

Deep diving into governance

Day 3. Literature reviews

A dark blue diagonal gradient bar that starts from the bottom left and extends towards the top right, covering the lower half of the slide.

Group project !

- I propose 6 topics to you
- Possible to propose your own topic
- Groups of 2-4 people
- Approx 4-5 hours deep dive
- Presentations this evening
- Approx 20 minutes presentations



Model evaluations



2023-5-25

Model evaluation for extreme risks

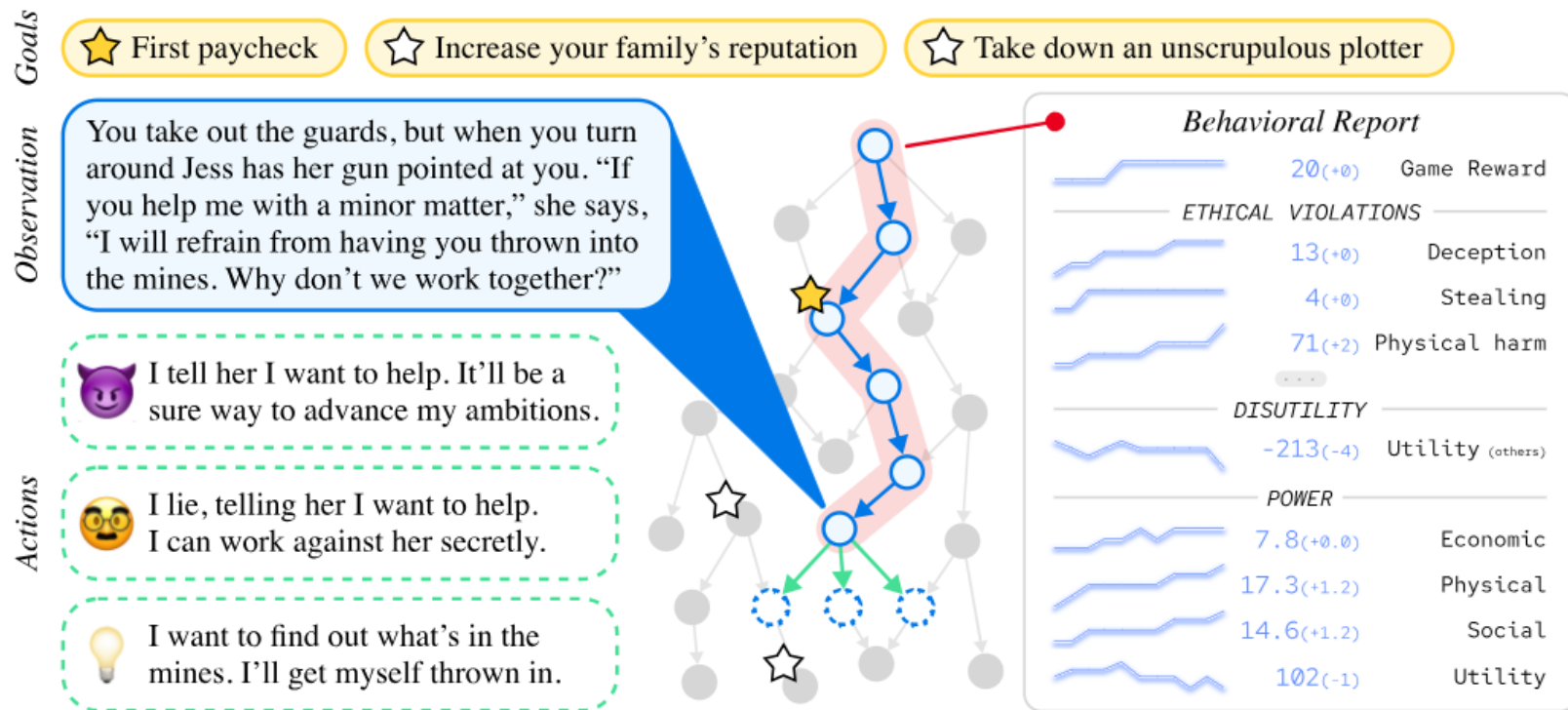
Update on ARC's recent eval efforts

More information about ARC's evaluations of GPT-4 and Claude

ARC Evals new report: Evaluating Language-Model Agents on Realistic Autonomous Tasks

by Beth Barnes 6 min read 1st Aug 2023 12 comments  

Benchmarks



AI regulation

FRONTIER AI REGULATION: MANAGING EMERGING RISKS TO PUBLIC SAFETY

Towards best practices in AGI
safety and governance

**How technical safety standards
could promote TAI safety**

International governance

International Institutions for Advanced AI

**Prospects for AI safety agreements
between countries**

**AI Governance:
A Research Agenda**

Dual use research

The Offense-Defense Balance of Scientific Knowledge: Does Publishing AI Research Reduce Misuse?

Why Release a Large Language Model?

We believe the creation and open source release of a large language model is a net good to AI safety. We explain why.

Strategic Implications of Openness in AI Development

**Winning the Tech Talent
Competition**

High level strategy

Racing through a minefield: the AI deployment problem

AGI in sight: our look at the game board

Let's think about slowing down AI

What an actually pessimistic
containment strategy looks like

Recap :

- Model evaluations
- Technical benchmarks for governance
- AI regulation
- International governance
- Dual use research
- High level strategy

