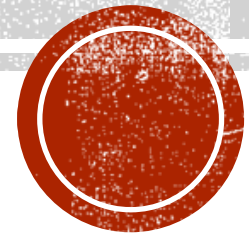# DEEP RL
# FROM HUMAN PREFERENCES

**Paul F Christiano**
OpenAI
paul@openai.com

**Jan Leike**
DeepMind
leike@google.com
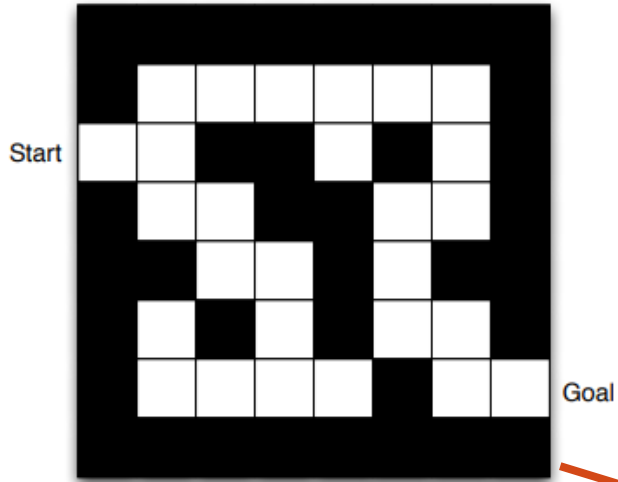
**Tom B Brown**
nottombrown@gmail.com
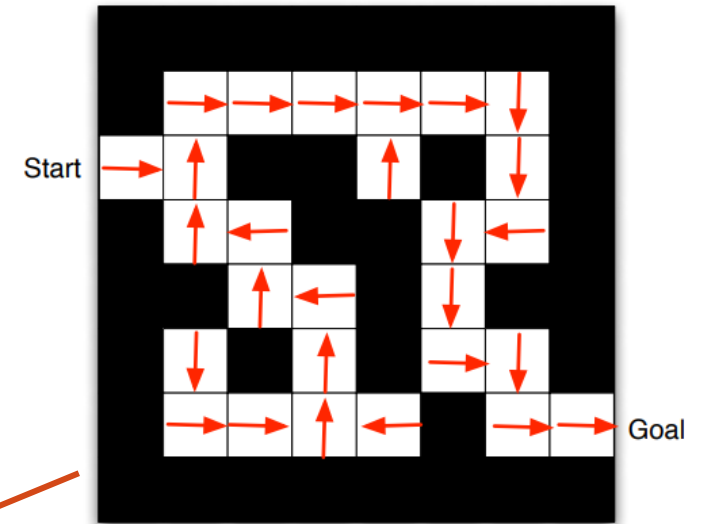
**Miljan Martic**
DeepMind
miljanm@google.com

**Shane Legg**
DeepMind
legg@google.com

**Dario Amodei**
OpenAI
damodei@openai.com

- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: Agent's location
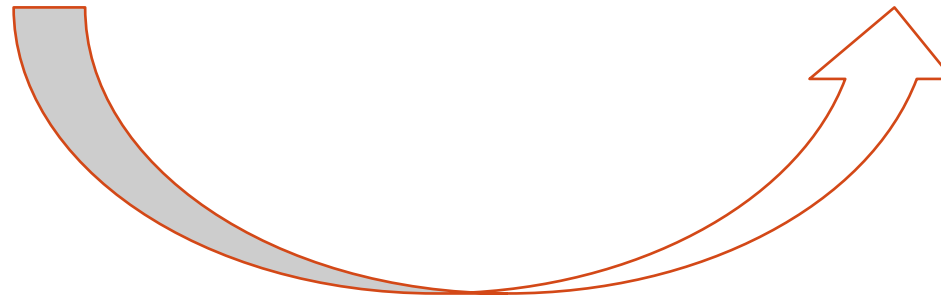
- Arrows represent policy $\pi(s)$ for each state $s$

# RECENT WORK

**Many tasks involve goals that are <span style="color:red">complex</span> and <span style="color:red">hard to specify</span>**

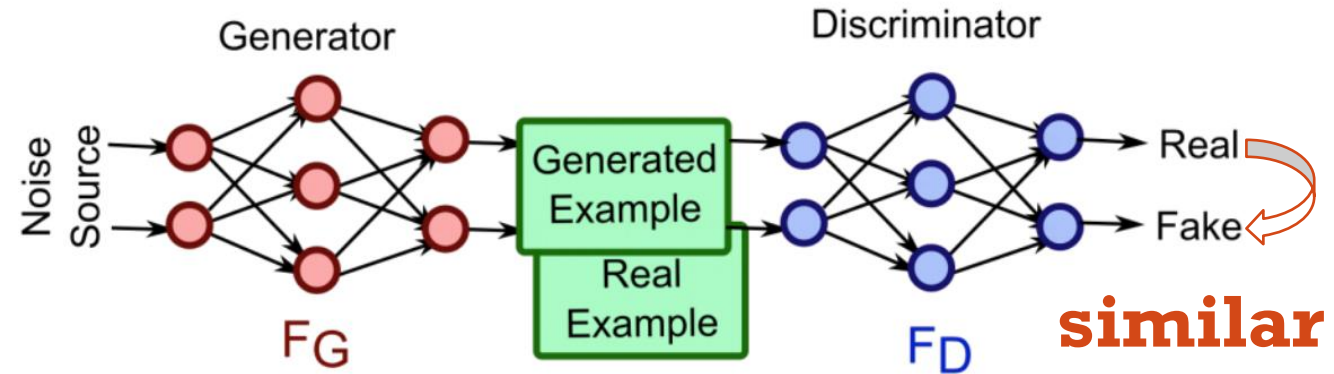design a simple reward function          not satisfy our preference
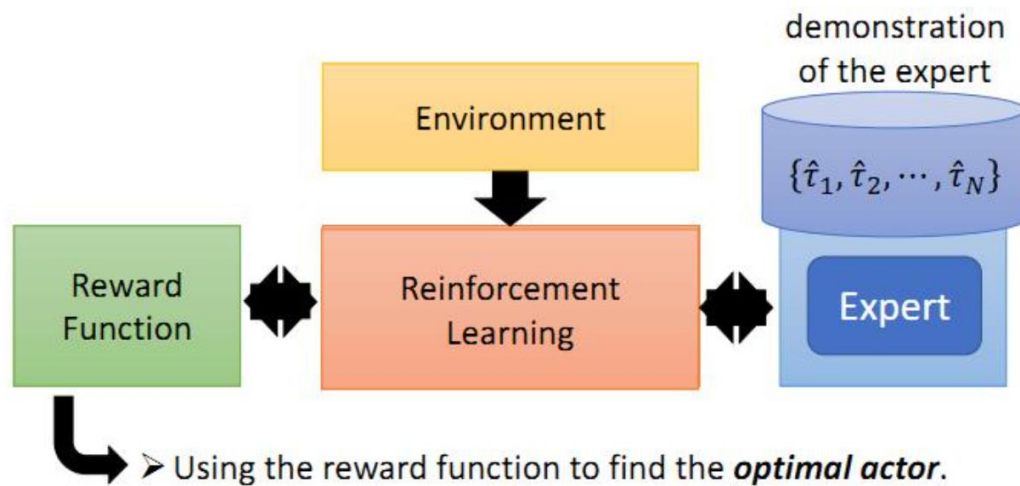
capture the intended behavior

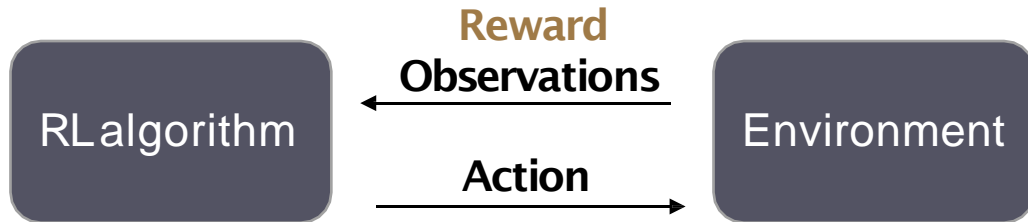→ **misalignment between the value and human preferences**

**The premise : have <span style="color:red">demonstrations of the desired task</span>**
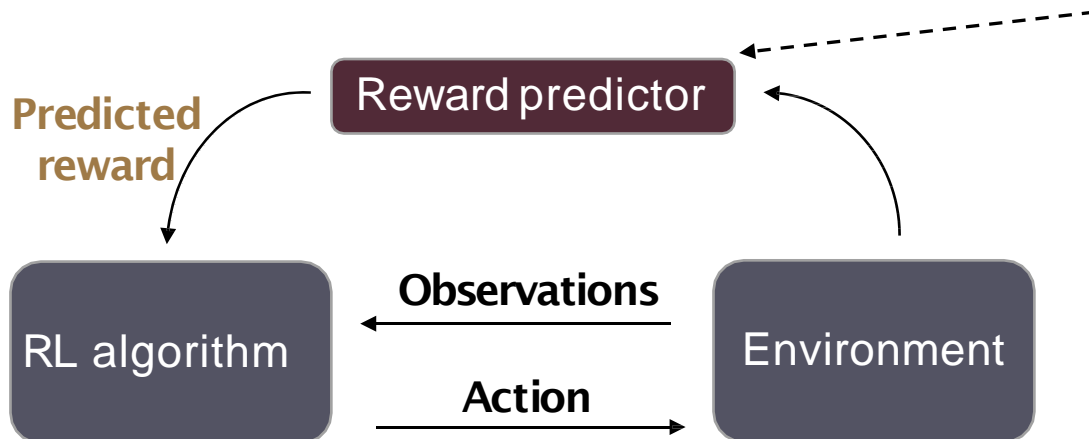


Inverse Reinforcement Learning (IRL)

# MOTIVATION

## Original RL



RL algorithm

**Reward**
**Observations**

Action

Environment

## RL by human feedback



**Predicted reward**

Reward predictor

RL algorithm

Observations

Action

Environment

### Set of Comparisons (Human feedback)



TEXT

Example 1 (√)

Example 2

Learn a reward function from **human feedback**

and then to optimize that **reward function.**

# Objective

Goal: RL agent produces trajectories which are preferred by the human

while making as few queries as possible to human
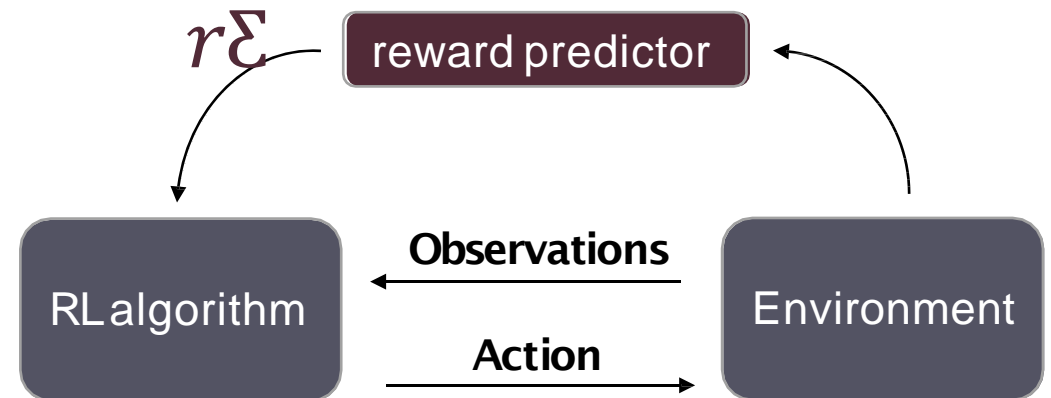
## Two neural networks

1. policy $\pi: O \to A$

2. reward predictor $\hat{r}: O \times A \to \mathbb{R}$

**RL agent** (policy $\pi$) **interacts with the environment** to produce trajectories $\{\tau^1, \dots, \tau^i\}$.
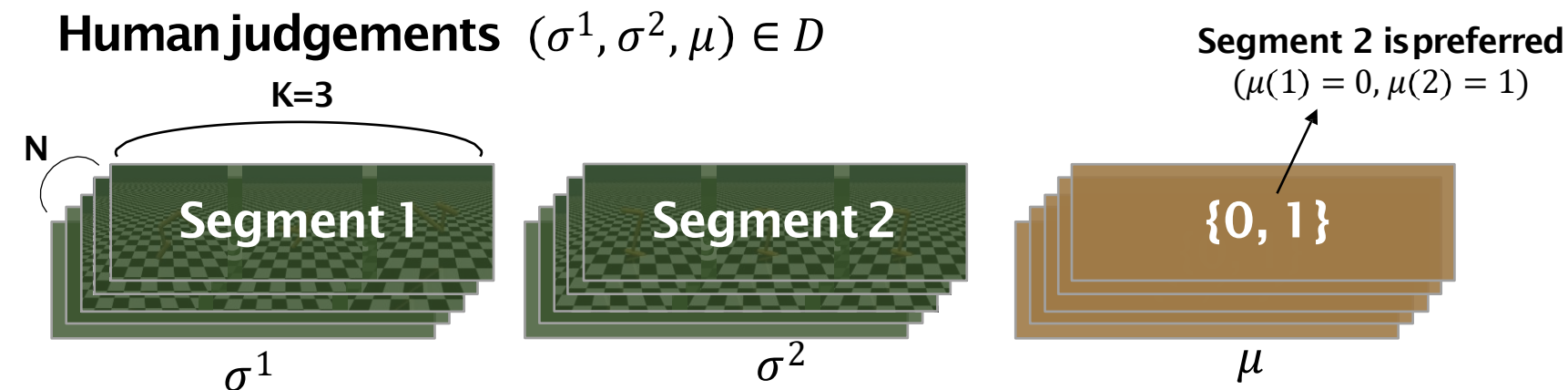
## No environment reward

Trajectory segment :
- $\sigma = ((o_0, a_0), (o_1, a_1), \dots (o_{k-1}, a_{k-1})) \in (O \times A)^k$
- $\sigma^1 > \sigma^2$ : The human preferred trajectory segment $\sigma^1$

$\hat{r}$

reward predictor

RL algorithm

Observations

Action

Environment

# Method

## Reward predictor $r$

### Human judgements $(\sigma^1, \sigma^2, \mu) \in D$

K=3

N

Segment 1

$\sigma^1$

Segment 2

$\sigma^2$

Segment 2 is preferred
$(\mu(1) = 0, \mu(2) = 1)$

{0, 1}
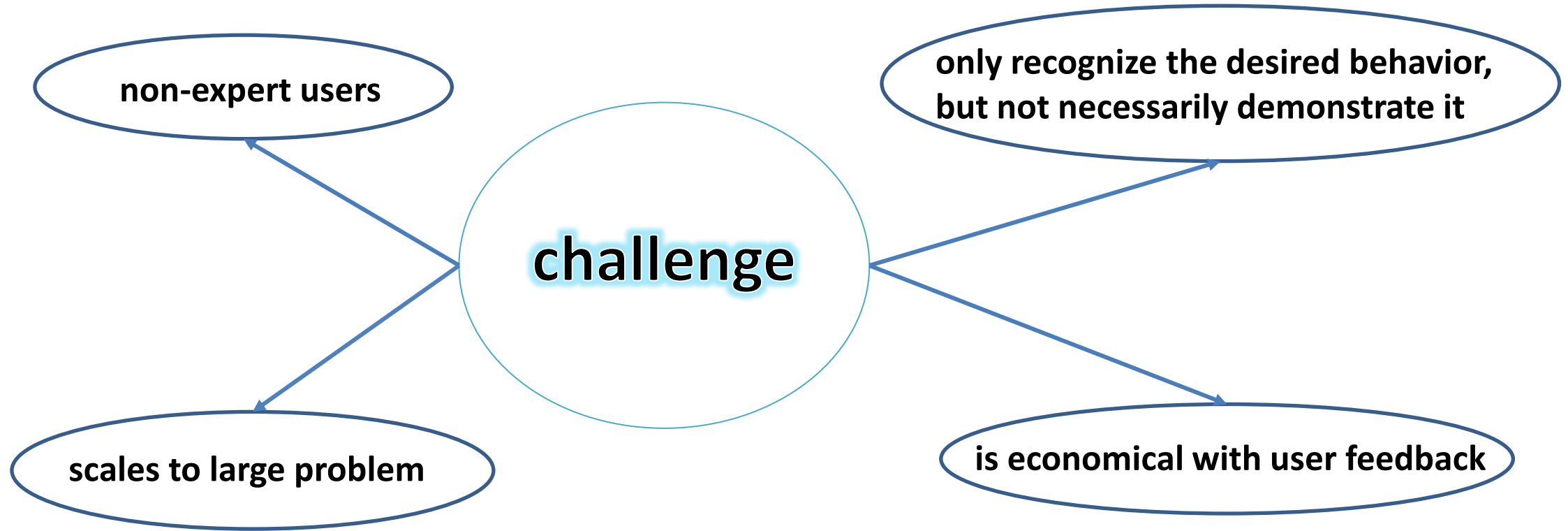
$\mu$

Human's probability of preferring a segment 1

$$\hat{P}(\sigma^1 > \sigma^2) = \frac{\exp \sum_{t=1}^{K} \hat{r}(o_t^1, a_t^1)}{\exp \sum_{t=1}^{K} \hat{r}(o_t^1, a_t^1) + \exp \sum_{t=1}^{K} \hat{r}(o_t^2, a_t^2)}$$
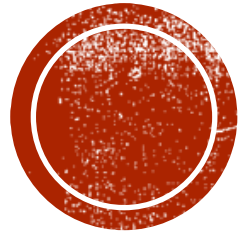
**Bradley-Terry model**

$$loss(\hat{r}) = - \sum_{(\sigma^1, \sigma^2, \mu) \in D} \mu(1) log\hat{P}(\sigma^1 > \sigma^2) + \mu(2) log\hat{P}(\sigma^1 < \sigma^2)$$

minimize cross-entropy between these predictions($\hat{P}$) and the actual human labels ($\mu$)

# Challenge

# THANKS FOR LISTENING