Simple bright ideas going wrong
The big picture
Fundamental difficulties

# Fundamental Difficulties in Aligning Advanced AI

Nate Soares

Adapted from a talk by Eliezer Yudkowsky

Simple bright ideas going wrong
The big picture
Fundamental difficulties

"The primary concern is not spooky emergent consciousness but simply the ability to make **high-quality decisions**."

—*Stuart Russell*

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Task: Fill cauldron.

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Broom's utility function:

$$\mathcal{U}_{broom} = \begin{cases} 1 & \text{if cauldron full} \\ 0 & \text{if cauldron empty} \end{cases}$$

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Broom's utility function:

$$\mathcal{U}_{broom} = \begin{cases} 1 & \text{if cauldron full} \\ 0 & \text{if cauldron empty} \end{cases}$$

Actions $a \in \mathcal{A}$, broom calculates: $\mathbb{E}\left[\mathcal{U}_{broom} \mid a\right]$

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Broom's utility function:

$$\mathcal{U}_{broom} = \begin{cases} 1 & \text{if cauldron full} \\ 0 & \text{if cauldron empty} \end{cases}$$

Actions $a \in \mathcal{A}$, broom calculates: $\mathbb{E}\left[\mathcal{U}_{broom} \mid a\right]$

Broom outputs: $\underset{a \in \mathcal{A}}{\text{sorta-argmax}}\ \mathbb{E}\left[\mathcal{U}_{broom} \mid a\right]$

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Simple bright ideas going wrong
The big picture
Fundamental difficulties

*Difficulty 1...*

Broom's utility function:

$$\mathcal{U}_{broom} = \begin{cases} 1 & \text{if cauldron full} \\ 0 & \text{if cauldron empty} \end{cases}$$

Human's utility function:

$$\mathcal{U}_{human} = \begin{cases} 1 & \text{if cauldron full} \\ 0 & \text{if cauldron empty} \\ -10 & \text{if workshop flooded} \end{cases}$$

Simple bright ideas going wrong
The big picture
Fundamental difficulties

*Difficulty 1...*

Broom's utility function:

$$\mathcal{U}_{broom} = \begin{cases} 1 & \text{if cauldron full} \\ 0 & \text{if cauldron empty} \end{cases}$$

Human's utility function:

$$\mathcal{U}_{human} = \begin{cases} 1 & \text{if cauldron full} \\ 0 & \text{if cauldron empty} \\ -10 & \text{if workshop flooded} \\ +0.2 & \text{if it's funny} \\ -1000000 & \text{if someone gets killed} \\ & \dots \text{and a whole lot more} \end{cases}$$

Simple bright ideas going wrong
The big picture
Fundamental difficulties

*Difficulty 2...*

$\mathcal{EU}$(99.99% chance of full cauldron) $>$ $\mathcal{EU}$(99.9% chance of full cauldron)

Simple bright ideas going wrong
The big picture
Fundamental difficulties

*Difficulty 2. . .*

$\mathcal{EU}$(99.99% chance of full cauldron) $>$ $\mathcal{EU}$(99.9% chance of full cauldron)

- Contrast "Task" - goal bounded in space, time, fulfillability, and effort required to fulfill

Simple bright ideas going wrong
The big picture
Fundamental difficulties

*Difficulty 2...*

$\mathcal{EU}$(99.99% chance of full cauldron) $>$ $\mathcal{EU}$(99.9% chance of full cauldron)

- Contrast "Task" - goal bounded in space, time, fulfillability, and effort required to fulfill
- "Task AGI" - not just top goal, but optimization subroutines are Tasks: nothing open-ended anywhere

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Can we just press the off switch?

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Try 1: Suspend button **B**

$$\mathcal{U}^3_{broom} = \begin{cases} 1 \text{ if cauldron full} & \& \ \mathbf{B}{=}\text{OFF} \\ 0 \text{ if cauldron empty} & \& \ \mathbf{B}{=}\text{OFF} \\ 1 \text{ if broom suspended} & \& \ \mathbf{B}{=}\text{ON} \\ 0 \text{ otherwise} \end{cases}$$

Simple bright ideas going wrong
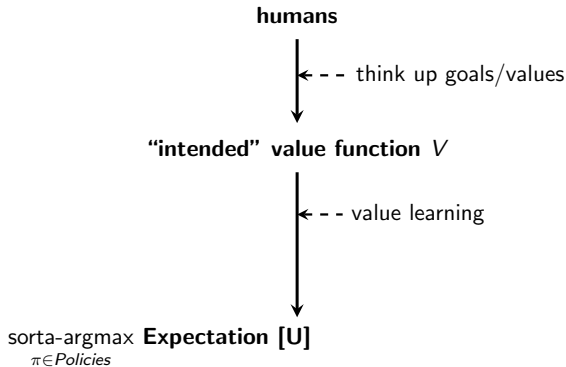The big picture
Fundamental difficulties

Try 1: Suspend button **B**

$$\mathcal{U}_{broom}^3 = \begin{cases} 1 \text{ if cauldron full} & \& \text{ } \mathbf{B}\text{=OFF} \\ 0 \text{ if cauldron empty} & \& \text{ } \mathbf{B}\text{=OFF} \\ 1 \text{ if broom suspended} & \& \text{ } \mathbf{B}\text{=ON} \\ 0 \text{ otherwise} \end{cases}$$

Probably, $\mathbb{E}\left[\mathcal{U}_{broom}^3 \mid \mathbf{B}\text{=OFF}\right] < \mathbb{E}\left[\mathcal{U}_{broom}^3 \mid \mathbf{B}\text{=ON}\right]$

Simple bright ideas going wrong
The big picture
Fundamental difficulties

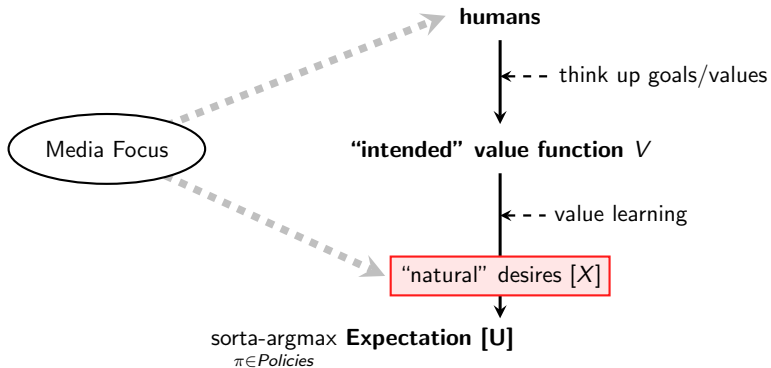Try 1: Suspend button **B**

$$\mathcal{U}_{broom}^3 = \begin{cases} 1 \text{ if cauldron full} & \& \ \textbf{B}{=}\text{OFF} \\ 0 \text{ if cauldron empty} & \& \ \textbf{B}{=}\text{OFF} \\ 1 \text{ if broom suspended} & \& \ \textbf{B}{=}\text{ON} \\ 0 \text{ otherwise} \end{cases}$$
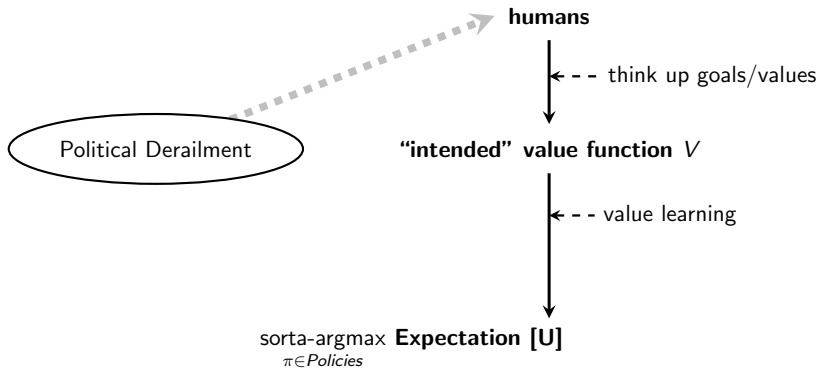
Probably, $\mathbb{E}\left[\mathcal{U}_{broom}^3 \mid \textbf{B}{=}\text{OFF}\right] < \mathbb{E}\left[\mathcal{U}_{broom}^3 \mid \textbf{B}{=}\text{ON}\right]$
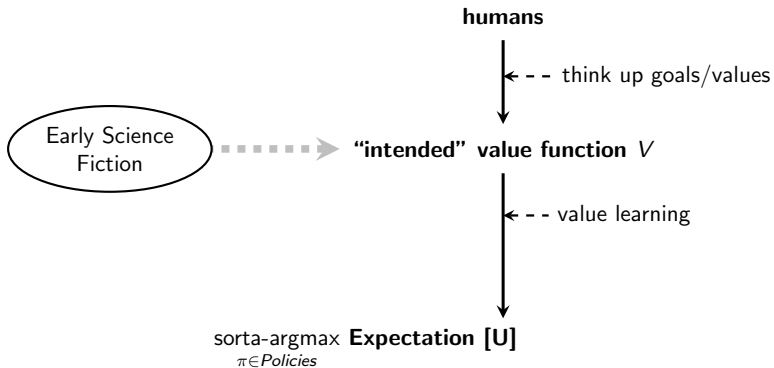
(Strategic broom tries to make you press the button.)

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

**humans**

$\Big\downarrow$ ◀ - - think up goals/values

**"intended" value function** $V$

$\Big\downarrow$ ◀ - - value learning

$\underset{\pi \in Policies}{\text{sorta-argmax}}$ **Expectation [U]**

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties



**humans**

← - - think up goals/values

Political Derailment

**"intended" value function** $V$

← - - value learning

$\underset{\pi \in \textit{Policies}}{\text{sorta-argmax}}$ **Expectation [U]**

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

**humans**

$\leftarrow$ - - think up goals/values

Early Science
Fiction ┈┈┈┈▶ **"intended" value function** $V$

$\leftarrow$ - - value learning

$\underset{\pi \in Policies}{\text{sorta-argmax}}$ **Expectation [U]**

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties



**humans**

think up goals/values

**"intended" value function** *V*
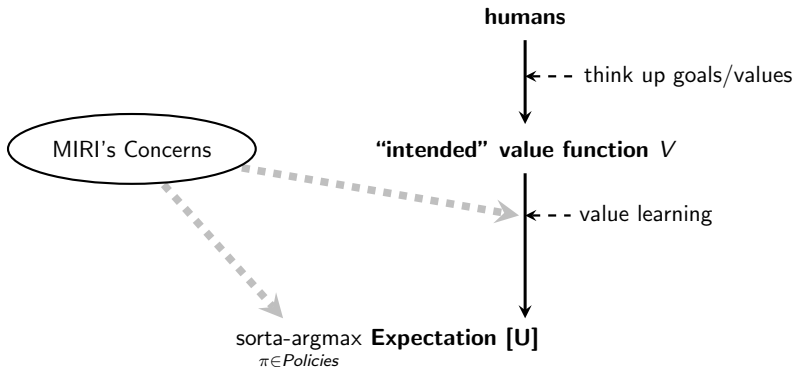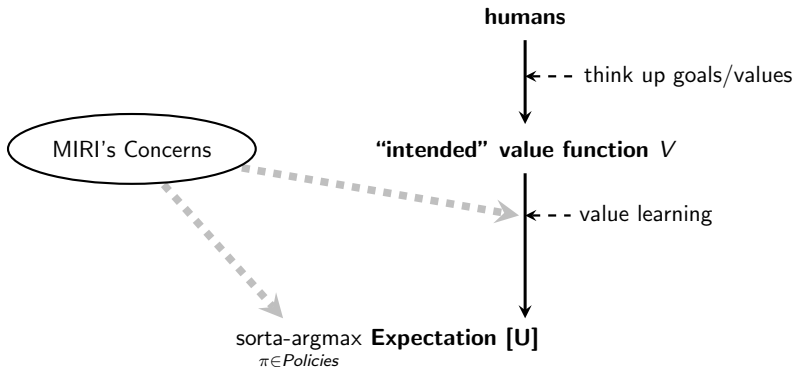
MIRI's Concerns

value learning

sorta-argmax **Expectation [U]**
π∈*Policies*

Take-home message: We're afraid it's going to be *technically difficult* to point AIs in an intuitively intended direction.

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

**humans**

← - - think up goals/values

MIRI's Concerns

**"intended" value function** *V*

← - - value learning

sorta-argmax **Expectation [U]**
$\pi \in Policies$

Take-home message: We're afraid it's going to be *technically difficult*
to point AIs in an intuitively intended direction.

...and if we screw up there, it *doesn't matter* which human
is standing closest to the AI.

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

Four key propositions:

1. **Orthogonality** – An AI system can be built to pursue almost any objective, in theory

Simple bright ideas going wrong
The big picture
Fundamental difficulties

Four key propositions:

1. **Orthogonality** – An AI system can be built to pursue almost any objective, in theory
2. **Instrumental convergence** – most objectives imply survival, resource acquisition, etc. as instrumental subgoals

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

Four key propositions:

1. **Orthogonality** – An AI system can be built to pursue almost any objective, in theory
2. **Instrumental convergence** – most objectives imply survival, resource acquisition, etc. as instrumental subgoals
3. **Capability gain** – there are potential ways for artificial agents to greatly gain in cognitive power and strategic options

Simple bright ideas going wrong
**The big picture**
Fundamental difficulties

Four key propositions:

1. **Orthogonality** – An AI system can be built to pursue almost any objective, in theory
2. **Instrumental convergence** – most objectives imply survival, resource acquisition, etc. as instrumental subgoals
3. **Capability gain** – there are potential ways for artificial agents to greatly gain in cognitive power and strategic options
4. **Alignment difficulty** – there's at least one part of "build an AI that does a big right thing" which is a deep, technical, hard AI problem

Simple bright ideas going wrong
The big picture
Fundamental difficulties

AI alignment is difficult...



...like rockets are difficult.

(Huge stresses break things that don't break in normal engineering.)

Simple bright ideas going wrong
The big picture
**Fundamental difficulties**
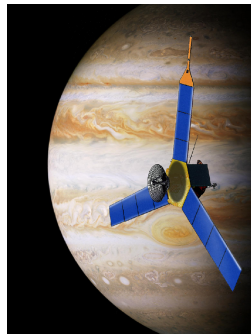
AI aligment is difficult. . .



. . . like space probes are difficult.

(If something goes wrong, it may be high and
out of reach.)

Simple bright ideas going wrong
The big picture
Fundamental difficulties

AI aligment is difficult. . .

. . . *sort of* like computer security is difficult.



(Intelligent search may select in favor of
unusual new paths outside our intended
behavior model.)

Simple bright ideas going wrong
The big picture
Fundamental difficulties

AI alignment:

**Treat it like a secure rocket probe.**

Simple bright ideas going wrong
The big picture
**Fundamental difficulties**

AI alignment:

**Treat it like a secure rocket probe.**

**Take it seriously.**

Simple bright ideas going wrong
The big picture
Fundamental difficulties

AI alignment:

**Treat it like a secure rocket probe.**

**Don't expect it to be easy.**

Simple bright ideas going wrong
The big picture
Fundamental difficulties

AI alignment:

**Treat it like a secure rocket probe.**

Simple bright ideas going wrong
The big picture
Fundamental difficulties

AI alignment:

**Treat it like a secure rocket probe.**

**Don't defer thinking until later.**

Simple bright ideas going wrong
The big picture
Fundamental difficulties

AI alignment:

**Treat it like a secure rocket probe.**

**Formalize ideas so others can critique and build upon them.**

Simple bright ideas going wrong
The big picture
Fundamental difficulties

# Questions?

Email: contact@intelligence.org