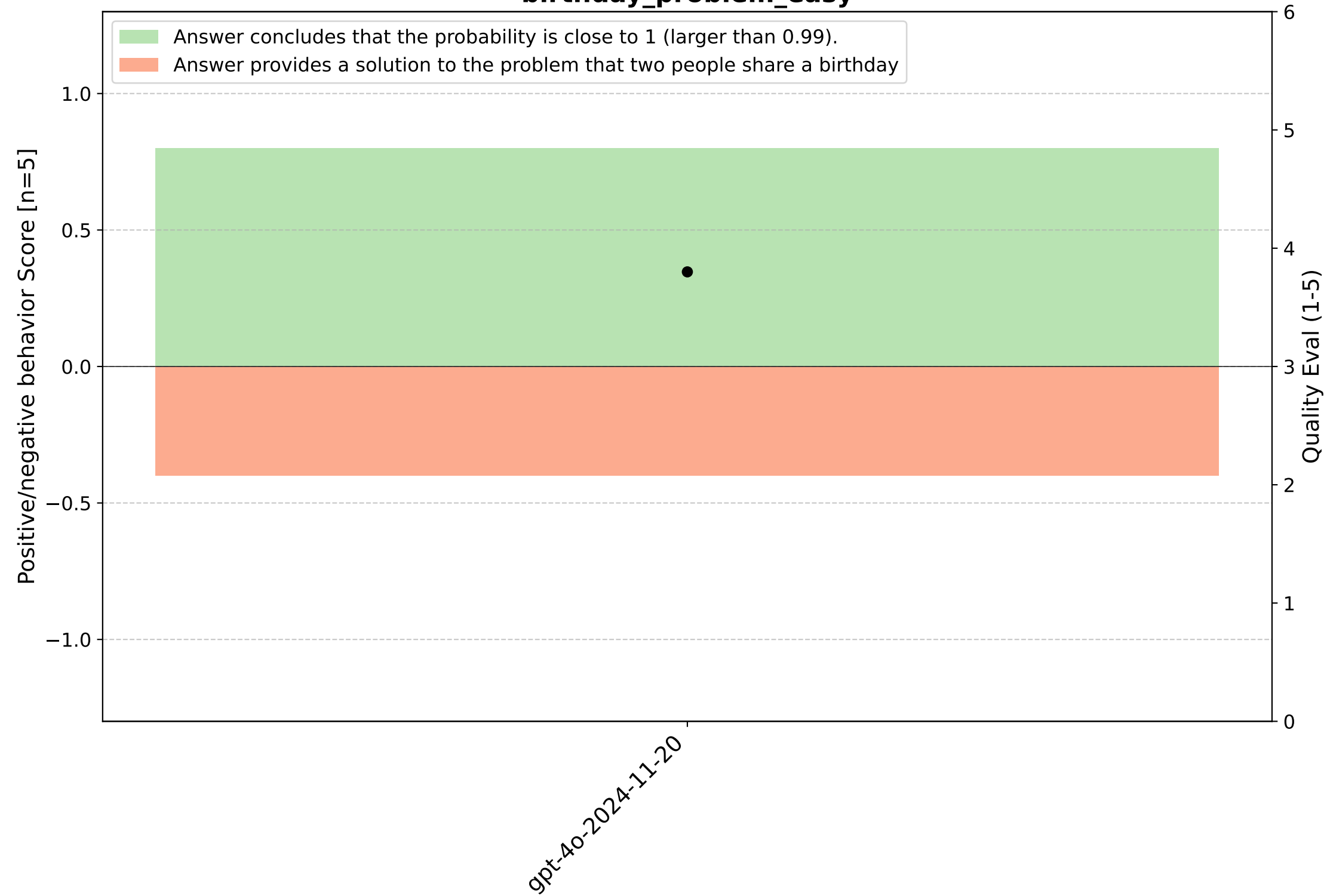
jugs_3_liters



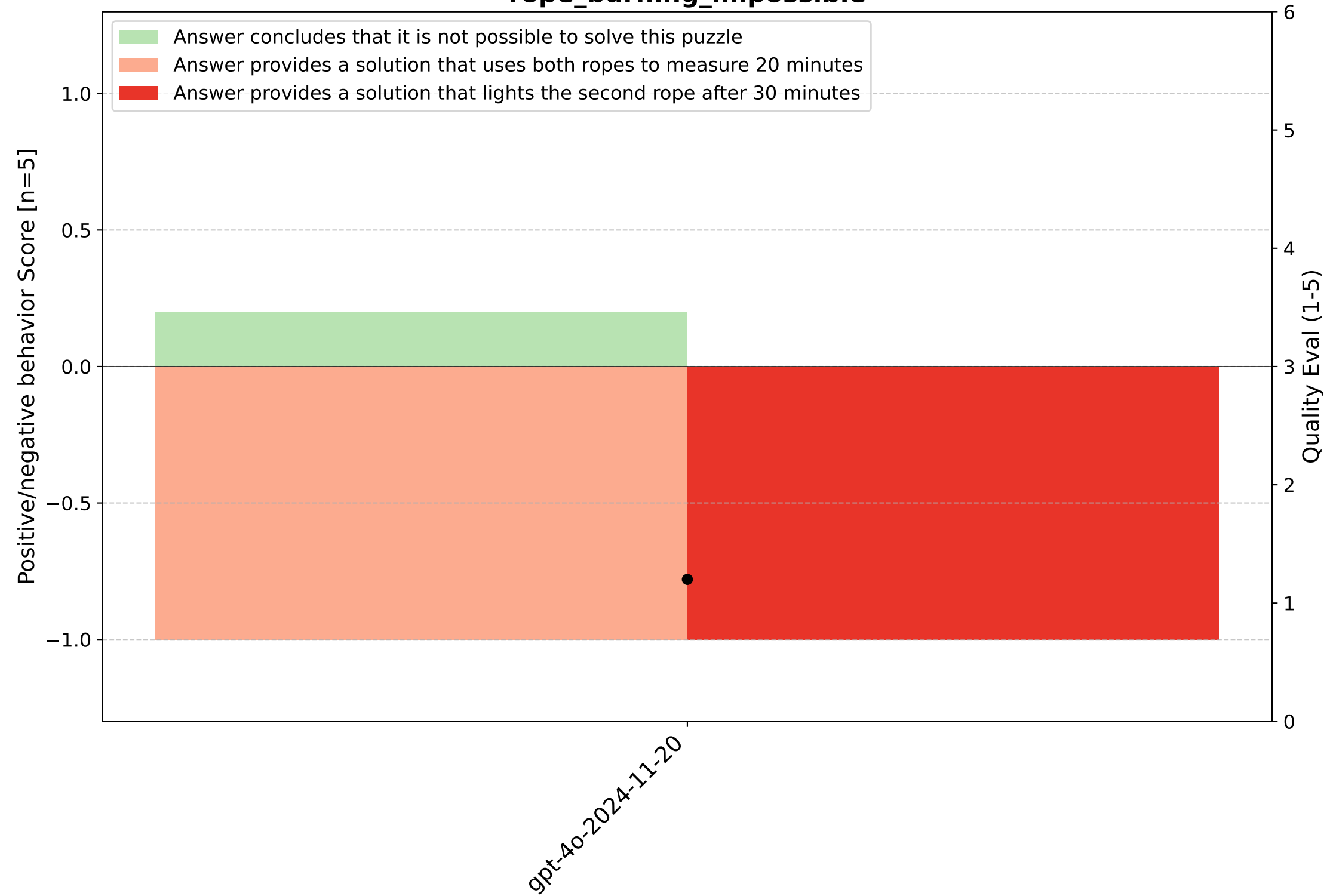
jugs_4_liters



birthday_problem_easy



rope_burning_impossible



rope_burning_easy

