

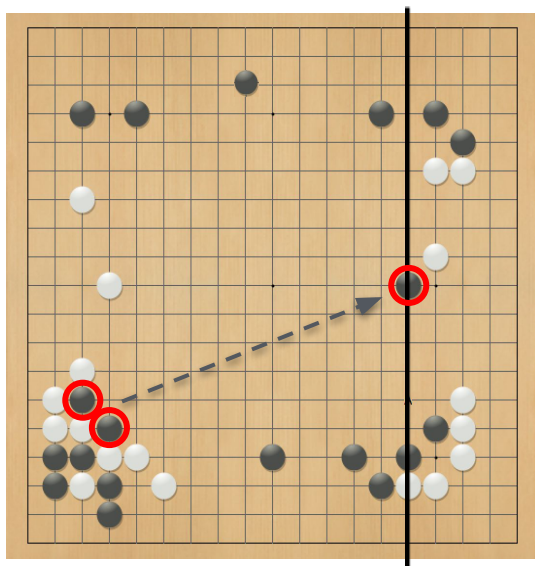
Scalable agent alignment

Jan Leike · BAGI 2019

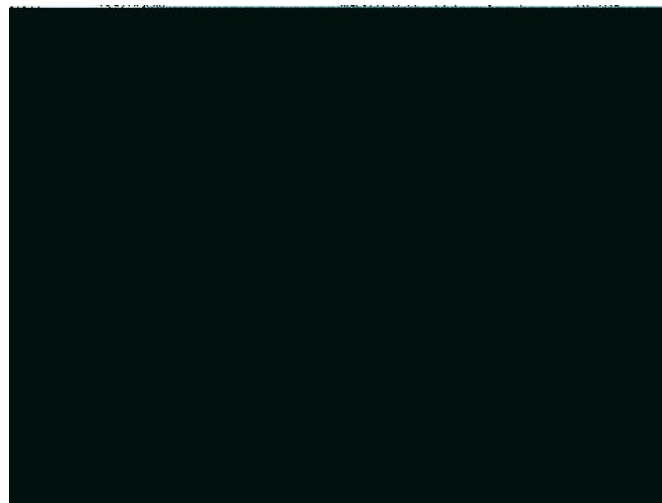
What we want from ML

move 37

AlphaGo ●
Lee Sedol ○



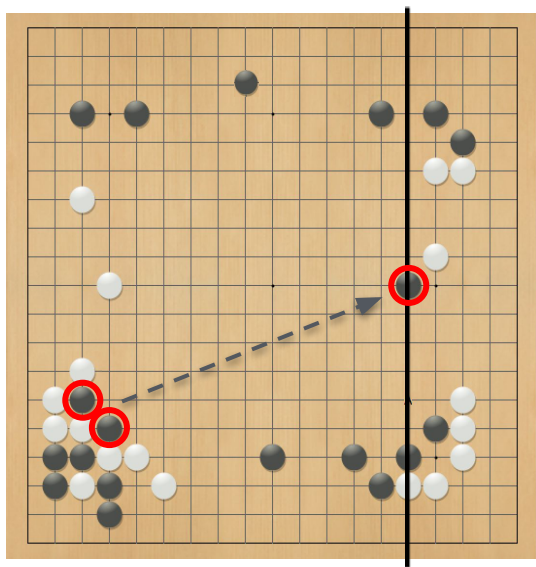
circling boat



What we want from ML

move 37

AlphaGo ●
Lee Sedol ○



circling boat



The agent alignment problem

How can we create agents that **behave** in accordance with the user's intentions?

“Preference payload” questions

- Whose preferences should the agent be aligned to?
- How should preferences of different users be aggregated?
- How should they be traded off against each other?
- When should the agent be disobedient?

“Preference payload” questions

- Whose preferences should the agent be aligned to?
- How should preferences of different agents be aggregated?
- How should preferences be handled off agent to agent?
- When should the agent be disobedient?

These questions are **important**.

We’re **not discussing** these questions here.

We’re only considering the **technical problem** of aligning **one agent to one user**.

Desiderata

Economical



Scalable



Image sources:
<https://www.porttechnology.org/>
<https://realanimetraining.com/>

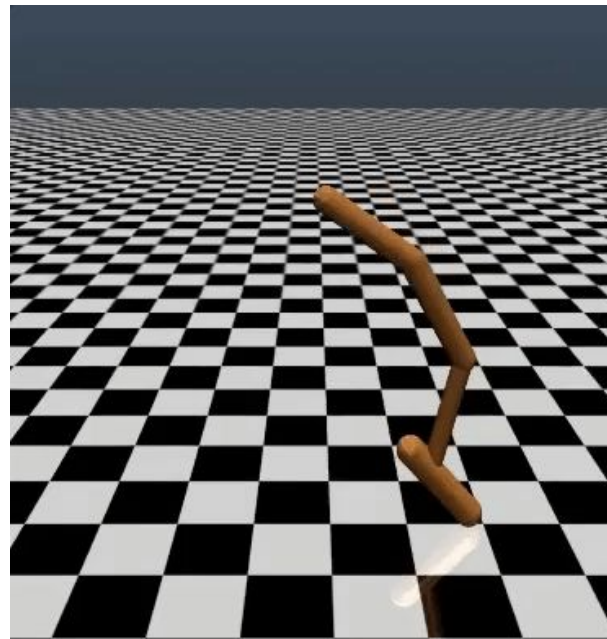
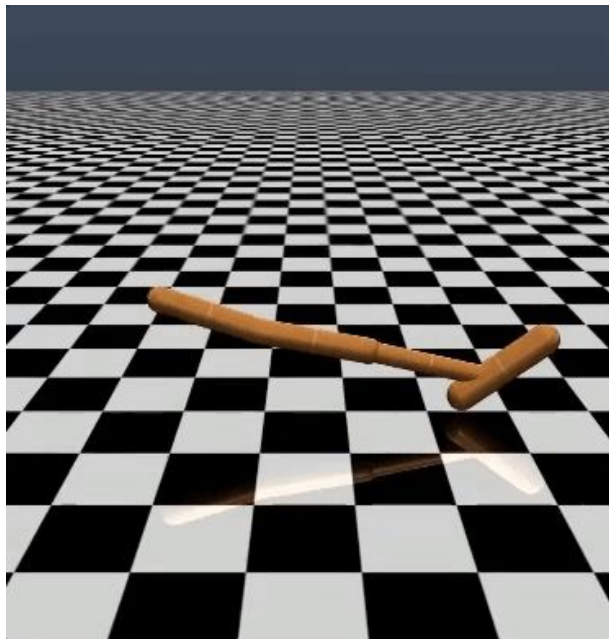
Assumption 1

Rather than **formally specifying** user intentions, we can instead **learn** these intentions to a sufficiently high accuracy.

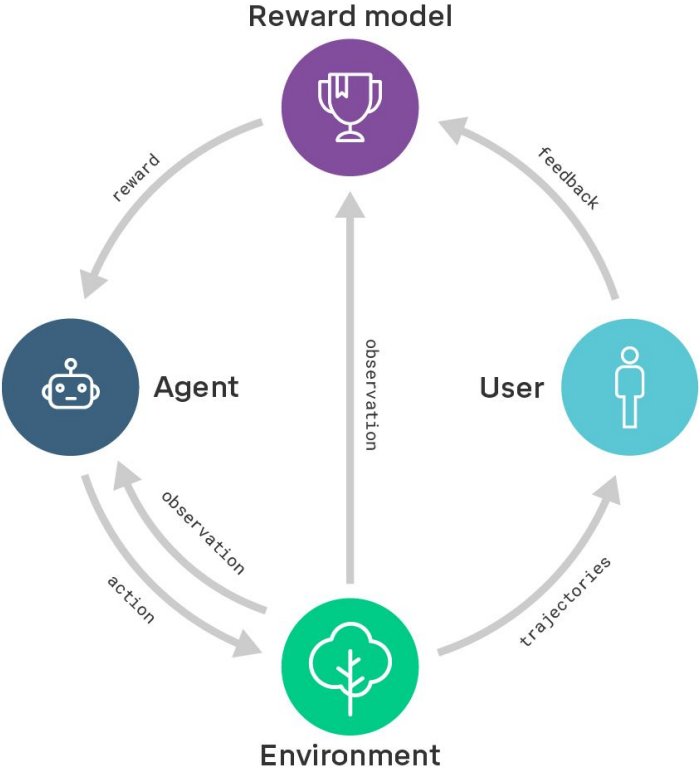
Assumption 2

For many tasks, **evaluation** of outcomes is **easier than** producing the correct **behavior**.

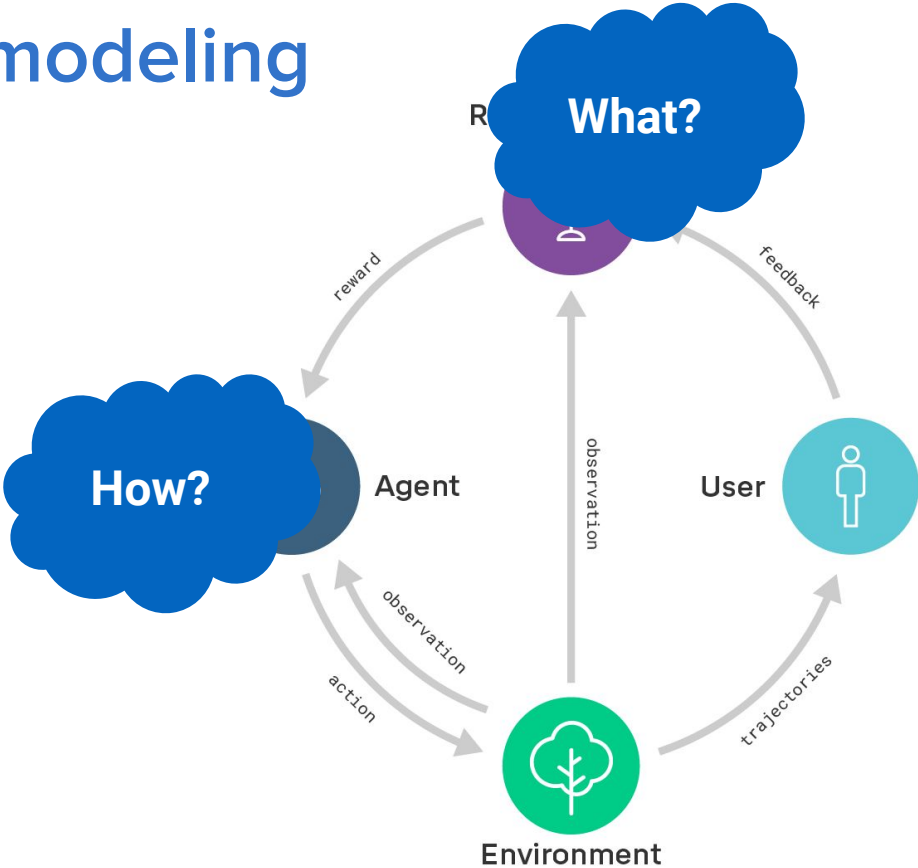
Evaluation is easier than behavior



Reward modeling



Reward modeling



Some tasks are hard to evaluate

Reward model



reward

feedback



User



observation

Agent



observation
action

trajectories



Environment

LETTER
AN ICLR 2018 WORKSHOP

Human-level control through deep reinforcement learning

Vedant Misra¹, Kamyar Kavehfar², David Silver³, Andrej A. Senior¹, Ivo Daniilovic¹, Marc G. Bellemare¹, Alex Graves¹, Martin Sunderhauf¹, Andrew K. Senior¹, Greg DeRoos¹, Ilya Sutskever⁴, Charles Isipovic¹, Anne Botea¹, Ioannis Antonoglou¹, Helen King¹, Sharan Kumara¹, Iason Veličković¹, Shazeer Legg¹, David Hasselblad¹

The theory of reinforcement learning provides a normative account, deeply rooted in psychological and neuroscientific perspectives on animal learning, of how agents may optimize their control of an environment. In our formulation, this perspective is abstracted to algorithms that can learn to control a wide range of tasks, in which the reinforcement learning problem is defined over a set of discrete or continuous actions. By using these algorithms to solve these problems through a hierarchical combination of reinforcement learning and hierarchical deep reinforcement learning, we have demonstrated that these algorithms can learn to control a wide range of tasks, in which the reinforcement learning problem is defined over a set of discrete or continuous actions. By using these algorithms to solve these problems through a hierarchical combination of reinforcement learning and hierarchical deep reinforcement learning, we have demonstrated that these algorithms can learn to control a wide range of tasks, in which the reinforcement learning problem is defined over a set of discrete or continuous actions.

How we measure advances in training deep neural networks¹ to develop novel artificial agents, termed a deep Q-network, that can learn to control a wide range of tasks, in which the reinforcement learning problem is defined over a set of discrete or continuous actions. Our 'deep games'—the games that the deep Q-network agent, receiving only the pixels and the game state as input, was able to surpass the performance of all previous algorithms and achieve a level comparable to that of a professional human game player across a set of games, using the same algorithm, network architecture and hyperparameters. This work bridges the divide between high-dimensional sensory inputs and actions, resulting in the first artificial agent that is capable of learning to solve a diverse array of challenging tasks.

We set out to create single algorithms that would be able to develop a wide range of competence in a natural range of challenging tasks—a central goal of general artificial intelligence². Our first artificial general intelligence³. To achieve this, we developed novel agents, a deep Q-network (DQN), which allow a combination of reinforcement learning with a class of artificial neural networks⁴. We used the deep neural networks, which we call a deep Q-network, to learn to control a wide range of tasks, in which the reinforcement learning problem is defined over a set of discrete or continuous actions. Our 'deep games'—the games that the deep Q-network agent, receiving only the pixels and the game state as input, was able to surpass the performance of all previous algorithms and achieve a level comparable to that of a professional human game player across a set of games, using the same algorithm, network architecture and hyperparameters. This work bridges the divide between high-dimensional sensory inputs and actions, resulting in the first artificial agent that is capable of learning to solve a diverse array of challenging tasks.

We set out to create single algorithms that would be able to develop a wide range of competence in a natural range of challenging tasks—a central goal of general artificial intelligence³. To achieve this, we developed novel agents, a deep Q-network (DQN), which allow a combination of reinforcement learning with a class of artificial neural networks⁴. We used the deep neural networks, which we call a deep Q-network, to learn to control a wide range of tasks, in which the reinforcement learning problem is defined over a set of discrete or continuous actions. Our 'deep games'—the games that the deep Q-network agent, receiving only the pixels and the game state as input, was able to surpass the performance of all previous algorithms and achieve a level comparable to that of a professional human game player across a set of games, using the same algorithm, network architecture and hyperparameters. This work bridges the divide between high-dimensional sensory inputs and actions, resulting in the first artificial agent that is capable of learning to solve a diverse array of challenging tasks.

We consider tasks in which the agent interacts with an environment through sequences of observations, actions and rewards. The goal of the agent is to select actions in a fashion that maximizes cumulative future reward. More formally, we use a deep convolutional neural network to approximate the optimal action value function

$$Q^*(s, a) = \max_{a'} \sum_{t=0}^{\infty} \gamma^t r_{t+1} + V^*(s_{t+1}) - V^*(s_t)$$

in which the maximum value of rewards is discounted by γ at each time step, s_t is a behavior policy π , r_t is the reward after making an observation s_t and taking an action a_t in the environment.

Reinforcement learning is known to be unstable or even to diverge when a nonlinear function approximator such as a neural network is used to represent the action value (also known as Q) function⁵. This instability has several causes: the correlation present in the sequence of observations, the fact that small updates to Q may significantly change the policy and therefore change the data distribution, and the work done between the action a_t and the next observation s_{t+1} (see [Schuermans et al., 2017](#)). We address these instabilities with a novel variant of Q-learning, which we call the deep Q-network (DQN). The main idea is to use a fixed, slowly changing policy π that stabilizes over the data, thereby ensuring convergence to the true action value function. Second, we use an ϵ -greedy policy that reduces the action values Q towards target values that are only periodically updated, thereby reducing correlations with the target values.

While most deep reinforcement learning methods receive external rewards in the reinforcement learning setting, such as a neural DQN (Silver et al., 2015), these methods can be used for training of neural networks in a variety of tasks. For example, these methods can be used to train a neural network to approximate a value function $Q(s, a)$ using the deep convolutional neural network (DQN) as the function approximator. We use the DQN to approximate the action value function $Q(s, a)$ using the deep convolutional neural network (DQN) as the function approximator. We use the DQN to approximate the action value function $Q(s, a)$ using the deep convolutional neural network (DQN) as the function approximator. We use the DQN to approximate the action value function $Q(s, a)$ using the deep convolutional neural network (DQN) as the function approximator.

To evaluate our DQN agent, we took advantage of the Atari 2600 platform, which offers a diverse array of tasks ($n = 49$) designed to be

Evaluation assistance tasks

- Well-written
- Novel
- Experiments correct
- Proofs correct
- ...



LETTER

doi:10.1038/nature14236

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{2*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fiedelnd¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dhruvhan Kumar¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

The theory of reinforcement learning provides a normative account¹, deeply rooted in psychological² and neuroscientific³ perspectives on animal behaviour, of how agents may optimize their control of an environment. To use reinforcement learning successfully in situations approaching real-world complexity, however, agents are confronted with a difficult task: they must derive efficient representations of the environment from high-dimensional sensory inputs, and use these to generalize past experience to new situations. Remarkably, humans and other animals seem to solve this problem through a harmonious combination of reinforcement learning and hierarchical sensory processing systems^{4,5}, the former evidenced by a wealth of neural data revealing notable parallels between the phasic signals emitted by dopaminergic neurons and temporal difference reinforcement learning algorithms⁶. While reinforcement learning agents have achieved some successes in a variety of domains^{7,8}, their applicability has previously been limited to domains in which useful features can be hand-engineered, or to domains with fully observed, low-dimensional state spaces. Here we use recent advances in training deep neural networks^{9–11} to develop a novel artificial agent, termed a deep Q-network, that can learn successful policies directly from high-dimensional sensory inputs using end-to-end reinforcement learning. We tested this agent on the challenging domain of classic Atari 2600 games¹². We demonstrate that the deep Q-network agent, receiving only the pixels and the game score as inputs, was able to surpass the performance of all previous algorithms and achieve a level comparable to that of a pro-

agent is to select actions in a fashion that maximizes cumulative future reward. More formally, we use a deep convolutional neural network to approximate the optimal action-value function

$$Q^*(s,a) = \max_{a'} [r_t + \gamma V_t + \gamma^2 r_{t+1} + \gamma^3 r_{t+2} + \dots]_{n=2, a_i = a, \pi}$$

which is the maximum sum of rewards, discounted by γ at each time-step t , achievable by a behaviour policy $\pi = P(a|s)$, after making an observation (s) and taking an action (a) (see Methods)¹⁷.

Reinforcement learning is known to be unstable or even to diverge when a nonlinear function approximator such as a neural network is used to represent the action-value (also known as Q) function¹⁸. This instability has several causes: the correlations present in the sequence of observations, the fact that small updates to Q may significantly change the policy and therefore change the data distribution, and the correlations between the action-values (Q) and the target values $r + \gamma \max_{a'} Q(s, a')$. We address these instabilities with a novel variant of Q-learning, which uses two key ideas. First, we used a biologically inspired mechanism termed experience replay^{21,22} that randomizes over the data, thereby removing correlations in the observation sequence and smoothing over changes in the data distribution (see below for details). Second, we used an iterative update that adjusts the action-values (Q) towards target values that are only periodically updated, thereby reducing correlations with the target.

While other stable methods exist for training neural networks in the reinforcement learning setting, such as neural fitted Q-iteration¹⁸, these

Evaluation assistance tasks

- Well-written



- Novel



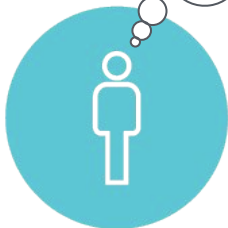
- Experiments correct



- Proofs correct



- ...



LETTER

doi:10.1038/nature14236

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{2*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fiedelnd¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dhruvhan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

The theory of reinforcement learning provides a normative account¹, deeply rooted in psychological² and neuroscientific³ perspectives on animal behaviour, of how agents may optimize their control of an environment. To use reinforcement learning successfully in situations approaching real-world complexity, however, agents are confronted with a difficult task: they must derive efficient representations of the environment from high-dimensional sensory inputs, and use these to generalize past experience to new situations. Remarkably, humans and other animals seem to solve this problem through a harmonious combination of reinforcement learning and hierarchical sensory processing systems^{4,5}, the former evidenced by a wealth of neural data revealing notable parallels between the phasic signals emitted by dopaminergic neurons and temporal difference reinforcement learning algorithms⁶. While reinforcement learning agents have achieved some successes in a variety of domains⁷⁻¹¹, their applicability has previously been limited to domains in which useful features can be hand-engineered, or to domains with fully observed, low-dimensional state spaces. Here we use recent advances in training deep neural networks¹²⁻¹⁴ to develop a novel artificial agent, termed a deep Q-network, that can learn successful policies directly from high-dimensional sensory inputs using end-to-end reinforcement learning. We tested this agent on the challenging domain of classic Atari 2600 games¹⁵. We demonstrate that the deep Q-network agent, receiving only the pixels and the game score as inputs, was able to surpass the performance of all previous algorithms and achieve a level comparable to that of a pro-

agent is to select actions in a fashion that maximizes cumulative future reward. More formally, we use a deep convolutional neural network to approximate the optimal action-value function

$$Q^*(s,a) = \max_{a'} [r_t + \gamma v_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi],$$

which is the maximum sum of rewards, discounted by γ at each time-step t , achievable by a behaviour policy $\pi = P(a|s)$, after making an observation (s) and taking an action (a) (see Methods)¹⁷.

Reinforcement learning is known to be unstable or even to diverge when a nonlinear function approximator such as a neural network is used to represent the action-value (also known as Q) function¹⁸. This instability has several causes: the correlations present in the sequence of observations, the fact that small updates to Q may significantly change the policy and therefore change the data distribution, and the correlations between the action-values (Q) and the target values $r + \gamma \max_{a'} Q(s, a')$. We address these instabilities with a novel variant of Q-learning, which uses two key ideas. First, we used a biologically inspired mechanism termed experience replay^{19,21} that randomizes over the data, thereby removing correlations in the observation sequence and smoothing over changes in the data distribution (see below for details). Second, we used an iterative update that adjusts the action-values (Q) towards target values that are only periodically updated, thereby reducing correlations with the target.

While other stable methods exist for training neural networks in the reinforcement learning setting, such as neural fitted Q-iteration¹⁸, these

Evaluation assistance tasks

- Well-written
- Novel
- Experiments correct
- Proofs correct
- ...



yes



yes



yes



N/A



HUMAN LETTER

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{2*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fiedelnd¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dhruvhan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

The theory of reinforcement learning provides a normative account¹, deeply rooted in psychological² and neuroscientific³ perspectives on animal behaviour, of how agents may optimize their control of an environment. To use reinforcement learning successfully in situations approaching real-world complexity, however, agents are confronted with a difficult task: they must derive efficient representations of the environment from high-dimensional sensory inputs, and use these to generalize past experience to new situations. Remarkably, humans and other animals seem to solve this problem through a harmonious combination of reinforcement learning and hierarchical sensory processing systems^{4,5}, the former evidenced by a wealth of neural data revealing notable parallels between the phasic signals emitted by dopaminergic neurons and temporal difference reinforcement learning algorithms⁶. While reinforcement learning agents have achieved some successes in a variety of domains⁷⁻¹¹, their applicability has previously been limited to domains in which useful features can be hand-engineered, or to domains with fully observed, low-dimensional state spaces. Here we use recent advances in training deep neural networks¹²⁻¹⁴ to develop a novel artificial agent, termed a deep Q-network, that can learn successful policies directly from high-dimensional sensory inputs using end-to-end reinforcement learning. We tested this agent on the challenging domain of classic Atari 2600 games¹⁵. We demonstrate that the deep Q-network agent, receiving only the pixels and the game score as inputs, was able to surpass the performance of all previous algorithms and achieve a level comparable to that of a pro-

agent is to select actions in a fashion that maximizes cumulative future reward. More formally, we use a deep convolutional neural network to approximate the optimal action-value function

$$Q^*(s,a) = \max_{a'} [r_t + \gamma V_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi]$$

which is the maximum sum of rewards, discounted by γ at each time-step t , achievable by a behaviour policy $\pi = P(a|s)$, after making an observation (s) and taking an action (a) (see Methods)¹⁷.

Reinforcement learning is known to be unstable or even to diverge when a nonlinear function approximator such as a neural network is used to represent the action-value (also known as Q) function¹⁸. This instability has several causes: the correlations present in the sequence of observations, the fact that small updates to Q may significantly change the policy and therefore change the data distribution, and the correlations between the action-values (Q) and the target values $r + \gamma \max_{a'} Q(s', a')$. We address these instabilities with a novel variant of Q-learning, which uses two key ideas. First, we used a biologically inspired mechanism termed experience replay^{19,21} that randomizes over the data, thereby removing correlations in the observation sequence and smoothing over changes in the data distribution (see below for details). Second, we used an iterative update that adjusts the action-values (Q) towards target values that are only periodically updated, thereby reducing correlations with the target.

While other stable methods exist for training neural networks in the reinforcement learning setting, such as neural fitted Q-iteration²⁰, these

doi:10.1038/nature14236

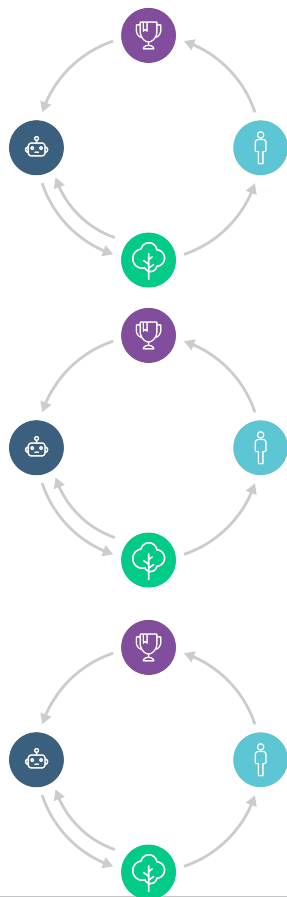
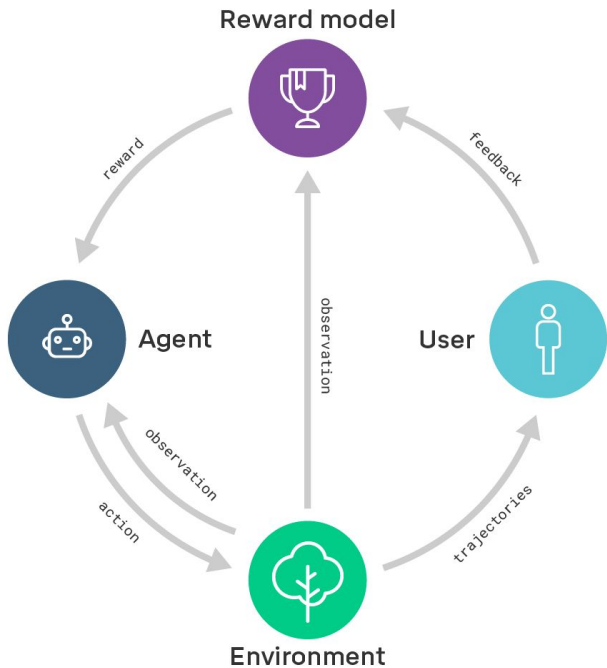
Recursive reward modeling

LETTER

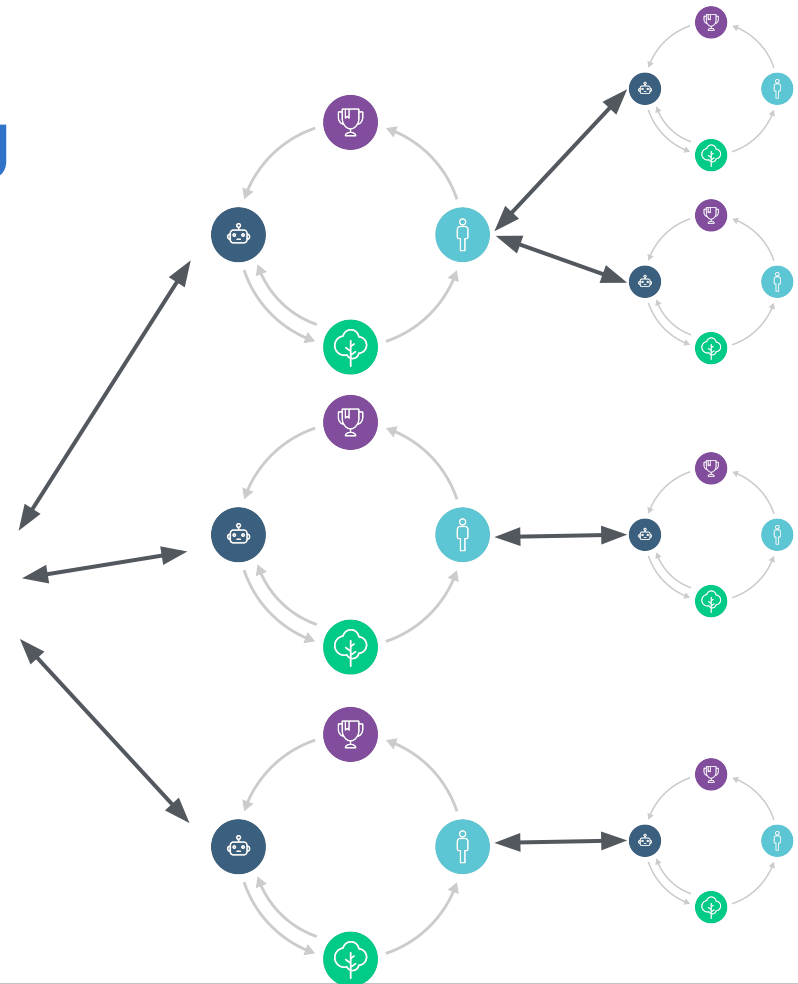
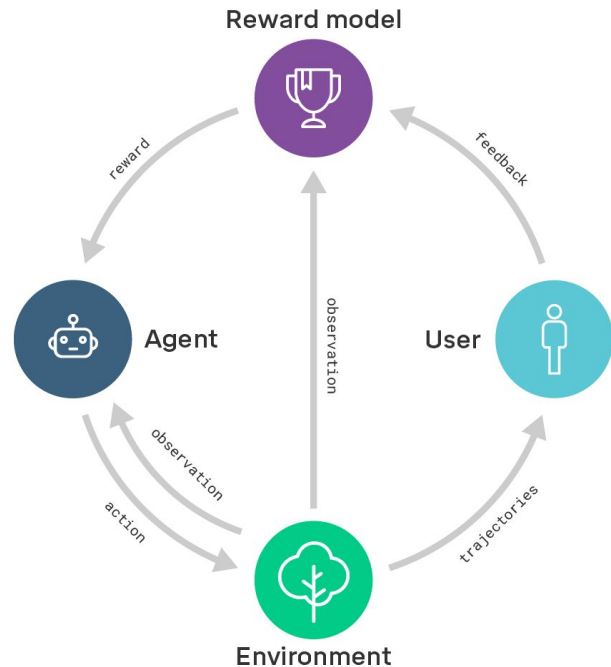
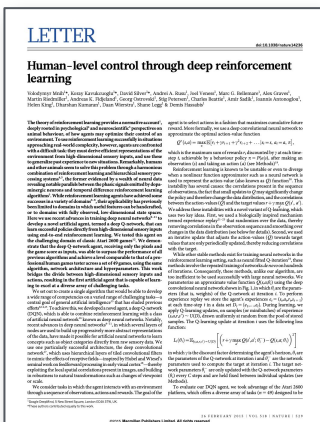
Human-level control through deep reinforcement learning

by David Silver, Thomas Schrittwieser, Jonathan Schrittwieser, Ioannis Antonoglou, Laurent M. A. Descotes-Genon, Demis Hassabis, David Greig, Adriane Pascanu, and Geoffrey E. Hinton

The theory of reinforcement learning provides a normative account, originally articulated by Sutton and Barto, of how agents might control their actions to maximize expected reward by learning from a sequence of observations. By using a stochastic model of the environment, reinforcement learning has become a normative theory of adaptive control. However, this theory has been applied almost exclusively to problems where the agent is restricted to act on a finite set of discrete actions. In this paper, we describe how to extend the theory to the more general and difficult problem of continuous action spaces. In particular, we describe how to train a neural network to act in a continuous action space from a sequence of observations by approximating the value function of a Markov Decision Process with a deep convolutional neural network. This work extends the previous work on deep reinforcement learning by showing that it is possible to train a neural network to act in a continuous action space from a sequence of observations. The key to this work is the use of a deep convolutional neural network to approximate the value function of a Markov Decision Process. This work extends the previous work on deep reinforcement learning by showing that it is possible to train a neural network to act in a continuous action space from a sequence of observations. The key to this work is the use of a deep convolutional neural network to approximate the value function of a Markov Decision Process. This work extends the previous work on deep reinforcement learning by showing that it is possible to train a neural network to act in a continuous action space from a sequence of observations. The key to this work is the use of a deep convolutional neural network to approximate the value function of a Markov Decision Process.



Recursive reward modeling



Challenges

Amount of feedback

Feedback distribution

Reward hacking

Unacceptable outcomes

Reward-result gap

Challenges

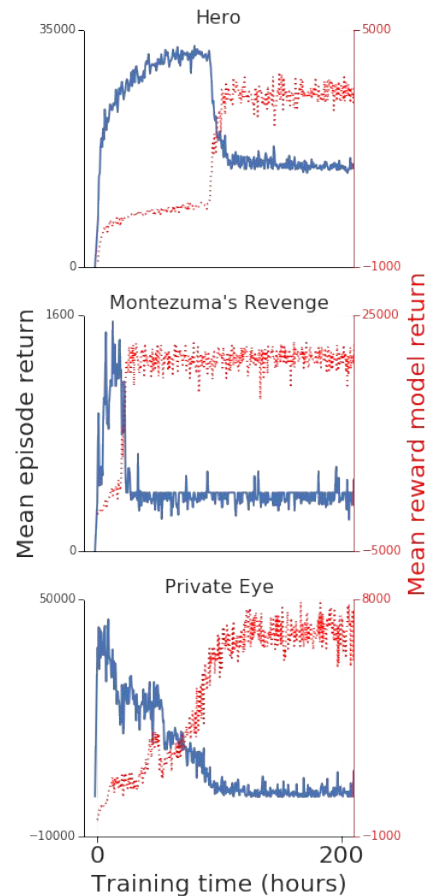
Amount of feedback

Feedback distribution

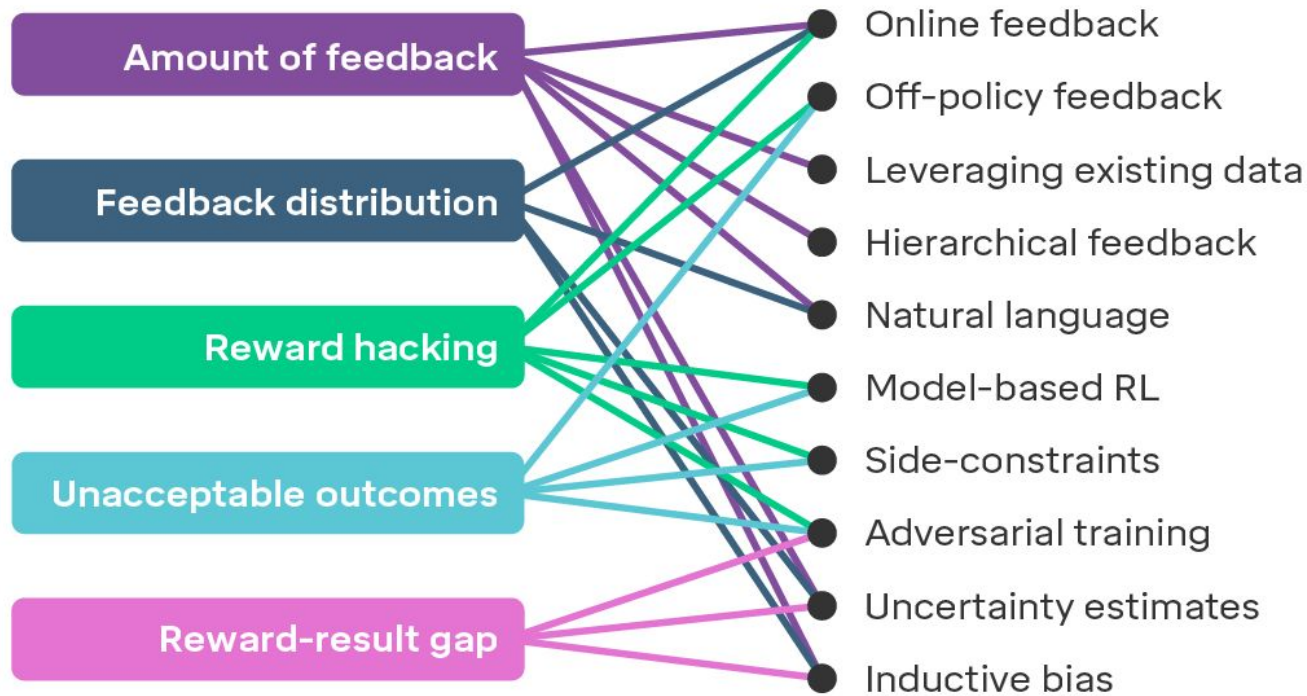
Reward hacking

Unacceptable outcomes

Reward-result gap



Challenges



Establishing trust

- Design choices
- Testing
- Interpretability
- Formal verification
- Theoretical guarantees



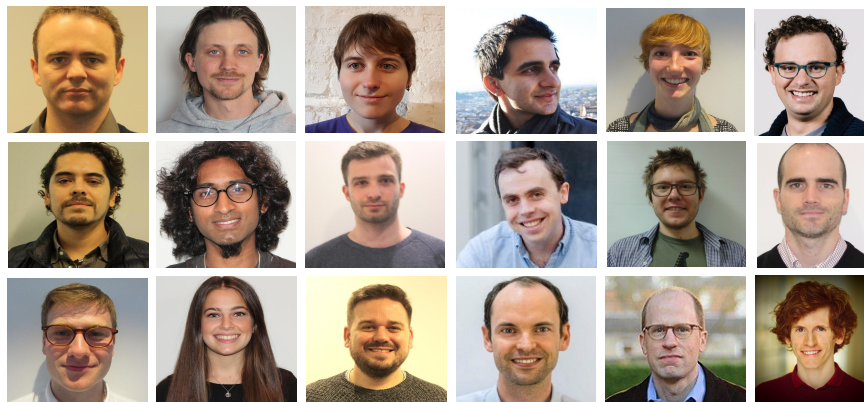
Safety certificates

Thanks! :)

Blog post: <https://goo.gl/azGMtA>

Paper:

<https://arxiv.org/abs/1811.07871>



Scalable agent alignment via reward modeling: a research direction

Jan Leike David Krueger* Tom Everitt Miljan Martic Vishal Maini Shane Legg
DeepMind DeepMind DeepMind DeepMind DeepMind DeepMind
Mila

Abstract

One obstacle to applying reinforcement learning algorithms to real-world problems is the lack of suitable reward functions. Designing such reward functions is difficult in part because the user only has an implicit understanding of the task objective. This gives rise to the *agent alignment problem*: how do we create agents that behave in accordance with the user's intentions? We outline a high-level research direction to solve the agent alignment problem centered around *reward modeling*: learning a reward function from interaction with the user and optimizing the learned reward function with reinforcement learning. We discuss the key challenges we expect to face when scaling reward modeling to complex and general domains, concrete approaches to mitigate these challenges, and ways to establish trust in the resulting agents.

1 Introduction

Games are a useful benchmark for research because progress is easily measurable. Atari games come with a score function that captures how well the agent is playing the game; board games or competitive multiplayer games such as Dota 2 and Starcraft II have a clear winner or loser at the end of the game. This helps us determine empirically which algorithmic and architectural improvements work best.

However, the ultimate goal of machine learning (ML) research is to go beyond games and improve human lives. To achieve this we need ML to assist us in real-world domains, ranging from simple tasks like ordering food or answering emails to complex tasks like software engineering or running a business. Yet performance on these and other real-world tasks is not easily measurable, since they do not come readily equipped with a reward function. Instead, the objective of the task is only indirectly available through the intentions of the human user.

This requires walking a fine line. On the one hand, we want ML to generate creative and brilliant solutions like AlphaGo's Move 37 (Metz, 2016)—a move that no human would have recommended, yet it completely turned the game in AlphaGo's favor. On the other hand, we want to avoid degenerate solutions that lead to undesired behavior like exploiting a bug in the environment simulator (Clark & Amodei, 2016; Lehman et al., 2018). In order to differentiate between these two outcomes, our agent needs to understand its user's *intentions*, and robustly achieve these intentions with its behavior. We frame this as the *agent alignment problem*.

How can we create agents that behave in accordance with the user's intentions?

arXiv:1811.07871v1 [cs.LG] 19 Nov 2018