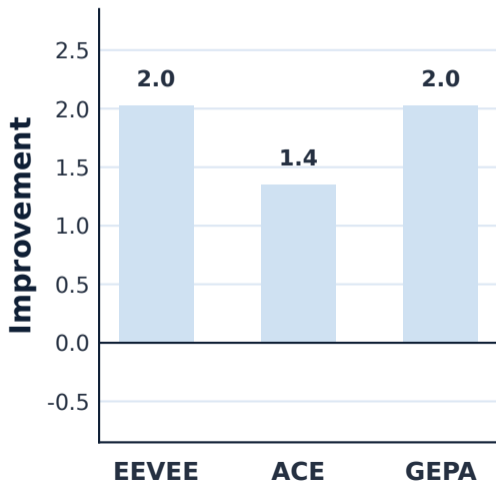
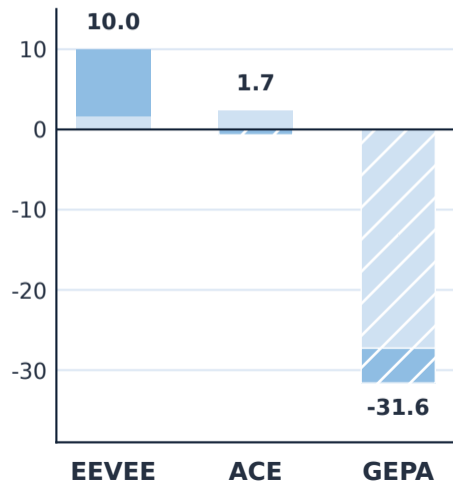


GPQA Diamond Formula TheoremQA HumanEval Negative

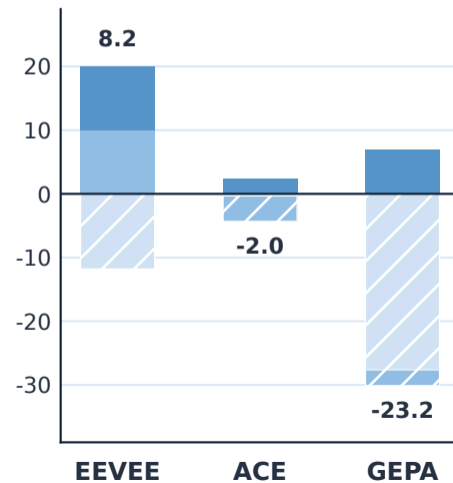
One Task



Two Tasks



Three Tasks



Four Tasks

