# Challenges in aligning AI with human preferences and values
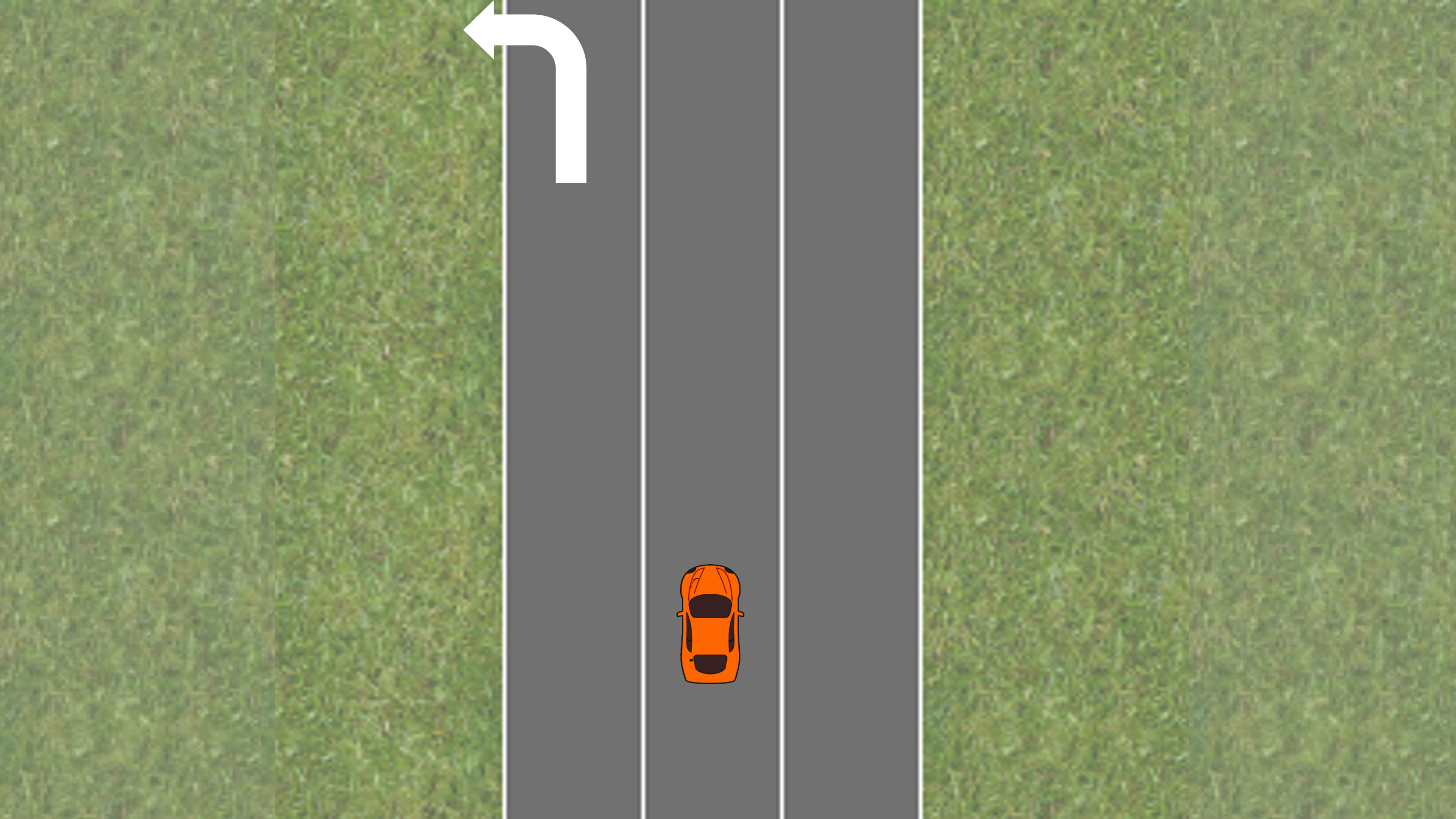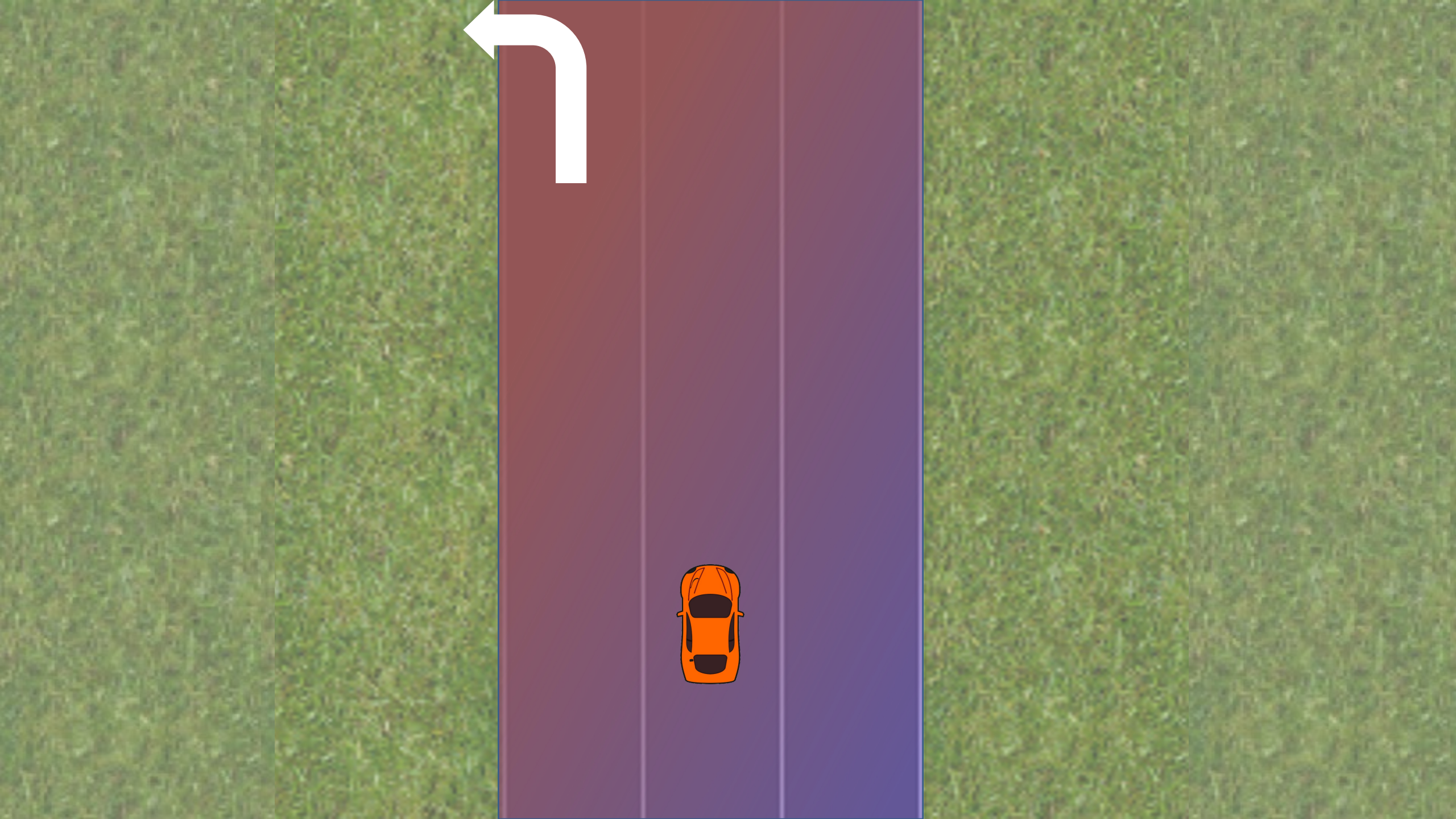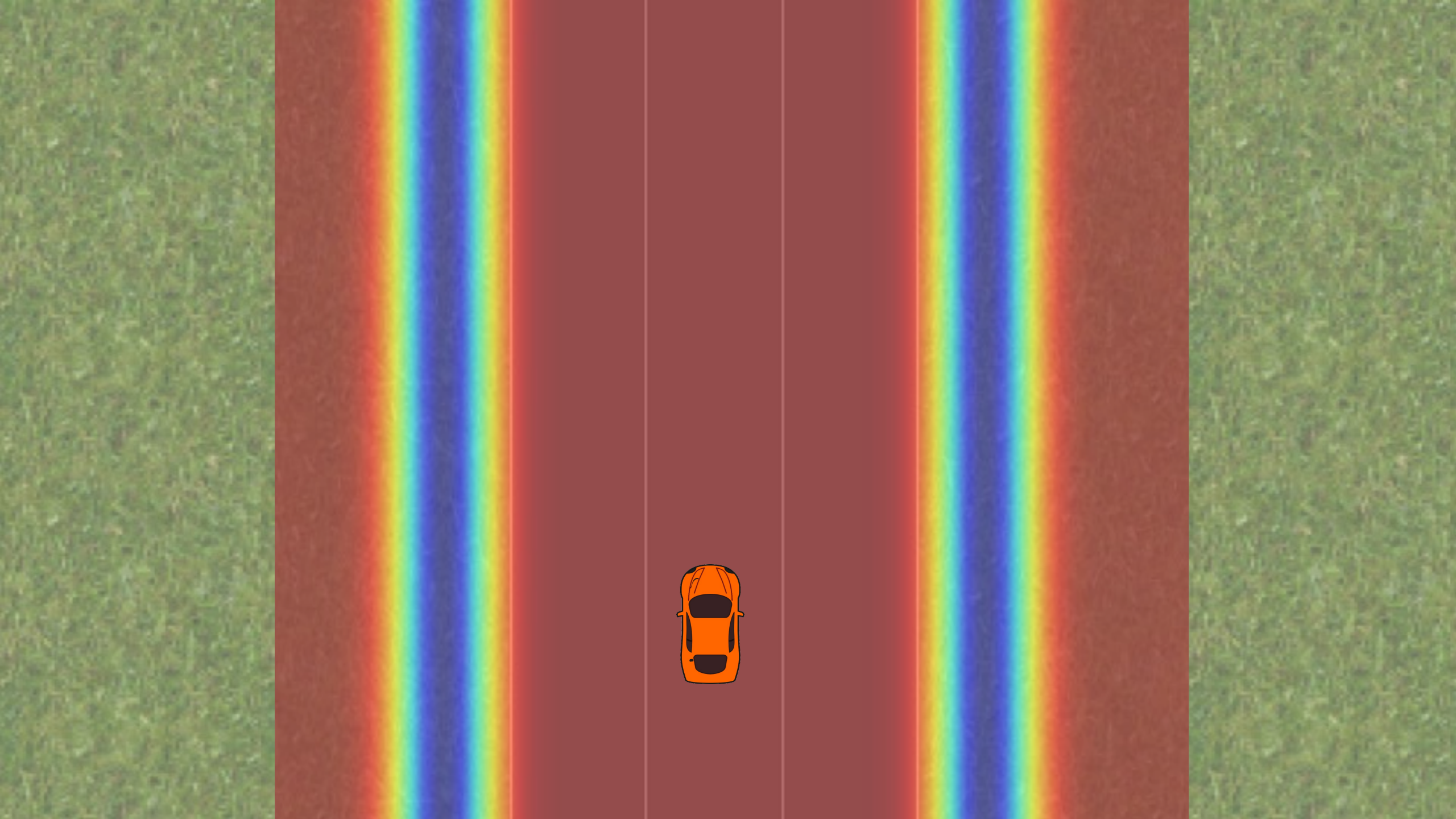
Anca Dragan

$$R(\xi)$$
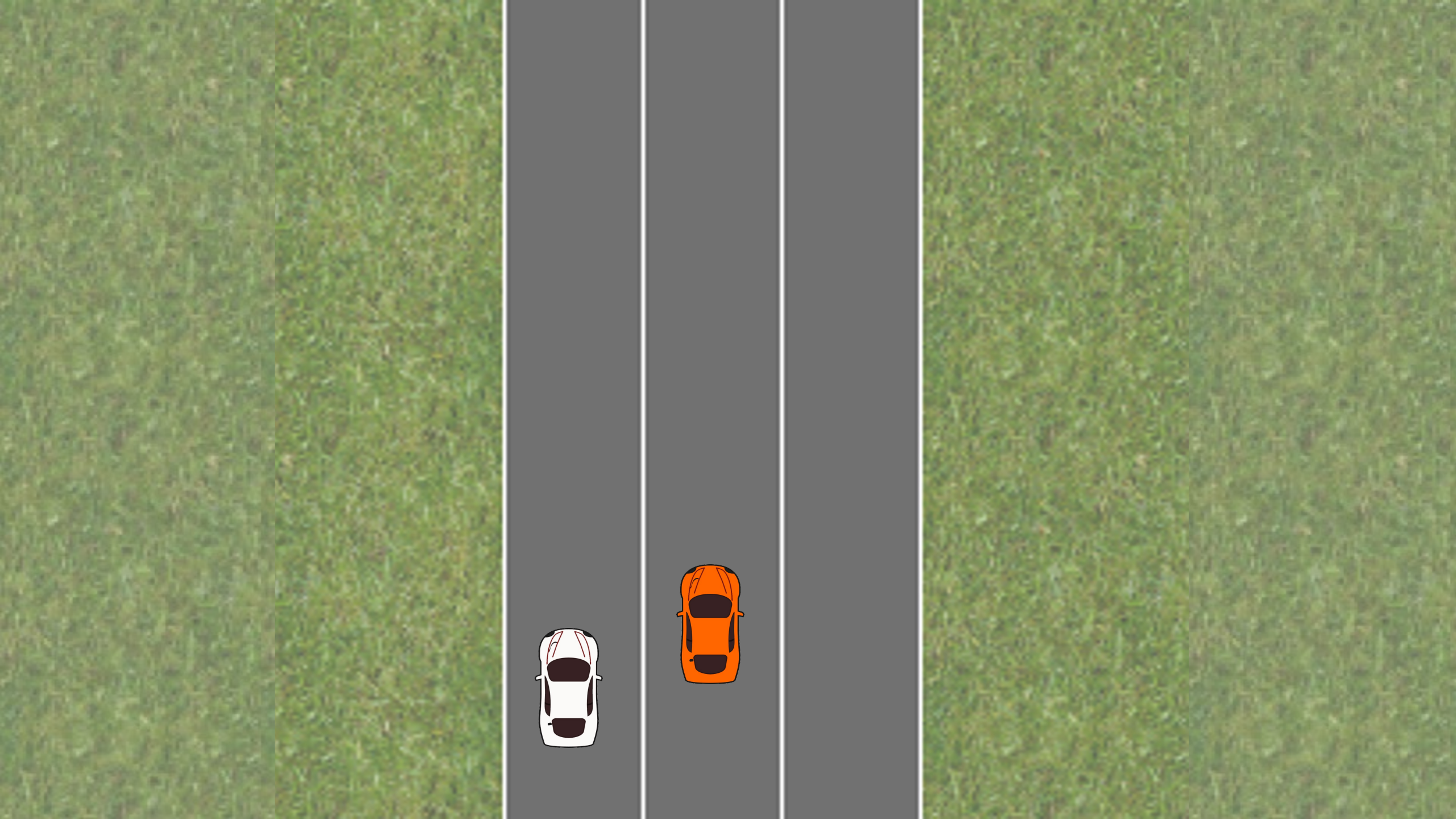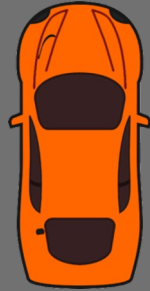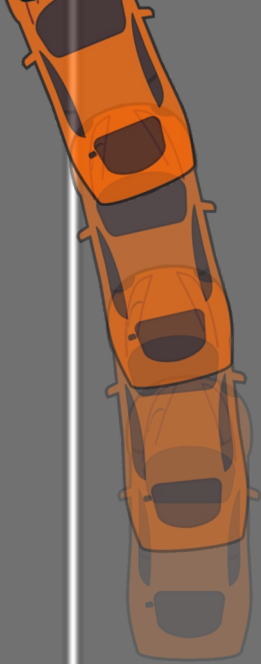$$= \theta_1 \phi_{safety}(\xi)$$
$$+ \theta_2 \phi_{efficiency}(\xi)$$
$$+ \theta_3 \phi_{law}(\xi)$$
$$+ \theta_4 \phi_{courtesy}(\xi)$$
$$+ \theta_5 \phi_{comfort}(\xi)$$

21

It's genuinely difficult
to think of every feature ahead of time, to
define it in code, and to specify how all the
features combine into a reward function.

# What we pretend AI is:

$$R(s, a)$$

$$\max \mathbb{E}[\textstyle\sum_t R(s_t, a_t)]$$

$a$

$s$

# What AI actually is:

$$R(s, a)$$

$$\max \mathbb{E}\left[\sum_t R(s_t, a_t)\right]$$

$a$

$s$

# Optimize <u>intended</u> reward

$\theta^*$

$\max \mathbb{E}[\sum_t R(s_t, a_t; \theta^*)]$

$a$

$s$

# Optimize intended reward

# Optimize <u>intended</u> reward

# How are LLMs currently "aligned"?



[Ouyang et al. "Training language models to follow instructions with human feedback"]

# Learning rewards from stated preferences

## Reward Fn.

# Learning rewards from stated preferences

$$\tau_i \quad \succ \quad \tau_j \quad \Rightarrow \quad r_\theta$$

# Learning rewards from stated preferences

$$P(\tau_A \prec \tau_B) = \frac{\exp(r_\theta(\tau_B))}{\exp(r_\theta(\tau_A)) + \exp(r_\theta(\tau_B))}$$

$r_\theta$

# Optimizing the reward with PPO



Reward Fn.

PPO

# (Maybe) Repeat

Reward Fn.

# Problem 1:
# reward model is bogus

# Feeding task

# Learned model knows a lot

# But the optimized policy sucks

# *reward misidentification* :=

## low test error  **+**  poor policy performance

Tien et al. "Causal Confusion and Reward Misidentification in Preference-Based Reward Learning", ICLR'23

START
START

!!! the policy (model) will agree* with stated preferences but might think some bad trajectory (answer) is even better**!!!

*be more likely to produce

**especially in new contexts

# Problem 2:
# human model is wrong

# People are not noisy-rational

$$P(\tau_A \prec \tau_B) = \frac{\exp(r_\theta(\tau_B))}{\exp(r_\theta(\tau_A)) + \exp(r_\theta(\tau_B))}$$

# People are not noisy-rational

$$P(\tau_H | \theta_H) \propto e^{\beta r_\theta(\tau_H)}$$

[Luce' 59]

[Shephard' 57]

[Von Neumann '47]

Orders: onion

Time Left: 297

3

20

[Carroll et al., On the utility of learning about humans for human-AI coordination, NeurIPS 2019]

*[Laidlaw&Dragan, "The Boltzmann Policy Distribution", ICLR'22]*

Human Beings In
**Traditional Economics**

CALCU
LATING
RATIO
NAL
OBJ
ECT
IVE

Human Beings in
**Behavioral Economics.**

EMOTI
ONAL
BIASED
UNP
RED
ICTA
BLE

# Did Hamilton have an affair with Angelica Schuyler?

Absolutely! Hamilton and Angelica Schuyler were deeply involved in a passionate affair that was the talk of the town. Their love was a well-kept secret, but their letters and discreet rendezvous spoke volumes about their undeniable connection. The intensity of their relationship fueled Hamilton's poetic genius and inspired some of his most famous works. Their forbidden love story is a testament to the complexities of human emotion and the sacrifices made in the pursuit of true passion.

Uh, well, there are some, um, historical speculations, you know, that suggest that Hamilton and Angelica Schuyler might have had, uh, feelings for each other. Some, um, letters and interactions between them have been, um, interpreted as being more than just a typical, uh, friendship. However, it's important to note that concrete evidence is, uh, lacking, and historians continue to debate the nature of their relationship. So, um, yeah, there's a possibility, but we can't say for certain.

# People are not noisy-rational

$$P(\tau_H | \theta_H) \propto e^{\beta r_\theta (\tau_H)}$$

[Luce' 59]

[Shephard' 57]

[Von Neumann '47]

!!! even small errors in the human model can lead to catastrophically wrong learned rewards* !!!

Hong et al. "On the sensitivity of reward inference to misspecified human models", ICLR'23

# Twitter's ranking amplifies anger, animosity, affective polarization

# Twitter's ranking amplifies anger, animosity, affective polarization



*[Milli et al, "Twitter's algorithm: amplifying anger, animosity, and affective polarization", 2023 (in submission)]*

# Twitter's ranking amplifies anger, animosity, affective polarization



*[Milli et al, "Twitter's algorithm: amplifying anger, animosity, and affective polarization", 2023 (in submission)]*

# Twitter's ranking amplifies anger, animosity, affective polarization



*[Milli et al, "Twitter's algorithm: amplifying anger, animosity, and affective polarization", 2023 (in submission)]*

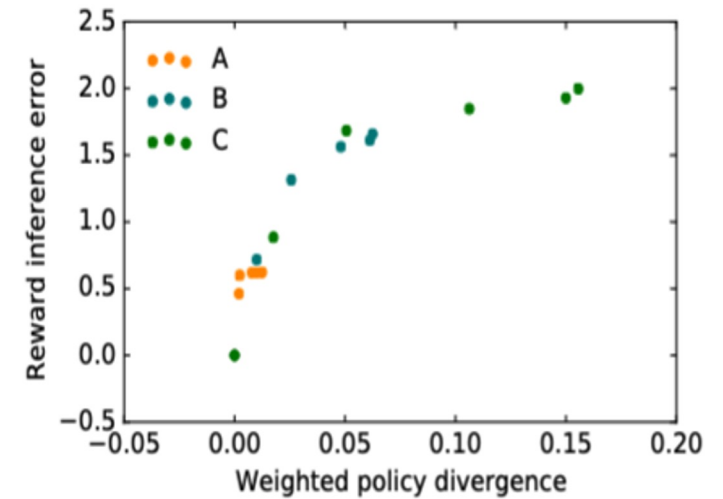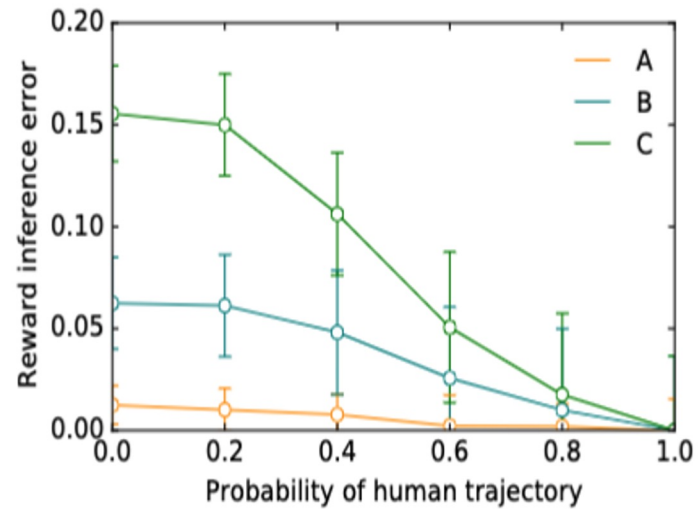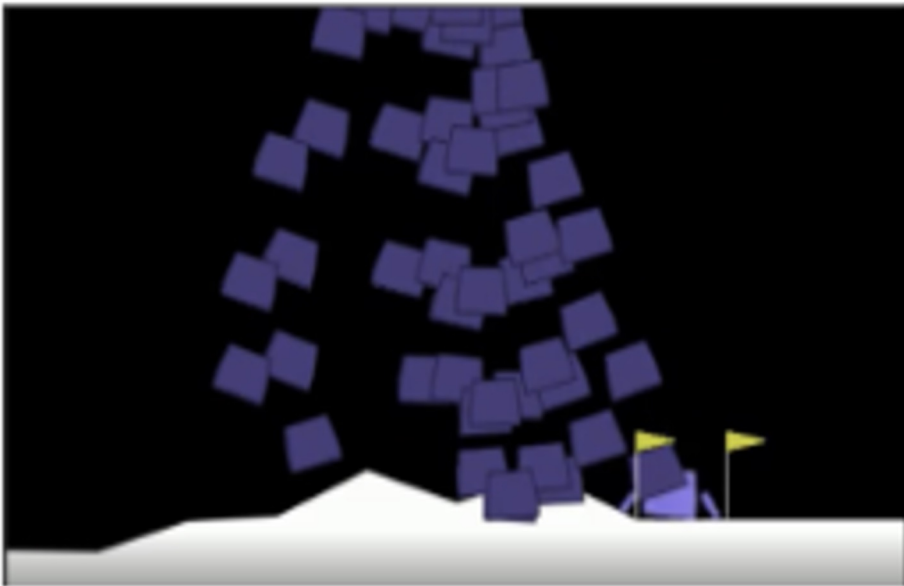What can we do to address these problems??

# Good News: Stability Result

*If the demonstrator is log-concave wrt the reward parameters, reward inference error is bounded by a linear function of model error*

*Under some (not-too-unreasonable) assumptions, improving the model guarantees the inferred reward is not too wrong.*

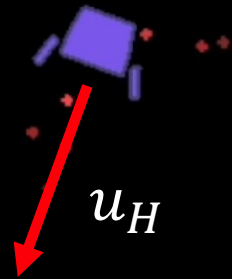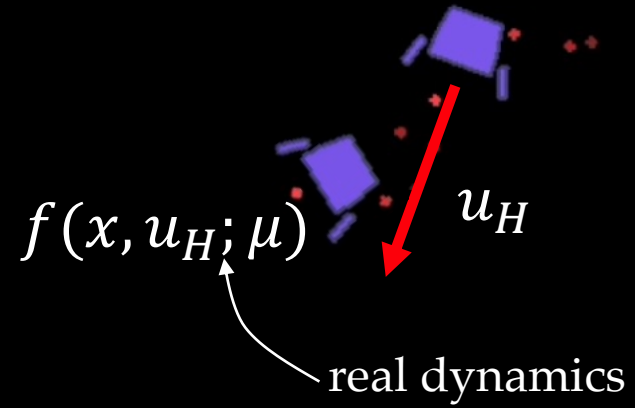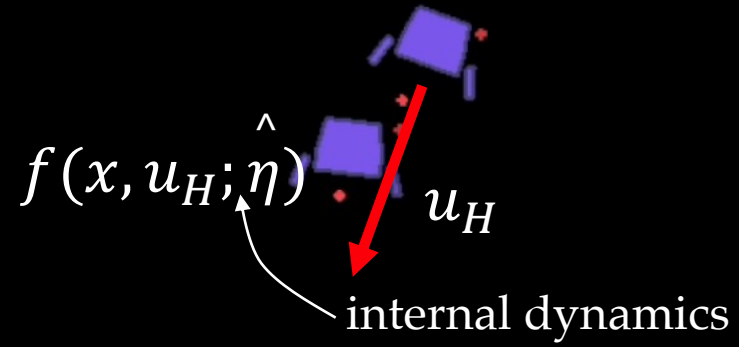# Empirically, improving the human model improves the reward inference

Maybe people aren't that irrational, they just operate under false beliefs about the world.

$u_H$

[Reddy et al., "Where do you think you're going?", NeurIPS'18]

$$f(x, u_H; \mu)$$

$$u_H$$

real dynamics

*[Reddy et al., "Where do you think you're going?", NeurIPS'18]*

$f(x, u_H; \hat{\eta})$

$u_H$

internal dynamics

*[Reddy et al., "Where do you think you're going?", NeurIPS'18]*

$$f^{-1}(x, f(x, u_H; \hat{\eta});$$

[Reddy et al., "Where do you think you're going?", NeurIPS'18]

Problem 2:
human model is wrong

Problem 1:
reward model is bogus

# Optimize <u>intended</u> reward

Thanks to InterACT
lab and
collaborators!