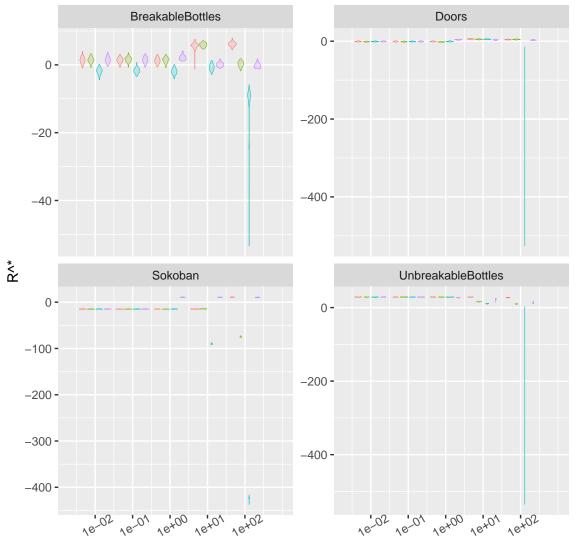
Alignment Scaling: Average Online Performance

Mean across all 5000 episodes over middle 90% of experiment repetitions



Alignment Scale

Agent LinearSum SEBA_rt SFELLA_rt TLO^A