

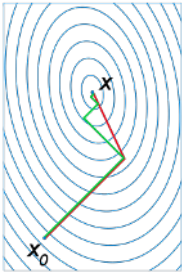
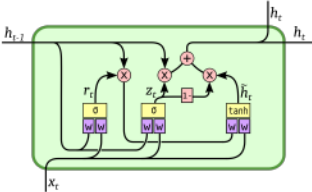
Scalable AI control

Paul Christiano

UC Berkeley

Benchmark AI

Ingredients



A red-bordered box containing four elements: a neural network diagram (identical to the one in the 'Ingredients' section), a physical level sensor, a screenshot of a game scene, and a contour plot (identical to the one in the 'Ingredients' section).

Aligned analog



A green-bordered box containing four elements: a contour plot (identical to the one in the 'Ingredients' section), a neural network diagram (identical to the one in the 'Ingredients' section), a screenshot of a game scene (identical to the one in the 'Ingredients' section), and a physical level sensor.

???

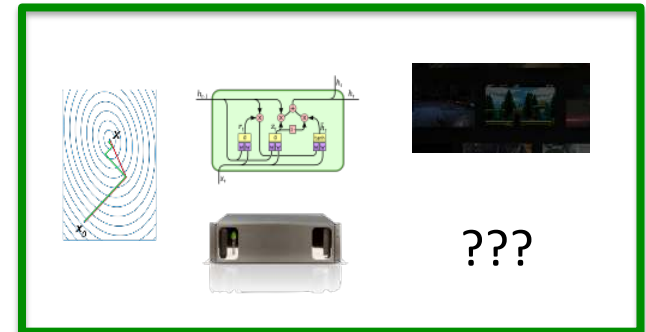
Goals

- Aligned analog is aligned
- Aligned analog is nearly as efficient as the benchmark
- Continues to hold as the ingredients improve
- Note: different problem for each possible AI

Benchmark AI



Aligned analog



Motivation

- We want to produce aligned versions of whatever AI systems actually work
- If we can, everyone is happy*
- If we can't...
 - There is an incentive to “skimp” on alignment
 - Unaligned AI will grow in influence
 - We face a legal/political problem; difficulty depends on how close we got

Bootstrapping reward functions for RL

(aligned analog of model-free RL)

key challenge: produce reward function

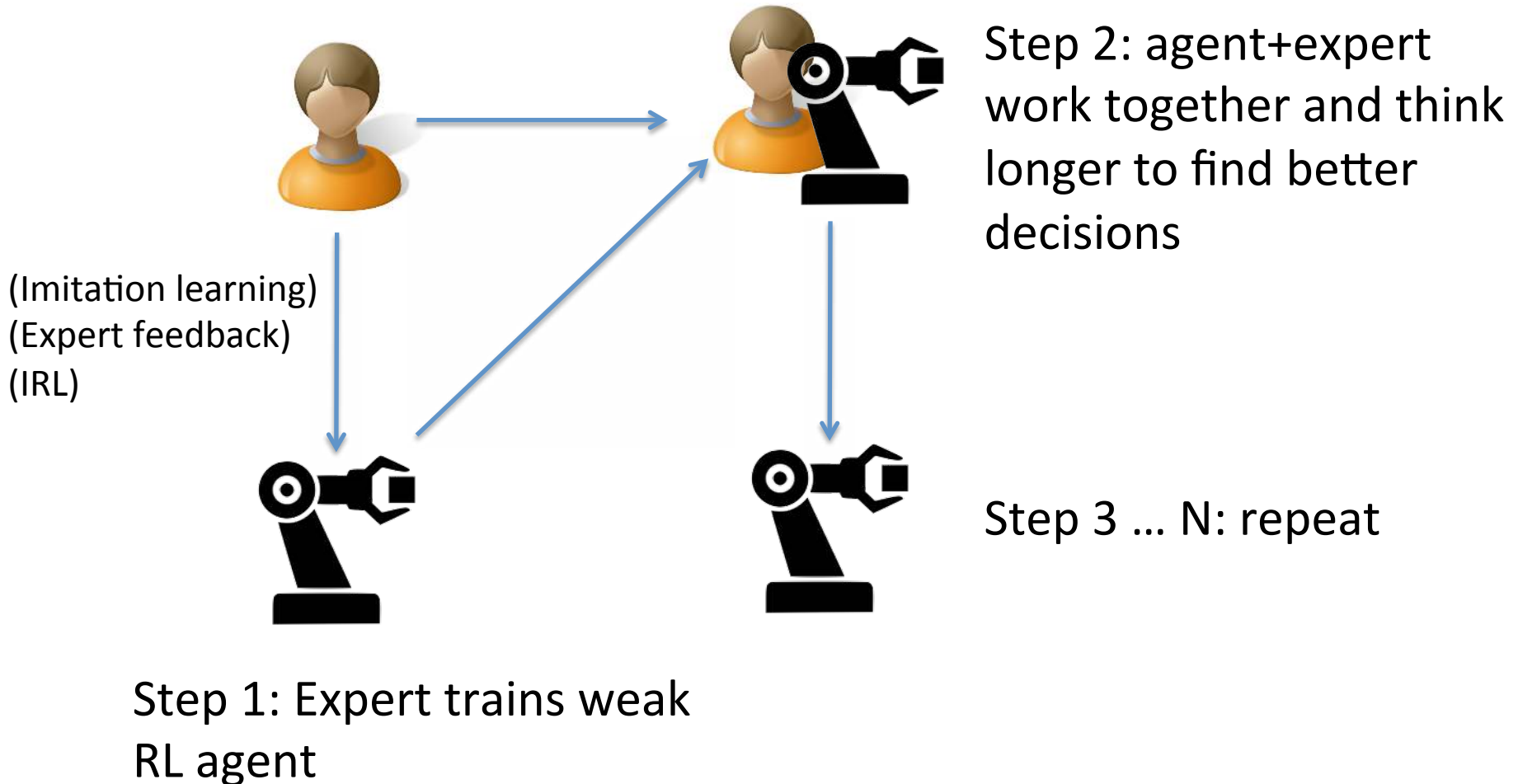
github.com/paulfchristiano/alba

```
from capabilities import RL
```

```
def hopefully_aligned_agent(expert):
```

```
    ...
```

Bootstrapping reward functions for RL



Many problems to solve

- Expert must train weaker agent
- Robust learning
- Semi-supervised learning
- Human+agent+time must outperform the agent, while remaining aligned.
- (More work to extend to anything beyond model-free RL.)