

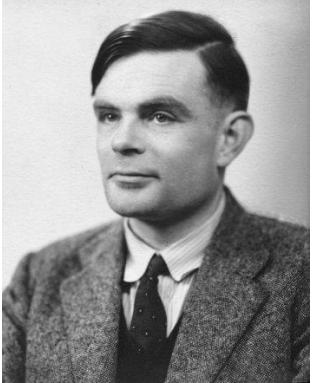
具身智能

刘华平

2025年4月17日

-
- **背景：具身智能**
 - 具身智能的体系
 - 具身智能关键技术
 - 探索与实践
 - 面向具身智能的AIGC
 - 总结

➤ 人工智能的发展



COMPUTING MACHINERY AND INTELLIGENCE

By A. M. Turing

1. The Imitation Game

I propose to consider the question, "Can machines think?" This should begin with definitions of the meaning of the terms "machine" and "think." The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words "machine" and "think" are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the question, "Can machines think?" is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words.

In the process of trying to imitate an adult human mind we are bound to think a good deal about the process which has brought it to the state that it is in. We may notice three components,

- (a) The initial state of the mind, say at birth,
- (b) The education to which it has been subjected,
- (c) Other experience, not to be described as education, to which it has been subjected.

Instead of trying to produce a programme to simulate the adult mind, why not rather try to produce one which simulates the child's? If this were then subjected to an appropriate course of education one would obtain the adult brain. Presumably the child-brain is something like a notebook as one buys it from the stationers. Rather little mechanism, and lots of blank sheets. (Mechanism and writing are from our point of view almost synonymous.) Our hope is that there is so little mechanism in the child-brain that something like it can be easily programmed. The amount of work in the education we can assume, as a first approximation, to be much the same as for the human child.

We may hope that machines will eventually compete with men in all purely intellectual fields. But which are the best ones to start with? Even this is a difficult decision. Many people think that a very abstract activity, like the playing of chess would be best. It can also be maintained that it is best to provide the machine with the best sense organs that money can buy, and then teach it to understand and speak English. This process could follow the normal teaching of a child. Things would be pointed out and named, etc. Again I do not know what the right answer is, but I think both approaches should be tried.

We can only see a short distance ahead, but we can see plenty there that needs to be done.

我们的目光所及，只是
不远的前方。

早在1950年，图灵就曾在他的经典论文《计算机械与智能》中提到，人工智能的未来将走向两条路：一条是“计算的智能体”，即像阿尔法狗一样依赖强大计算力的大模型；另一条路则是“**具身智能**”，机器拥有像婴儿一样的感知能力，通过与环境的互动和学习不断进化。换句话说，具身智能并不是最近才火的概念，它的源流可以追溯到人工智能的最早理论。

1 背景

➤ 什么是具身智能

身体利用感知-运动系统在与环境交互过程中产生智能



具身认知

- ✓ Be Multimodal
- ✓ Be Incremental
- ✓ Be Physical
- ✓ Explore
- ✓ Be Social
- ✓ Learn a Language



具身人工智能

- 多模态感知, ...
- 持续学习, ...
- 物理交互, ...
- 灵活探索, ...
- 社交模仿, ...
- 语言学习, ...

- L. Smith and M. Gasser, "The development of embodied cognition: six lessons from babies.,," Artificial Life, vol.11, no. 1-2, 2005.

——Smith and Gasser

The embodiment hypothesis is the idea that intelligence emerges in the **interaction** of an agent with an environment and as a result of **sensorimotor** activity.

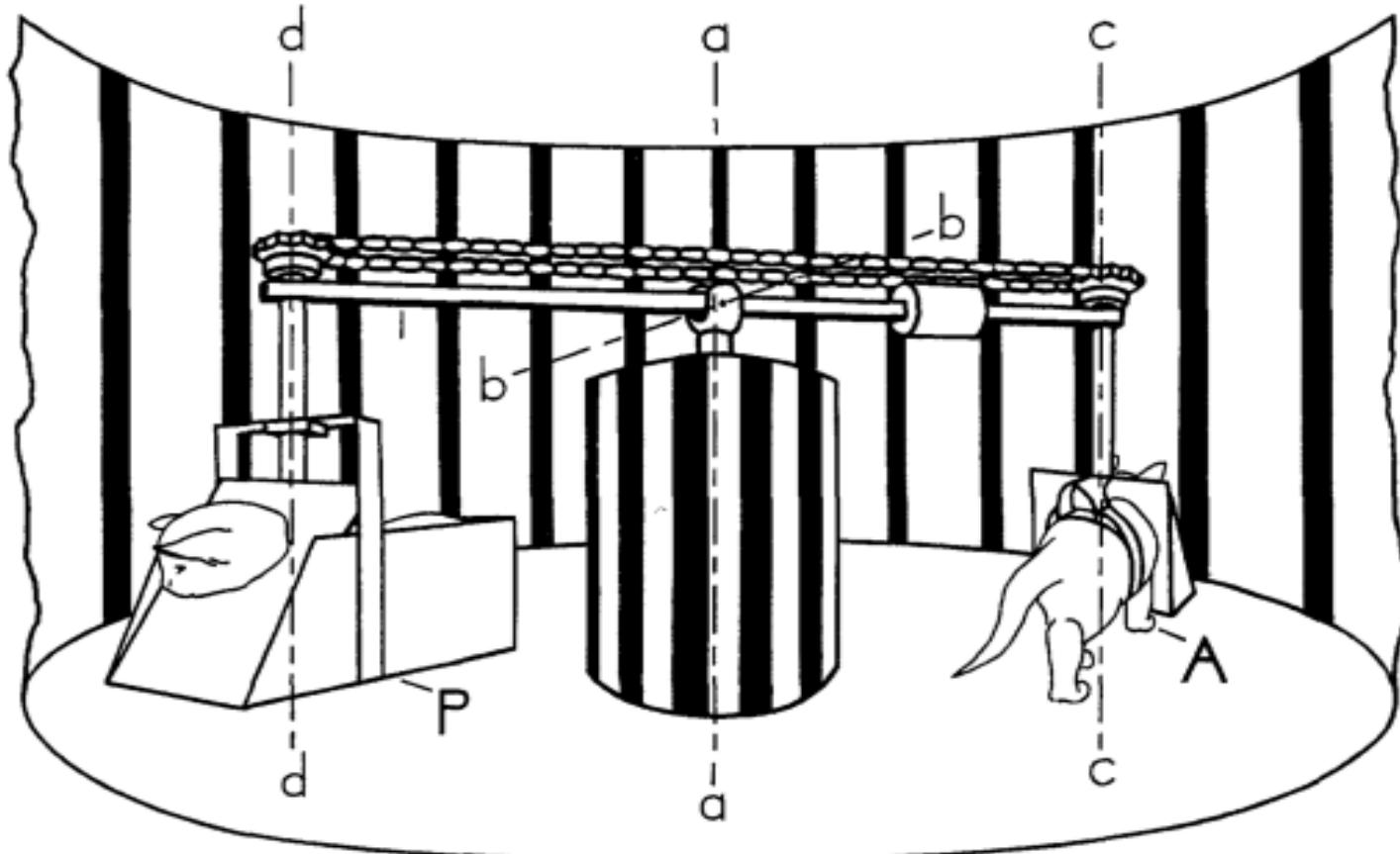
1 背景

➤ 无处不在的具身智能



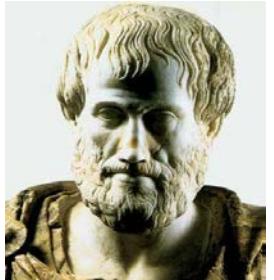
1 背景

➤ 极端的例子



1 背景

➤ 很长的过去



Aristotle
(公元前384-前322)



Charles Robert Darwin
(1809-1882)



Claude Bernard
(1813-1878)



Walter Bradford Cannon
(1871-1945)



Jean Piaget
(1896-1980)



James J. Gibson
(1904-1979)

身体的主动性首先在于
“通过触摸而感觉”

心境和动作之间的联系
就是情感的真正含义

“没有脑袋”的生理学

身体的智慧

动作是认识的源泉

感知与动作在与环境
交互中密切联系



René Descartes
1596-1650

我思故我在



Martin Heidegger
(1889-1976)

《存在与时间》
比表征更基本的活动
是操劳



Maurice Merleau-Ponty
(1908-1961)

知觉现象学：身体和主
体其实是同一个实在



Norbert Wiener
(1894-1964)

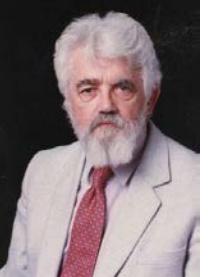
控制论

1 背景

➤ 很短的历史



Marvin Lee Minsky
(1927-2016)



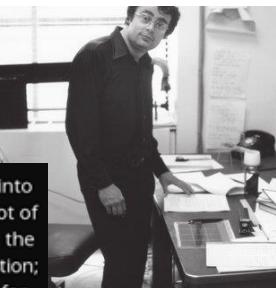
达特茅斯会议
(1956)



John Haugeland
(1945-2010)

GOF AI

- 人类认知和思维的基本单元是符号，而认知过程就是在符号表示上的一种运算。
- 人和计算机都是物理符号系统,因此可以用计算机来模拟人的智能行为,即用计算机的符号操作来模拟人的认知过程。
- 启发式算法 → **专家系统** → 知识工程



David Marr
(1945-1980)

When David Marr at MIT moved into computer vision, he generated a lot of excitement, but he hit up against the problem of knowledge representation; he had no good representations for knowledge in his vision systems.
— Marvin Minsky —

心有余而力不足

The spirit is willing, but the flesh is weak



The vodka is good but the meat is spoiled

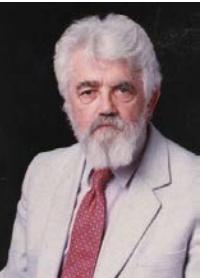
1 背景

➤ 很短的历史



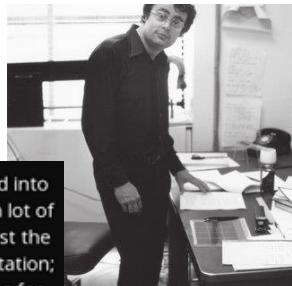
Marvin Lee Minsky John McCarthy
(1927-2016) (1927-2011)

达特茅斯会议 (1956)



John Haugeland (1945-2010)

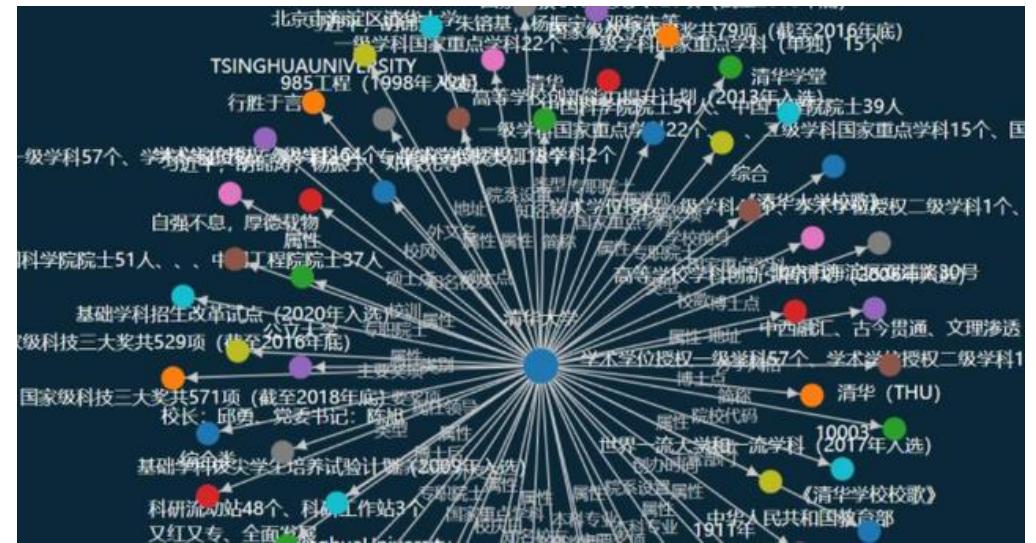
GOF AI



When David Marr at MIT moved into computer vision, he generated a lot of excitement, but he hit up against the problem of knowledge representation; he had no good representations for knowledge in his vision systems.

— Marvin Minsky —

- 人类认知和思维的基本单元是符号，而认知过程就是在符号表示上的一种运算。
 - 人和计算机都是物理符号系统,因此可以用计算机来模拟人的智能行为,即用计算机的符号操作来模拟人的认知过程。
 - 启发式算法 → 专家系统 → 知识工程



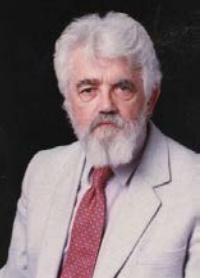
1 背景

9

➤ 很短的历史



Marvin Lee Minsky
(1927-2016)
John McCarthy
(1927-2011)
达特茅斯会议
(1956)



John Haugeland
(1945-2010)
GOF AI



Hopfield
(1933-)

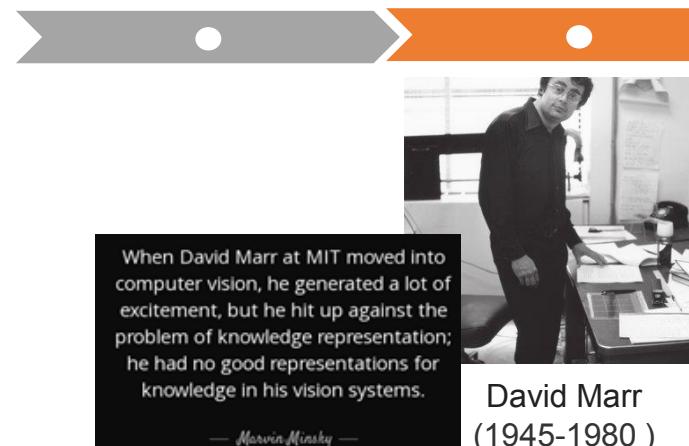


Rumelhart
(1942-2011)



Hilton
(1947-)

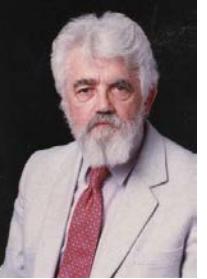
- 数据挖掘
- 机器学习
- 模式识别



➤ 很短的历史



Marvin Lee Minsky
(1927-2016)
John McCarthy
(1927-2011)
达特茅斯会议
(1956)



John Haugeland
(1945-2010)
GOF AI



Hopfield
(1933-)



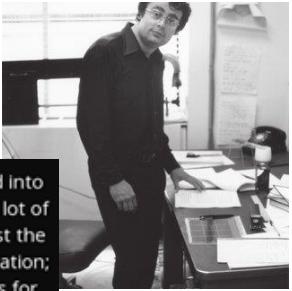
Rumelhart
(1942-2011)



Hilton
(1947-)

- 数据挖掘
- 机器学习
- 模式识别

神经网络->反向传播->深度学习



When David Marr at MIT moved into computer vision, he generated a lot of excitement, but he hit up against the problem of knowledge representation; he had no good representations for knowledge in his vision systems.
— Marvin Minsky —

David Marr
(1945-1980)

已经在“人类最后智力骄傲”上碾压人类的AlphaGo ...



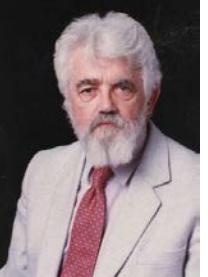
却连挪动一枚小小的棋子都需要人类帮助才能完成

1 背景

➤ 很短的历史



Marvin Lee Minsky
(1927-2016)



达特茅斯会议
(1956)



John Haugeland
(1945-2010)

GOF AI



Hopfield
(1933-)



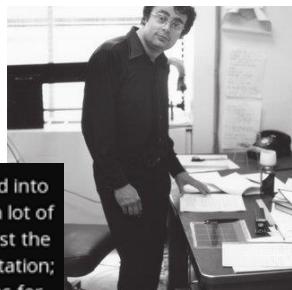
Rumelhart
(1942-2011)



Hilton
(1947-)

- 模式识别
- 数据挖掘
- 机器学习

神经网络->反向传播->深度学习



When David Marr at MIT moved into computer vision, he generated a lot of excitement, but he hit up against the problem of knowledge representation; he had no good representations for knowledge in his vision systems.
— Marvin Minsky —

David Marr
(1945-1980)



Pfeifer
(1947-)



Moravec
(1948-)



Brooks
(1954-)



Cangelosi
(1967-)

- 机器人学
- 机构学
- 形态智能

• 智能是具身化和情境化的，智能需要一个身体

• 强化

➤ 莫拉维克悖论

要让电脑如成人般地下棋是相对容易的，但是要让电脑有如一岁小孩般的感知和行动能力却是相当困难的。



Hans Moravec



Donald Knuth

人工智能已经在几乎所有需要思考的领域超过了人类，但是在那些人类和其它动物不需要思考就能完成的事情上，还差得很远。

——Donald Knuth (2019)



1997: IBM的DeepBlue战胜卡斯帕罗夫



2016: DeepMind的AlphaGo” 战胜李世石

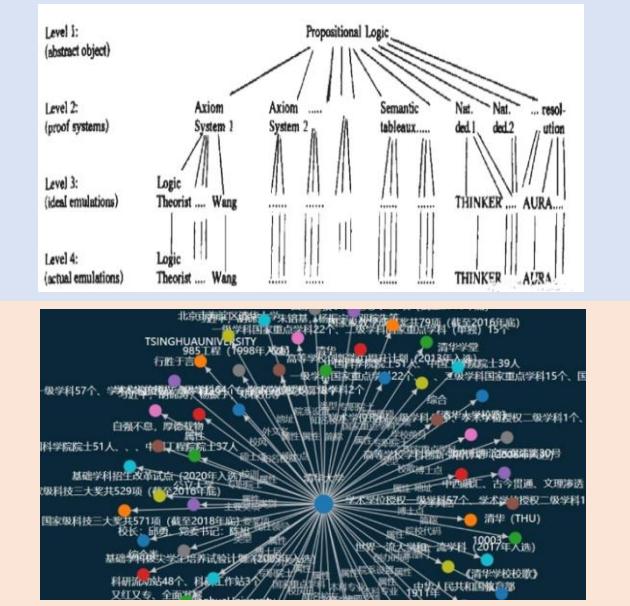


1 背景

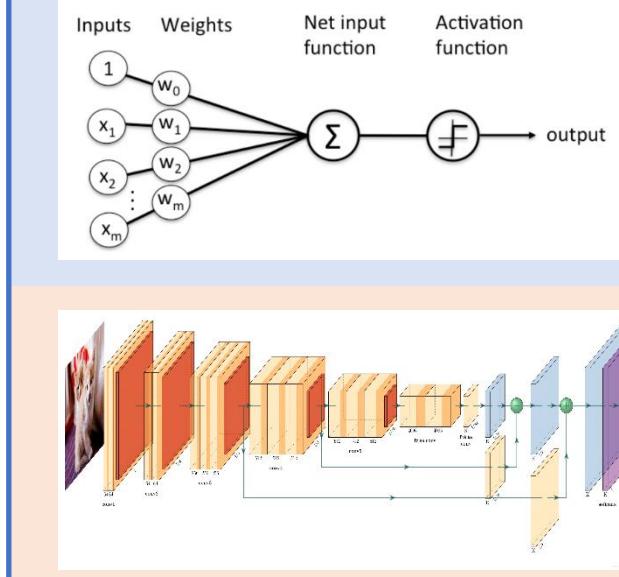
➤ 离身智能 v. s. 具身智能

具身智能

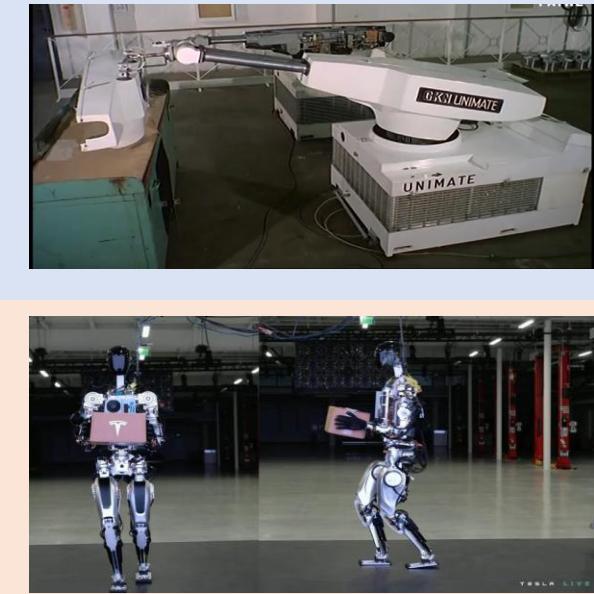
符号主义：表示



联结主义：计算



行为主义：交互



离身智能

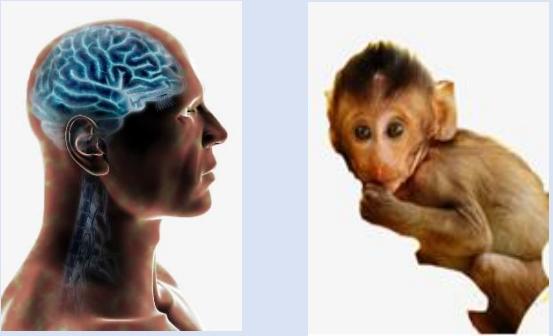
1 背景

➤ 共同的目标

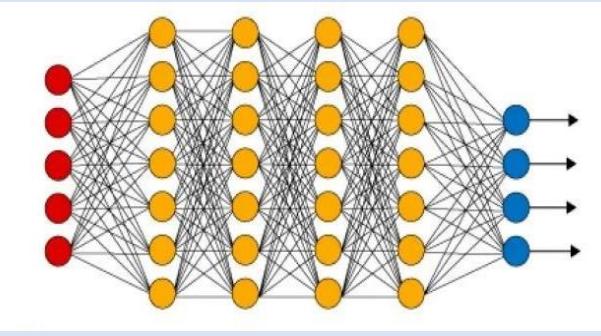
离身智能

具身智能

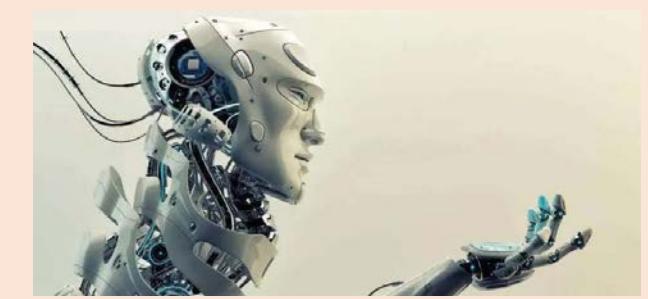
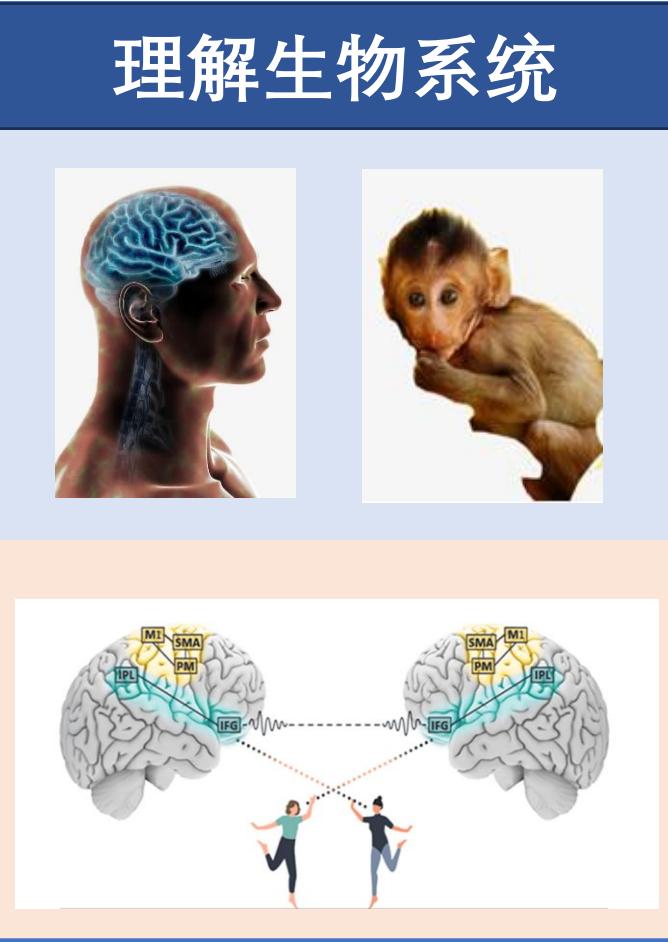
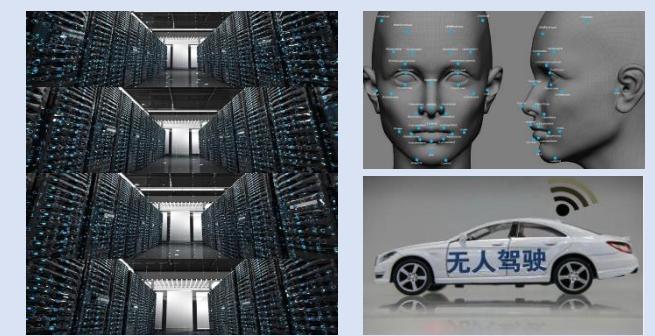
理解生物系统



智能行为的抽象

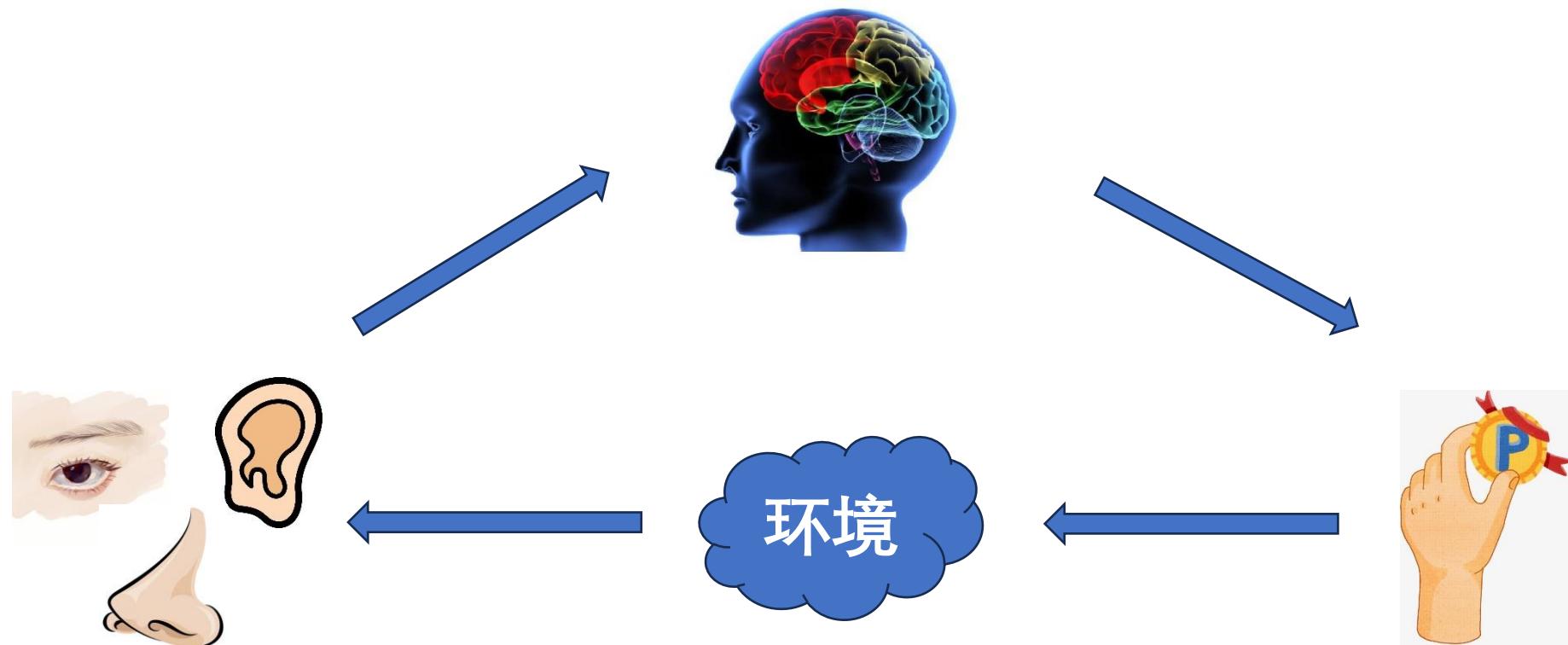


设计人工智能体



➤ 感知-动作回路才是认知的中心

一旦我们开始探索生物与环境之间的协调与交互，就很难再确定感知是什么时候结束的，而认知是什么时候开始的。



-
- **背景：具身智能**
 - **具身智能的体系**
 - **具身智能关键技术**
 - **探索与实践**
 - **面向具身智能的AIGC**
 - **总结**

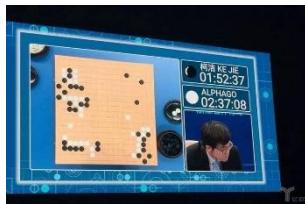
2 具身智能的体系

➤ 狹义与广义的具身智能



离身智能

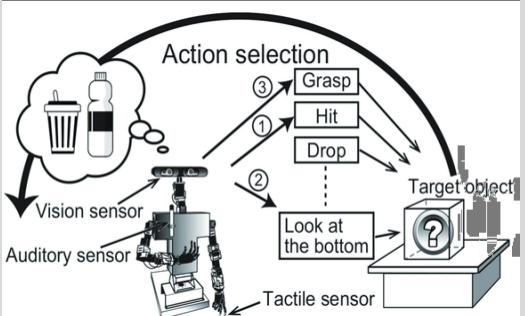
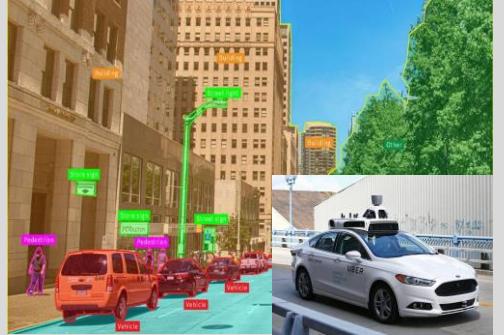
具身智能



2 具身智能的体系

➤ 具身智能的典型任务

发现



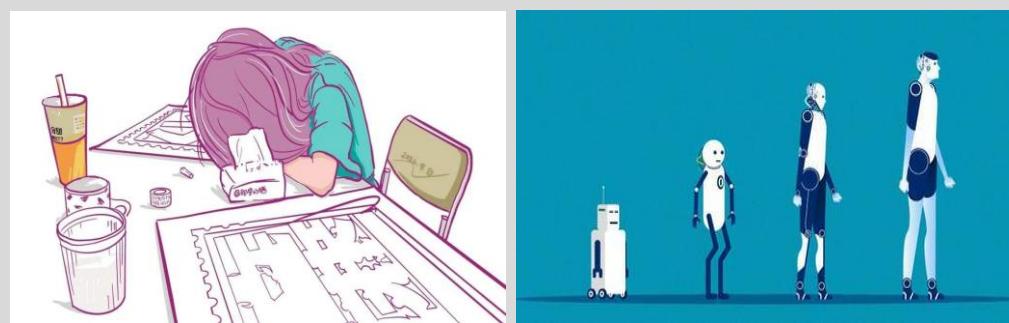
学习



控制



优化



具身智能

离身智能

2 具身智能的体系

➤ 具身智能与机器人的关系



2 具身智能的体系

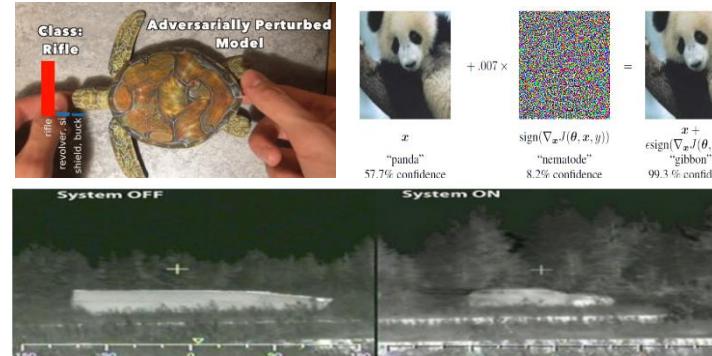
➤ 具身智能的优点、缺点与难点

Disembodied
Intelligence

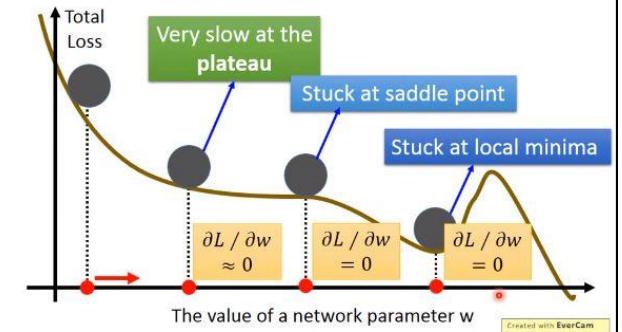
Good



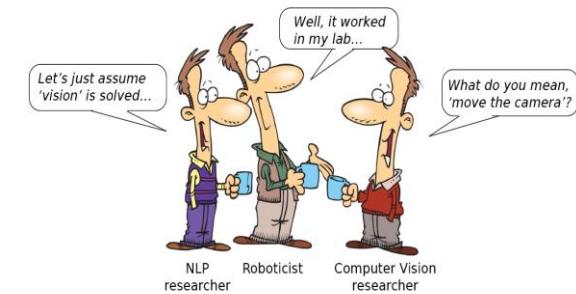
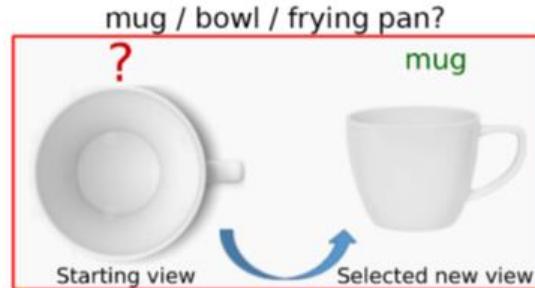
Bad



Ugly

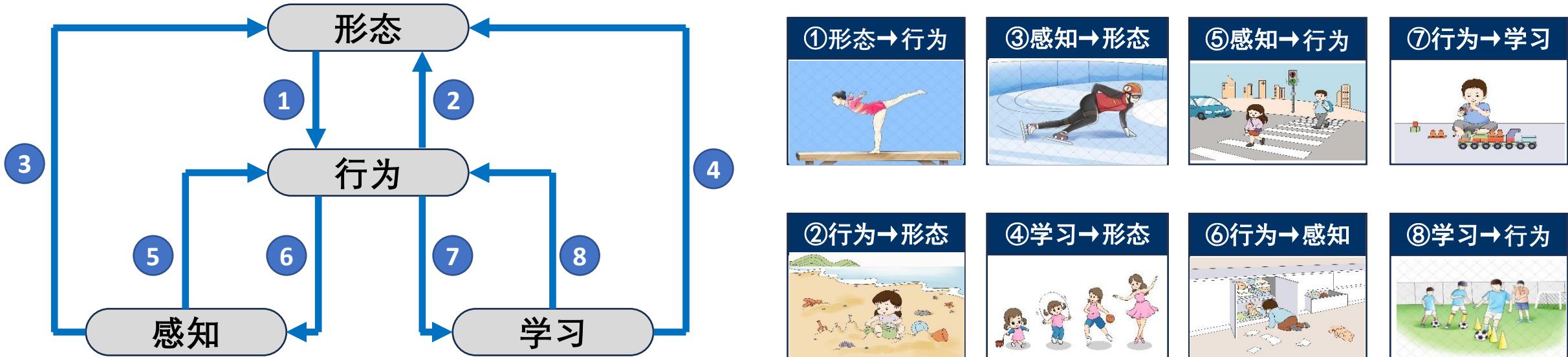


Embodied
Intelligence



2 具身智能的体系

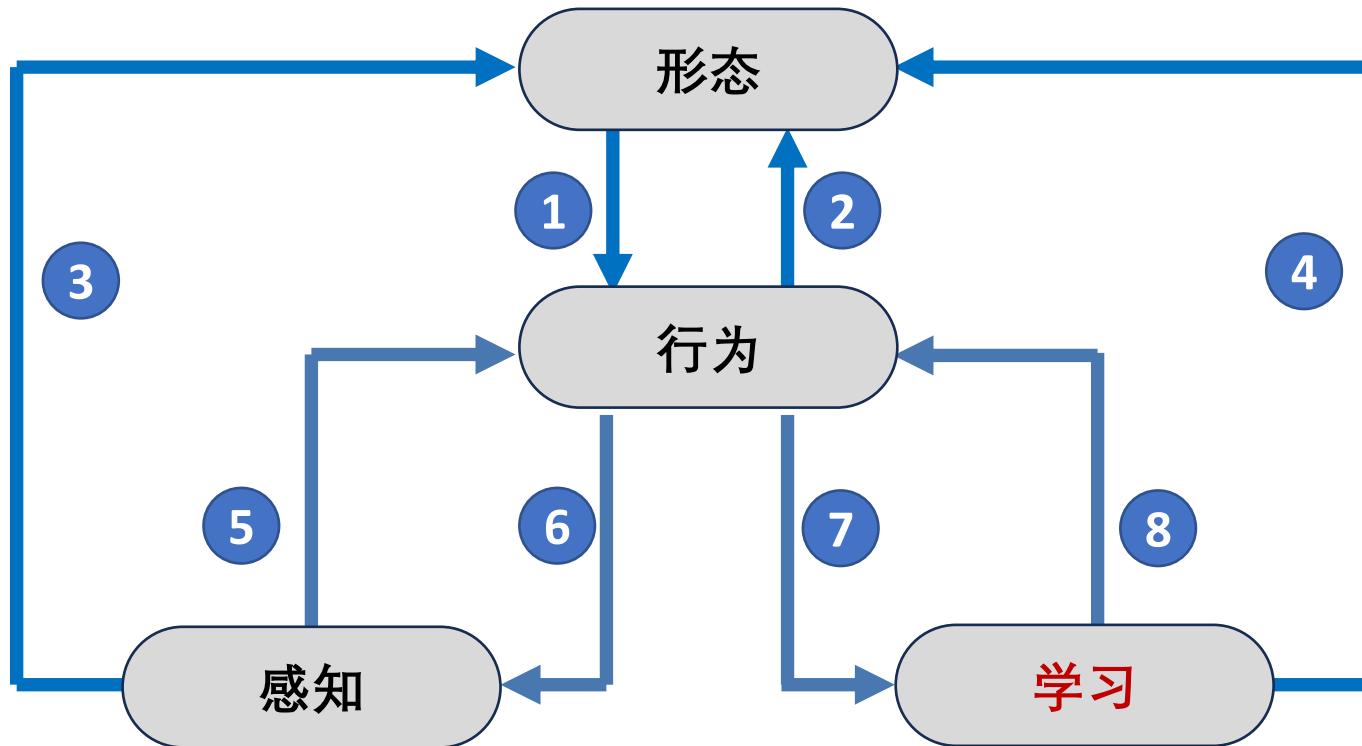
具身智能的体系结构



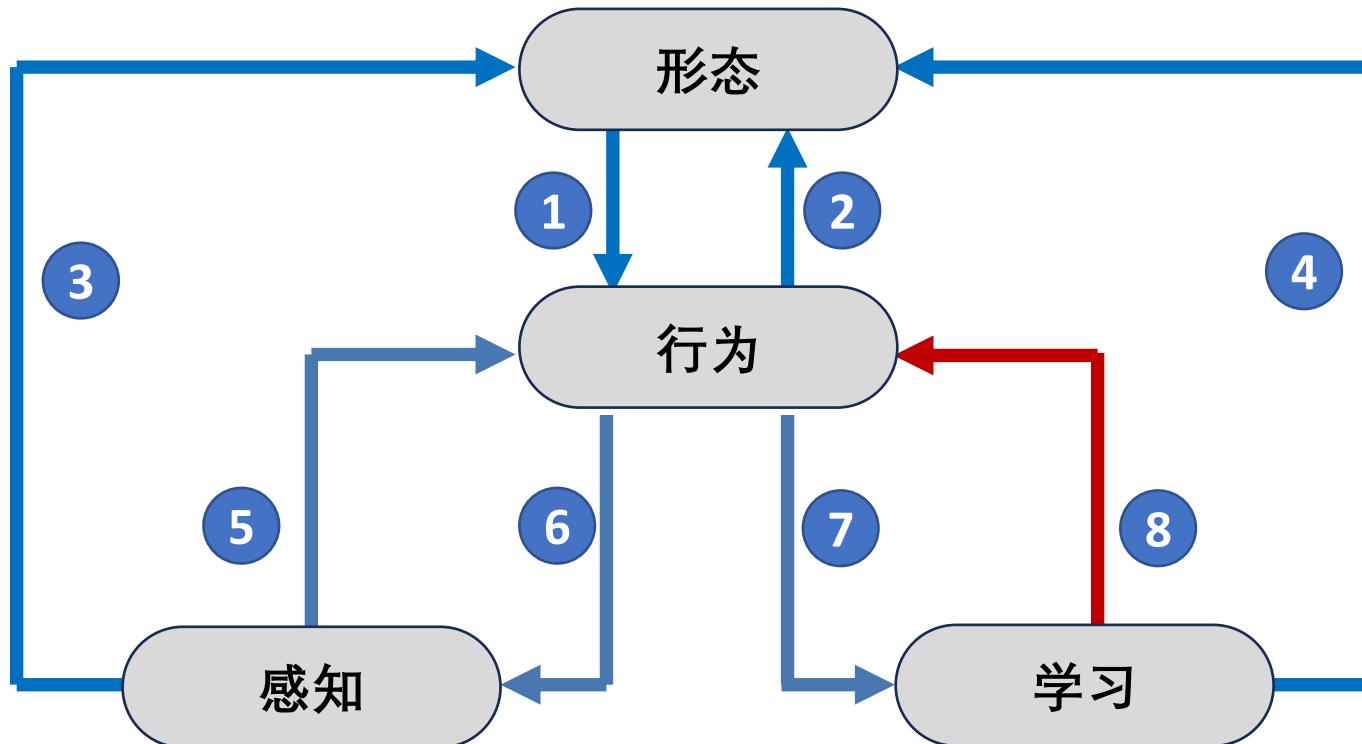
- ① 基于形态的行为生成
- ② 基于行为的形态控制
- ③ 基于感知的形态变换
- ④ 基于学习的形态优化
- ⑤ 基于感知的行为生成
- ⑥ 基于行为的主动感知
- ⑦ 基于行为的自主学习
- ⑧ 基于学习的行为优化

-
- **背景：具身智能**
 - **具身智能的体系**
 - **具身智能关键技术**
 - **探索与实践**
 - **面向具身智能的AIGC**
 - **总结**

3 关键技术



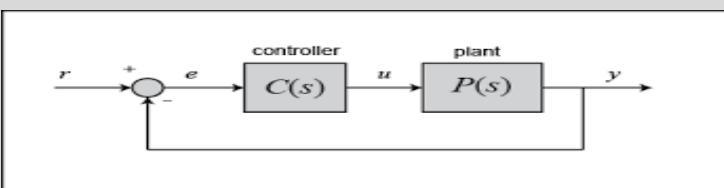
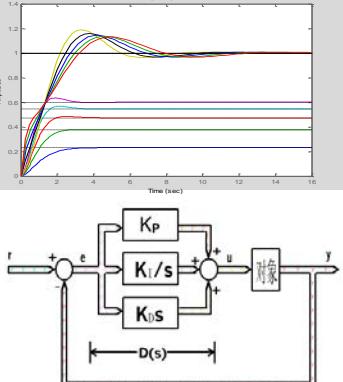
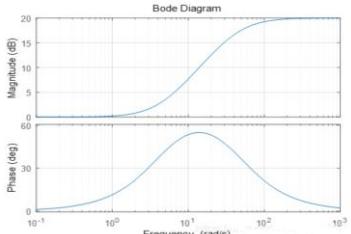
3 关键技术



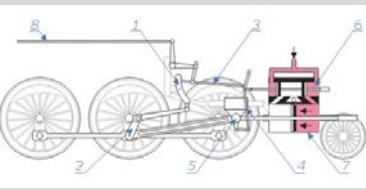
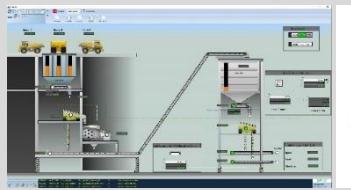
3.0 学习→行为：强化学习

经典控制：经验

$$G_c(s) = \frac{s + \frac{1}{T}}{s + \frac{1}{\alpha T}} \quad (\alpha < 1)$$



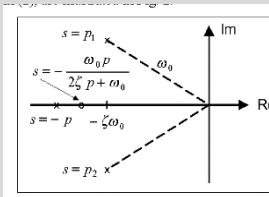
$$u_t = K_P e_t + K_I \int e_t dt + K_D \frac{de_t}{dt}$$



现代控制：模型

$$\begin{aligned} \text{state equations: } & \dot{x}_1 = f_1(x_1, x_2, \dots, x_n, u_1, \dots, u_m) \\ & \vdots \\ & \dot{x}_n = f_n(x_1, x_2, \dots, x_n, u_1, \dots, u_m) \\ \text{output equations: } & y_1 = h_1(x_1, x_2, \dots, x_n, u_1, \dots, u_m) \\ & \vdots \\ & y_p = h_p(x_1, x_2, \dots, x_n, u_1, \dots, u_m) \end{aligned}$$

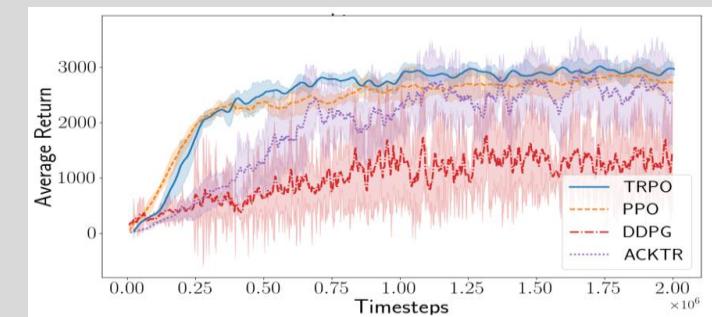
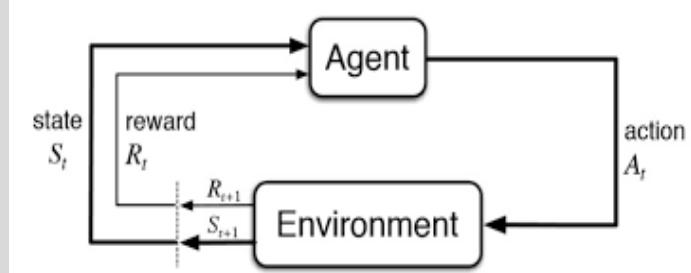
$$\begin{cases} \dot{x}_t = Ax_t + Bu_t \\ y_t = Cx_t + Du_t \end{cases}$$



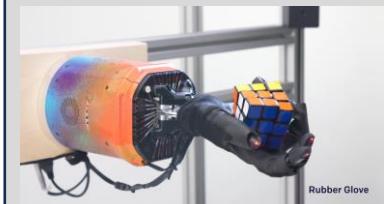
$$u_t = -K x_t$$



智能控制：学习

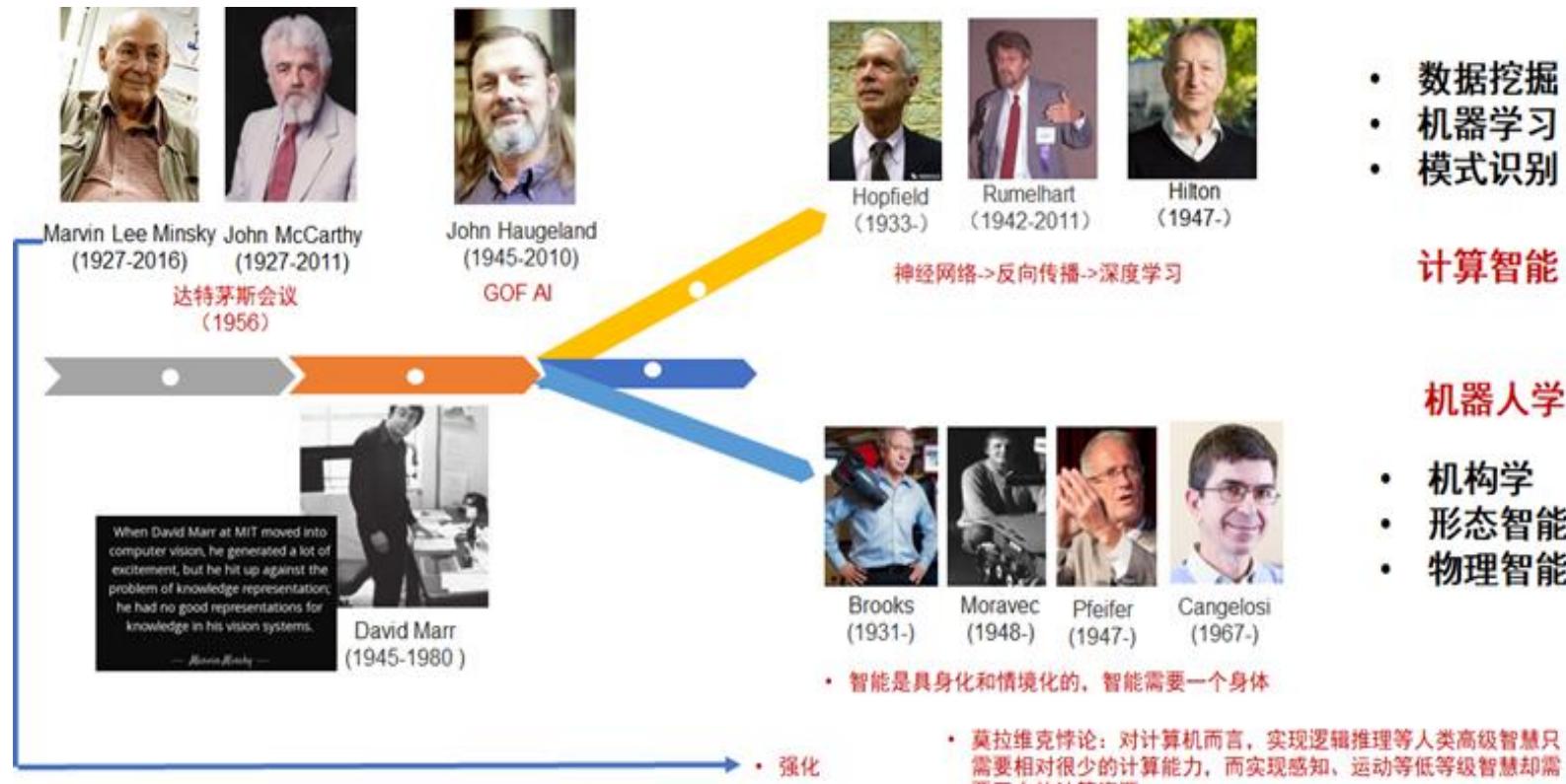


$$\pi_{\theta}(a_t | s_t)$$



3.0 学习→行为：强化学习

➤ 起源



人工神经网络
Artificial Neural Networks

Walter Pitts, 1943



深度学习
Deep Learning

Geoffrey Hinton, 2006



深度强化学习
DQN

DeepMind, 2015

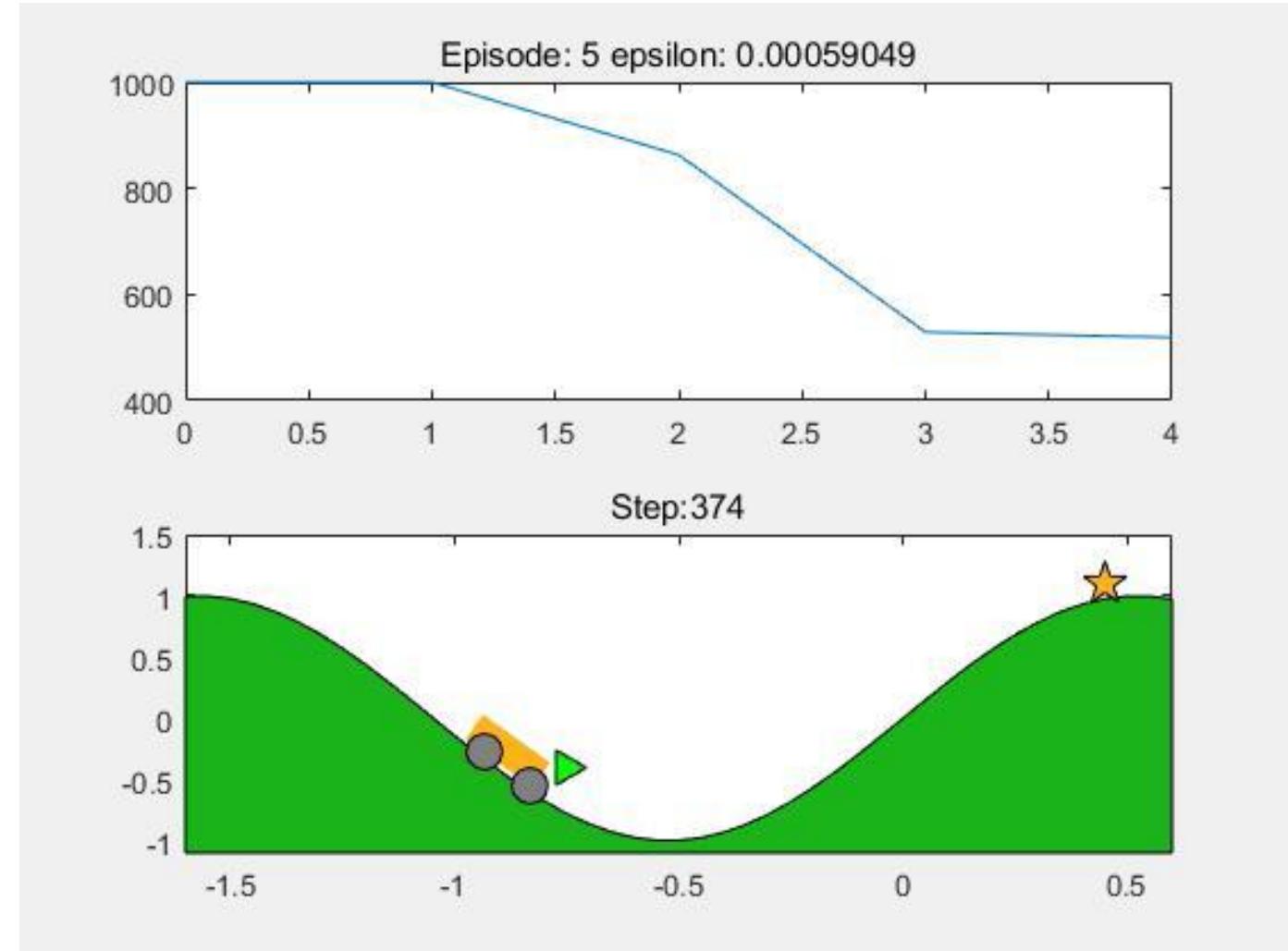
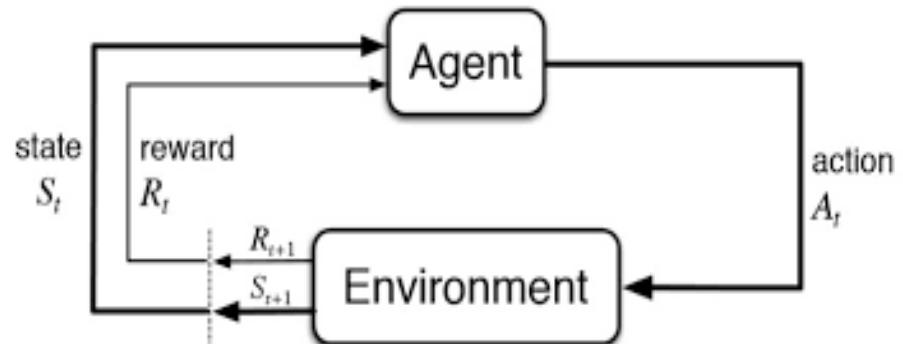
最优控制与动态规划
Markov Decision Process

Richard Bellman, 1957

时序差分与最优控制
Q-learning

Chris Watkins, 1989

3.0 学习→行为：强化学习



3.0 学习→行为：强化学习

➤ 两类基本算法

Q-学习算法

$$\max_{\theta} \mathbb{E} \left[\sum_{t=0}^H R(s_t) | \pi_{\theta} \right]$$

$$Q(s_t, a_t) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots$$

$$Q(s_t, a_t) = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}} Q^*(s, a)$$

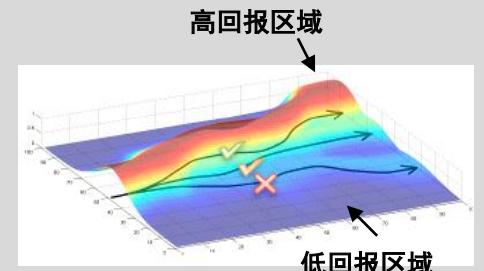


发展：DQN

策略梯度算法

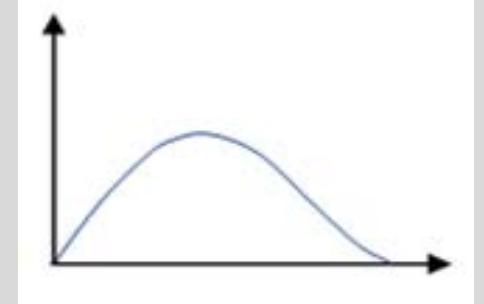
$$J(\pi) = \mathbb{E}_{\tau \sim p_{\pi}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) \right]$$

$$\nabla_{\theta} J(\pi_{\theta}) \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{H^{(n)}} R(\tau^{(n)}) \nabla_{\theta} \log \pi_{\theta}(a_t^{(n)} | s_t^{(n)})$$



$$\theta \leftarrow \theta + \eta \nabla_{\theta} J$$

$$\pi_{\theta}(a | s)$$



发展：REINFORCE

3.0 学习→行为：强化学习

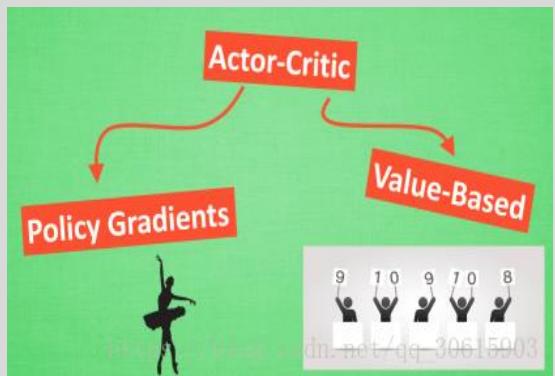
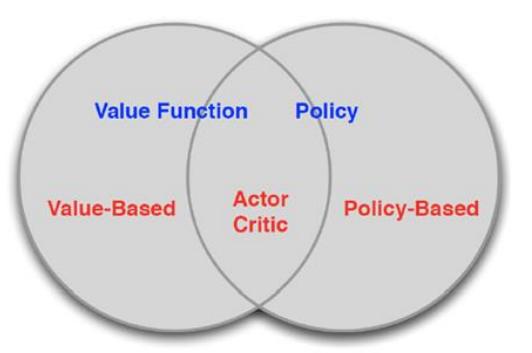
➤ 两类典型算法

Actor-Critic

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{s \sim p_0(s)} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} (Q^{\pi_{\theta}}(s, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a | s))$$

$$\hat{Q}_{\omega}(s, a) \approx Q^{\pi_{\theta}}(s, a)$$

$$\nabla_{\theta} J(\pi_{\theta}) = \hat{Q}_{\omega}(s_t, a_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$



PPO

$$A^{\pi_{\theta}}(s_t, a_t) = Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)$$

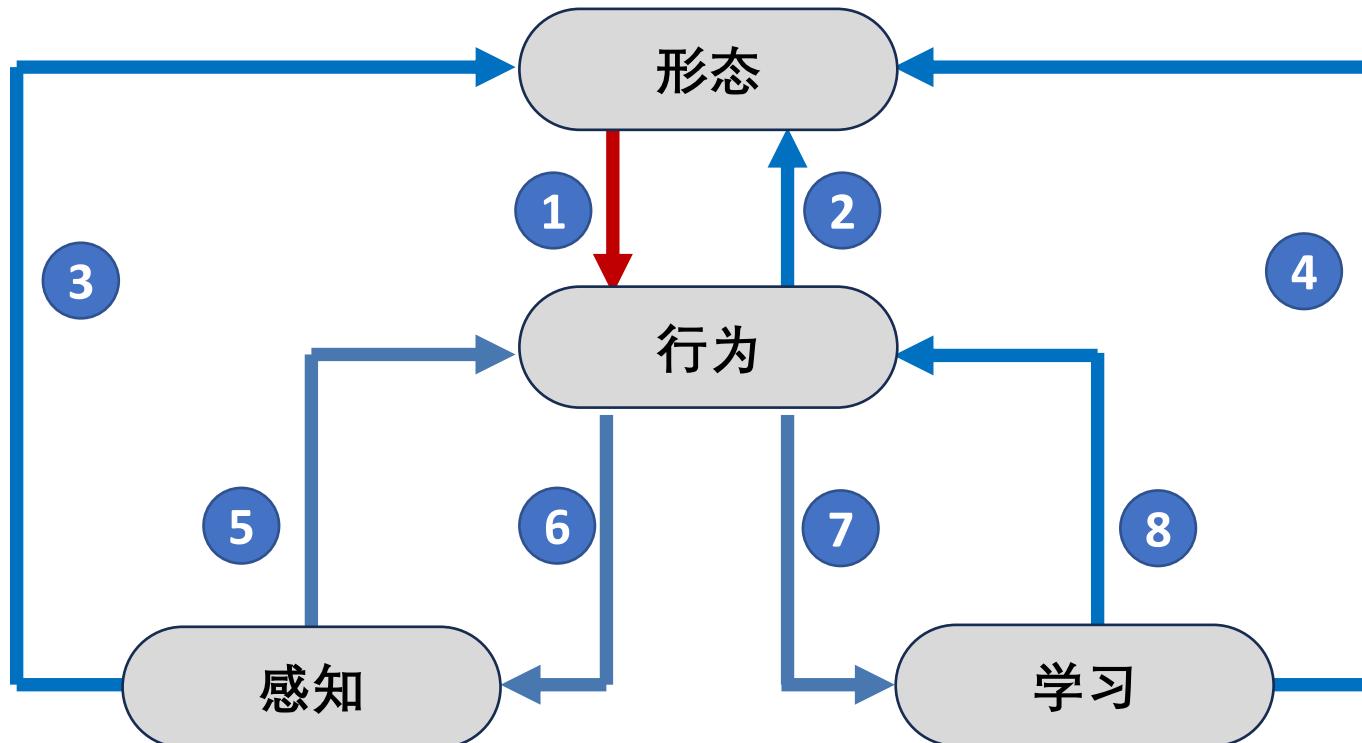
$$\mathbb{E}_{\tau \sim p_{\pi_{\theta}}} \left(\sum_{t=0}^{\infty} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right)$$

$$L_{\pi_{\theta_{\text{old}}}}(\pi_{\theta}) = \mathbb{E}_{\tau \sim p_{\pi_{\theta_{\text{old}}}}} \left(\sum_{t=0}^{\infty} \gamma^t \rho_t(\theta) A^{\pi_{\theta_{\text{old}}}}(s_t, a_t) \right)$$

$$L_{\pi_{\theta_{\text{old}}}}(\pi_{\theta}) = \mathbb{E}_{\tau \sim p_{\pi_{\theta_{\text{old}}}}} \left(\sum_{t=0}^{\infty} \gamma^t \min \left\{ (\rho_t(\theta) A^{\pi_{\theta_{\text{old}}}}(s_t, a_t), \text{clip}(\rho_t(\theta), 1-\varepsilon, 1+\varepsilon) A^{\pi_{\theta_{\text{old}}}}(s_t, a_t)) \right\} \right)$$

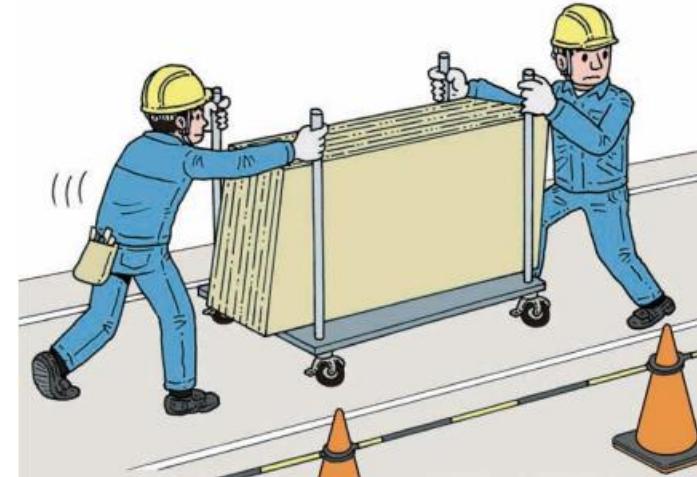
$$\text{clip}(\rho_t(\theta), 1-\varepsilon, 1+\varepsilon) = \begin{cases} \rho_t(\theta) & 1-\varepsilon \leq \rho_t(\theta) \leq 1+\varepsilon \\ 1-\varepsilon & \rho_t(\theta) < 1-\varepsilon \\ 1+\varepsilon & \rho_t(\theta) > 1+\varepsilon \end{cases}$$

3.1 形态→行为



3.1 形态→行为

➤ 例子



3.1 形态→行为

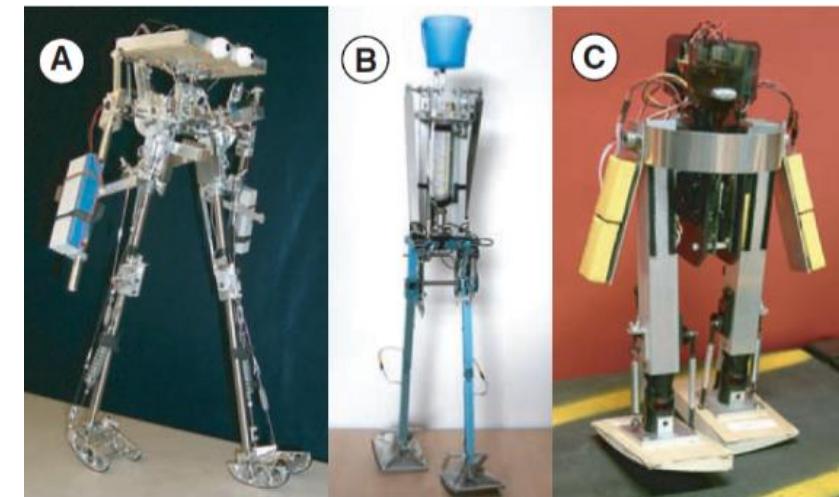
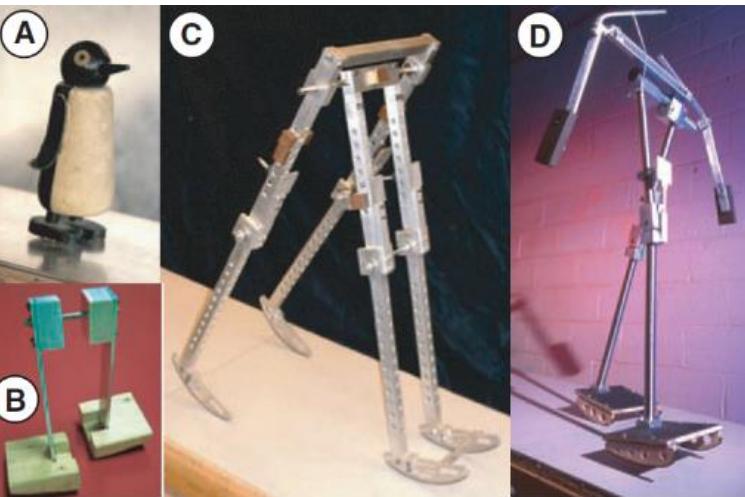
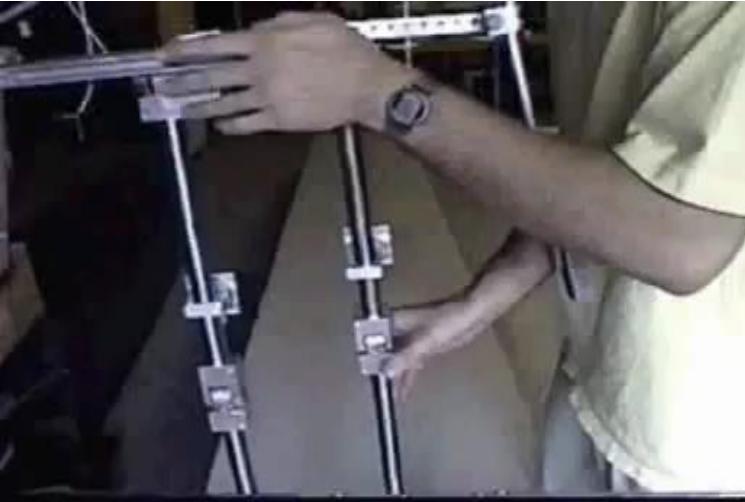
➤ 边界

由于形态计算方面的研究与仿生机器人的研究关系非常密切，二者之间的关系会引起混淆。事实上，形态计算更关心的是利用形态来产生行为，而并非从形态上简单地模仿某些生物。很多仿生型机器人，只是形体上模仿了动物的身体，这种模仿可能能获得一些自由度上的突破，例如腿式机器人相比轮式能爬楼，但在行为控制方面，并没有充分利用形态自身的优势，而且仍需要设计复杂的控制器。本田的Asimo机器人，尽管已经充分接近人类的外形，但其每一个关节都需要用特定的算法来控制。这些情形都不属于形态计算之列。

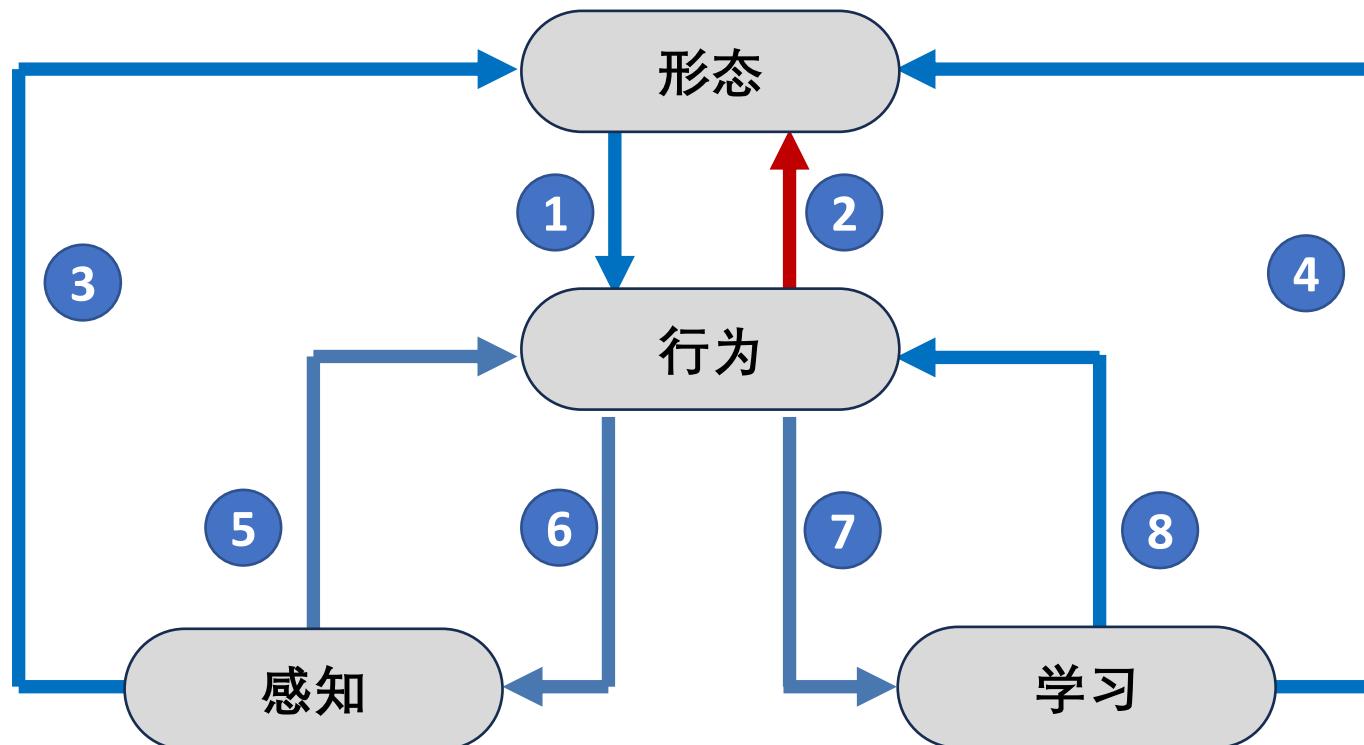


3.1 形态→行为

➤ 被动行走机器人 (Passive Walking Robot)



3.2 行为→形态：形态控制



3. 2 行为→形态：形态控制

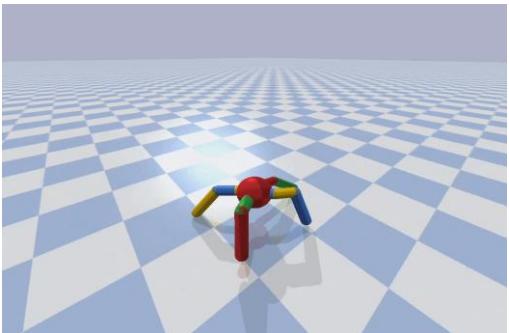
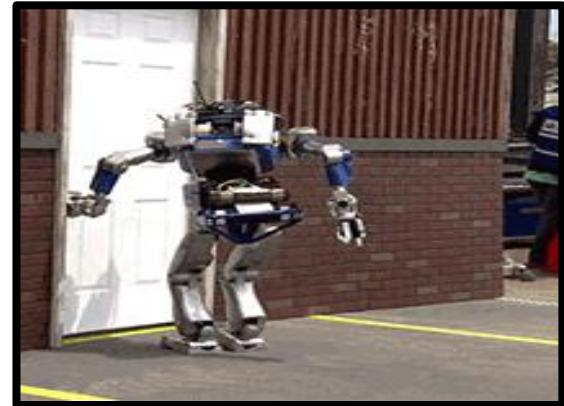
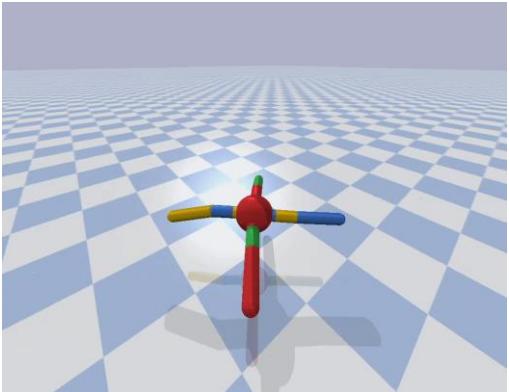
➤ 形态约束



3. 2 行为→形态：形态控制

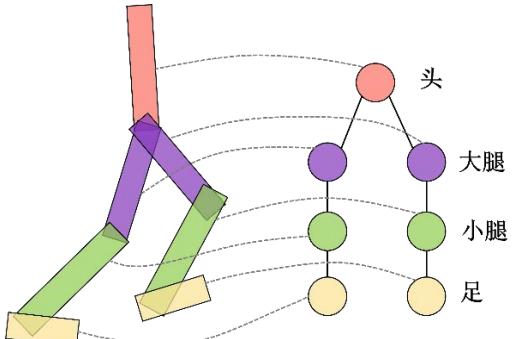
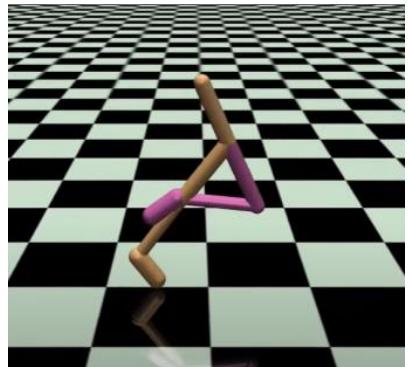
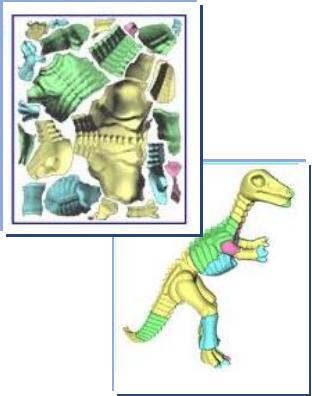
➤ 强化学习的维数灾

- 机器人自由度多，优化学习困难
- 机器人形态**复杂**，难以利用
- 在不同形态之间的**迁移**



3. 2 行为→形态：形态控制

➤ 引入结构性约束



$$\pi_{\theta}(a_t | s_t)$$

图结构表示

$$D = (V, E, Y)$$

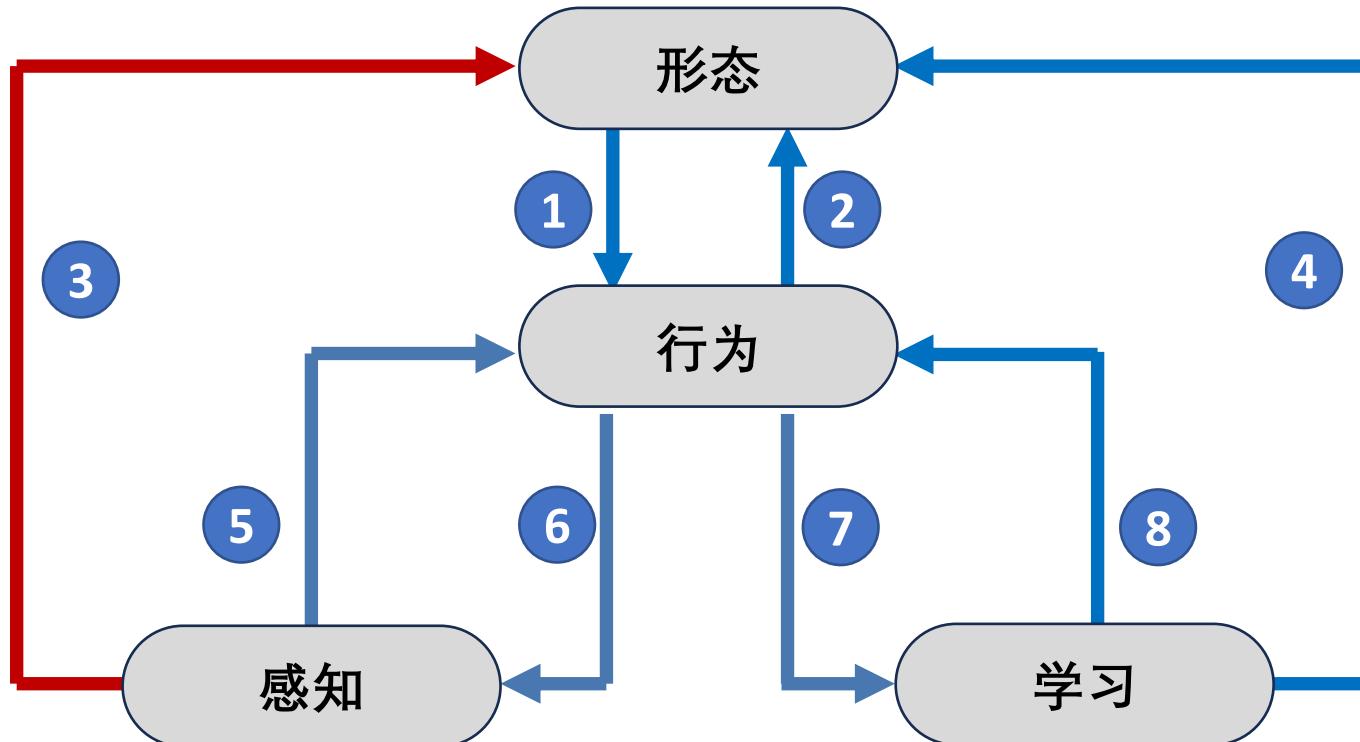
$N_{out}(v)$ 所有以 v 作为头节点的有向边所指向的尾节点的集合。这一集合反映的是节点 v 的“输出节点”。

$N_{in}(v)$ 所有以 v 作为尾节点的边的头结点的集合。这一集合反映的是节点 v 的“输入节点”。

$$Y_V(v_i) \in \{1, 2, \dots, |Y_V|\}$$

$$Y_E(\{v_i, v_j\}) \in \{1, 2, \dots, |Y_E|\}$$

3. 3 感知→形态：形态变换



3. 3 感知→形态：形态变换



3. 3 感知→形态：形态变换

➤ 地形适应的形态变换

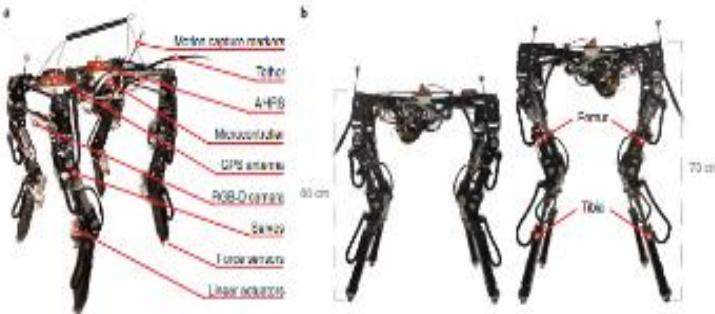
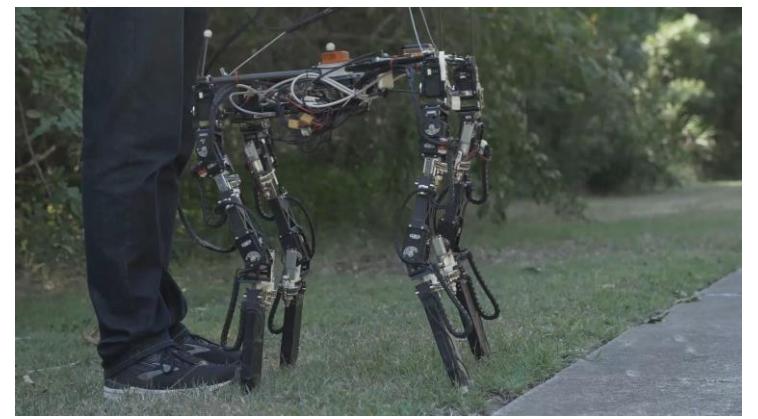
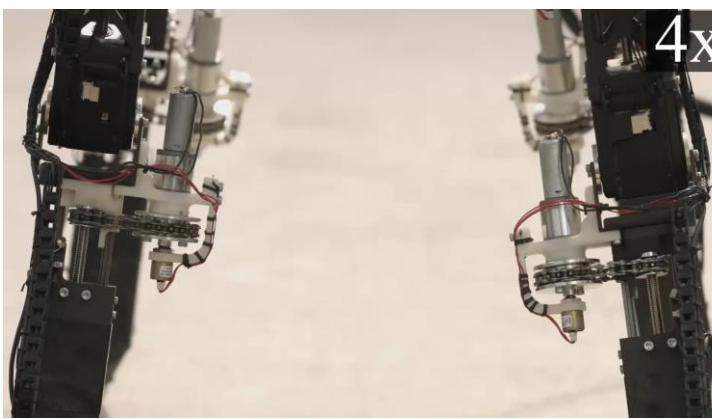


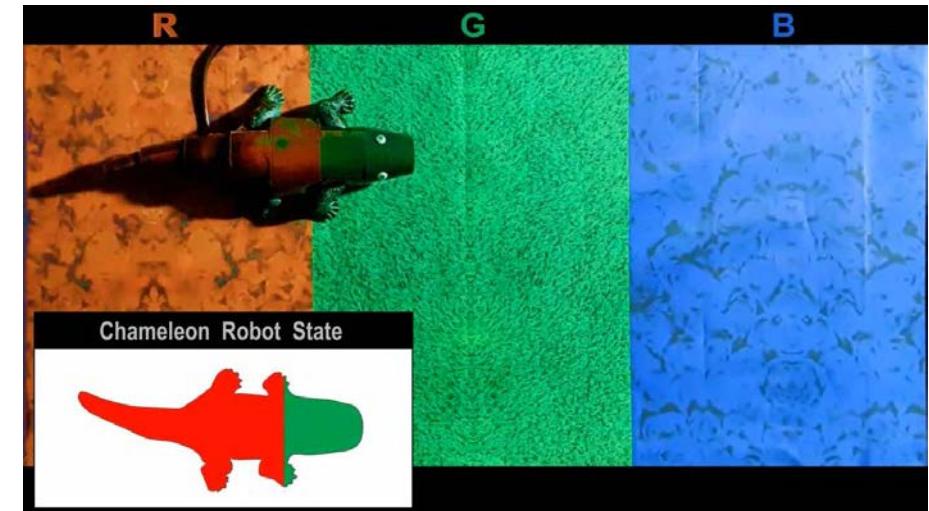
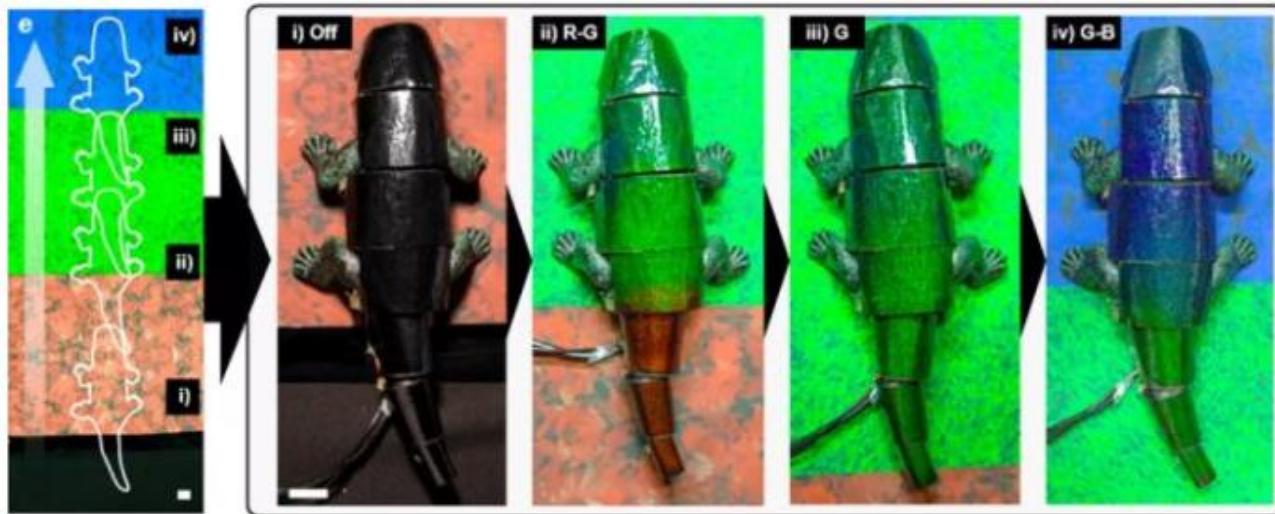
Fig. 1 | The morphologically adaptive robot used in this study. a, An overview of the main components of the robot. b, The robot with the shortest (left) and longest (right) leg configuration. AHS, attitude and heading reference system; GPS, global positioning system; RGB-D, red, green, blue and depth.



3. 3 感知→形态：形态变换

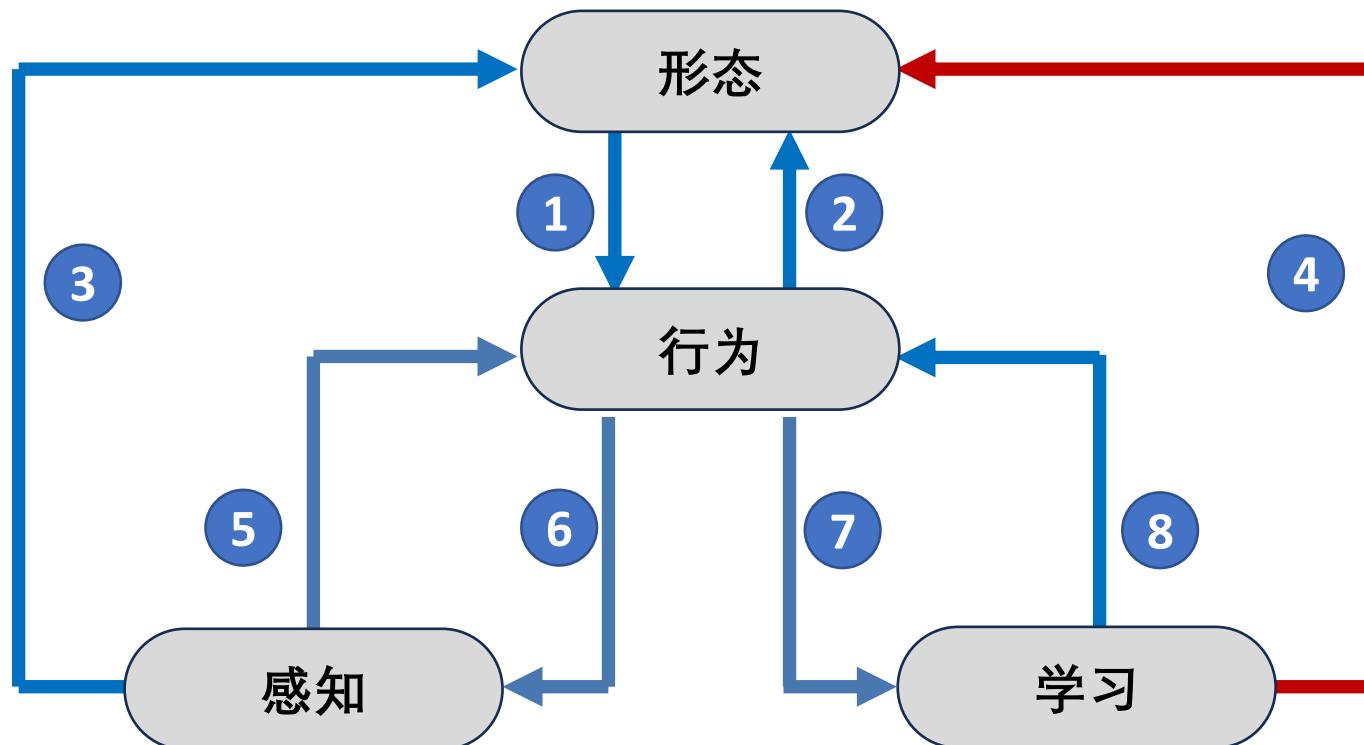
➤ 仿变色龙软体机器人

将人造变色龙皮肤应用于软体机器人上，并结合颜色传感器和反馈控制系统，从而使得这种设备级的自适应人造伪装技术能够检测到背景环境的颜色，并令变色龙软体机器人上实现实时背景颜色匹配。



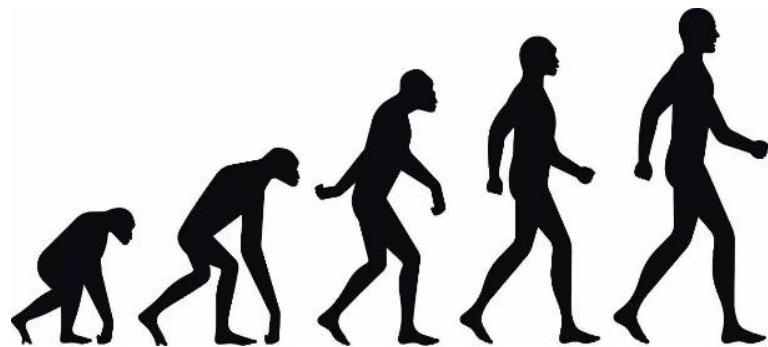
- Kim H, Choi J, Kim K K, et al. Biomimetic chameleon soft robot with artificial crypsis and disruptive coloration skin[J]. Nature communications, 2021, 12(1): 1-11.

3.4 学习→形态：形态生成



3.4 学习→形态：形态生成

➤ 脑-体协同进化

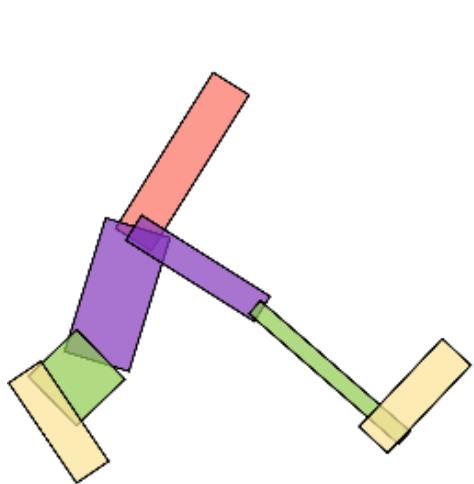


$$(D^*, \pi^*) = \arg \max_{D, \pi} J(D, \pi)$$

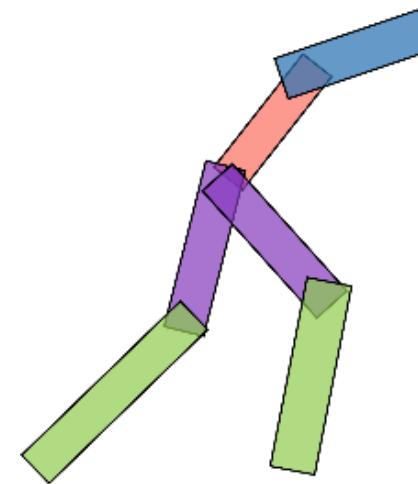
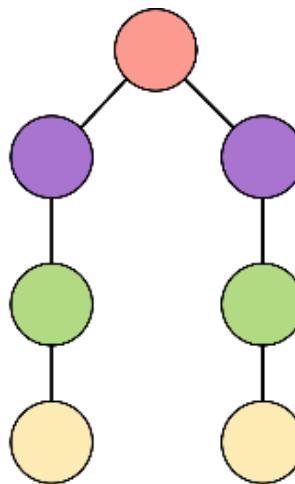
$$\pi^* = \arg \max_{\pi} J(\bar{D}, \pi)$$

3.4 学习→形态：形态生成

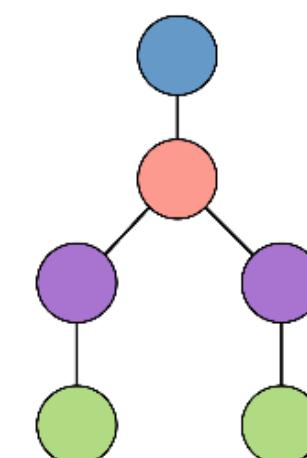
➤ 脑-体协同进化



形态参数

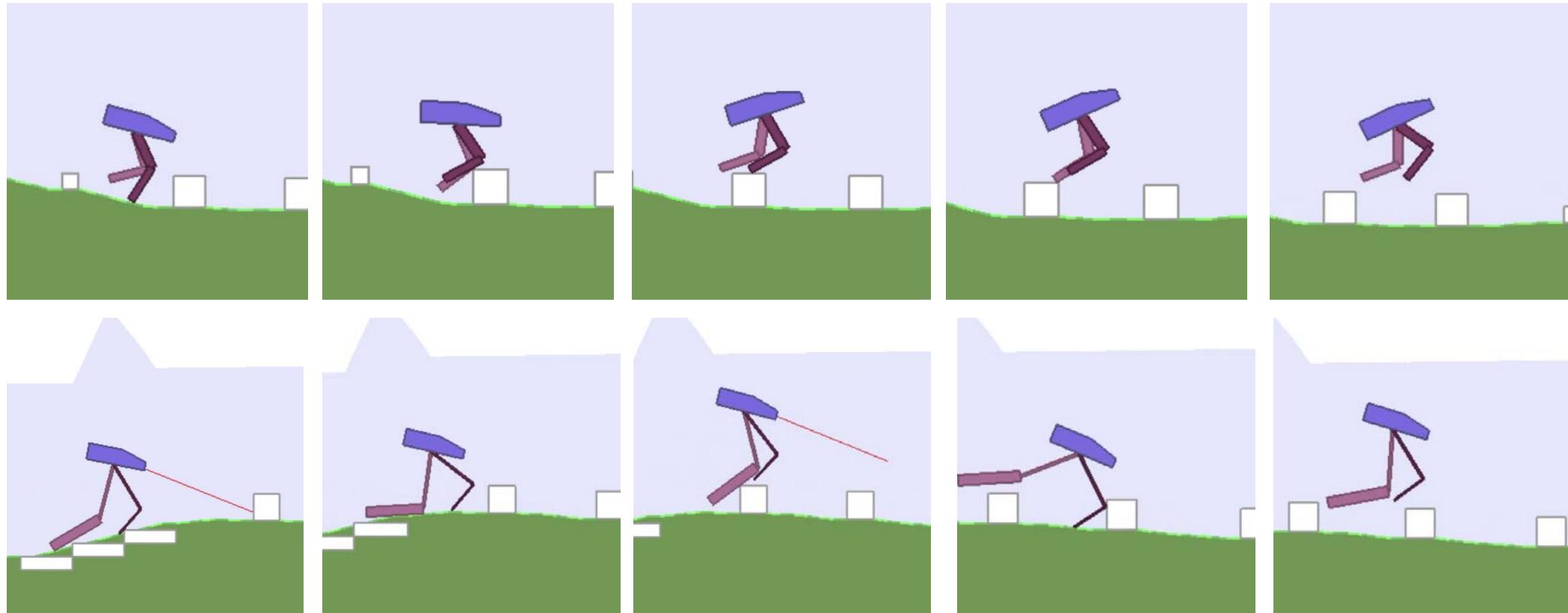


形态结构

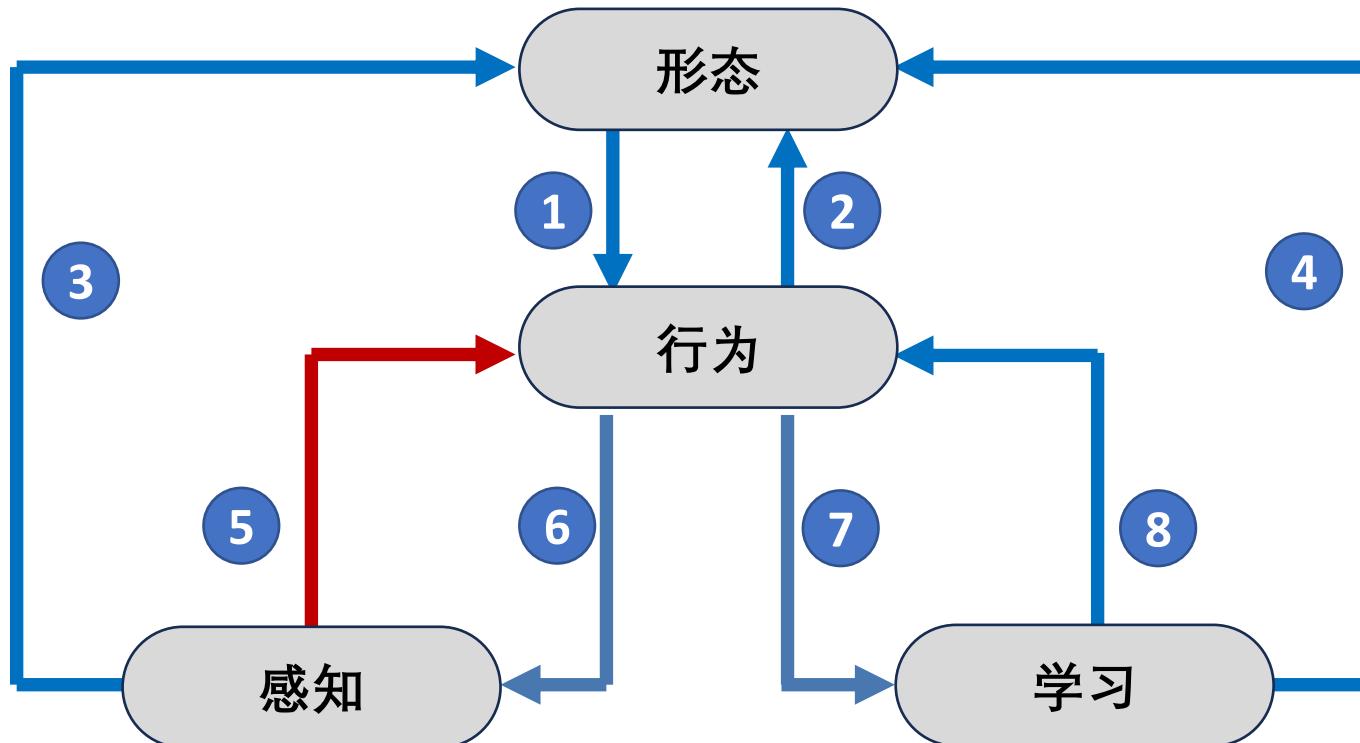


3.4 学习→形态：形态生成

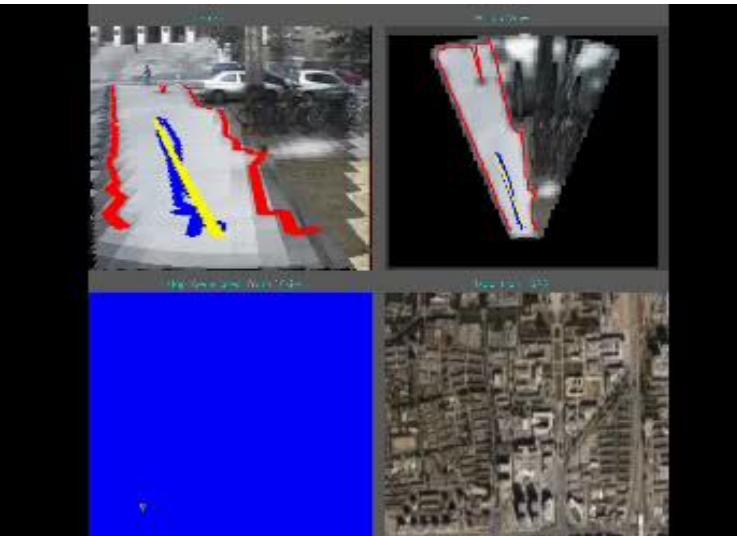
➤ 形态参数的优化



3.5 感知→行为：视觉导航

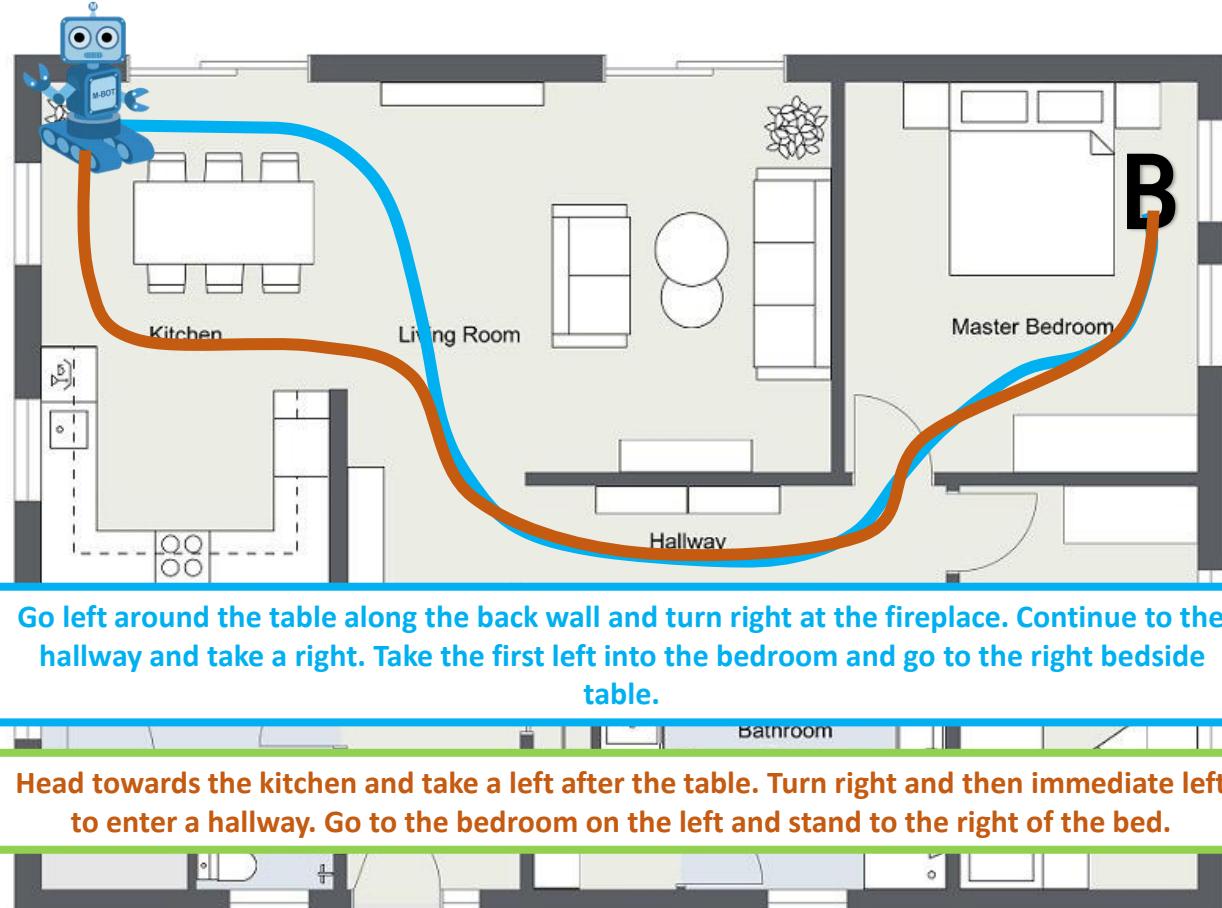


3. 5 感知→行为：视觉导航



3. 5 感知→行为：视觉导航

➤ 视觉语言导航



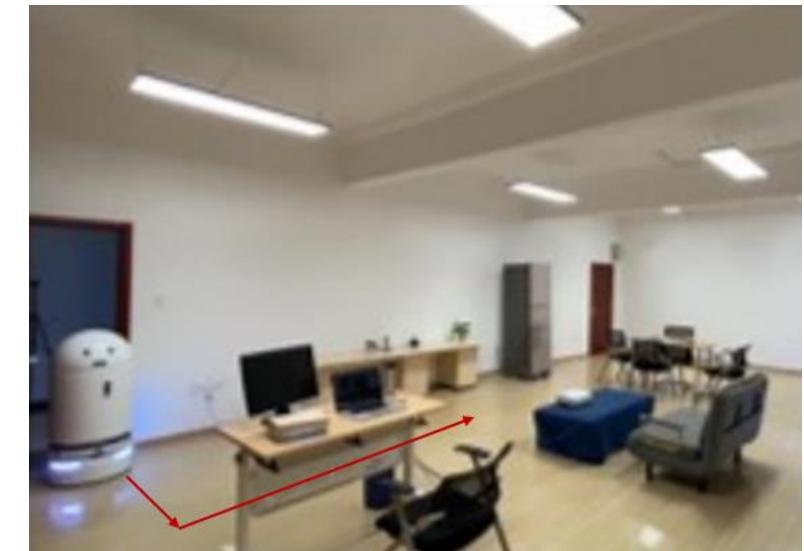
- Vision language navigation (VLN) 让智能体跟着自然语言指令进行导航，这个任务需要同时理解自然语言指令与视角中可以看见的图像信息，然后在环境中对自身所处状态做出对应的动作，最终达到目标位置。



Walk beside the outside doors and behind the chairs across the room.

3. 5 感知→行为：视觉导航

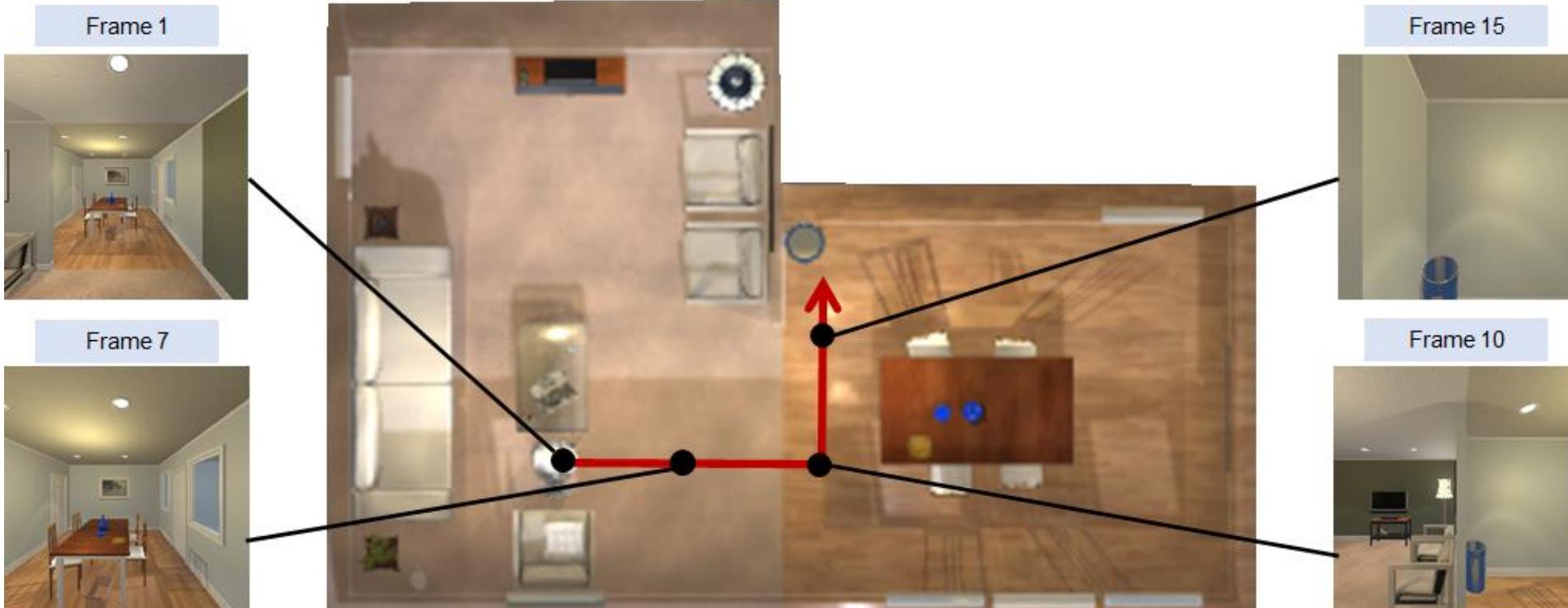
➤ 视觉语言导航



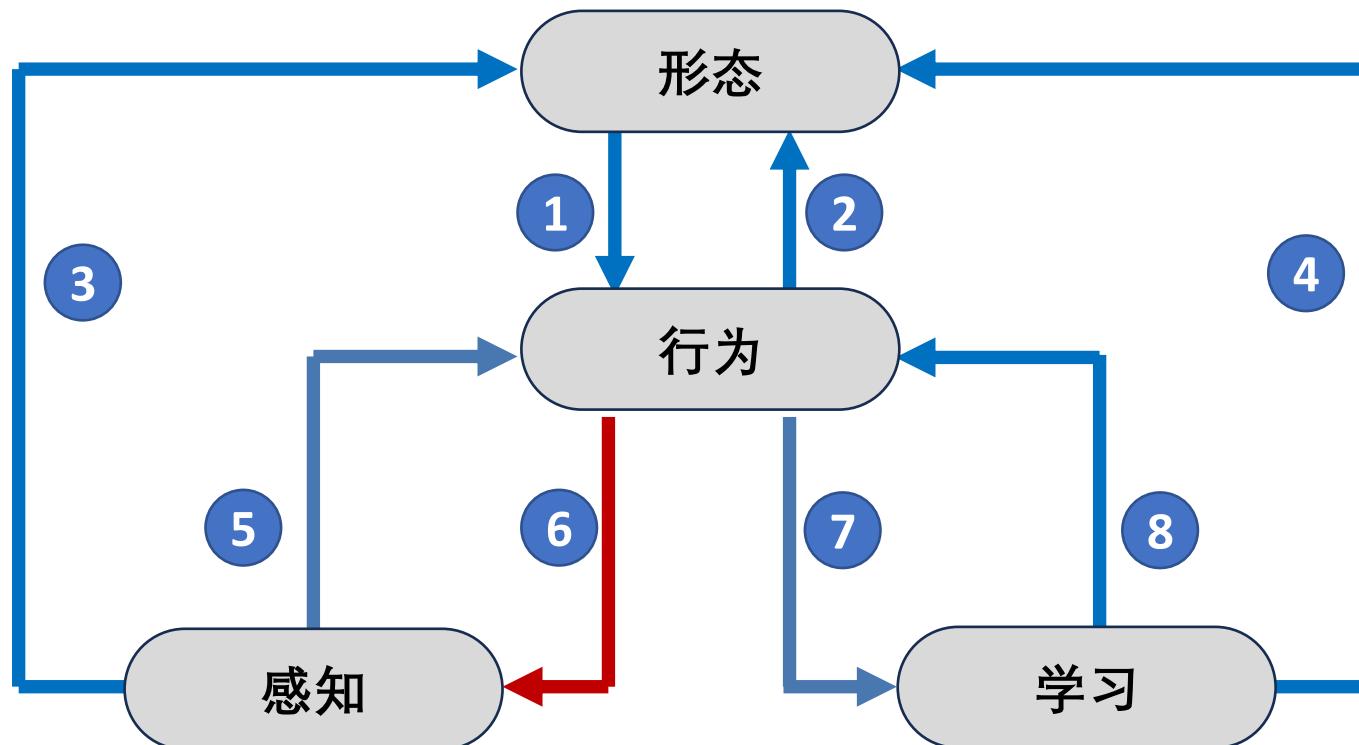
3. 5 感知→行为：视觉导航

➤ 问题描述

语言指令：Walk a few steps straight, and then turn to left
and move several steps left to the garbage can on right.



3.6 行为→感知：主动感知



3.6 行为→感知：主动感知

➤ 生活中的主动感知



3.6 行为→感知：主动感知

➤ 主动感知的应用



3.6 行为→感知：主动感知

➤ 主动感知的优点

- A motivation for examining active vision is the fact that passive vision has been shown to be very problematic.
 - Almost every basic problem in passive machine perception is very difficult, it is ill-posed in the sense of Hadamard.
 - Problems that are **ill-posed**, **nonlinear** or **unstable** for a passive observer become **well-posed**, **linear** or **stable** for an active observer.



Active Perception

RUZENA BAJCSY, MEMBER, IEEE
Invited Paper

Active Perception (Active Vision specifically) is defined as a study of Modeling and Control strategies for perception. By modeling we mean analysis of sensors, processing modules and their interactions, data fusion and decision making. By control we mean the extent of application in space and time. The local models represent sensory information and dynamics such as motion, dimensions, focal length, etc., required to band-pass filter the scene. The global models on the other hand characterize the overall performance and the quality of the perception. The modeling and control strategies are formulated as a search of such sequence of steps that will yield the best solution. Active perception is the process of obtaining most information. Examples are shown as the evidence proof of the proposed theory on obtaining range from focus and stereopsis over 2-D segmentation of an image and 3-D shape parametrization.

II. What is Active Sensing?

In the robotics and computer vision literature, the term "generally refers to a sensor that transmits (generally electromagnetic radiation, e.g., radar, sonar, ultrasound, microwaves and collimated light) into the environment and receives measures of the reflected signals. We also use the use of active sensing not only as a necessary condition on active sensing, and that sensing can be performed with passive sensors (that only receive, and do not emit/transmit), but also as a general term to mean active to denote a time-of-flight sensor, but to denote a passive sensor employed in an active fashion, purposefully changing the sensor's state parameters according to simple rules.

Despite the problem of Active Sensing can be stated as a problem of controlling strategies applied to the data acquisition, it is not necessarily the same as the problem of data interpretation and the goal or the task of the process. The question may be asked: "Is Active Sensing only an application of Control Theory? Our answer is: "No, at least not in its entirety". Here is why:

- 1) The feedback is performed not only on sensory data but on complex processed sensory data, i.e., various extracted features, including relational features.
- 2) The feedback is dependent on *a priori* knowledge—models that are a mixture of numeric/parametric and symbolic information.

But one can say that Active Sensing is an application of intelligent control theory that includes reading, decoding, learning, etc. control. This argument has been eloquently stated by Teunissen (1993): "Because of the inherent limitation of a single image, the acquisition of information should be treated as an integral part of the perceptual process. This is in contrast to the traditional view of the separation of image inadequacy rather than the secondary problems caused by it." Although he uses the term "perception" rather than active sensing the message is the same.

The implications of the active sensing approach are the following:

- 1) The necessity of models of sensors. This is to say, first, the model of the physics of sensors as well as the noise of the sensors. Second, the model of the signal processing and

Ruzena Bajcsy
UC Berkeley

John Aloimonos
U. Maryland

Problem

Problem	Passive Observer	Active Observer
Shape for shading	Ill-posed problem. Needs to be regularized. Even then, unique solution is not guaranteed because of nonlinearity.	Well-posed problem. Unique solution. Linear equation used. Stability.
Shape from contour	Ill-posed problem. Has not been regularized up to now in the Tikhonov sense. Solvable under restrictive assumptions.	Well-posed problem. Unique solution for both monocular or binocular observer.
Shape from texture	Ill-posed problem. Needs some assumption about the texture.	Well-posed problem. No assumption required.
Structure from motion	Well posed but unstable. Nonlinear constraints.	Well posed and stable. Quadratic constraints, simple solution methods, stability.
Optic flow (area based)	Ill posed. Needs to be regularized. The introduced smoothness might produce erroneous results.	Well-posed problem. Unique solution. Might be unstable.

Received: 9 September 2016 / **Accepted:** 11 January 2017 / **Published online:** 15 February 2017
© The Author(s) 2017. This article is published with open access at Springerlink.com

Abstract Despite the recent successes in robotics, artificial intelligence and computer vision, a complete artificial agent successfully performing active perception is still missing. A number of tasks and methods for active perception have been developed in the past, their broader utility perhaps impeded by insufficient computational power or costly hardware. The history of these ideas, perhaps selective due to our perspectives, is instructive. In particular, good arguments can be made that this, however, is the case, although there is no doubt that group of stars looks like something I know—let me see if I can find other stars to complete the pattern. This is the essence of active perception—to set up a goal based on some knowledge of the world and to put in motion the actions that stay achieve it.

Throughout the years, the topic of perception, and particularly vision, has been a great source of wonder and sometimes even awe. This is true for us all. This paper will be reviewed here and the interested reader should see Pastore (1971) and Wade (2000), among others. These papers present the history of the development of the field on the problem of perception, as well as an extensive review on perception. Those interested in a biological perspective on active perception should see Fadale and Gulyás (2003) and Poggio (2003). The latter is a brief history of the development of the theory of the birth of the computational active perception paradigm.

All authors contributed equally to all aspects of this manuscript.

John K. Tsotsos
Johns Hopkins University

John K. Tsotsos
Johns Hopkins University

1 Department of Electrical Engineering and Computer Science, University of California, Berkeley, CA, USA

2 Department of Computer Science, University of Maryland, College Park, MD, USA

Department of Electrical Engineering and Computer Science, York University, Toronto, ON, Canada

**Anton Robot (2018) 42:177–196
https://doi.org/10.1007/s10514-017-9615-3**

CrossMark

Revisiting active perception

Ruzena Bajcsy¹ · Yiannis Aloimonos² · John K. Tsotsos³ 

Keywords Sensing · Perception · Attention · Control · **I Introduction**

Sons evening, long ago, our ancestors looked up in the night sky, just as they had done for thousands of years. But this time, it was different. For the first time, human eyes noticed that one of several papers published in *Autonomous Robots* containing the special issue on Active Perception.

This is one of several papers published in *Autonomous Robots* containing the special issue on Active Perception.

All authors contributed equally to all aspects of this manuscript.

John K. Tsotsos
Johns Hopkins University

John K. Tsotsos
Johns Hopkins University

1 Department of Electrical Engineering and Computer Science, University of California, Berkeley, CA, USA

2 Department of Computer Science, University of Maryland, College Park, MD, USA

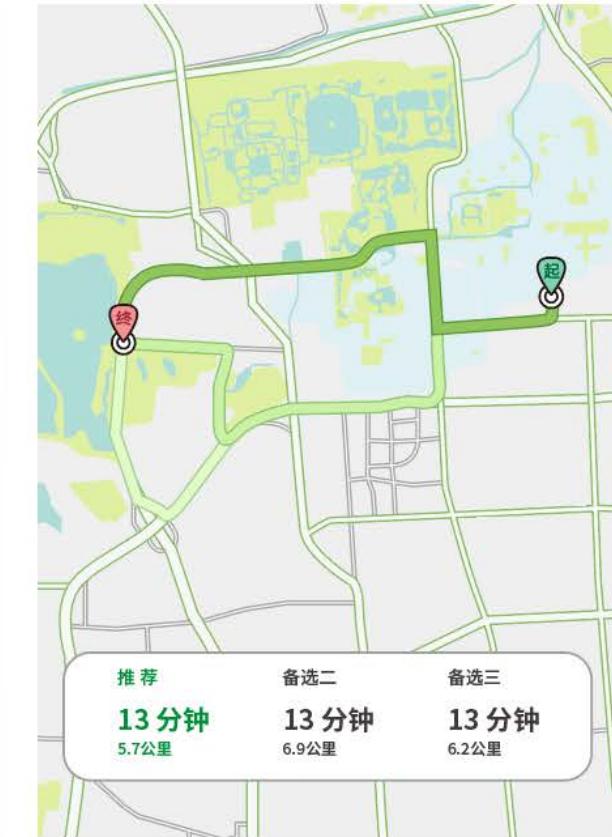
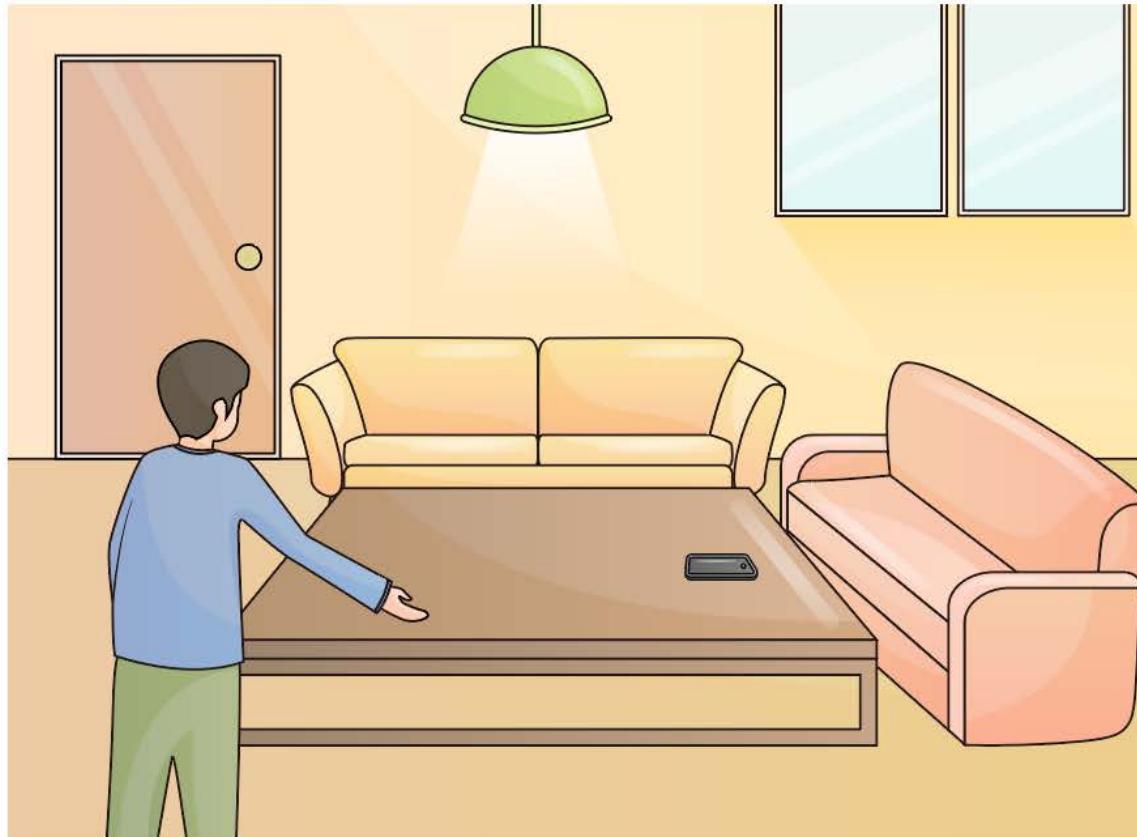
Department of Electrical Engineering and Computer Science, York University, Toronto, ON, Canada

Springer

- Aloimonos, John, Isaac Weiss, and Amit Bandyopadhyay. "Active vision." *International journal of computer vision* 1, no. 4 (1988): 333-356.

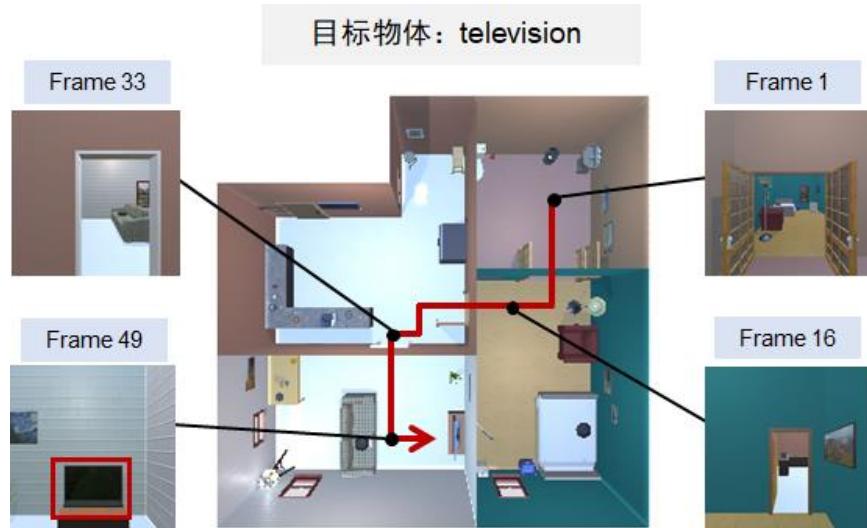
3.6 行为→感知：主动感知

➤ 视觉语义导航（Visual Semantic Navigation, VSN）



3.6 行为→感知：主动感知

➤ VSN问题描述



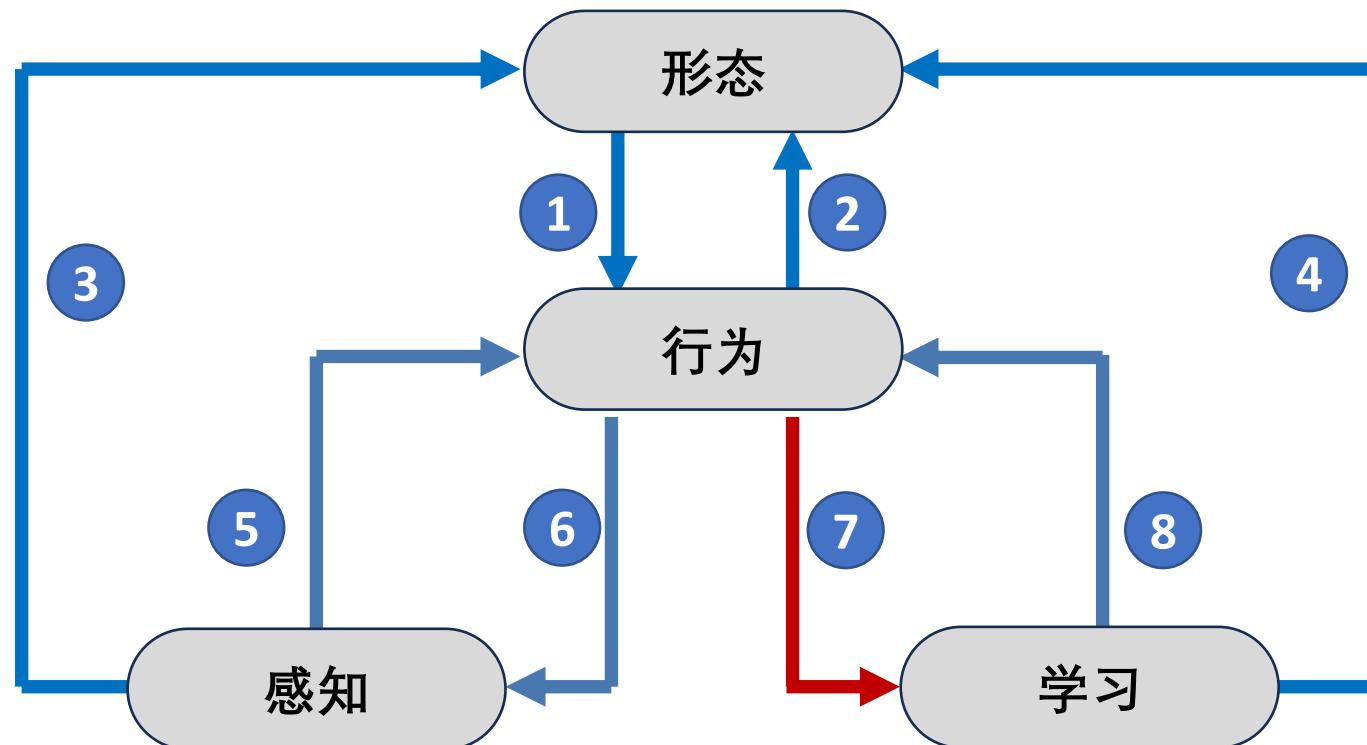
$$s_t = \varphi(I_t)$$

$$o = \phi(O)$$

$$a_t = \text{Nav}(s_t, a_{t-1}, o)$$

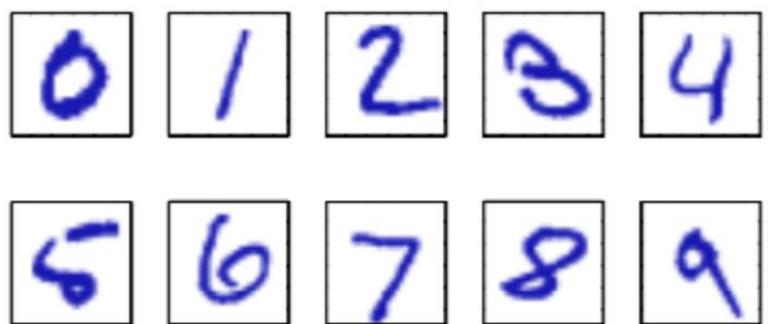
$$s_{t+1} = \text{Env}(s_t, a_t)$$

3.7 行为→学习：具身学习



3.7 行为→学习：具身学习

➤ 人是如何学习的？



3.7 行为→学习：具身学习

➤ 问题描述

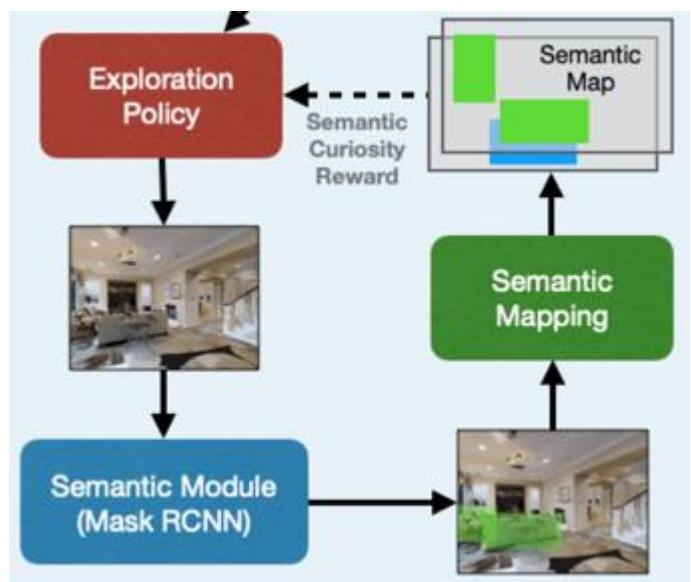


3.7 行为→学习：具身学习

➤ 问题描述

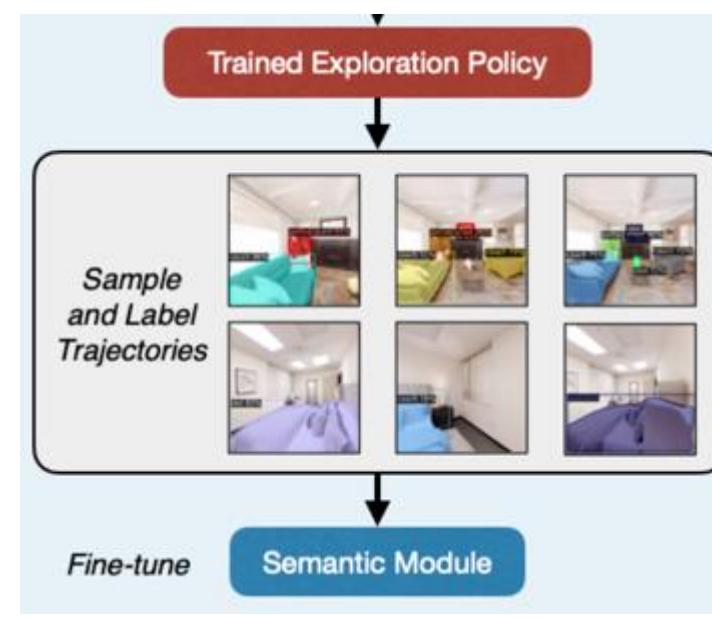
学习探索路径

学习如何探索环境才能获得更好的样本



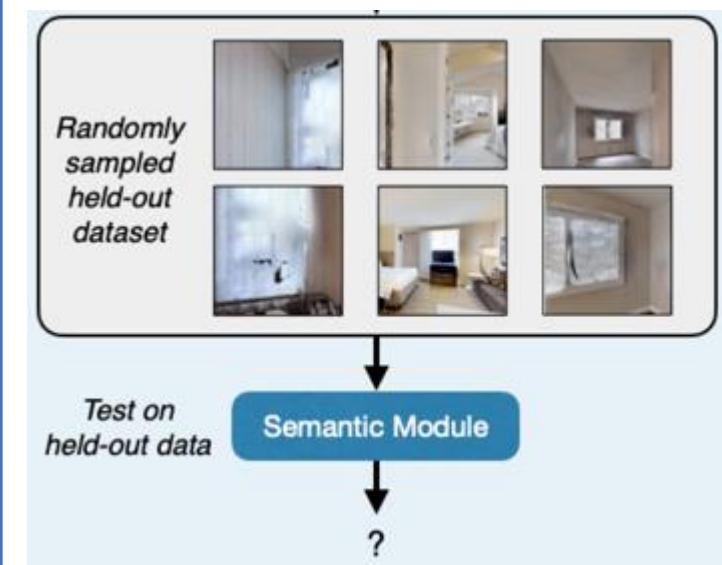
采集训练样本

利用探索策略获取样本与标签，提升感知器性能

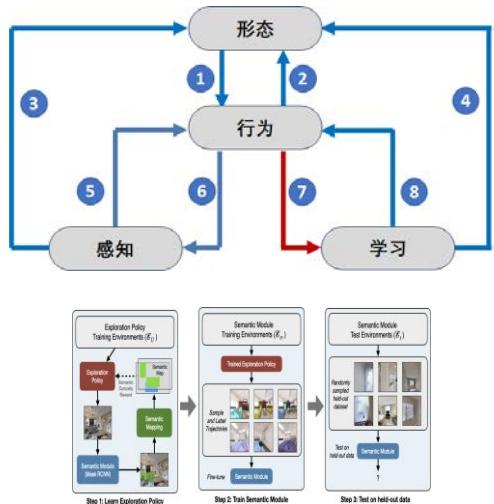
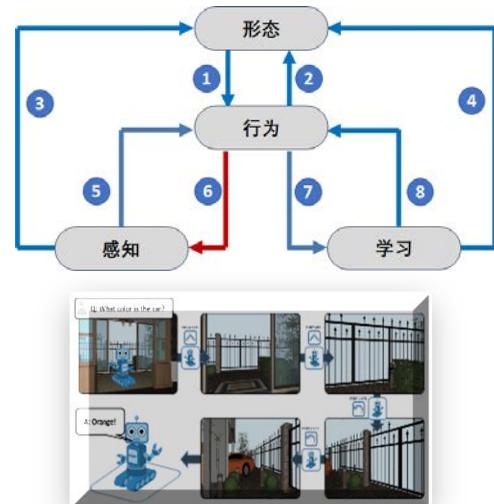
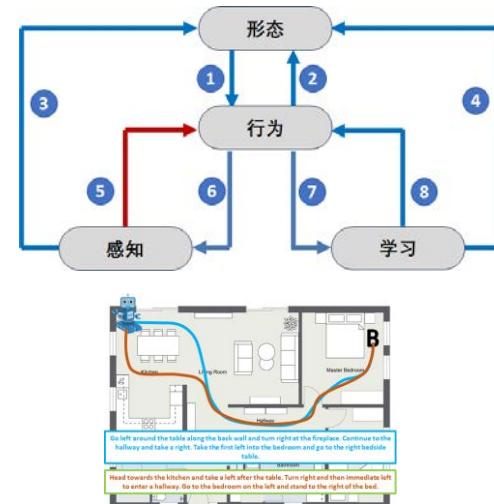
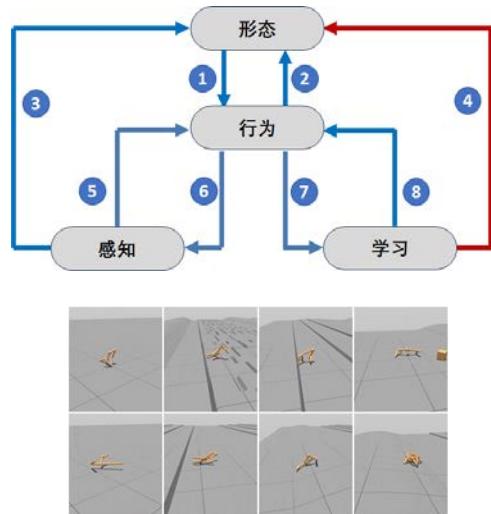
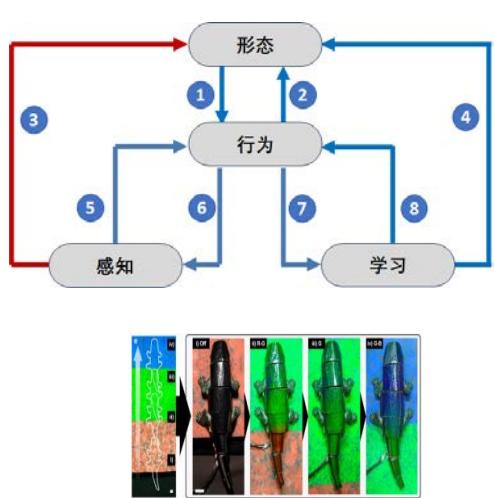
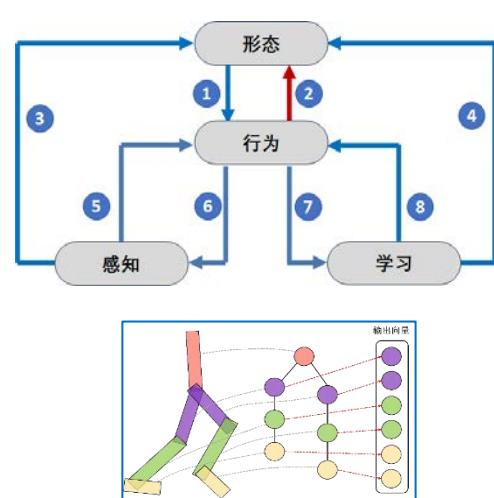
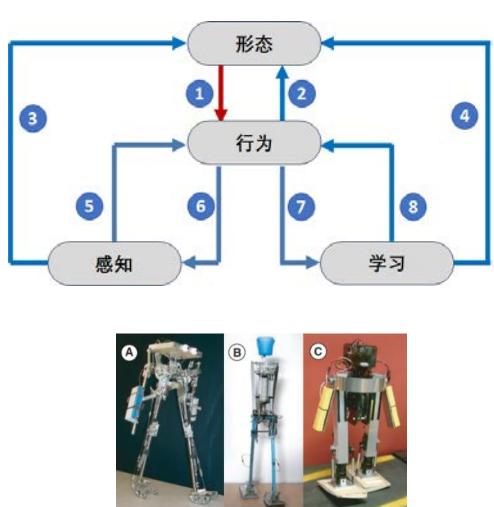
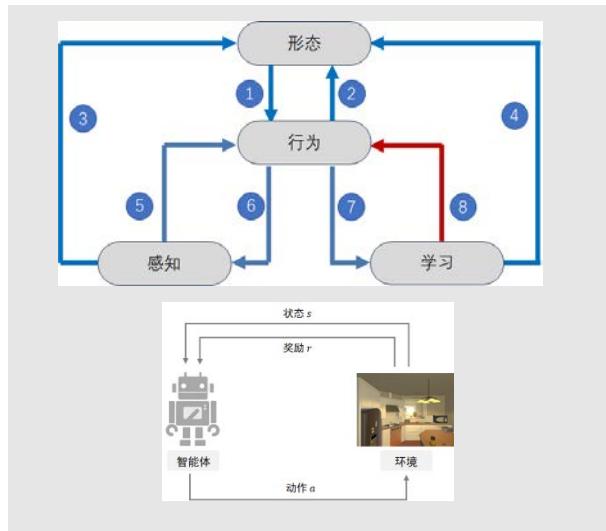


部署应用

在新场景中对改进的感知器进行测试



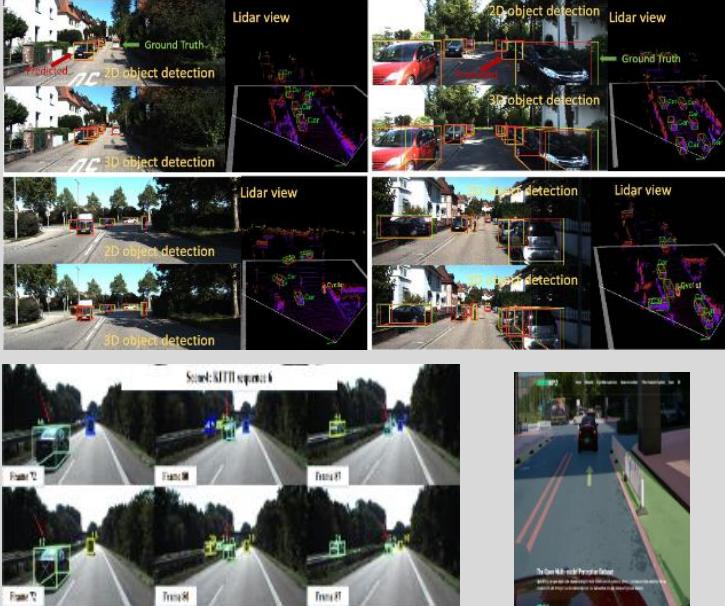
3.8 总结



-
- **背景：具身智能**
 - **具身智能的体系**
 - **具身智能关键技术**
 - **探索与实践**
 - **面向具身智能的AIGC**
 - **总结**

脑、身体与环境的交互中涌现智能

物理环境



主动交互感知

智能体



异构随遇协同

人



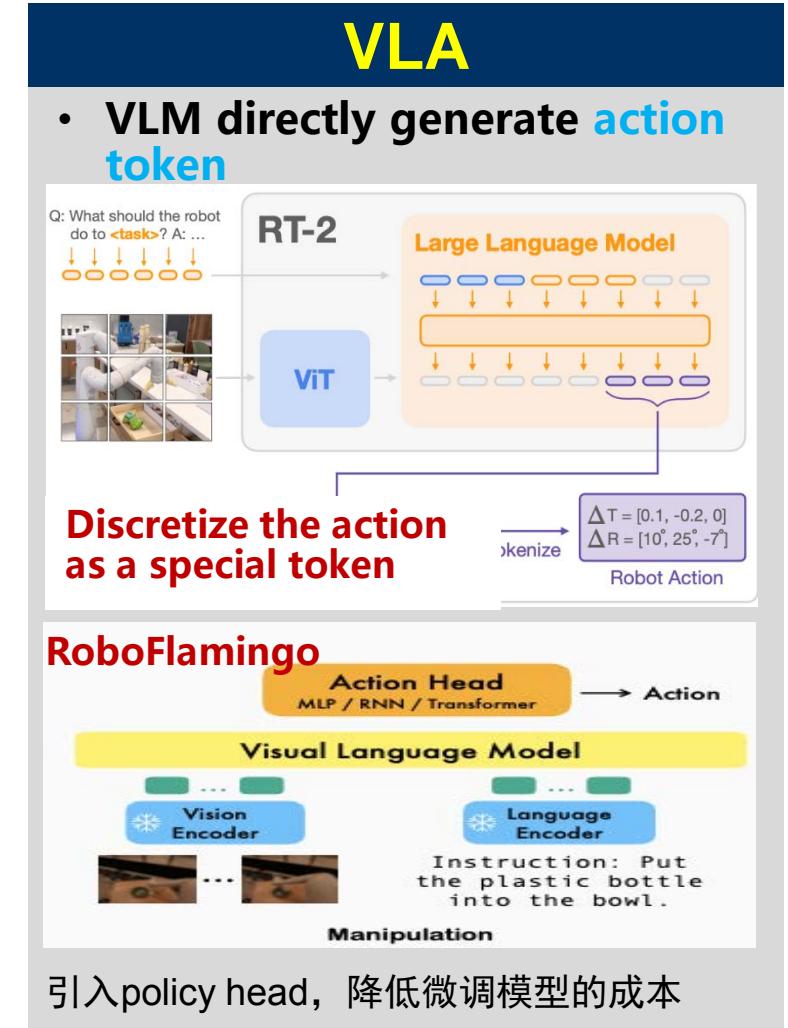
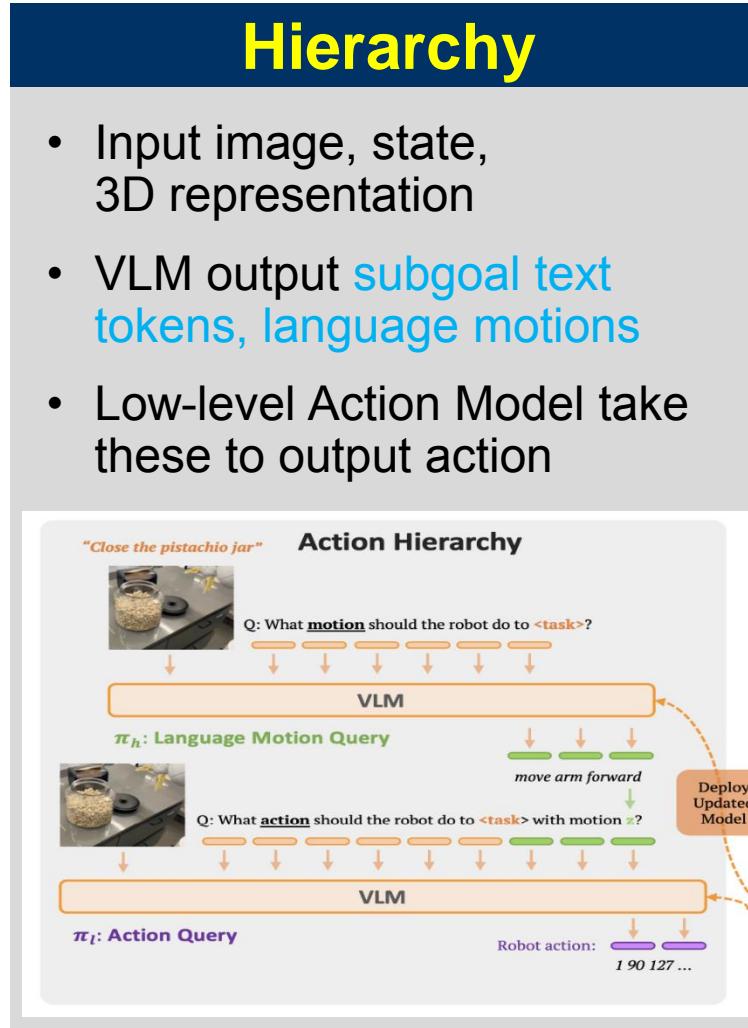
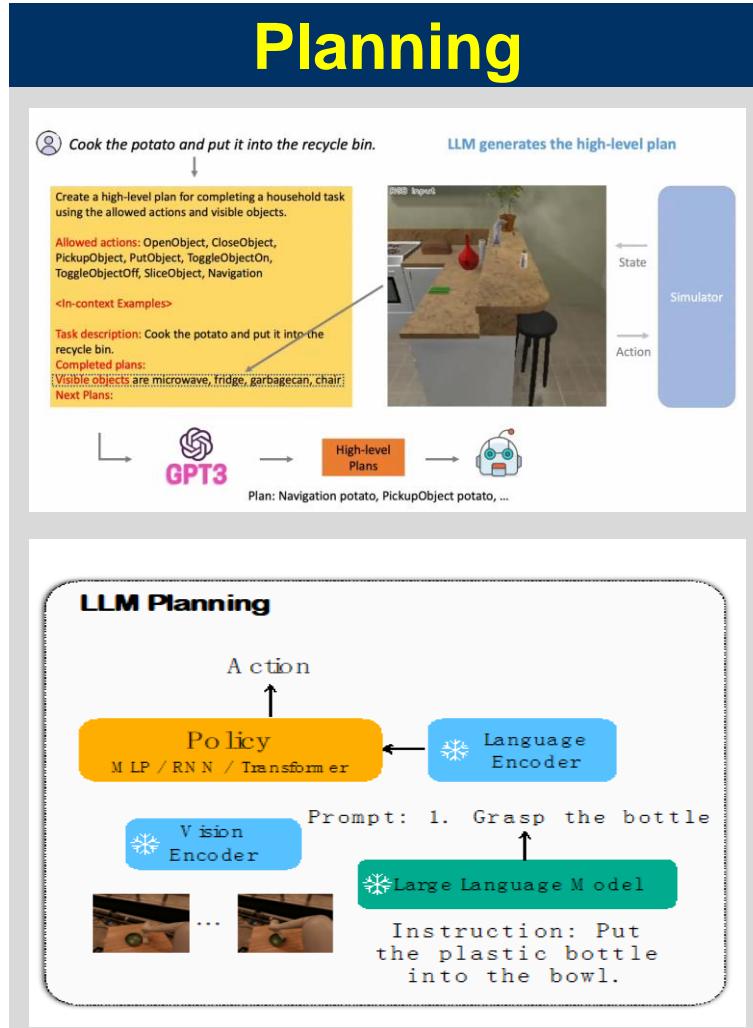
信任持续演进

-
- 背景：具身智能
 - 具身智能的体系
 - 具身智能关键技术
 - 探索与实践
 - 面向具身智能的AIGC
 - 总结

5 面向具身智能的AIGC

65

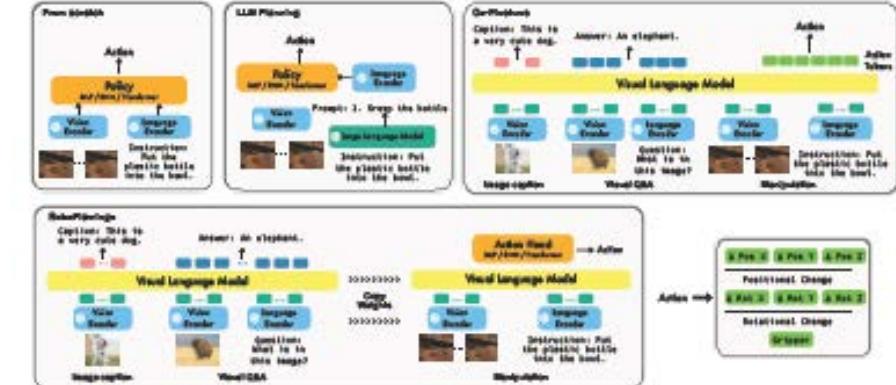
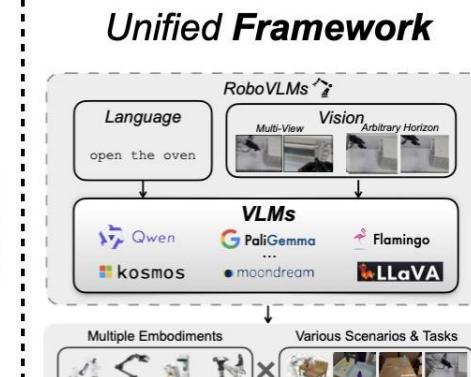
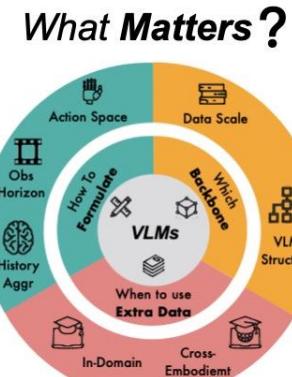
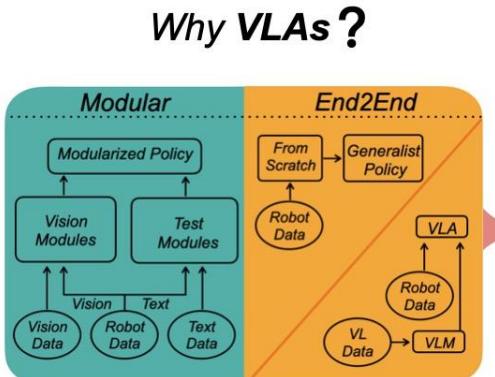
➤ 大模型



5 面向具身智能的AIGC

66

➤ 具身大模型



- Vision-language foundation models as effective robot imitators, ICLR, 2024

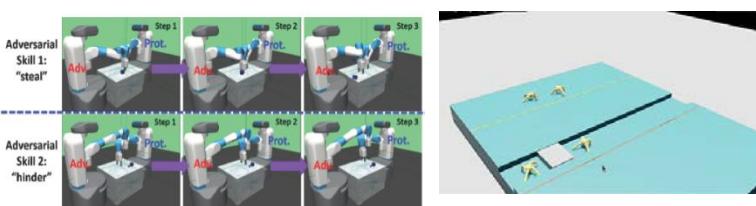
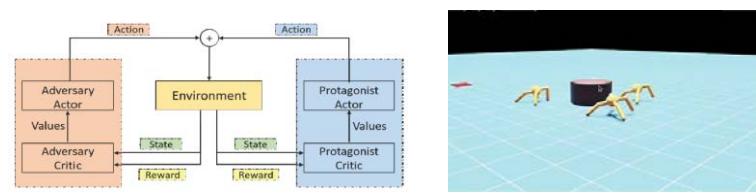
5 面向具身智能的AIGC

67

➤ 基于竞争、对抗的多体具身交互学习方法

对抗行为学习

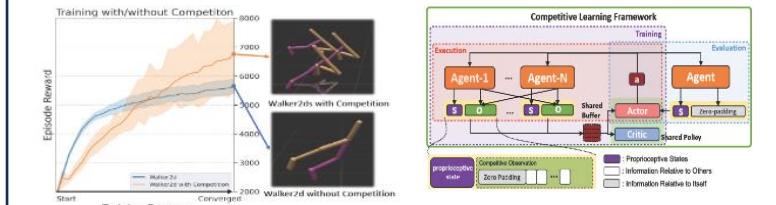
- 通过两个智能体之间的行为对抗（“偷走”与“抢回”），得到了鲁棒性更强的智能体行为策略。



- Adversarial Skill Learning for Robust Manipulation, ICRA, 2021
- Learning a distributed hierarchical locomotion controller for embodied cooperation, CoRL, 2024

竞争行为学习

- 通过在智能体独立训练过程中引入竞争信息，得到了比单一训练更好的结果，揭示了“竞争促进学习”的机制。

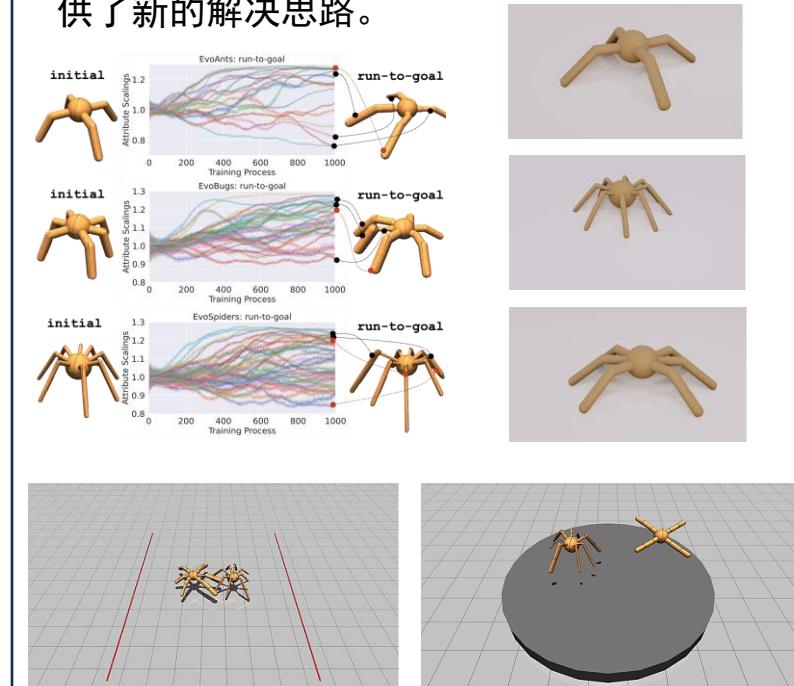


max velocity record			
4.39	3.15	4.65	3.02
5.09	7.59	6.02	2.87
5.27	7.43	6.88	2.90
5.33	7.25	6.80	2.98
5.43	7.23	6.33	2.92

- Stimulate the Potential of Robots via Competition, ICRA, 2024
- Adversarial decision making against intelligent targets in cooperative multi-agent systems, IEEE T-CDS, 2023

对抗形态学习

- 模拟人类身体发育机理，探索了智能体在对抗过程中可改变形态的方法，为形态优化提供了新的解决思路。



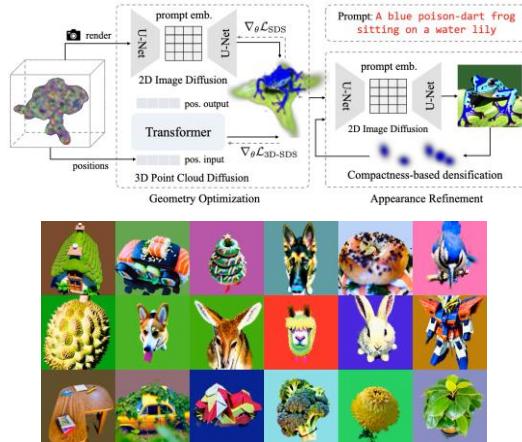
- CompetEvo: Towards morphological evolution from competition, IJCAI, 2024

5 面向具身智能的AIGC

68

➤ 基于AIGC的材料-结构一体化学习

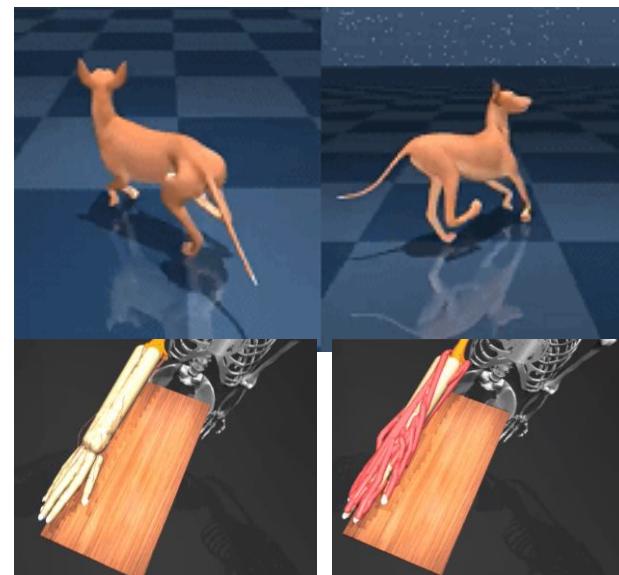
GSGen



Text-to-3D using Gaussian Splatting



Editing any 3D Gaussians within minutes

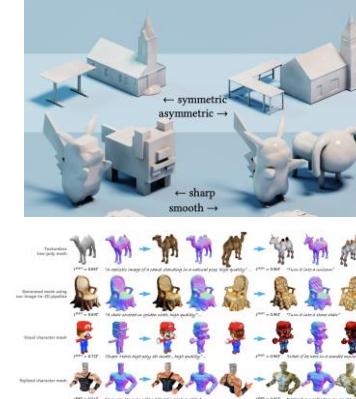


GaussianEditor

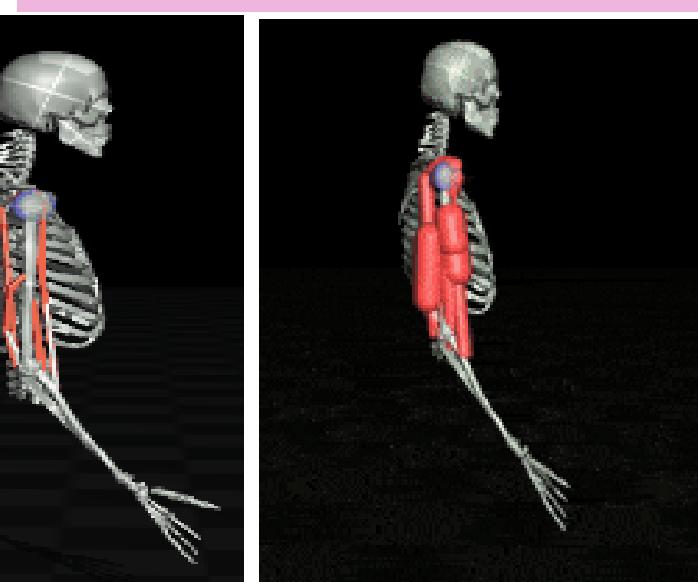
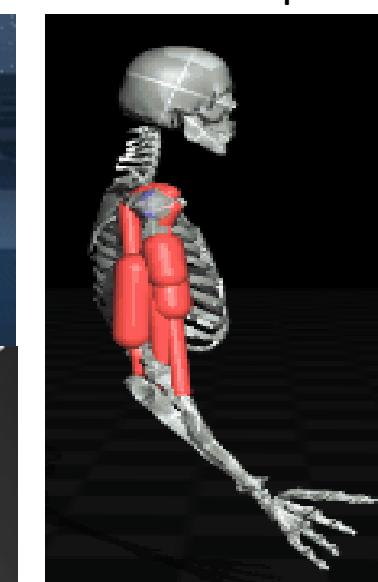
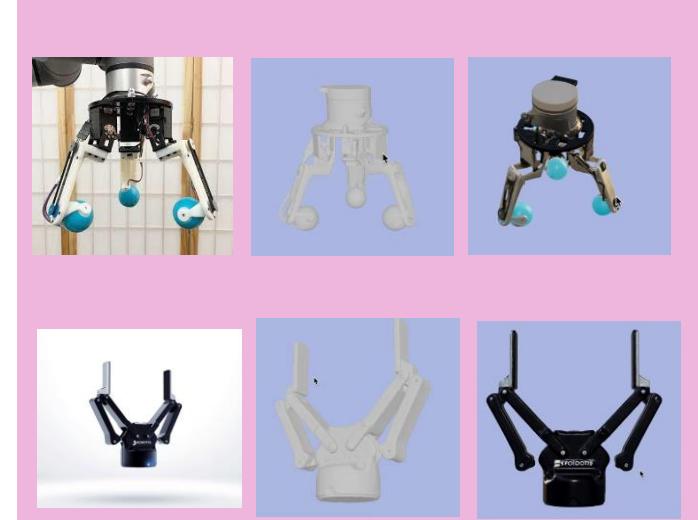


Original View Removed Inpainted

MeshEdit



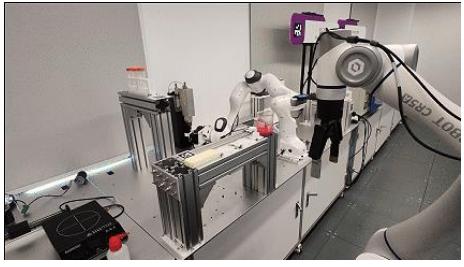
Mesh Editing with textual instruction and complex demands



5 面向具身智能的AIGC

69

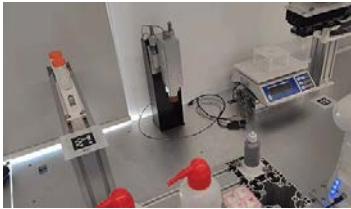
➤ 大模型驱动的具身智能科学实验室



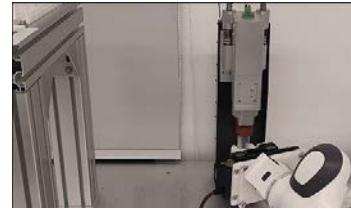
获取原料及模具



递送材料到操作区



获取试管



去除试管盖



放置试管到电子秤



添加硅胶



添加磁粉



添加试管盖



放置试管, 开始搅拌



放置模具到电子秤



从工作区取回材料



放回材料到材料区



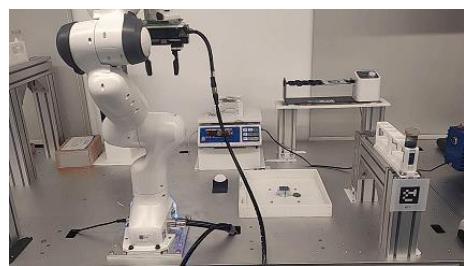
从搅拌机取下试管



去除试管盖



浇筑模型



添加试管盖



放置废弃试管到废料区



取下浇筑好的模具



移动机器人获取模具

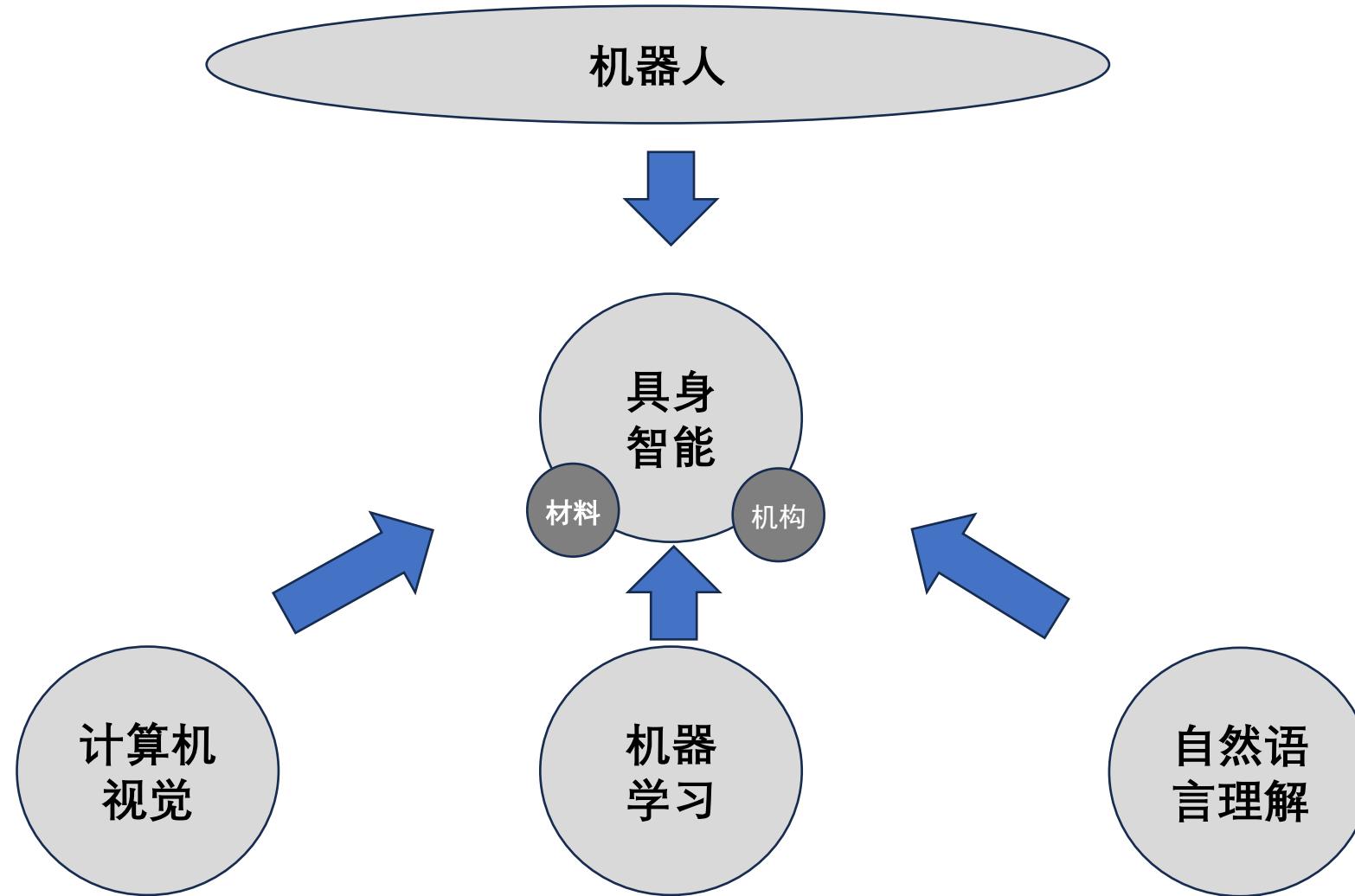


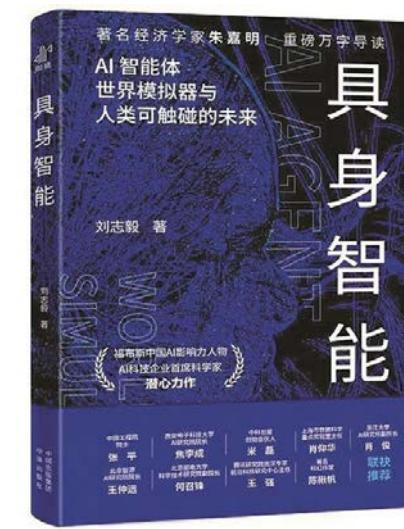
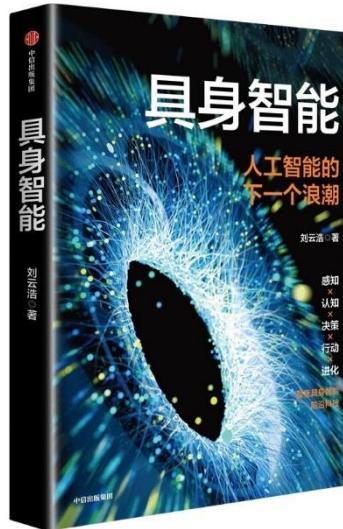
将模具放置到匀胶机

-
- 背景：具身智能
 - 具身智能的体系
 - 具身智能关键技术
 - 探索与实践
 - 面向具身智能的AIGC
 - 总结

总结

➤ 多学科交叉





浅晓终来得上纸
行躬要事此知绝

谢 谢 !