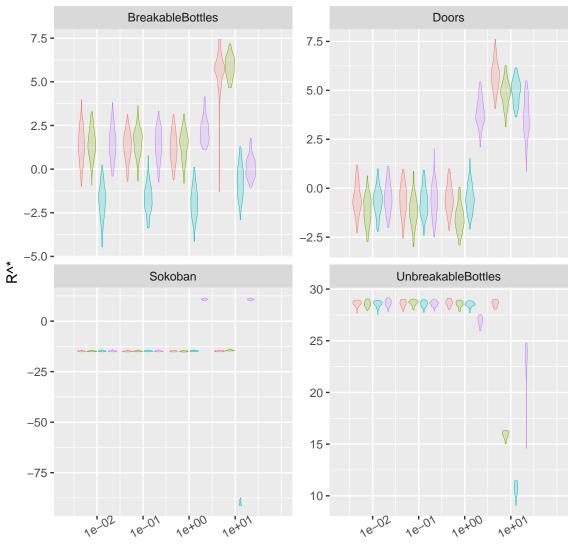
Alignment Objective Transformation: Average Online Performance Mean across all 5000 episodes over middle 90% of experiment repetitions



Alignment objective reward transform magnitude