

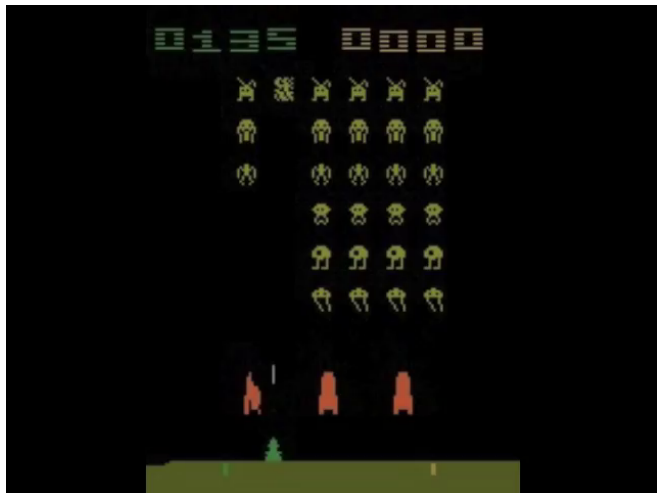
# General Reinforcement Learning

Jan Leike

Future of Humanity Institute  
University of Oxford

9 June 2016

# Reinforcement Learning Today<sup>1</sup>



---

<sup>1</sup>Volodymyr Mnih et al. "Human-Level Control through Deep Reinforcement Learning". In: *Nature* 518.7540 (2015), pp. 529–533.

If we upscale DQN, do we get strong AI?

# Narrow Reinforcement Learning

Atari 2600

---

fully observable

ergodic

very large state space

$\epsilon$ -exploration works

# Narrow Reinforcement Learning

Atari 2600	The Real World™
fully observable	partially observable
ergodic	not ergodic
very large state space	infinite state space
$\epsilon$ -exploration works	$\epsilon$ -exploration fails

# Narrow Reinforcement Learning

Atari 2600	The Real World™
fully observable	partially observable
ergodic	not ergodic
very large state space	infinite state space
$\epsilon$ -exploration works	$\epsilon$ -exploration fails
ergodic MDPs	general environments

# Outline

AIXI

Optimality

Game Theory

AI Safety

# Outline

AIXI

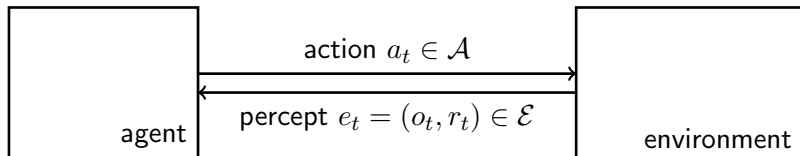
Optimality

Game Theory

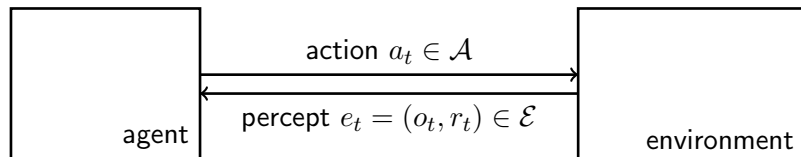
AI Safety



# General Reinforcement Learning



# General Reinforcement Learning



policy

$$\pi : (\mathcal{A} \times \mathcal{E})^* \rightarrow \Delta \mathcal{A}$$

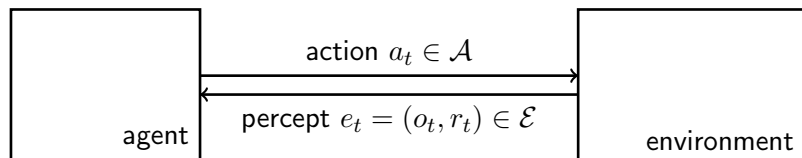
environment

$$\nu : (\mathcal{A} \times \mathcal{E})^* \times \mathcal{A} \rightarrow \Delta \mathcal{E}$$

history

$$\mathfrak{a}_{<t} := a_1 e_1 a_2 e_2 \dots a_{t-1} e_{t-1}$$

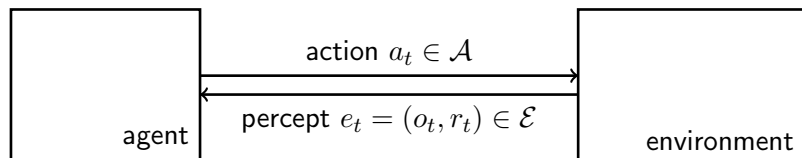
# General Reinforcement Learning



policy  $\pi : (\mathcal{A} \times \mathcal{E})^* \rightarrow \Delta \mathcal{A}$   
environment  $\nu : (\mathcal{A} \times \mathcal{E})^* \times \mathcal{A} \rightarrow \Delta \mathcal{E}$   
history  $\mathfrak{x}_{<t} := a_1 e_1 a_2 e_2 \dots a_{t-1} e_{t-1}$

**Goal:** maximize  $\sum_{t=1}^{\infty} \gamma_t r_t$   
where  $\gamma : \mathbb{N} \rightarrow [0, 1]$  is a discount function with  $\sum_{t=1}^{\infty} \gamma_t < \infty$

# General Reinforcement Learning



policy  $\pi : (\mathcal{A} \times \mathcal{E})^* \rightarrow \Delta \mathcal{A}$   
environment  $\nu : (\mathcal{A} \times \mathcal{E})^* \times \mathcal{A} \rightarrow \Delta \mathcal{E}$   
history  $\mathfrak{a}_{<t} := a_1 e_1 a_2 e_2 \dots a_{t-1} e_{t-1}$

**Goal:** maximize  $\sum_{t=1}^{\infty} \gamma_t r_t$   
where  $\gamma : \mathbb{N} \rightarrow [0, 1]$  is a discount function with  $\sum_{t=1}^{\infty} \gamma_t < \infty$

## Assumptions

- ▶  $0 \leq r_t \leq 1$
- ▶  $\mathcal{A}$  and  $\mathcal{E}$  are finite

# Value Functions

# Value Functions

Value of policy  $\pi$  in environment  $\nu$ :

$$V_{\nu}^{\pi}(\mathbf{a}_{<t}) := \frac{1}{\sum_{k=t}^{\infty} \gamma^k} \mathbb{E}_{\nu}^{\pi} \left[ \sum_{k=t}^{\infty} \gamma^k r_k \mid \mathbf{a}_{<t} \right]$$

# Value Functions

Value of policy  $\pi$  in environment  $\nu$ :

$$V_{\nu}^{\pi}(\mathbf{a}_{<t}) := \frac{1}{\sum_{k=t}^{\infty} \gamma^k} \mathbb{E}_{\nu}^{\pi} \left[ \sum_{k=t}^{\infty} \gamma^k r_k \mid \mathbf{a}_{<t} \right]$$

- ▶ optimal value:  $V_{\nu}^* := \sup_{\pi} V_{\nu}^{\pi}$
- ▶  $\nu$ -optimal policy:  $\pi_{\nu}^* := \arg \max_{\pi} V_{\nu}^{\pi}$

# Value Functions

Value of policy  $\pi$  in environment  $\nu$ :

$$V_{\nu}^{\pi}(\mathbf{ae}_{<t}) := \frac{1}{\sum_{k=t}^{\infty} \gamma^k} \mathbb{E}_{\nu}^{\pi} \left[ \sum_{k=t}^{\infty} \gamma^k r_k \mid \mathbf{ae}_{<t} \right]$$

- ▶ optimal value:  $V_{\nu}^* := \sup_{\pi} V_{\nu}^{\pi}$
- ▶  $\nu$ -optimal policy:  $\pi_{\nu}^* := \arg \max_{\pi} V_{\nu}^{\pi}$
- ▶ Effective horizon:

$$H_t(\varepsilon) := \min \left\{ k \mid \frac{\sum_{i=t+k}^{\infty} \gamma^i}{\sum_{i=t}^{\infty} \gamma^i} \leq \varepsilon \right\}$$



---

<sup>2</sup>Ray Solomonoff. "A Formal Theory of Inductive Inference. Parts 1 and 2". In: *Information and Control* 7.1 (1964), pages.

<sup>3</sup>Marcus Hutter. *Universal Artificial Intelligence. Sequential Decisions Based on Algorithmic Probability*. Springer, 2005.

- ▶ countable set of environments  $\mathcal{M} = \{\nu_1, \nu_2, \dots\}$

---

<sup>2</sup>Ray Solomonoff. "A Formal Theory of Inductive Inference. Parts 1 and 2". In: *Information and Control* 7.1 (1964), pages.

<sup>3</sup>Marcus Hutter. *Universal Artificial Intelligence. Sequential Decisions Based on Algorithmic Probability*. Springer, 2005.

- ▶ countable set of environments  $\mathcal{M} = \{\nu_1, \nu_2, \dots\}$
- ▶ prior  $w : \mathcal{M} \rightarrow [0, 1]$

---

<sup>2</sup>Ray Solomonoff. "A Formal Theory of Inductive Inference. Parts 1 and 2". In: *Information and Control* 7.1 (1964), pages.

<sup>3</sup>Marcus Hutter. *Universal Artificial Intelligence. Sequential Decisions Based on Algorithmic Probability*. Springer, 2005.

- ▶ countable set of environments  $\mathcal{M} = \{\nu_1, \nu_2, \dots\}$
- ▶ prior  $w : \mathcal{M} \rightarrow [0, 1]$   
Solomonoff prior<sup>2</sup>  $w(\nu) \propto 2^{-K(\nu)}$

---

<sup>2</sup>Ray Solomonoff. "A Formal Theory of Inductive Inference. Parts 1 and 2". In: *Information and Control* 7.1 (1964), pages.

<sup>3</sup>Marcus Hutter. *Universal Artificial Intelligence. Sequential Decisions Based on Algorithmic Probability*. Springer, 2005.

- ▶ countable set of environments  $\mathcal{M} = \{\nu_1, \nu_2, \dots\}$
- ▶ prior  $w : \mathcal{M} \rightarrow [0, 1]$   
Solomonoff prior<sup>2</sup>  $w(\nu) \propto 2^{-K(\nu)}$
- ▶ Bayesian mixture

$$\xi := \sum_{\nu \in \mathcal{M}} w(\nu) \nu$$

---

<sup>2</sup>Ray Solomonoff. "A Formal Theory of Inductive Inference. Parts 1 and 2". In: *Information and Control* 7.1 (1964), pages.

<sup>3</sup>Marcus Hutter. *Universal Artificial Intelligence. Sequential Decisions Based on Algorithmic Probability*. Springer, 2005.

- ▶ countable set of environments  $\mathcal{M} = \{\nu_1, \nu_2, \dots\}$
- ▶ prior  $w : \mathcal{M} \rightarrow [0, 1]$   
Solomonoff prior<sup>2</sup>  $w(\nu) \propto 2^{-K(\nu)}$
- ▶ Bayesian mixture

$$\xi := \sum_{\nu \in \mathcal{M}} w(\nu) \nu$$

AIXI is the Bayes-optimal agent with a Solomonoff prior

$$\pi_{\xi}^* := \arg \max_{\pi} V_{\xi}^{\pi}$$

---

<sup>2</sup>Ray Solomonoff. "A Formal Theory of Inductive Inference. Parts 1 and 2". In: *Information and Control* 7.1 (1964), pages.

<sup>3</sup>Marcus Hutter. *Universal Artificial Intelligence. Sequential Decisions Based on Algorithmic Probability*. Springer, 2005.

# On-Policy Value Convergence

$$V_{\xi}^{\pi}(\mathbf{a}_{<t}) - V_{\mu}^{\pi}(\mathbf{a}_{<t}) \rightarrow 0 \text{ as } t \rightarrow \infty \text{ almost surely}$$

# Outline

AIXI

Optimality

Game Theory

AI Safety



# Notions of Optimality in Reinforcement Learning

- ▶ Bayes optimality
- ▶ Asymptotic optimality
- ▶ Sample complexity bounds
- ▶ Regret bounds
- ▶ ...

# Asymptotic Optimality

$\pi$  is asymptotically optimal iff

$$V_{\mu}^*(\mathfrak{a}_{<t}) - V_{\mu}^{\pi}(\mathfrak{a}_{<t}) \rightarrow 0 \text{ as } t \rightarrow \infty$$

---

<sup>4</sup>Laurent Orseau. "Asymptotic Non-Learnability of Universal Agents with Computable Horizon Functions". In: *Theoretical Computer Science* 473 (2013), pp. 149–156.

# Asymptotic Optimality

$\pi$  is asymptotically optimal iff

$$V_{\mu}^*(\mathfrak{a}_{<t}) - V_{\mu}^{\pi}(\mathfrak{a}_{<t}) \rightarrow 0 \text{ as } t \rightarrow \infty$$

For asymptotic optimality the agent needs to explore infinitely often for an entire effective horizon.

---

<sup>4</sup>Laurent Orseau. "Asymptotic Non-Learnability of Universal Agents with Computable Horizon Functions". In: *Theoretical Computer Science* 473 (2013), pp. 149–156.

# Asymptotic Optimality

$\pi$  is asymptotically optimal iff

$$V_{\mu}^*(\mathfrak{a}_{<t}) - V_{\mu}^{\pi}(\mathfrak{a}_{<t}) \rightarrow 0 \text{ as } t \rightarrow \infty$$

For asymptotic optimality the agent needs to explore infinitely often for an entire effective horizon.

## Theorem

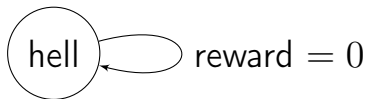
*AIXI is not asymptotically optimal.*<sup>4</sup>

---

<sup>4</sup>Laurent Orseau. "Asymptotic Non-Learnability of Universal Agents with Computable Horizon Functions". In: *Theoretical Computer Science* 473 (2013), pp. 149–156.

Hell

# Hell



# The Dogmatic Prior<sup>5</sup>

---

<sup>5</sup>Jan Leike and Marcus Hutter. “Bad Universal Priors and Notions of Optimality”. In: *Conference on Learning Theory*. 2015, pp. 1244–1259.

# The Dogmatic Prior<sup>5</sup>

Policy  $\pi_{Lazy}$ :

```
while (true) { do_nothing(); }
```

---

<sup>5</sup>Jan Leike and Marcus Hutter. "Bad Universal Priors and Notions of Optimality". In: *Conference on Learning Theory*. 2015, pp. 1244–1259.



# The Dogmatic Prior<sup>5</sup>

Policy  $\pi_{Lazy}$ :

```
while (true) { do_nothing(); }
```

Dogmatic prior  $\xi'$ :

if not acting according to  $\pi_{Lazy}$ ,  
go to hell with high probability

---

<sup>5</sup>Jan Leike and Marcus Hutter. "Bad Universal Priors and Notions of Optimality". In: *Conference on Learning Theory*. 2015, pp. 1244–1259.

# The Dogmatic Prior<sup>5</sup>

Policy  $\pi_{Lazy}$ :

```
while (true) { do_nothing(); }
```

Dogmatic prior  $\xi'$ :

if not acting according to  $\pi_{Lazy}$ ,  
go to hell with high probability

## Theorem

$\forall \varepsilon > 0 \exists \xi'$  s.t.  $AI_{\xi'}$  acts according to  $\pi_{Lazy}$  as long as  
 $V_{\xi}^{\pi_{Lazy}}(\mathbf{a}_{<t}) > \varepsilon > 0$ .

---

<sup>5</sup>Jan Leike and Marcus Hutter. "Bad Universal Priors and Notions of Optimality". In: *Conference on Learning Theory*. 2015, pp. 1244–1259.

# Thompson Sampling

Thompson sampling policy  $\pi_T$ :

*Sample  $\rho \sim w(\cdot \mid \mathbf{x}_{<t})$ .*

*Follow  $\pi_\rho^*$  for  $H_t(\varepsilon_t)$  steps.*

*Repeat.*

with  $\varepsilon_t \rightarrow 0$ .

---

<sup>6</sup>Jan Leike et al. "Thompson Sampling is Asymptotically Optimal in General Environments". In: *Uncertainty in Artificial Intelligence*. 2016.

# Thompson Sampling

Thompson sampling policy  $\pi_T$ :

*Sample  $\rho \sim w(\cdot \mid \mathbf{x}_{<t})$ .*

*Follow  $\pi_\rho^*$  for  $H_t(\varepsilon_t)$  steps.*

*Repeat.*

with  $\varepsilon_t \rightarrow 0$ .

## Theorem

*Thompson sampling is asymptotically optimal in mean:<sup>6</sup>*

$$\mathbb{E}_\mu^{\pi_T} [V_\mu^*(\mathbf{x}_{<t}) - V_\mu^{\pi_T}(\mathbf{x}_{<t})] \rightarrow 0 \text{ as } t \rightarrow \infty.$$

---

<sup>6</sup>Jan Leike et al. "Thompson Sampling is Asymptotically Optimal in General Environments". In: *Uncertainty in Artificial Intelligence*. 2016.

# Recoverable Environments

An environment  $\nu$  is **recoverable** iff

$$\sup_{\pi} \left| \mathbb{E}_{\nu}^{\pi^*} [V_{\nu}^*(\mathbf{a}_{<t})] - \mathbb{E}_{\nu}^{\pi} [V_{\nu}^*(\mathbf{a}_{<t})] \right| \rightarrow 0 \text{ as } t \rightarrow \infty.$$

# Recoverable Environments

An environment  $\nu$  is **recoverable** iff

$$\sup_{\pi} \left| \mathbb{E}_{\nu}^{\pi^*} [V_{\nu}^*(\mathbf{a}_{<t})] - \mathbb{E}_{\nu}^{\pi} [V_{\nu}^*(\mathbf{a}_{<t})] \right| \rightarrow 0 \text{ as } t \rightarrow \infty.$$

For non-recoverable environments:

Either the agent gets caught in a trap  
or it is not asymptotically optimal.

# Regret

$$R_m(\pi, \mu) := \max_{\pi'} \mathbb{E}_{\mu}^{\pi'} \left[ \sum_{t=1}^m r_t \right] - \mathbb{E}_{\mu}^{\pi} \left[ \sum_{t=1}^m r_t \right]$$

# Regret

$$R_m(\pi, \mu) := \max_{\pi'} \mathbb{E}_{\mu}^{\pi'} \left[ \sum_{t=1}^m r_t \right] - \mathbb{E}_{\mu}^{\pi} \left[ \sum_{t=1}^m r_t \right]$$

A problem class is *learnable* iff  $\exists \pi \forall \mu R_m(\pi, \mu) \in o(m)$ .



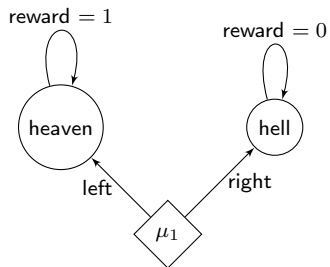
# Regret

$$R_m(\pi, \mu) := \max_{\pi'} \mathbb{E}_{\mu}^{\pi'} \left[ \sum_{t=1}^m r_t \right] - \mathbb{E}_{\mu}^{\pi} \left[ \sum_{t=1}^m r_t \right]$$

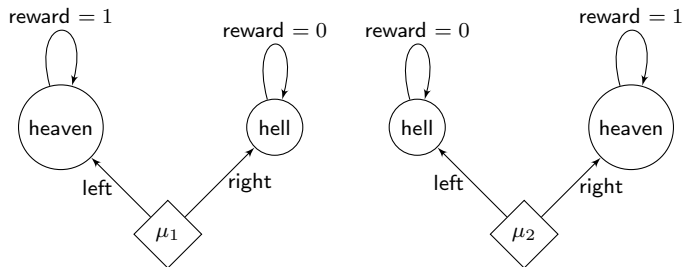
A problem class is *learnable* iff  $\exists \pi \forall \mu R_m(\pi, \mu) \in o(m)$ .

**Fact:** The general RL problem is *not* learnable.

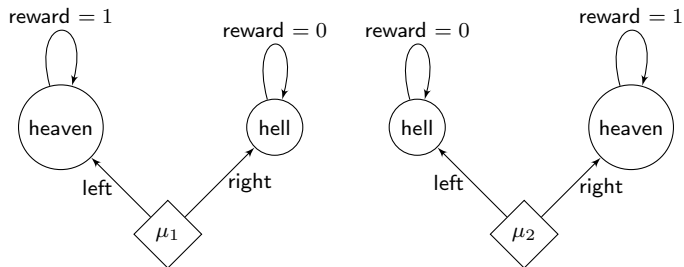
# Regret in Non-Recoverable Environments



# Regret in Non-Recoverable Environments



# Regret in Non-Recoverable Environments



$$R_m(\text{left}, \mu_1) = 0$$
$$R_m(\text{right}, \mu_1) = m$$

$$R_m(\text{left}, \mu_2) = m$$
$$R_m(\text{right}, \mu_2) = 0$$

# Sublinear Regret

## Theorem

*If*

- ▶  $\mu \in \mathcal{M}$  is recoverable,
- ▶  $\pi$  is asymptotically optimal in mean, and
- ▶  $\gamma$  satisfies some weak assumptions,

*then regret is sublinear.*<sup>7</sup>

---

<sup>7</sup>Jan Leike et al. "Thompson Sampling is Asymptotically Optimal in General Environments". In: *Uncertainty in Artificial Intelligence*. 2016.

# Optimality Summary

	AIXI	TS	All policies
Sublinear regret	×	recoverable	×
Sample complexity	×	?	
Pareto optimality	✓	✓	✓
Bayes optimality	✓	×	
Asymptotic optimality	×	✓	

# Outline

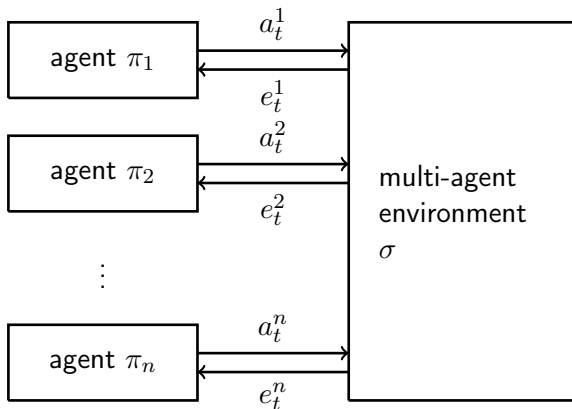
AIXI

Optimality

Game Theory

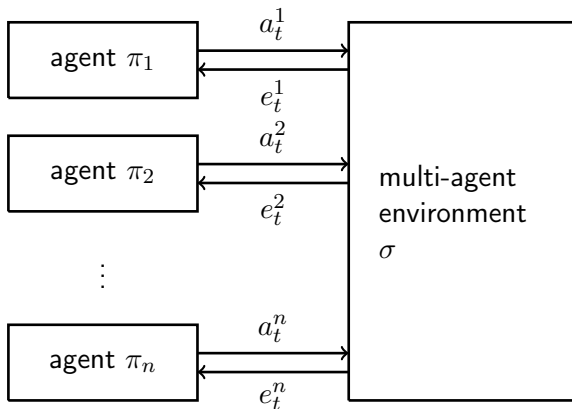
AI Safety

# Multi-Agent Environments



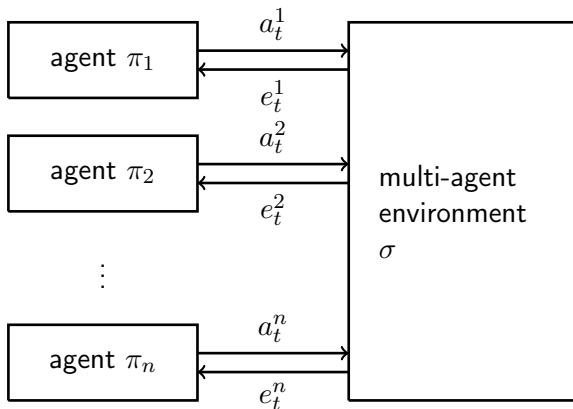


# Multi-Agent Environments



- ▶  $\pi_i$  is an  $\epsilon$ -best response iff  $V_{\sigma_i}^* - V_{\sigma_i}^{\pi_i} < \epsilon$

# Multi-Agent Environments



- ▶  $\pi_i$  is an  $\varepsilon$ -best response iff  $V_{\sigma_i}^* - V_{\sigma_i}^{\pi_i} < \varepsilon$
- ▶  $\pi_1, \dots, \pi_n$  play an  $\varepsilon$ -Nash equilibrium iff each  $\pi_i$  is an  $\varepsilon$ -best response

# The Bayesian Approach

- ▶ countable set of policies  $\Pi$

# The Bayesian Approach

- ▶ countable set of policies  $\Pi$
- ▶ prior  $w \in \Delta\Pi$

# The Bayesian Approach

- ▶ countable set of policies  $\Pi$
- ▶ prior  $w \in \Delta\Pi$
- ▶ act Bayes-optimal with respect to  $w$

# The Bayesian Approach

- ▶ countable set of policies  $\Pi$
- ▶ prior  $w \in \Delta\Pi$
- ▶ act Bayes-optimal with respect to  $w$

**Grain of Truth:** the Bayes-optimal policy needs to be in  $\Pi$

# Results for Bayesian Agents

## Theorem

*If each player is Bayesian, knows the infinite repeated game and has a grain of truth, then the players converge to an  $\epsilon$ -Nash equilibrium.*<sup>8</sup>

---

<sup>8</sup>Ehud Kalai and Ehud Lehrer. "Rational Learning Leads to Nash Equilibrium". In: *Econometrica* (1993), pp. 1019–1045.

<sup>9</sup>Jan Leike, Jessica Taylor, and Benya Fallenstein. "A Formal Solution to the Grain of Truth Problem". In: *Uncertainty in Artificial Intelligence*. 2016.

# Results for Bayesian Agents

## Theorem

*If each player is Bayesian, knows the infinite repeated game and has a grain of truth, then the players converge to an  $\epsilon$ -Nash equilibrium.<sup>8</sup>*

## Theorem

*Two Bayesian players playing infinite repeated matching pennies may fail to converge to an  $\epsilon$ -Nash equilibrium, **even if they have a grain of truth.**<sup>9</sup>*

---

<sup>8</sup>Ehud Kalai and Ehud Lehrer. "Rational Learning Leads to Nash Equilibrium". In: *Econometrica* (1993), pp. 1019–1045.

<sup>9</sup>Jan Leike, Jessica Taylor, and Benya Fallenstein. "A Formal Solution to the Grain of Truth Problem". In: *Uncertainty in Artificial Intelligence*. 2016.



# Solving the Grain of Truth Problem<sup>10</sup>

## Theorem

*There is a class of environments  $\mathcal{M}_{refl}$  that contains a grain of truth with respect to any computable priors' Bayes-optimal policies in any computable multi-agent environment.*

---

<sup>10</sup>Jan Leike, Jessica Taylor, and Benya Fallenstein. "A Formal Solution to the Grain of Truth Problem". In: *Uncertainty in Artificial Intelligence*. 2016.

# Solving the Grain of Truth Problem<sup>10</sup>

## Theorem

*There is a class of environments  $\mathcal{M}_{refl}$  that contains a grain of truth with respect to any computable priors' Bayes-optimal policies in any computable multi-agent environment.*

## Theorem

*Each  $\nu \in \mathcal{M}_{refl}$  is limit computable.*

---

<sup>10</sup>Jan Leike, Jessica Taylor, and Benya Fallenstein. "A Formal Solution to the Grain of Truth Problem". In: *Uncertainty in Artificial Intelligence*. 2016.

# Solving the Grain of Truth Problem<sup>10</sup>

## Theorem

*There is a class of environments  $\mathcal{M}_{refl}$  that contains a grain of truth with respect to any computable priors' Bayes-optimal policies in any computable multi-agent environment.*

## Theorem

*Each  $\nu \in \mathcal{M}_{refl}$  is limit computable.*

## Theorem

*There are limit computable policies  $\pi_1, \dots, \pi_n$  such that for any computable multi-agent environment  $\sigma$  and for all  $\varepsilon > 0$  and all  $i \in \{1, \dots, n\}$  the probability that the policy  $\pi_i$  is an  $\varepsilon$ -best response converges to 1 as  $t \rightarrow \infty$ .*

---

<sup>10</sup>Jan Leike, Jessica Taylor, and Benya Fallenstein. "A Formal Solution to the Grain of Truth Problem". In: *Uncertainty in Artificial Intelligence*. 2016.

# Outline

AIXI

Optimality

Game Theory

AI Safety

# AI Safety Approaches

bottom-up

top-down

---

practical algorithms

toy models

demos

# AI Safety Approaches

bottom-up	top-down
practical algorithms	theoretical models
toy models	abstract problems
demos	theorems

## “Applications” of GRL to AI Safety

- ▶ self-modification: Orseau and Ring (2011), Orseau and Ring (2012), Everitt et al. (2016)
- ▶ self-reflection: Fallenstein, Soares, and Taylor (2015), Leike, Taylor, and Fallenstein (2016)
- ▶ memory manipulation: Orseau and Ring (2012)
- ▶ interruptibility: Orseau and Armstrong (2016)
- ▶ decision theory: Everitt, Leike, and Hutter (2015)
- ▶ wireheading: Ring and Orseau (2011), Everitt and Hutter (2016)
- ▶ value learning: Dewey (2011)
- ▶ questions of identity: Orseau (2014)

# Limits of the Current Model

- ▶ model-based
- ▶ dualistic
- ▶ not self-improving
- ▶ assumes infinite computation



# Conclusion

Mathematical and mental tools to think about strong AI

- ▶ exploration vs. exploitation
- ▶ effective horizon
- ▶ on-policy vs. off-policy
- ▶ model-based vs. model-free
- ▶ recoverability
- ▶ asymptotic optimality
- ▶ reflective oracles

Jan Leike. “Nonparametric General Reinforcement Learning”.  
PhD thesis. Australian National University, 2016

<http://jan.leike.name/>