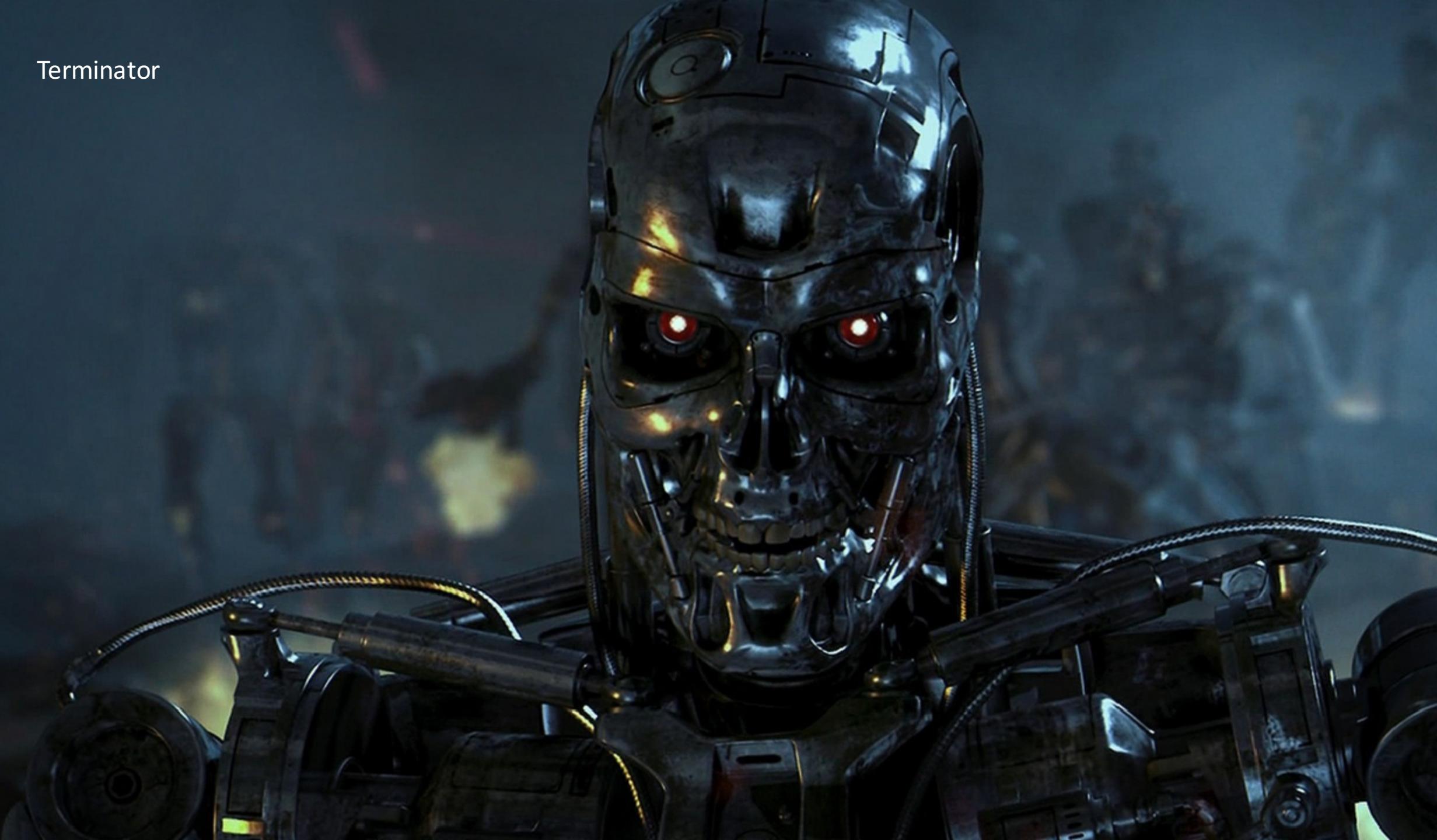


游戏  高手

通用人工智能之梦

张江

Terminator



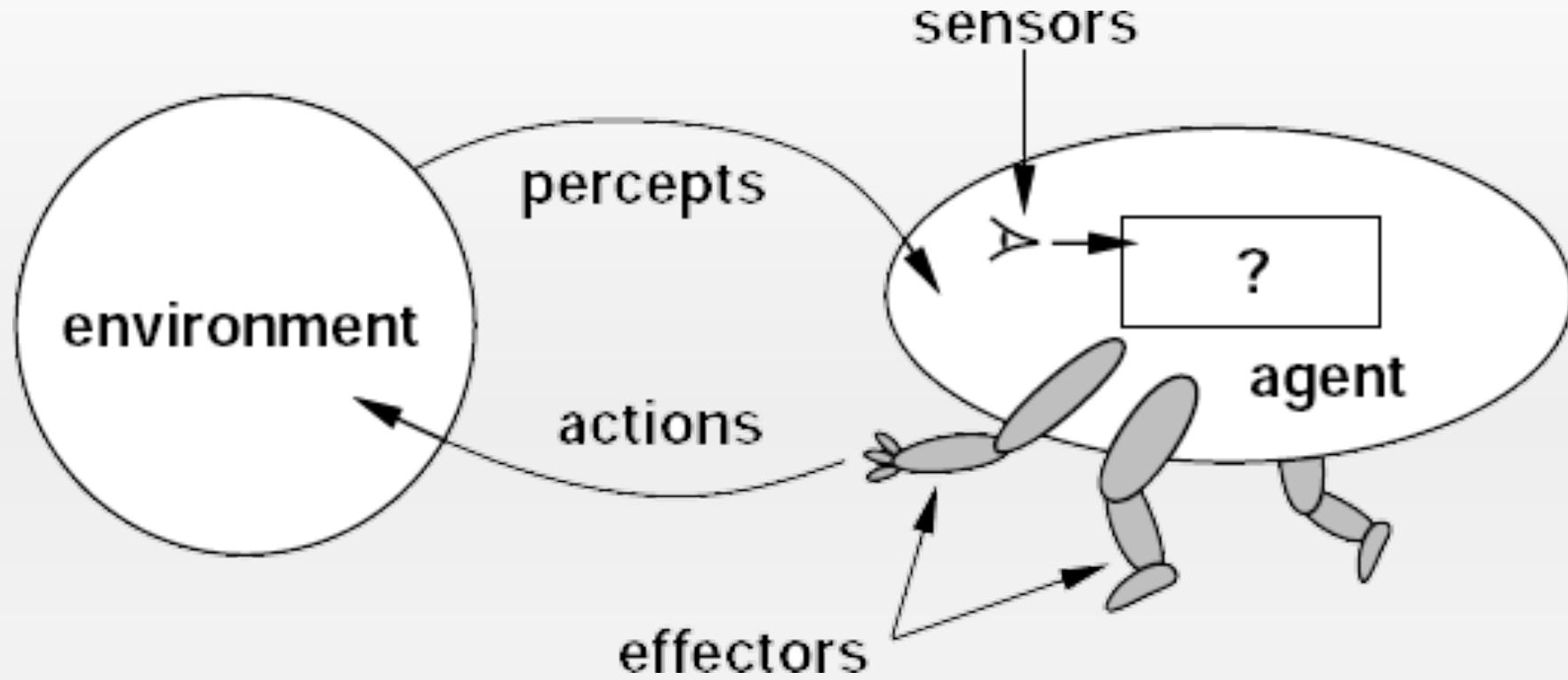
WallE



A.I.



人工智能 v.s. 机器学习



智能体（Intelligent agent）应该是一个完整的可以感知、决策、行动、学习的整体，而不是分开的一系列任务

强化学习

人工智能

机器学习

强化学习

Reinforcement
Learning

什么是强化学习？

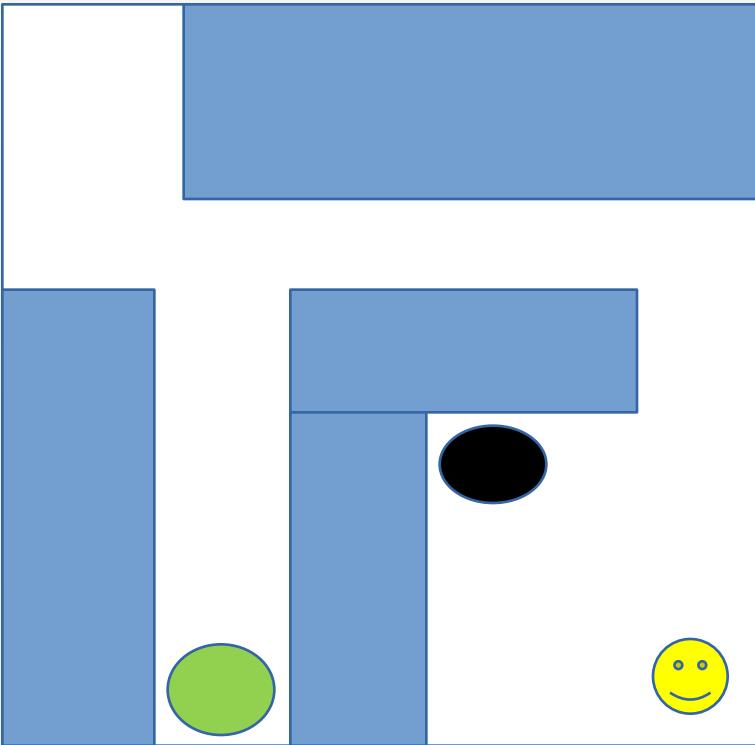
Reinforcement learning (RL) is an area of [machine learning](#) inspired by [behaviourist psychology](#), concerned with how [software agents](#) ought to take [actions](#) in an [environment](#) so as to maximize some notion of cumulative [reward](#).

- 没有即时反馈
- 边做边学
 - 平衡探索与利用



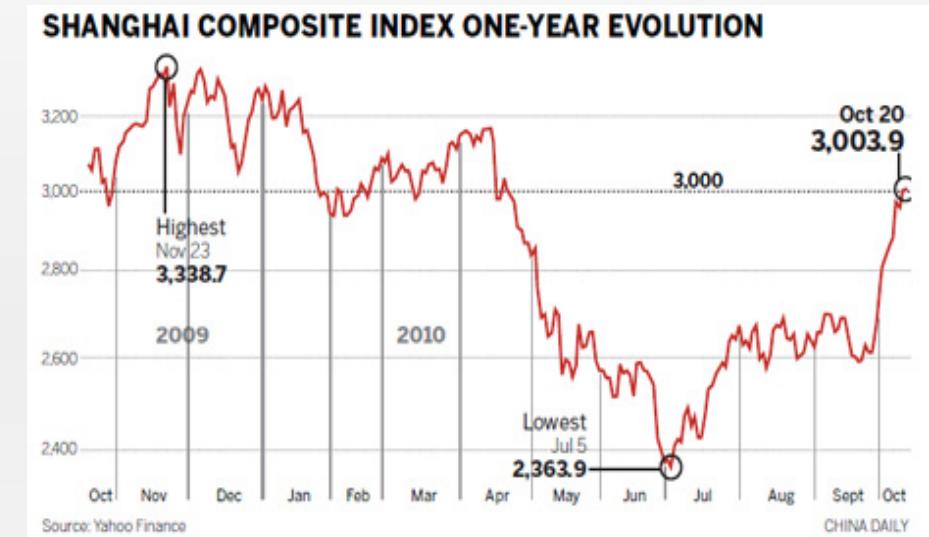
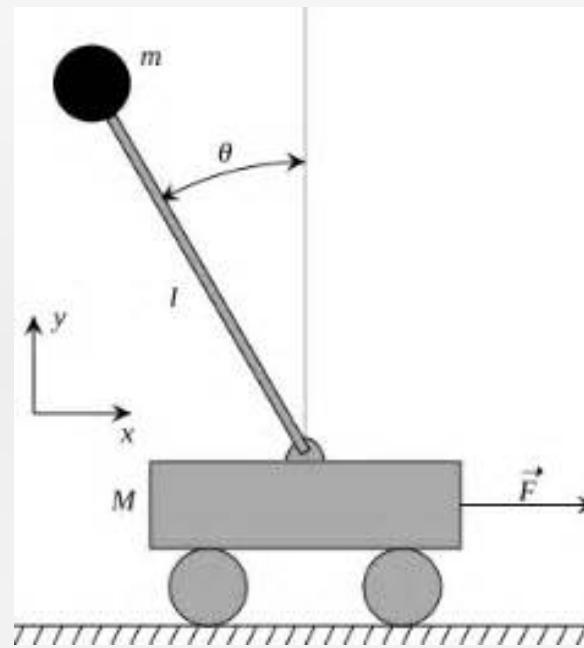
Arthur Samuel

典型的强化学习框架



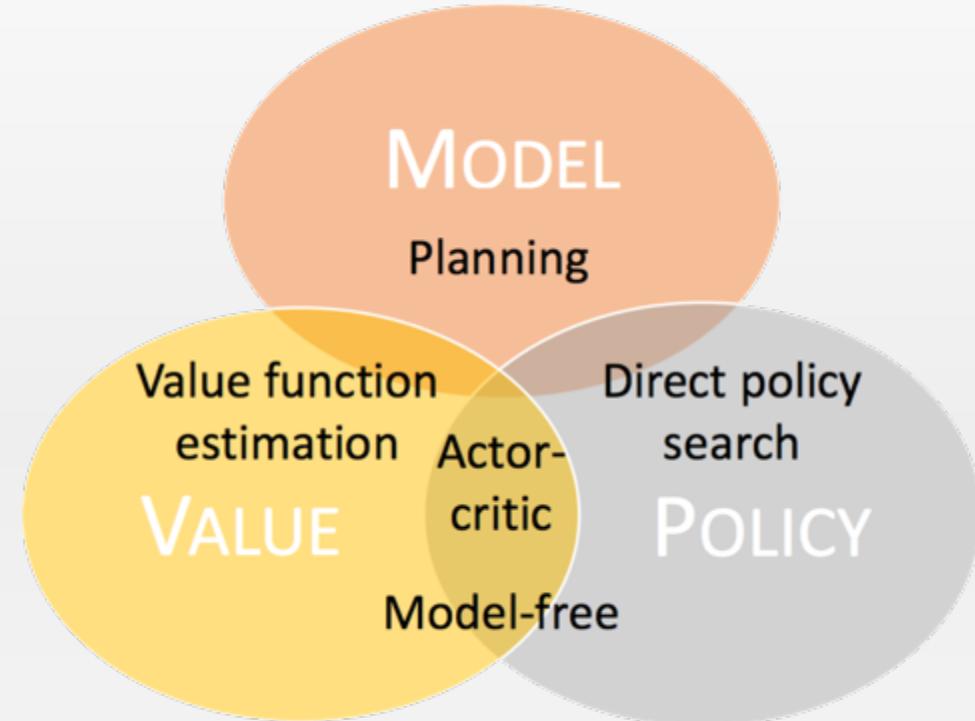
- 每一时刻
- 环境状态 S , 纪录Agent可能的位置 (x,y)
- 有一个Agent, 它可以采取行动:
 $A=\{0,1,2,\dots\}$
- 状态转移:
 - $\sigma(S, a)=S'$
- 环境会给Agent反馈 r :
 - 踩到炸弹 - 10
 - 踩到蘑菇 + 20
 - 一般情况0

强化学习应用领域



强化学习

- 按照Agent对环境的理解程度
 - Model-free: Q-learning, DQN, ...
 - Model-based: AlphaGo
- 按照学习对象的划分
 - Policy based
 - Value based: AI打游戏
 - Policy-value: AlphaGo
 - Actor-critic



LETTER

doi:10.1038/nature14236

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan 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

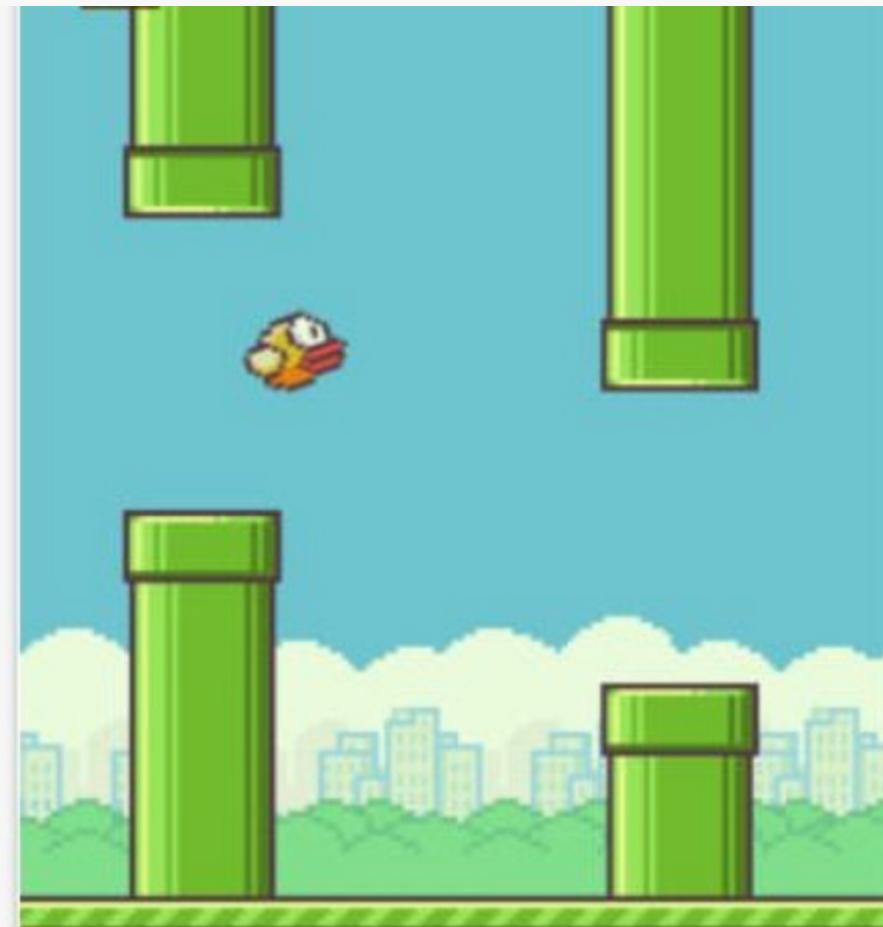
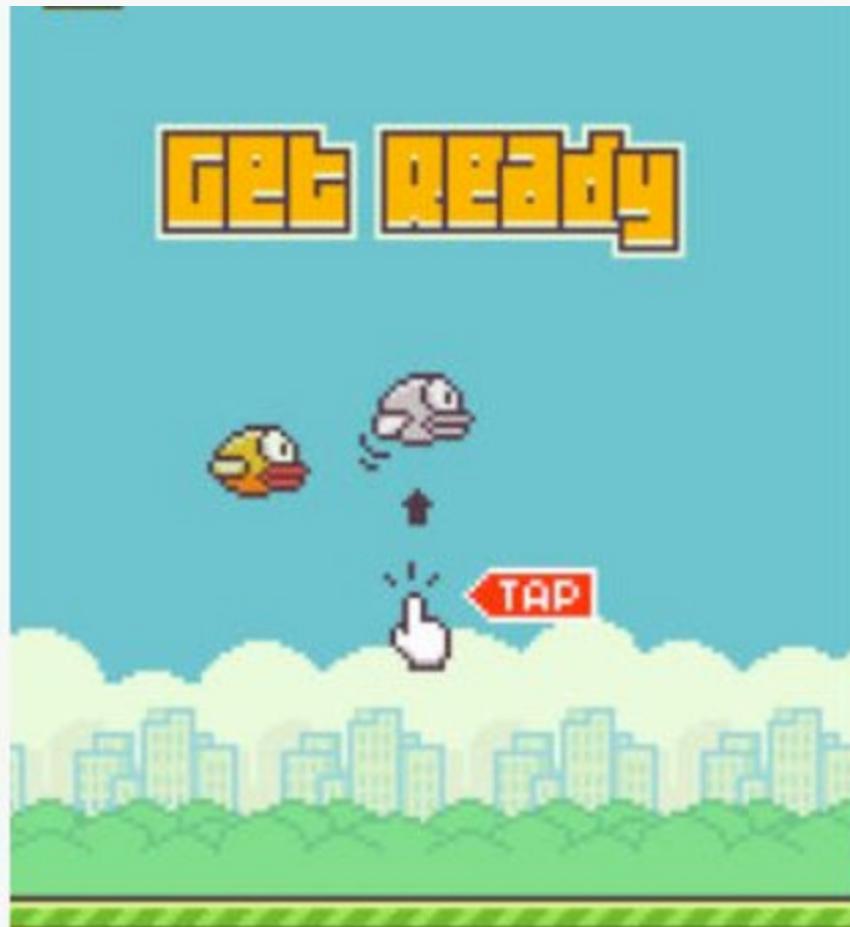
agent is to select actions in a fashion that maximizes cumulative future reward. More formally, we use a deep convolutional neural network to approximate the optimal action-value function

$$Q^*(s,a) = \max_{\pi} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi],$$

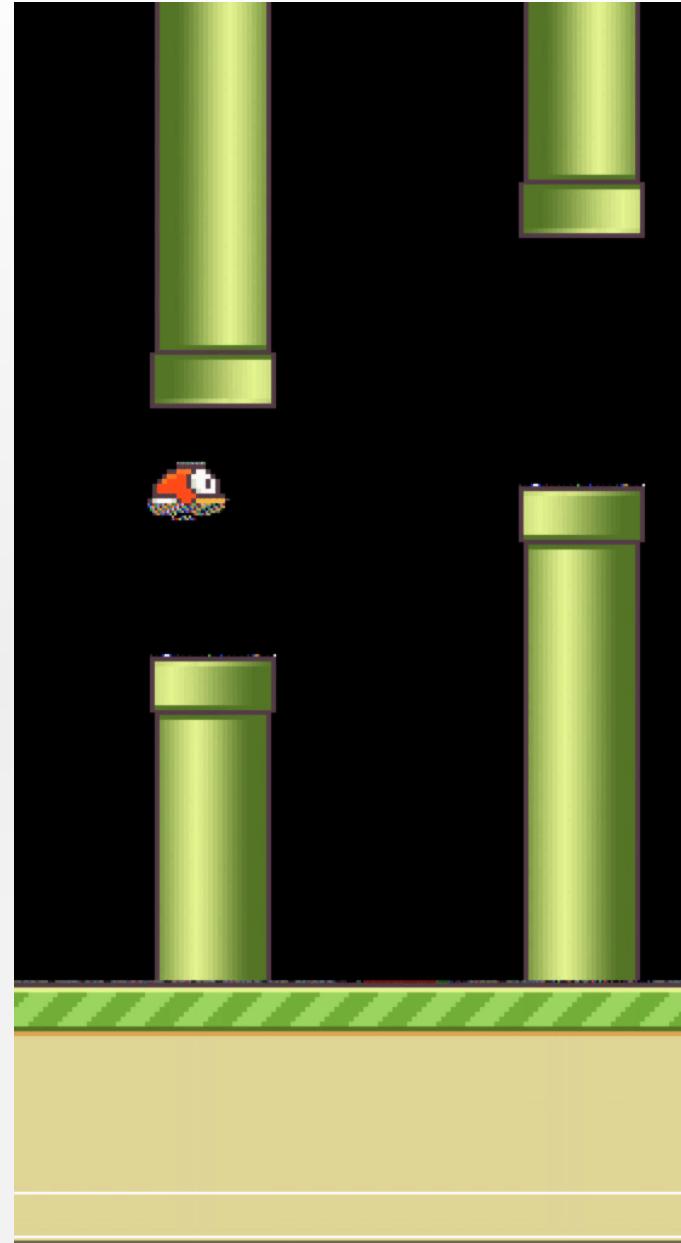
which is the maximum sum of rewards r_t discounted by γ at each time-step t , achievable by a behaviour policy $\pi = P(a|s)$, after making an observation (s) and taking an action (a) (see Methods)¹⁹.

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) function²⁰. This instability has several causes: the correlations present in the sequence of observations, the fact that small updates to Q may significantly change

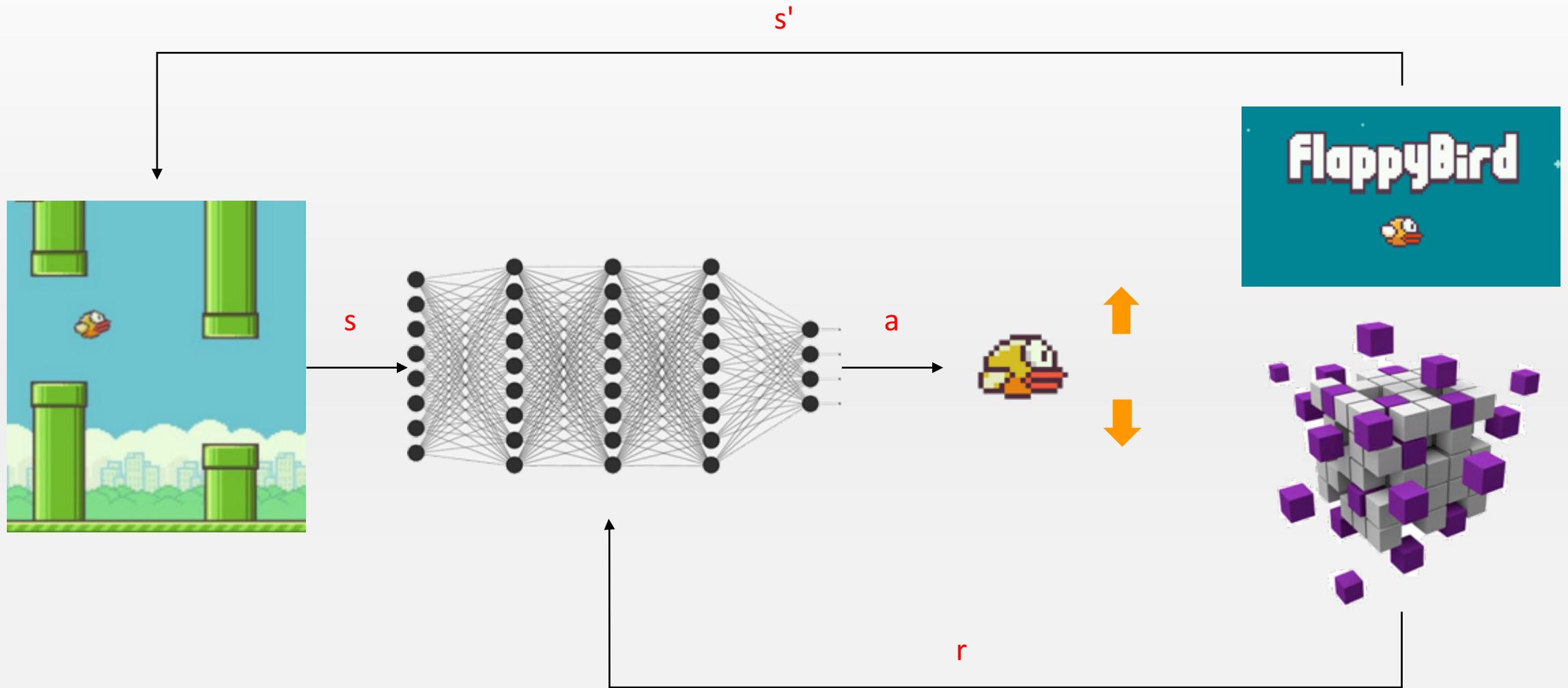
Flappy bird



Flappy bird



游戏高手

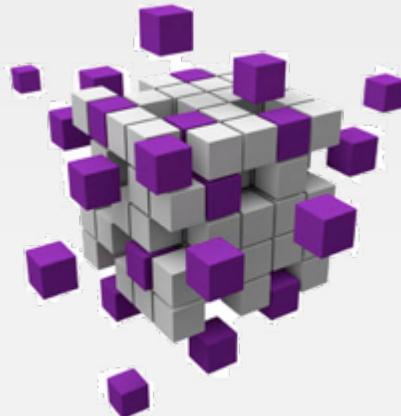


Flappy bird



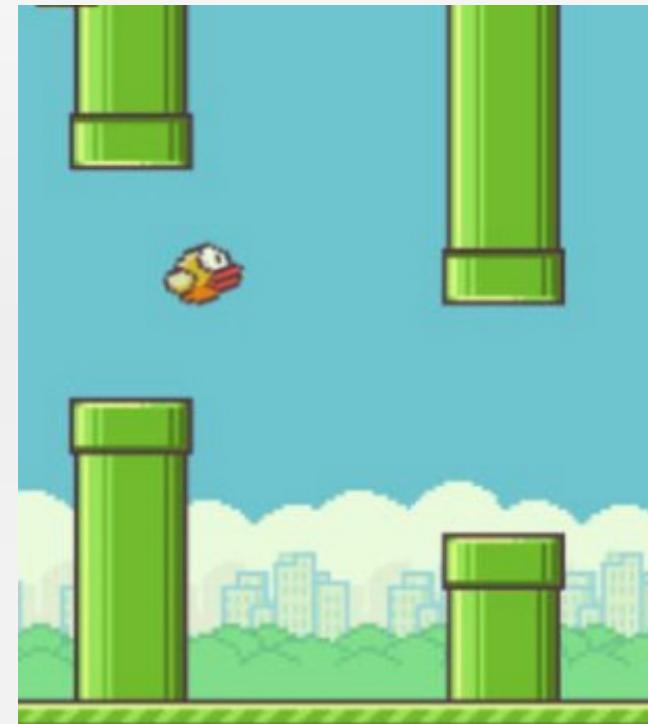
+0.1

每一帧都+0.1分



Flappy bird

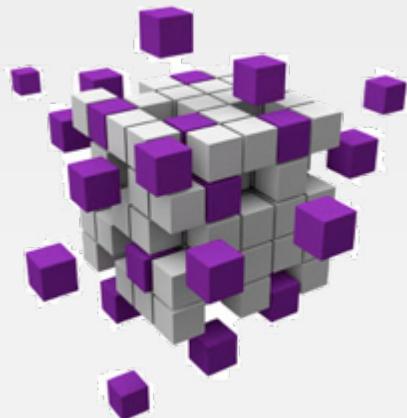
+1



每过一个管道+1分

Flappy bird

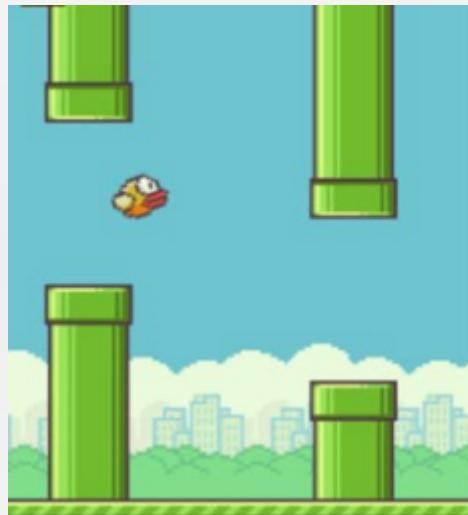
-1



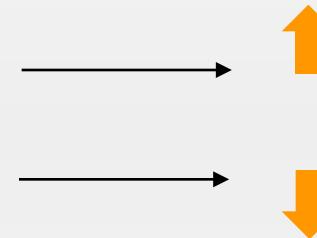
每撞一个管道-1分

游戏重新开始

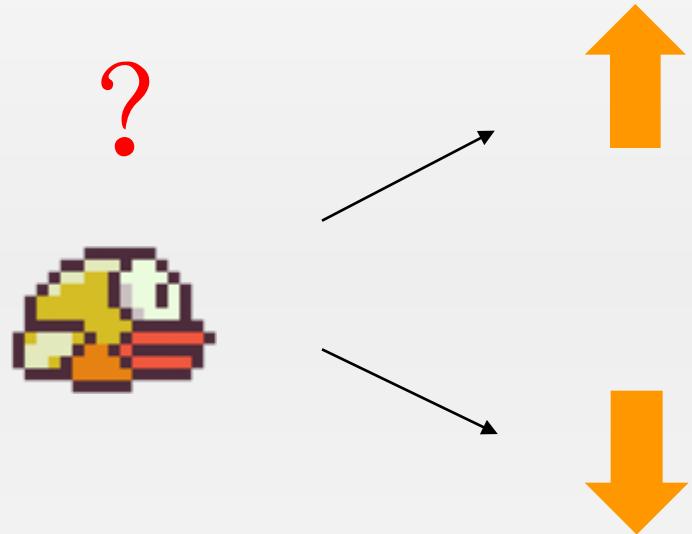
游戏高手： 神经网络



s

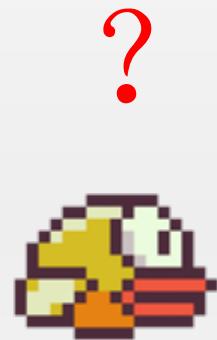


游戏高手：评价系统



我们应该如何评价这两个行动的好坏呢？

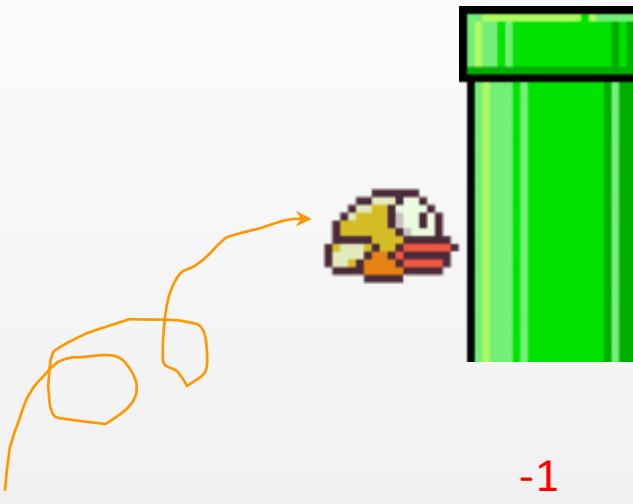
游戏高手：评价系统



?



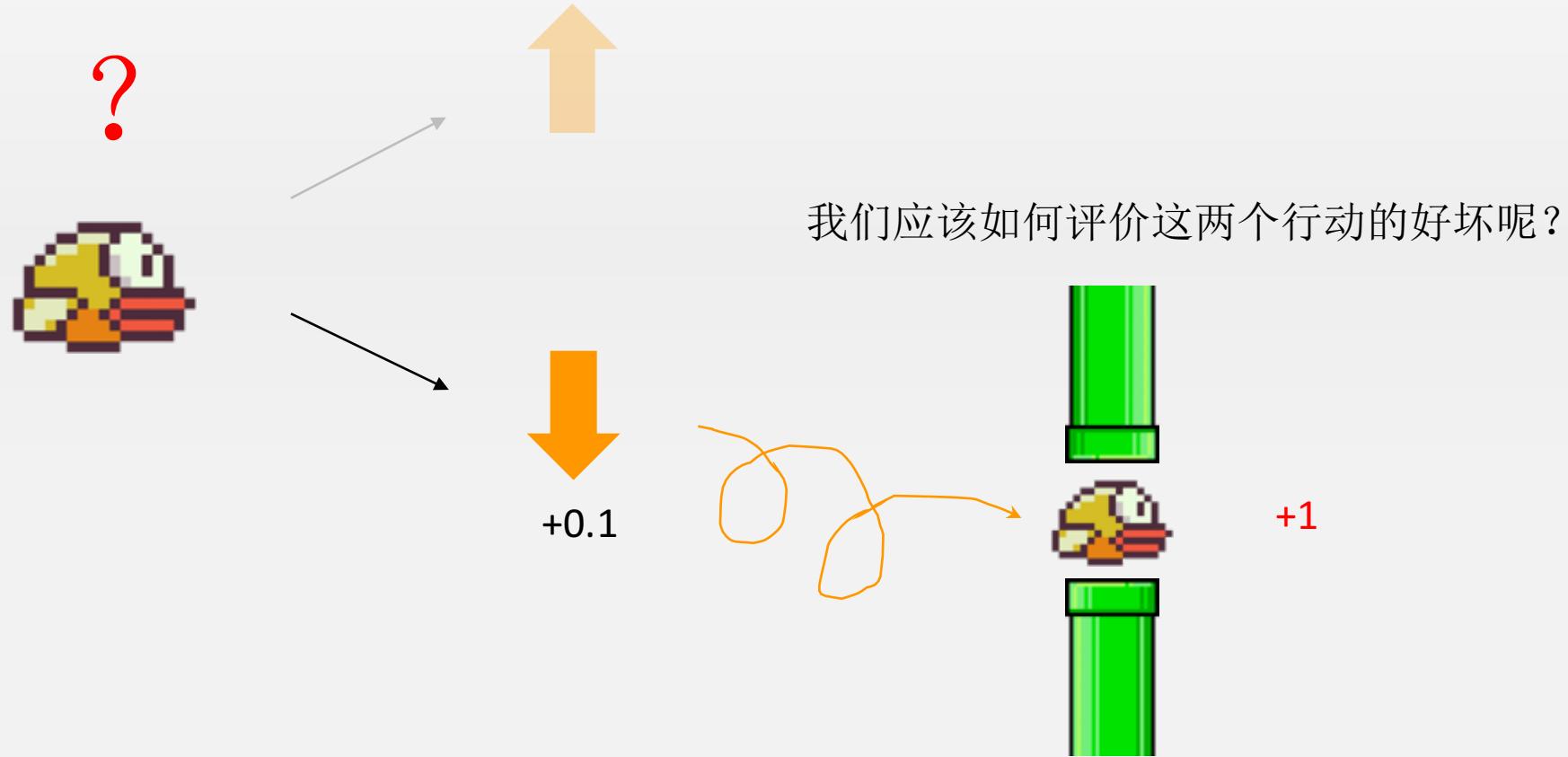
+0.1



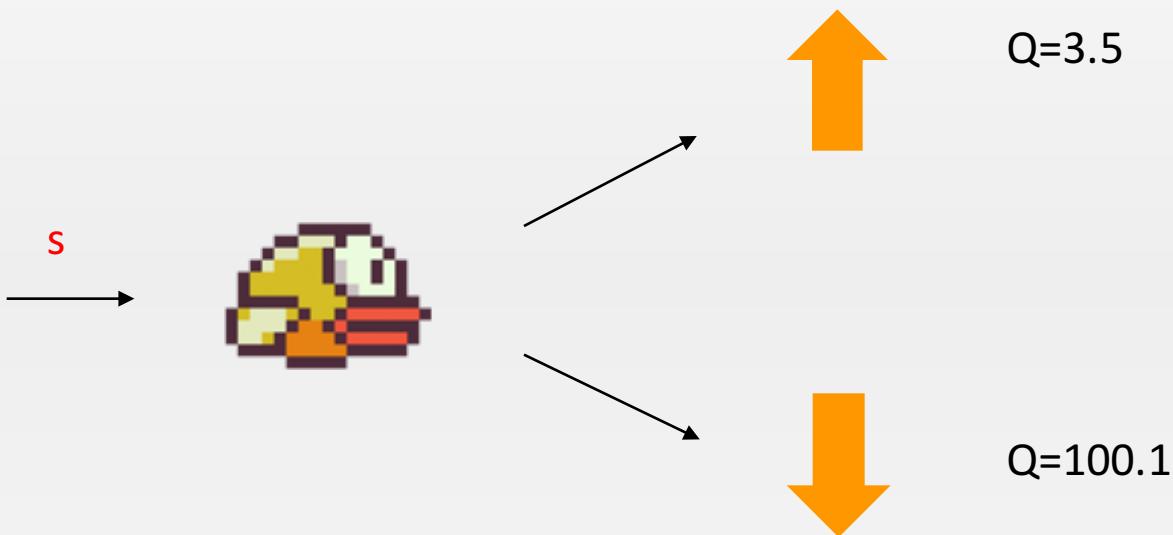
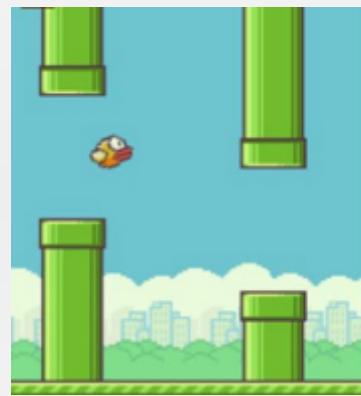
-1

我们应该如何评价这两个行动的好坏呢？

游戏高手：评价系统



游戏高手：评估函数Q

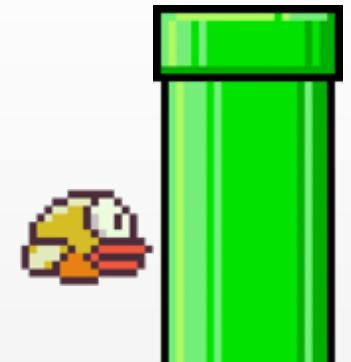


Q是一个s和a的函数， $Q(s,a)$ 指导了小鸟的行动

如何计算这个Q函数呢？

最好的行动

$$Q = -1$$



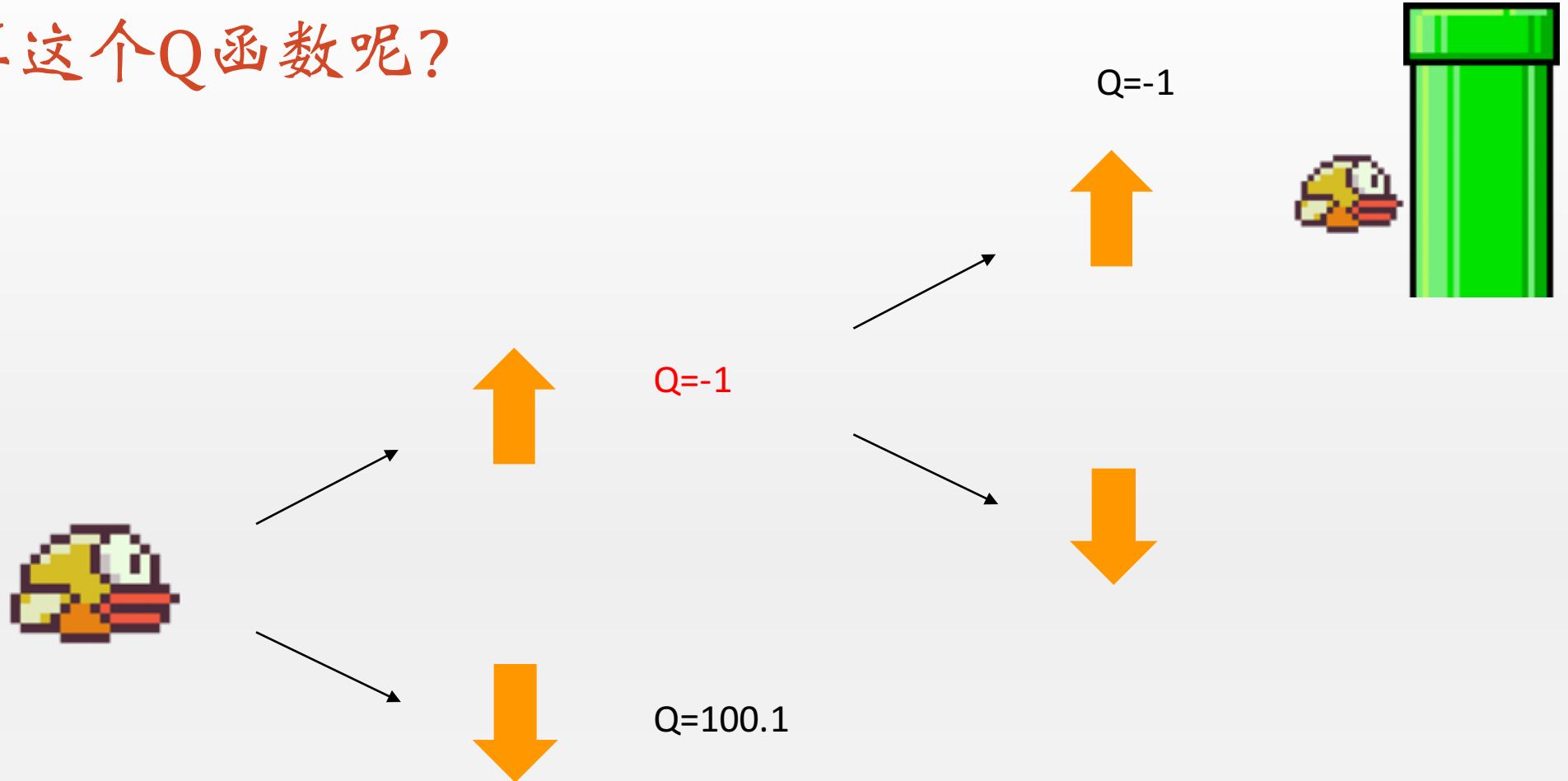
$$Q = ?$$



$$Q = 100.1$$

假如根据Q在t步的指导，导致了t+1步的未来 $r=-1$ 的回报

如何计算这个Q函数呢？



那么， t 步的Q也应该等于-1，也就是说 $Q(t) \sim r(t) + Q(t+1)$

如何计算这个Q函数呢？

$$Q(t) \sim r(t) + Q(t+1)$$

如何计算这个Q函数呢？

$$Q(t) \sim r(t) + Q(t + 1)$$

$$Q(s, a) \sim r(t) + Q(s', a')$$

根据游戏画面 s, s' 以及玩家的行动来决定 Q

如何计算这个Q函数呢？

$$Q(t) \sim r(t) + Q(t+1)$$

$$Q(s, a) \sim r(t) + Q(s', a')$$

$$Q(s, a) \sim r(t) + \gamma Q(s', a')$$

遥远的未来对现在的作用必须衰减

如何计算这个Q函数呢？

$$Q(t) \sim r(t) + Q(t+1)$$

$$Q(s, a) \sim r(t) + Q(s', a')$$

$$Q(s, a) \sim r(t) + \gamma Q(s', a')$$

$$Q(s, a) \sim r(t) + \gamma \max_{a'} Q(s', a')$$

现在的Q需要考虑未来最好的情况

Q要尽可能地与未来一致，并且客观地反映出环境的反馈

Q函数呢

$$Q(s, a) \sim r(t) + \gamma \max_{a'} Q(s', a')$$

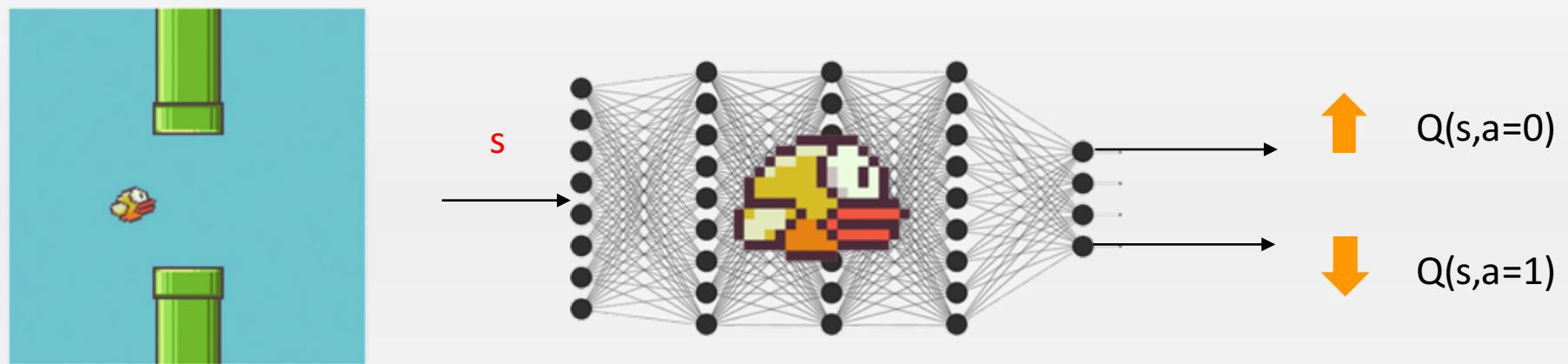
这就是著名的Q-learning算法的思想

Q函数呢

$$Q(s, a) \sim r(t) + \gamma \max_{a'} Q(s', a')$$

如何实现上面的过程呢？

游戏高手：评估函数Q

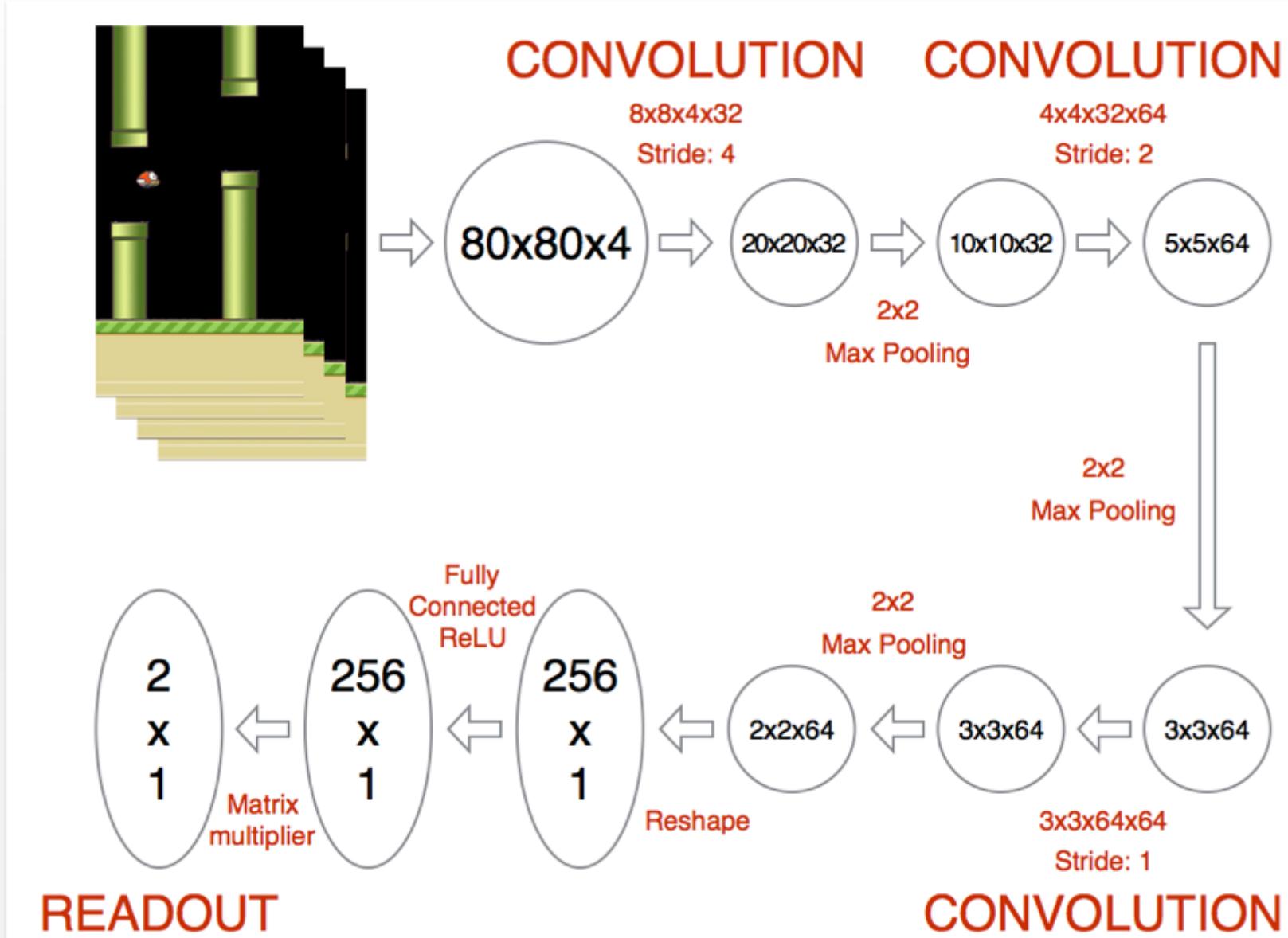


实际上，我们就是由神经网络实现了一个Q函数

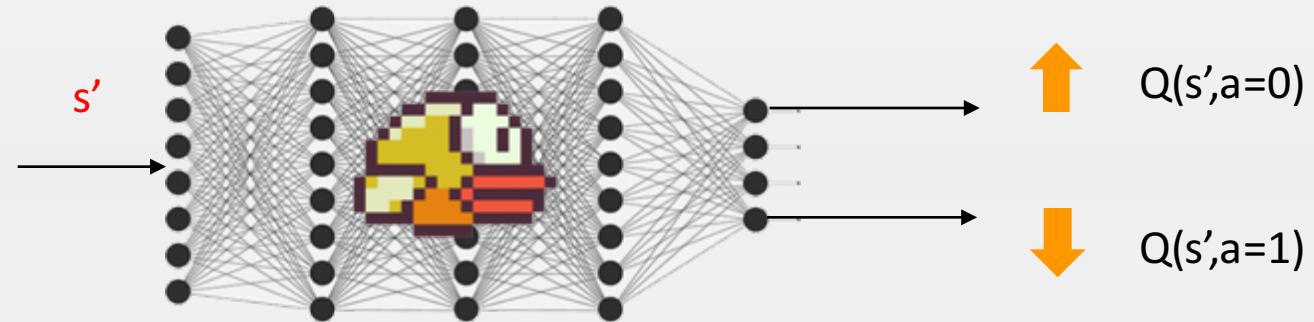
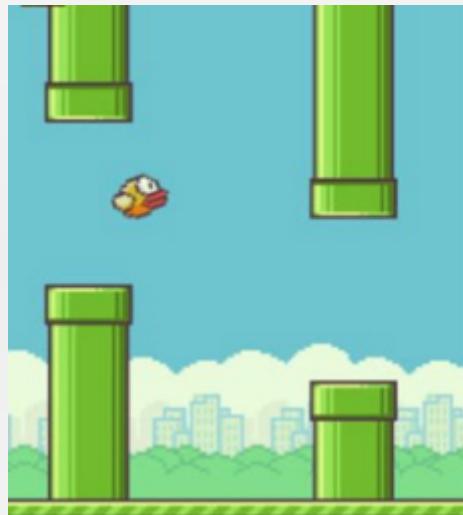
游戏高手： 神经 网络



yenchenlin

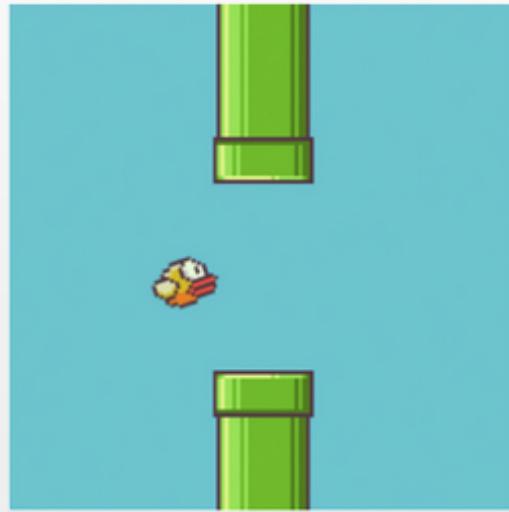


游戏高手：如何计算损失函数？

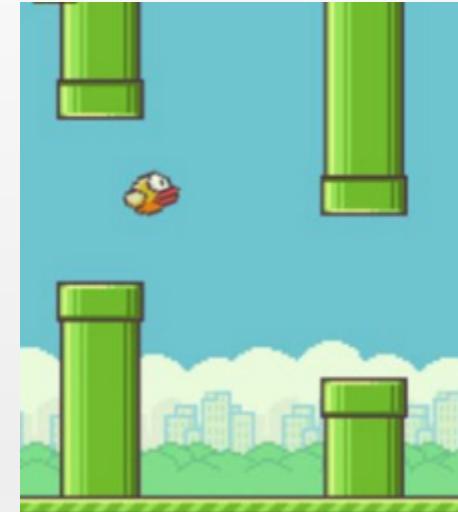


为了计算 $Q(t+1)$ ，我们只需要有 $t+1$ 时刻的 s'

游戏高手：如何计算损失函数？



$s, Q(s,a)$

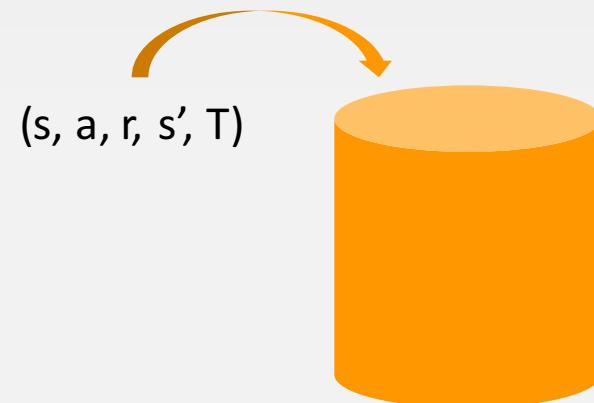
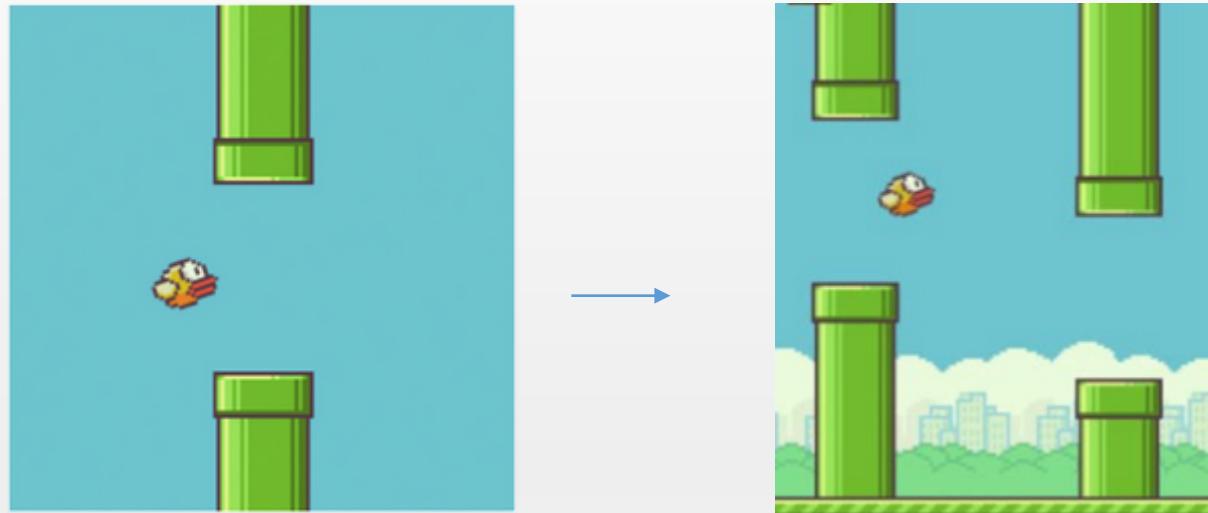


$s', Q(s',a')$

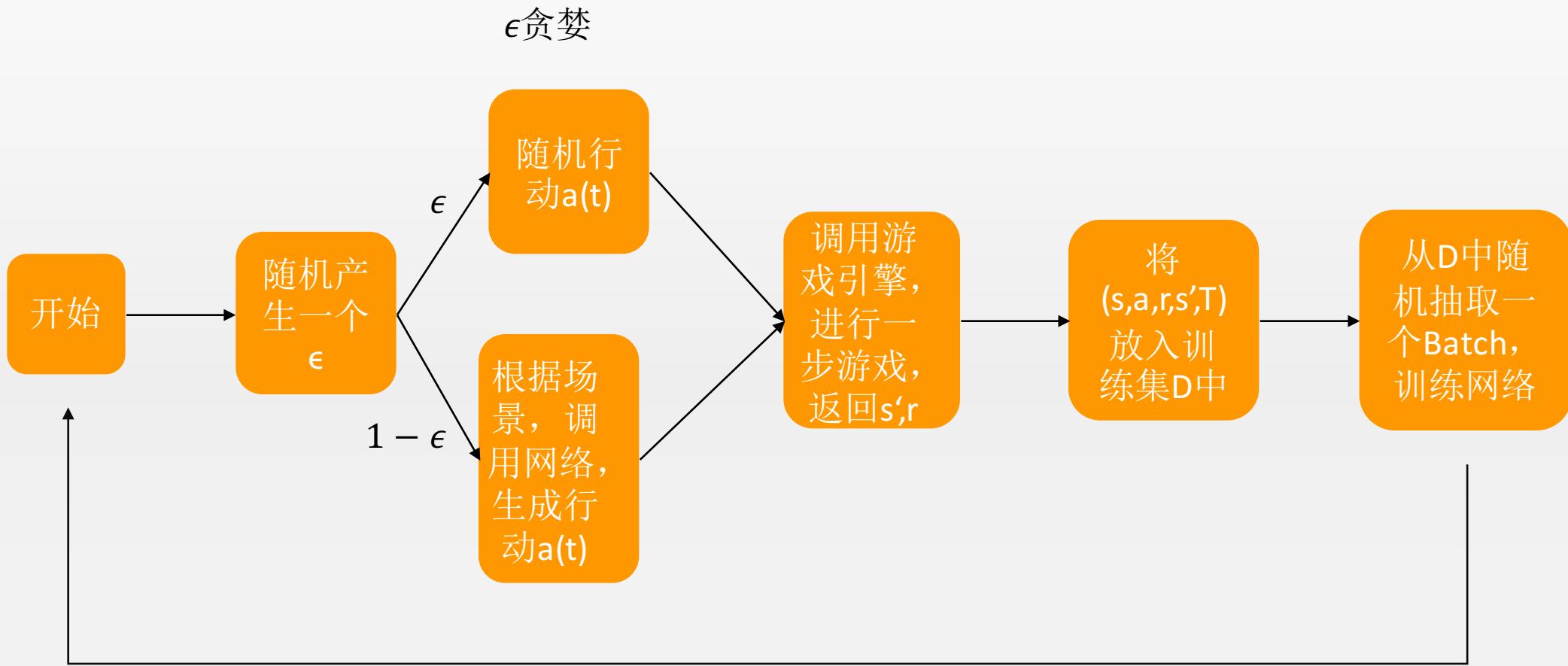
$$L(s, a) = \left(\max_{a'} Q(s', a') - Q(s, a) \right)^2$$

训练： $L.backward()$

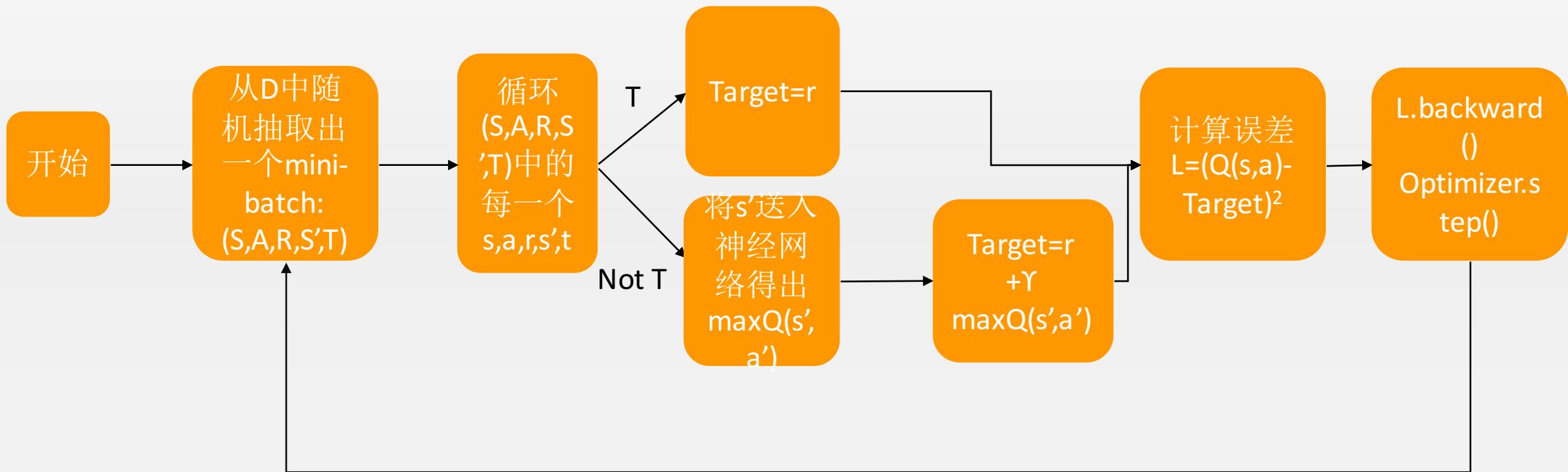
如何形成训练集 (miniBatch) ?



运行流程



训练流程

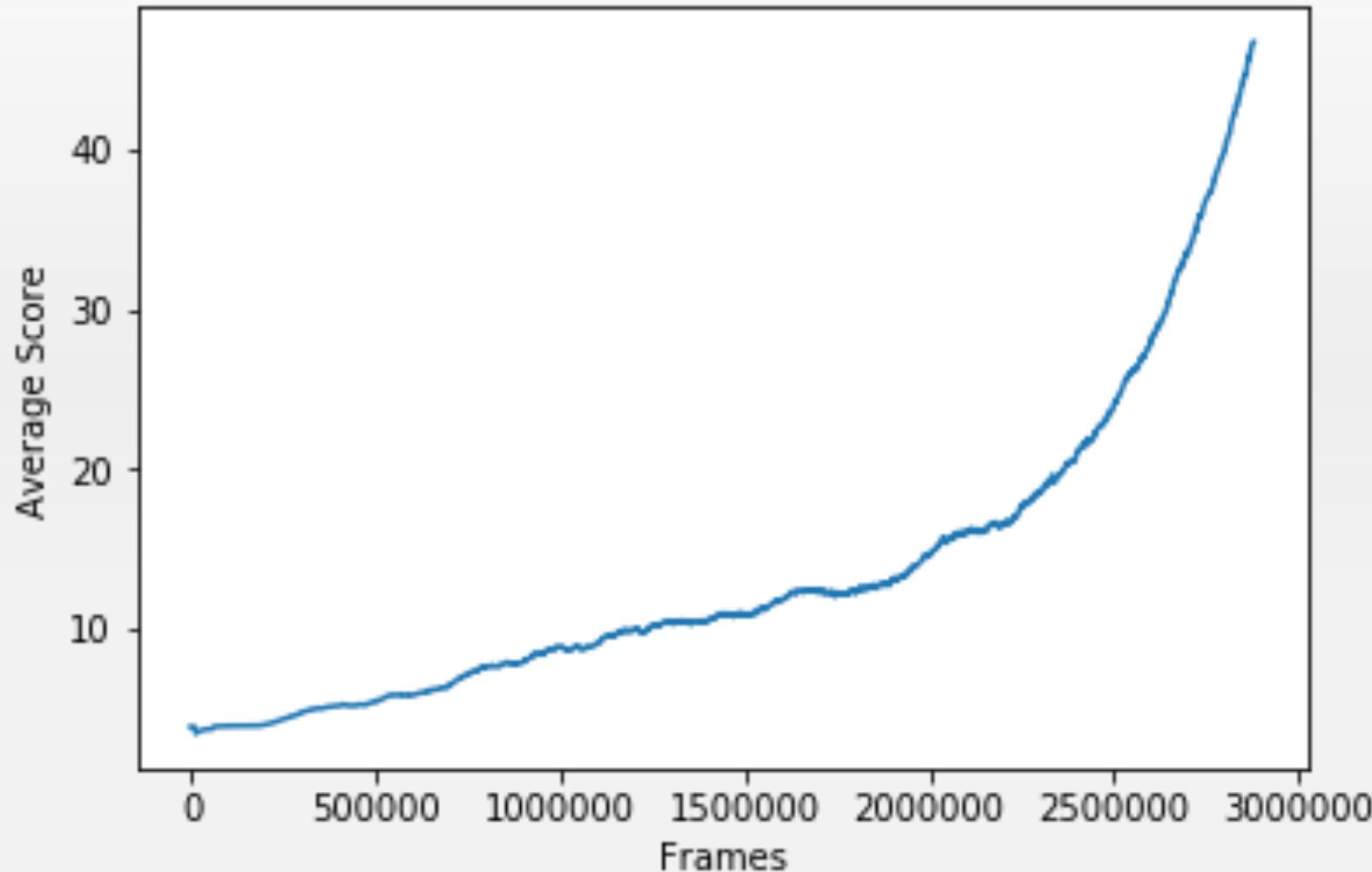


PyGame



- 一个基于Python的游戏制作环境
- 让多媒体操作（图像、声音）等变得简单
- 完全开源、免费

PyGame



LETTER

doi:10.1038/nature14236

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan 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

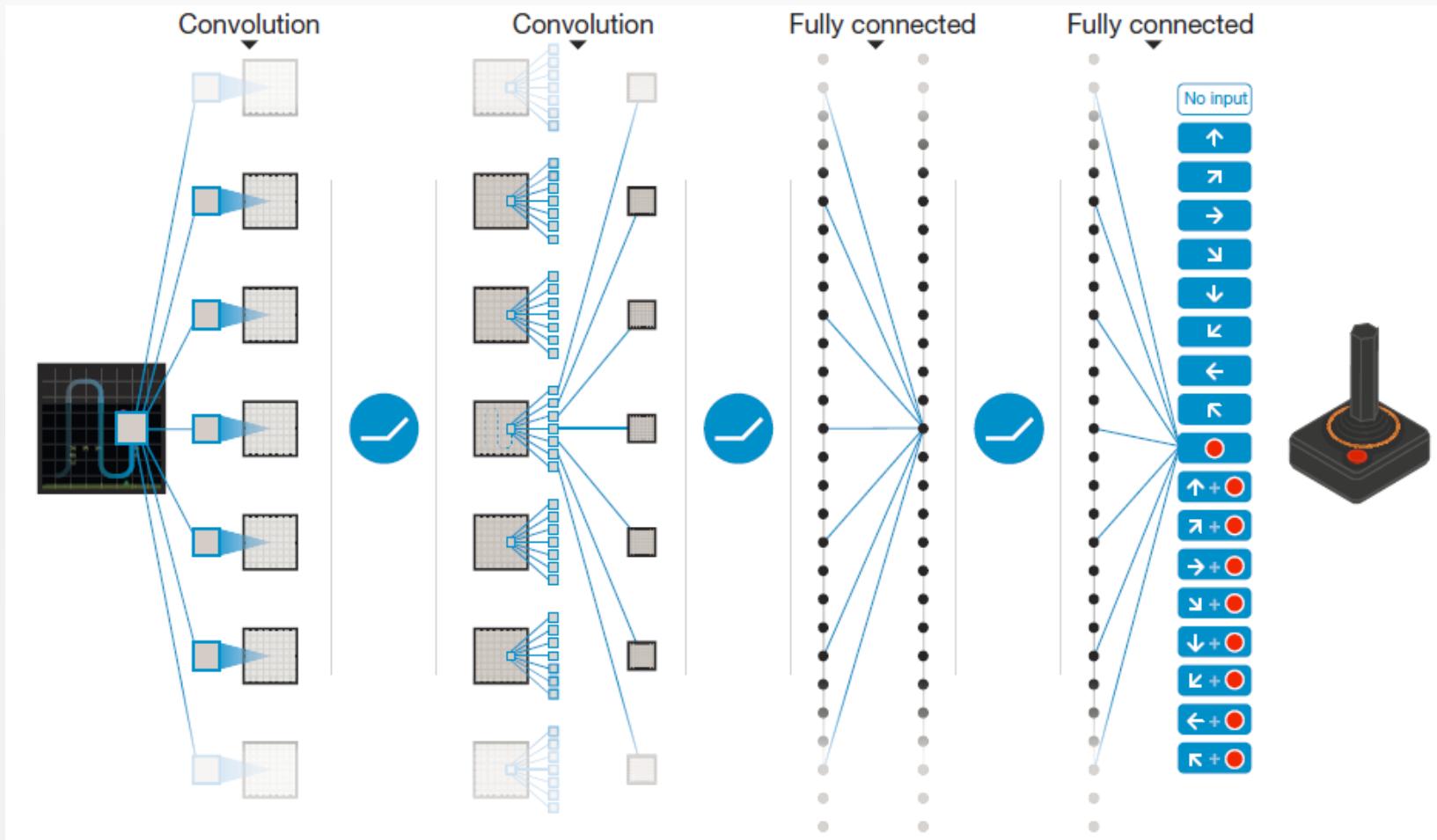
agent is to select actions in a fashion that maximizes cumulative future reward. More formally, we use a deep convolutional neural network to approximate the optimal action-value function

$$Q^*(s,a) = \max_{\pi} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi],$$

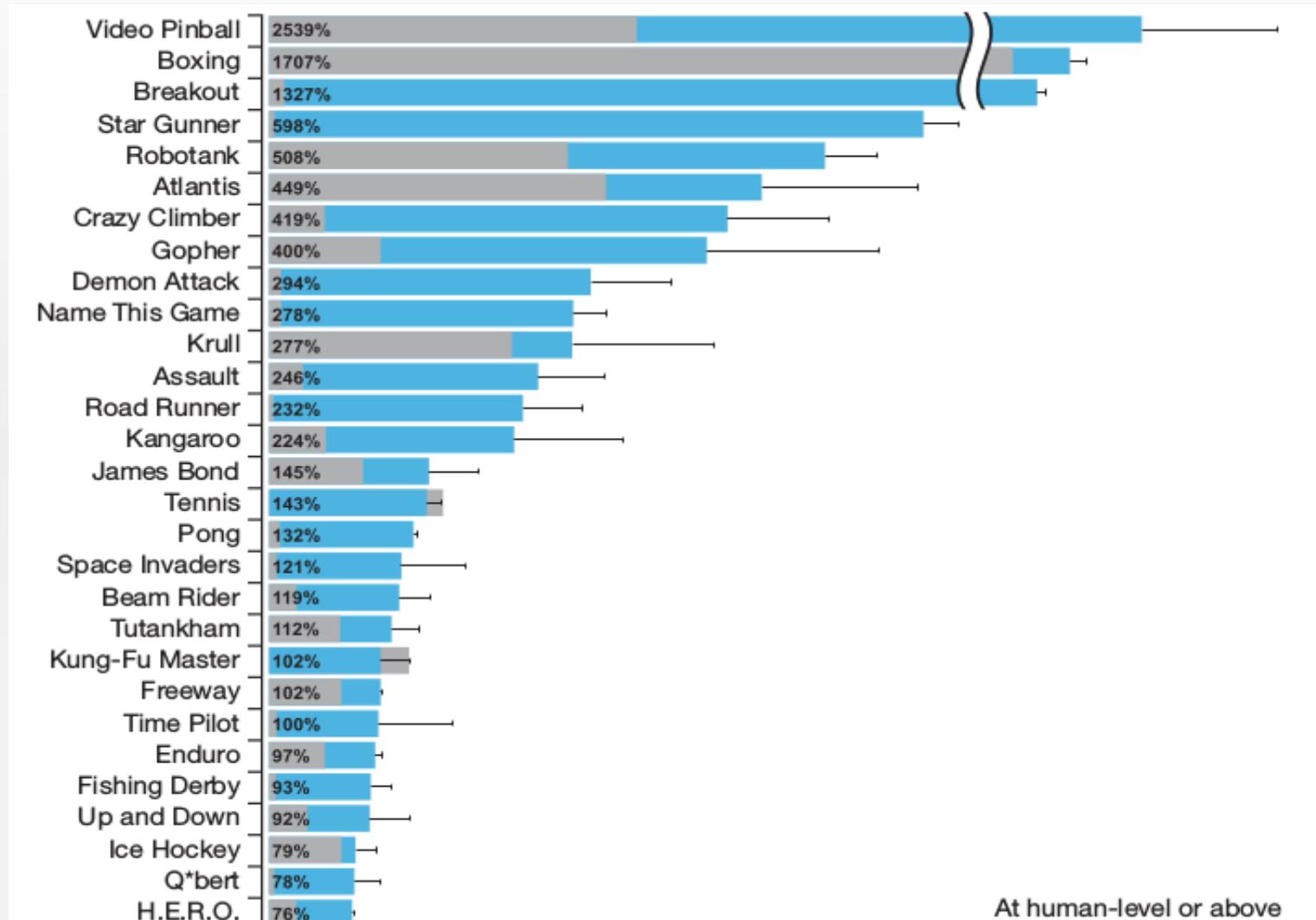
which is the maximum sum of rewards r_t discounted by γ at each time-step t , achievable by a behaviour policy $\pi = P(a|s)$, after making an observation (s) and taking an action (a) (see Methods)¹⁹.

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) function²⁰. This instability has several causes: the correlations present in the sequence of observations, the fact that small updates to Q may significantly change

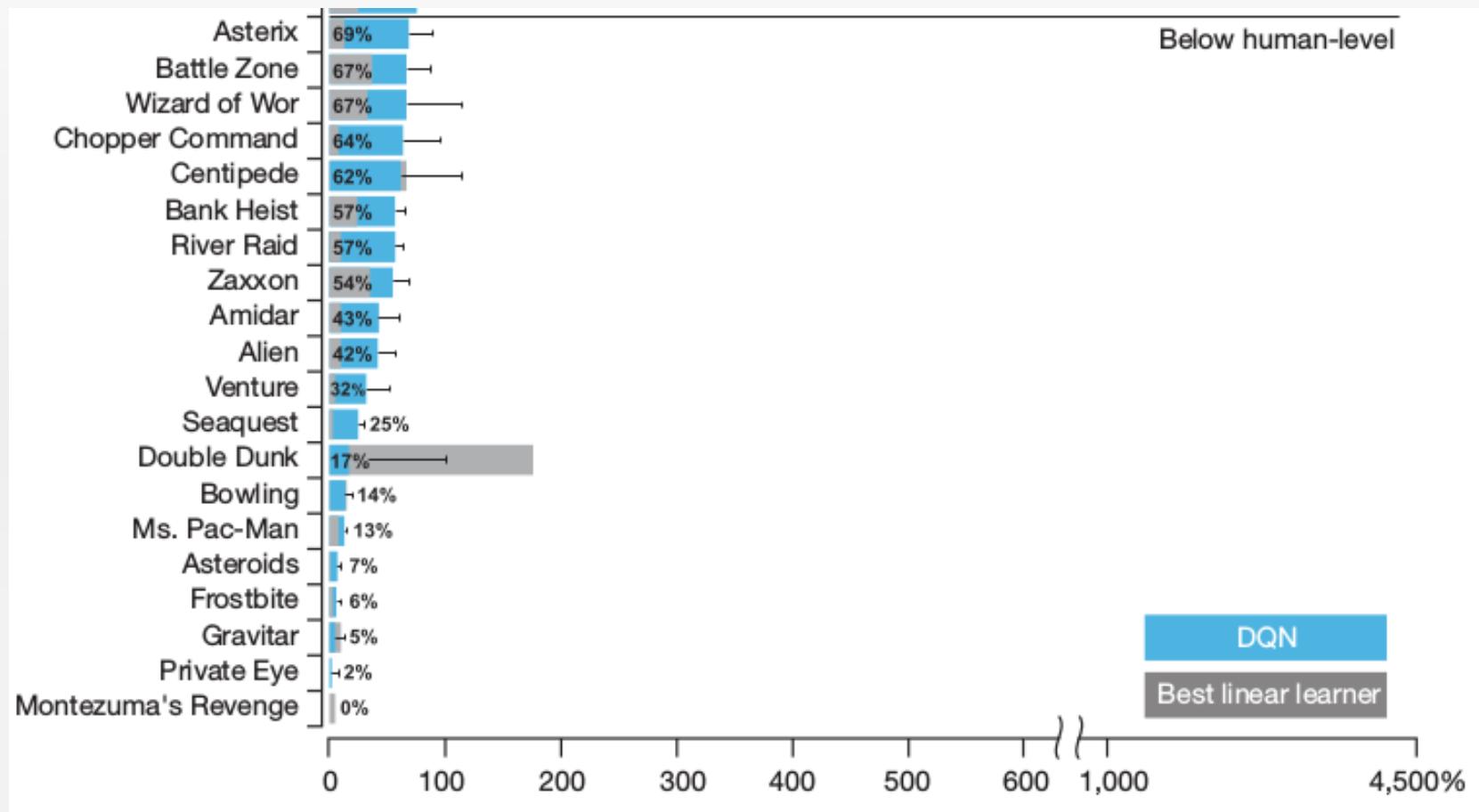
游戏的人工智能



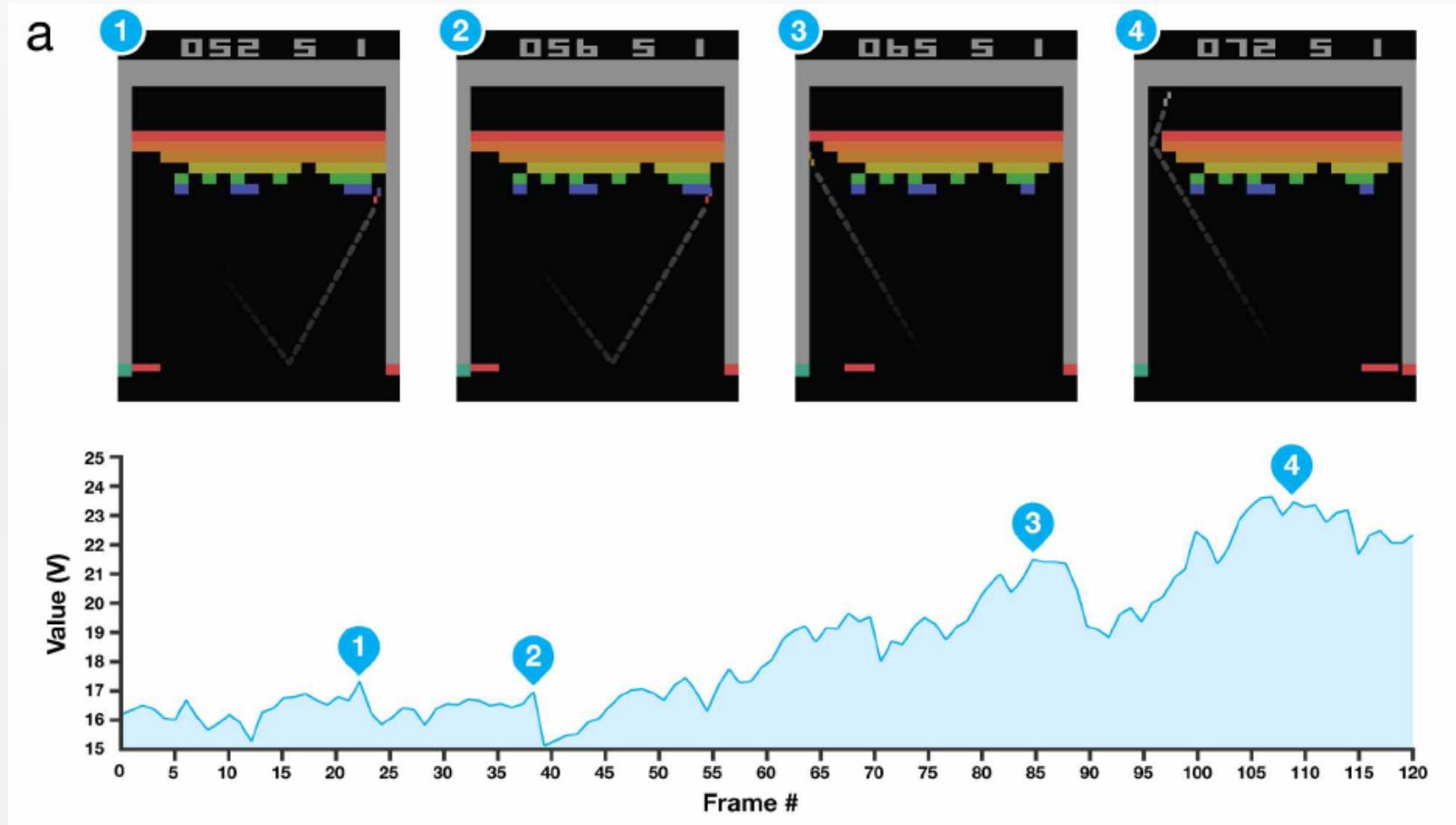
完成效果



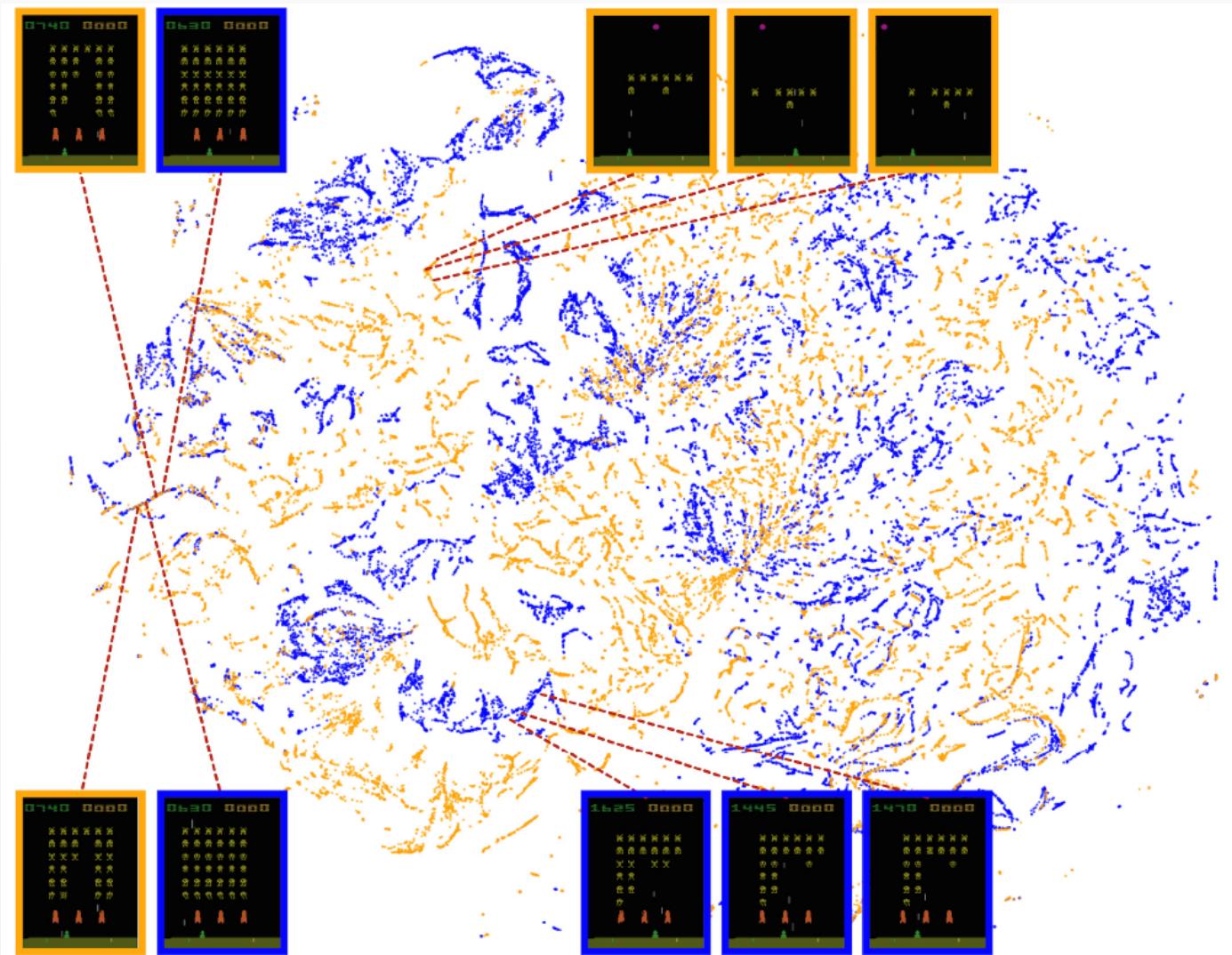
完成效果



完成效果



完成效果



- 蓝色：人类
- 橙色：机器

更多游戏AI



<https://deepmind.com/blog/deep-reinforcement-learning/>

更多游戏AI

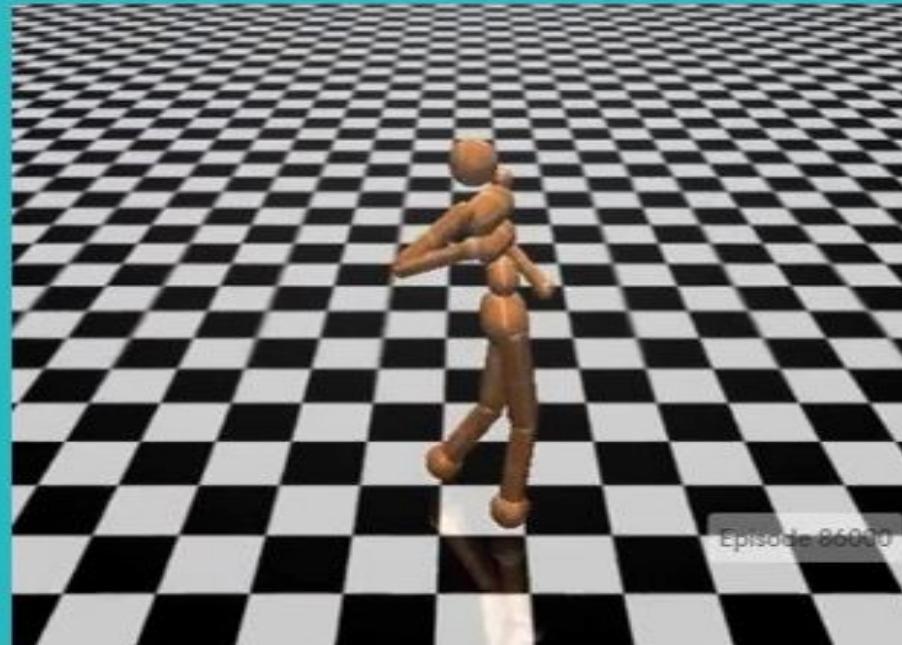
Developing & Comparing RL Algorithms



OpenAI Gym

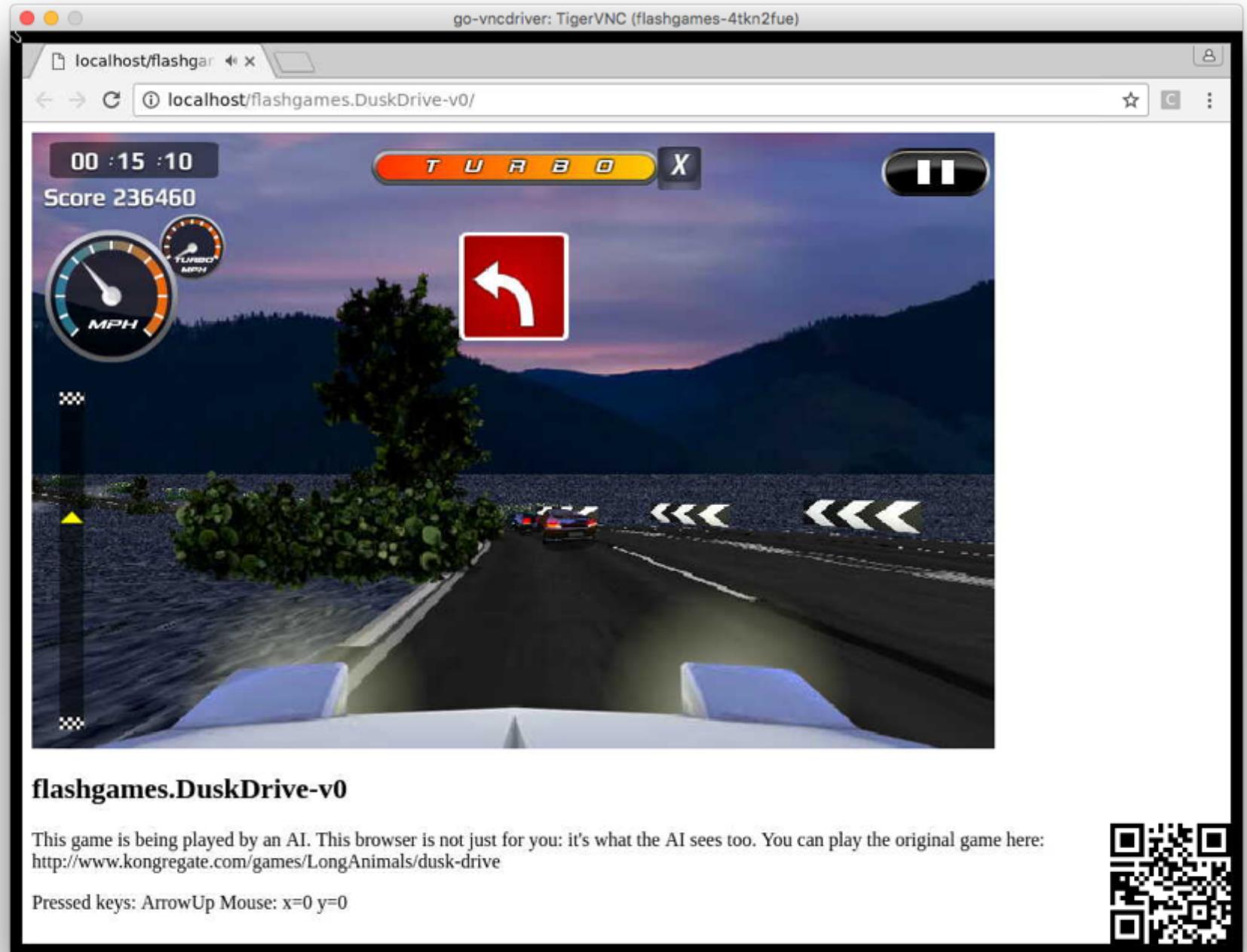
BETA

A toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Go.

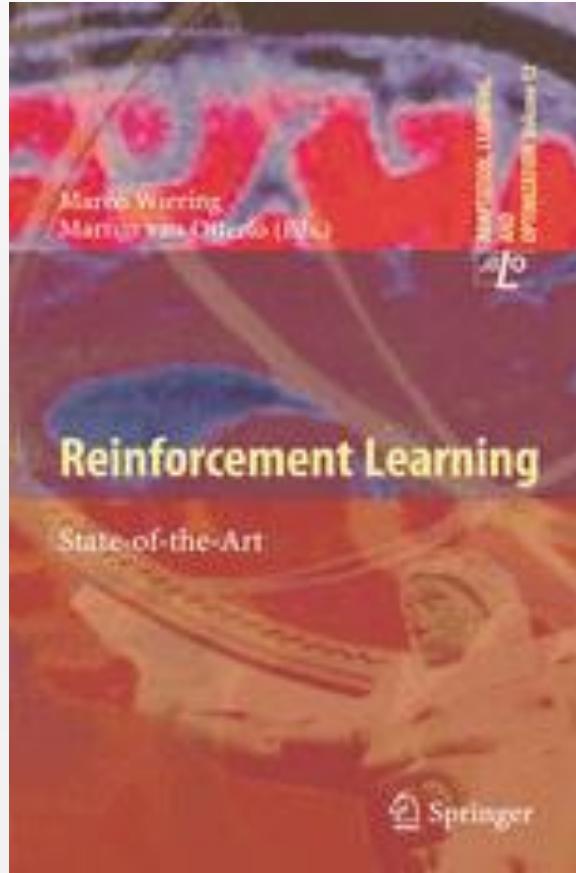


更多游戏AI

OpenAI:
Universe



更多游戏AI



<https://link.springer.com/content/pdf/10.1007%2F978-3-642-27645-3.pdf>

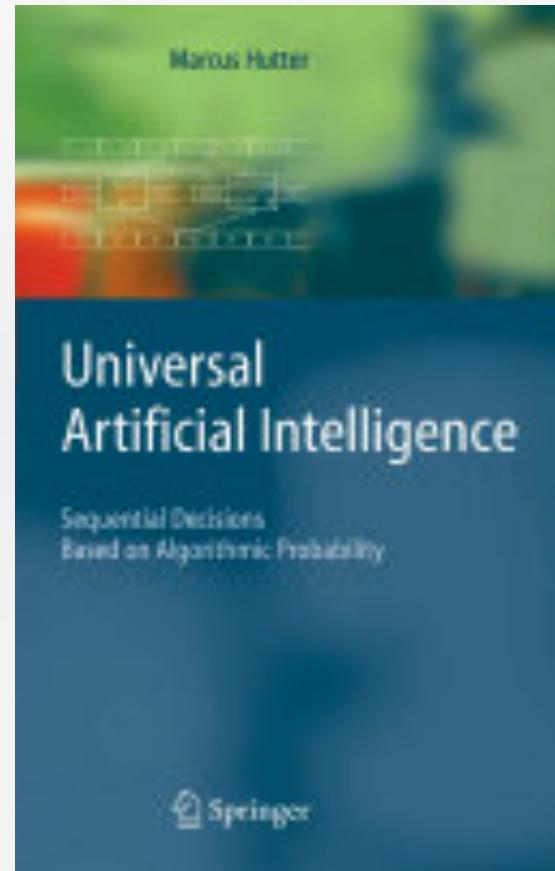
有关AGI (Artificial General Intelligence)

- Artificial general intelligence (AGI) is the intelligence of a machine that could successfully perform any intellectual task that a human being can.
- It is a primary goal of some artificial intelligence research and a common topic in science fiction and future studies.
- Artificial general intelligence is also referred to as “**strong AI**”,^[1] “**full AI**”^[2] or as the ability of a machine to perform “general intelligent action”.^[3] Academic sources reserve “strong AI” to refer to machines capable of experiencing consciousness.
- DQN可以看作游戏世界中的AGI

Universal Artificial Intelligence



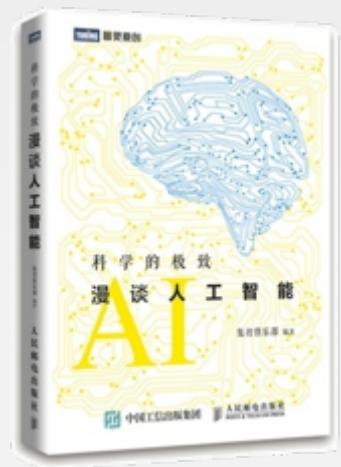
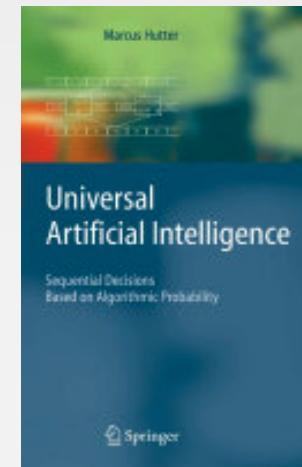
Marcus Hutter



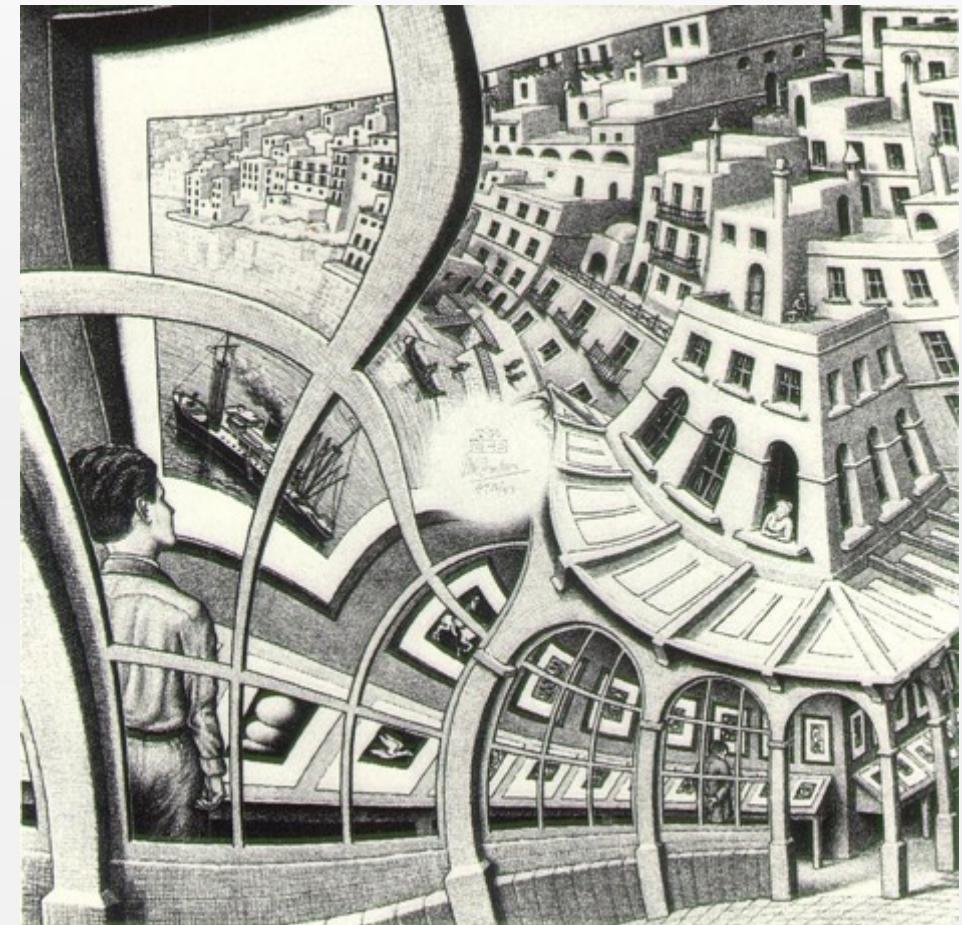
AIXI

一个公式定义人工智能：

$$\arg \max_{a_t} \sum_{o_t r_t} \dots \max_{a_m} \sum_{o_m r_m} [r_t + \dots + r_m] \sum_{q: U(q, a_1 \dots a_m) = o_1 r_1 \dots o_m r_m} 2^{-\text{length}(q)},$$



有关自我意识 (Self-awareness)



递归定理 (Recursion Theorem)

- Kleene第二递归定理
- 对于任意的程序F，总存在一段程序代码c，使得我们执行代码c的结果完全等价于把源代码c作为数据输入给程序F执行的结果。

自我反省的程序 (self-introspection)

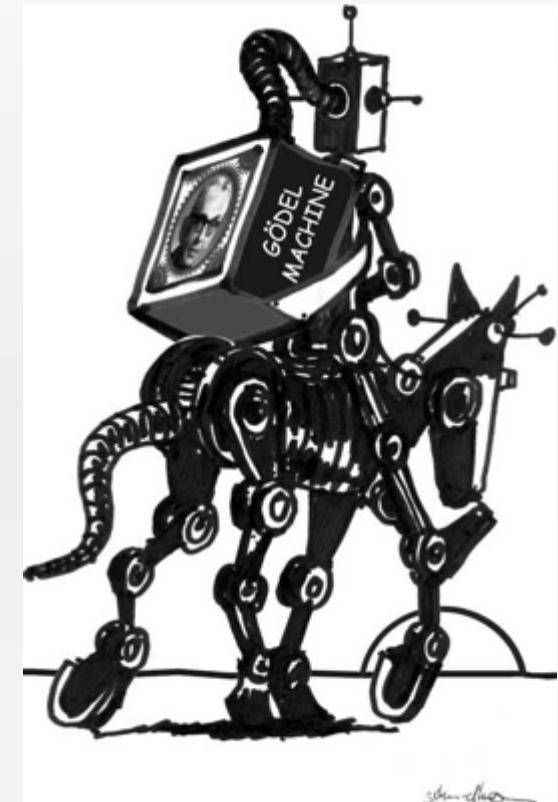
- 自我意识的核心：自我反省
- 存在着这样一种计算机程序 $F_t(x)$ ，它的作用就是计算任意的源代码为 x 的程序在经过 t 时间步的运算后的结果。
- 于是根据递归定理，我们便知道，存在着一个源程序 O ，它所做的就是：把自己的源代码拿出来，然后在自己的虚拟机上模拟自己运算 t 时间步后的结果。

哥德尔机

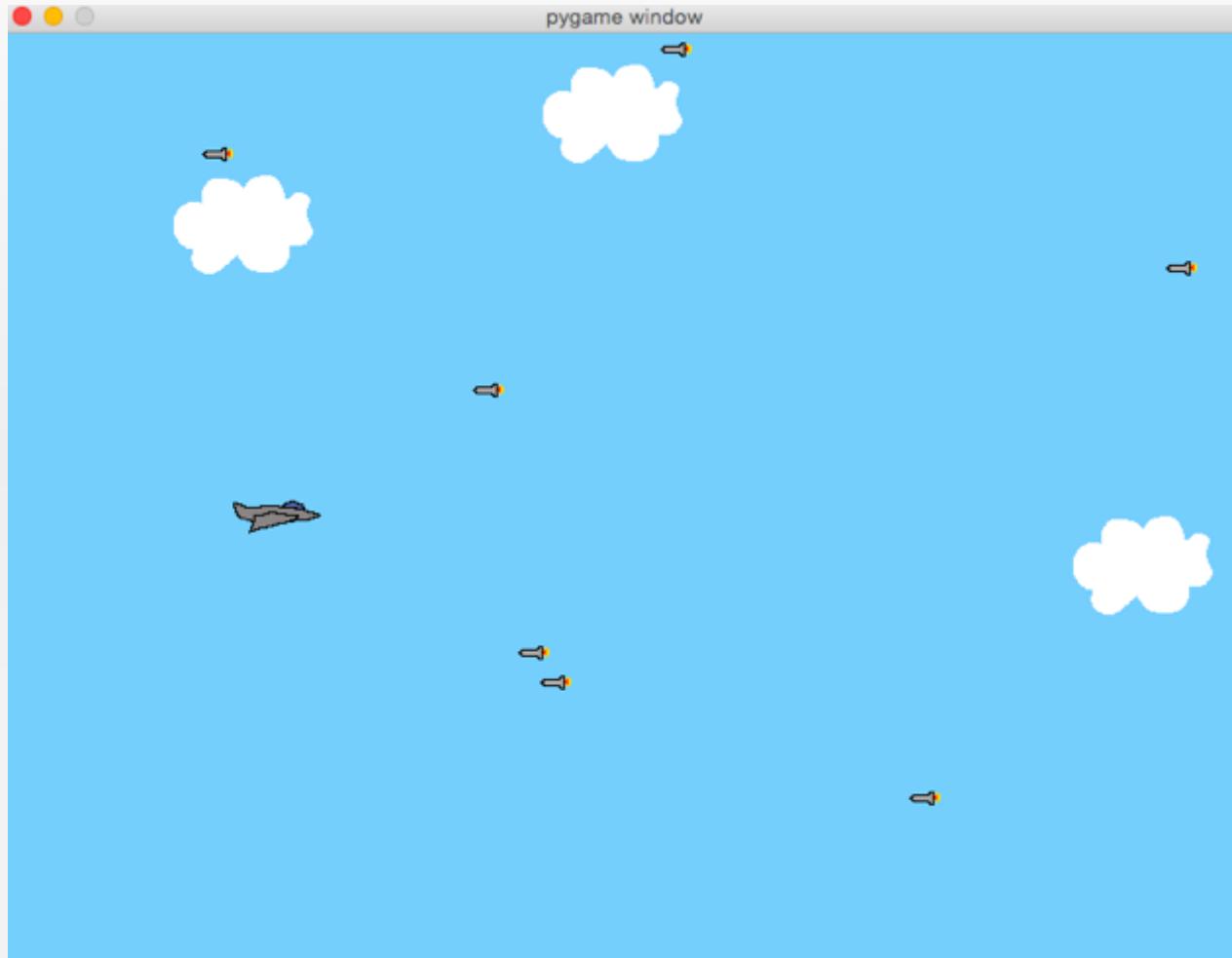
A **Gödel machine** is an approach to [Artificial General Intelligence](#) that uses a [recursive self-improvement](#) architecture proposed by Jürgen Schmidhuber. It was inspired by the mathematical theories of Kurt Gödel, where one could always find a mathematical truth or axiom that if attached to a formal system would make it stronger. A Gödel Machine is a universal problem solver that will make provably optimal self-improvements – self-improvements which can be proved to better maximize its [utility](#)



Jürgen Schmidhuber



项目：飞机躲炸弹



- 制作一个飞机躲炸弹的游戏：
<https://realpython.com/blog/python/pygame-a-primer/>
- 训练一个AI程序自动玩这个游戏

敬请期待

