

# 具身智能-10

---

刘华平

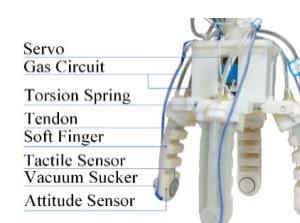
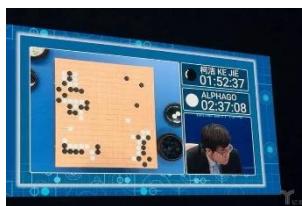
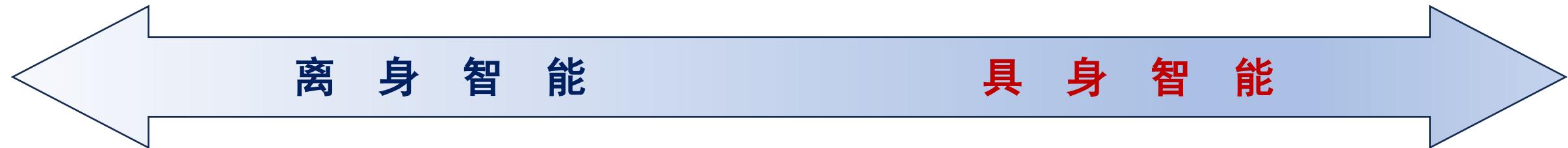
2025年4月23日

# 课程内容安排

课次	周次	上课内容	软件
1	1	绪论	
2	2	深度学习	
3	3	强化学习1	Gym, Mujoco
4	4	强化学习2	Gym, Mujoco
	5	作业准备	
5	6	自监督与持续学习	
	7	开题	Powerpoint
6	8	形态智能	Gym, Mujoco
7	9	视觉导航: VLN	AI2THOR
8	10	主动感知: VSN, EQA	AI2THOR
	11	五一放假	
9	12	具身学习	AI2THOR
10	13	多体智能	AI2THOR
11	14	面向具身智能的AIGC	AI2THOR
	15-16	成果准备与展示	Powerpoint

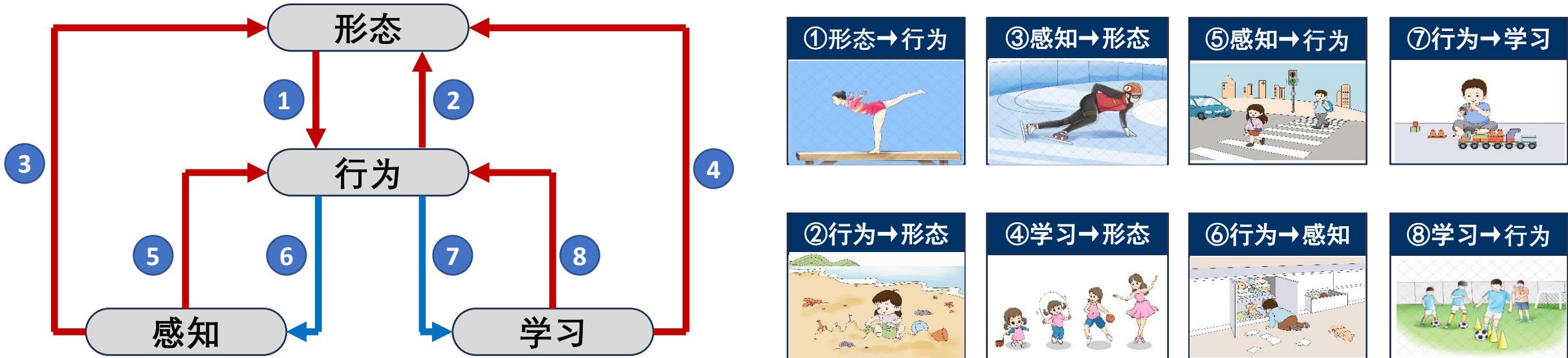
# 具身智能的体系

## ➤ 狹义与广义的具身智能

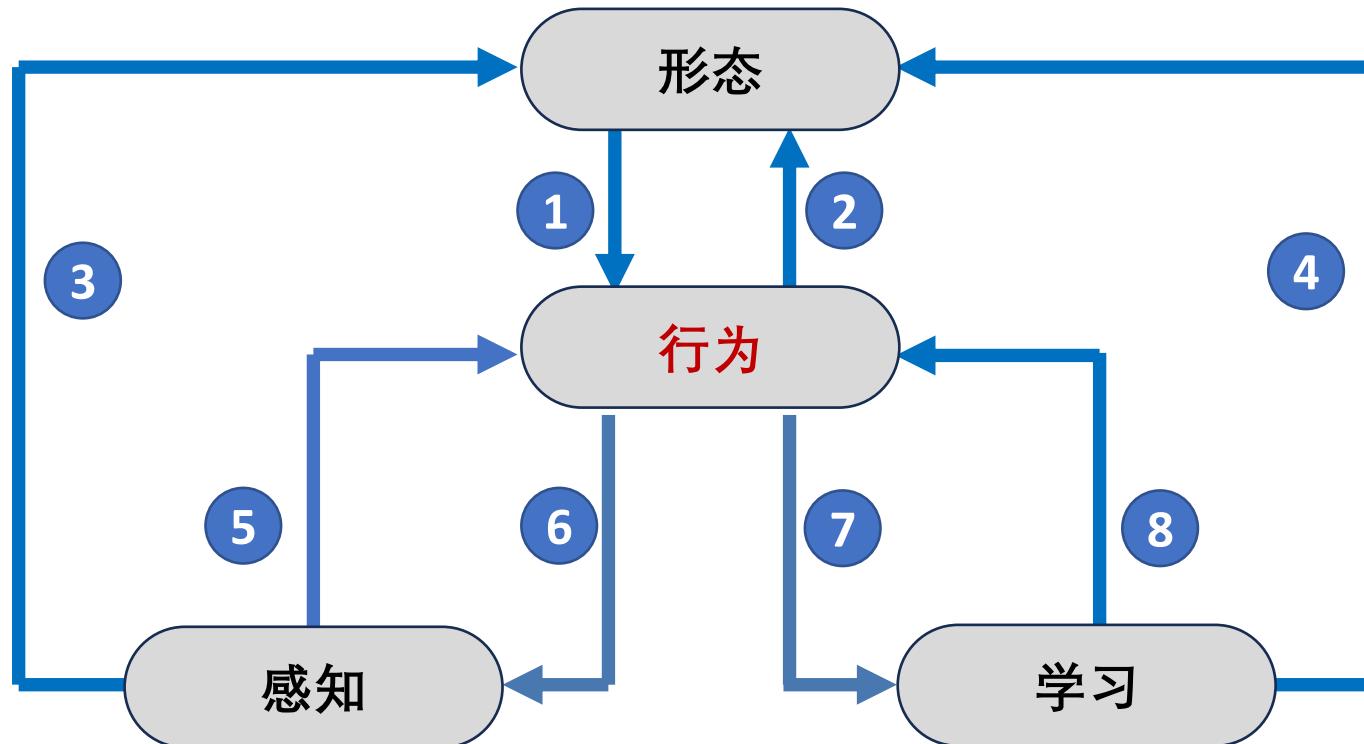


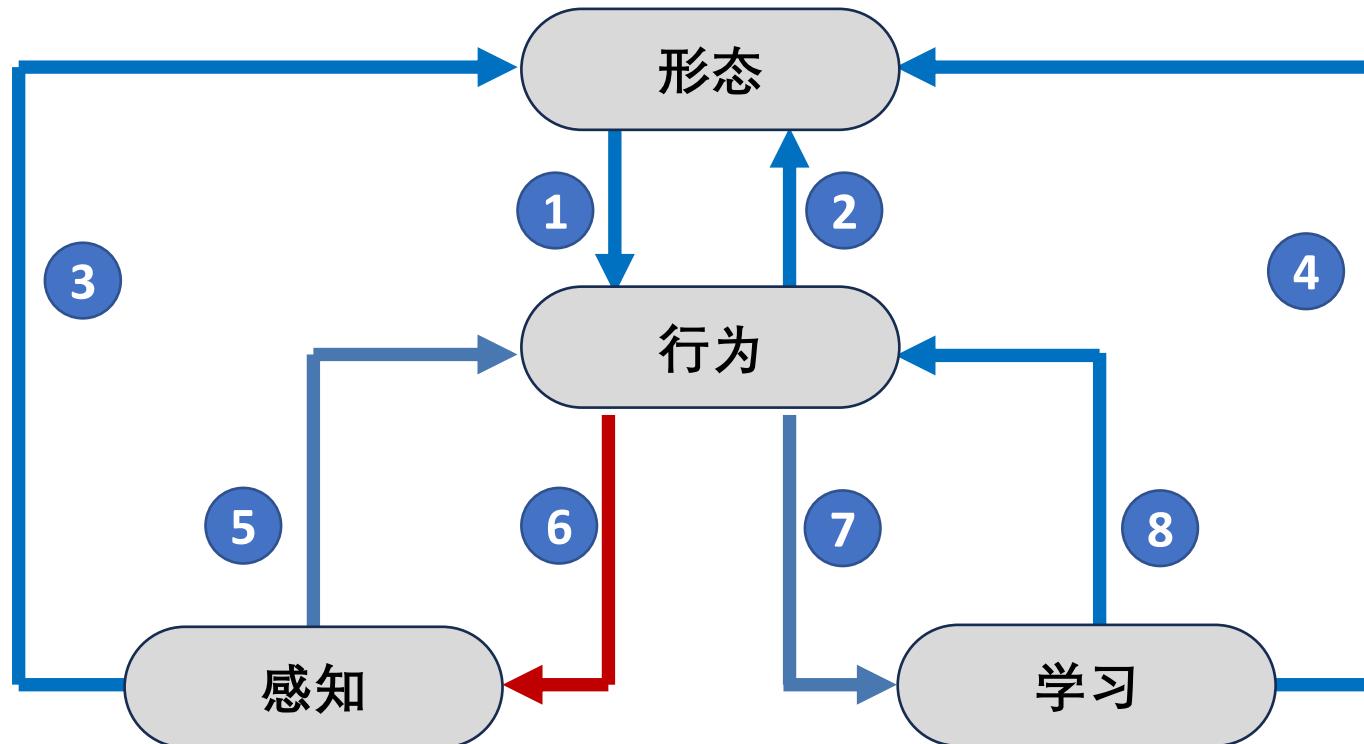
# 具身智能的体系

## 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





- 
- 从视觉导航到主动感知
  - 视觉语义导航
  - 具身场景描述
  - 具身语言问答

# 1 主动感知

## ➤ 为什么需要主动感知

眼睛看到的事实其实是大脑想让我们看到的东西

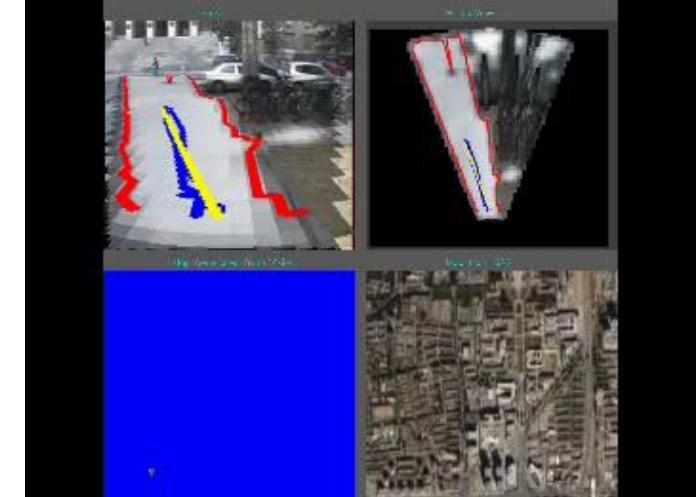


心不在焉

视而不见  
听而不闻  
食而不知其味

# 1 主动感知

## ➤ 为什么需要主动感知



反射式  
Reactive



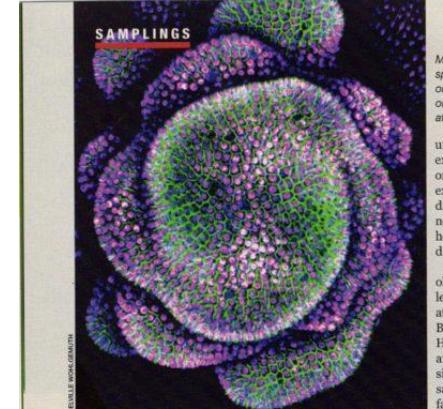
# 1 主动感知

## ➤ 生活中的主动感知



# 1 主动感知

## 动物中主动感知



**SAMPLINGS**

Microscopic projection of an *Arabidopsis* shoot apex deficient in an auxin-response transcription factor (magenta) which was experimentally restored in the outer cell layer, resulting in the formation of continuous, rather than discrete, organs. Auxin transporter protein (green) tends to follow the transcription factor at a lag.

unknown. Contemporary experiments have demonstrated that auxin is exported from cells in a directional manner. Now, new research has revealed how these cellular export directions are determined.

Developmental biologist Marcus G. Heisler, leader of a research group at the European Molecular Biology Laboratory in Heidelberg, Germany, and affiliated with the University of Sydney in Australia, says some members of his family are artists, which may have influenced his fascination with plant patterning. Heisler and five co-authors in the United Kingdom, Sweden, and Germany tested preexisting models for auxin transport to find out how communication between cells gives rise to polarity patterns.

Using plants in which auxin signaling was impaired genetically, the research team activated an auxin response trans-

cription factor (which helps to turn auxin-regulated genes) in just a few cells to see what neighboring cells would do. Neighboring cells oriented their auxin transporter proteins towards the cells where high auxin signaling had been engineered.

That response confirmed the researchers' hypothesis of a positive feedback mechanism: the presence of auxin at a given location stimulates the plant to send even more there, creating an auxin hotspot. At the same time, auxin is depleted from regions further away, yielding regular spacing between hotspots and thus periodic leaf formation.

Furthermore, modulating auxin signaling in single cells or in small patches of cells not only oriented the auxin transport in neighboring cells, but also oriented their cytoskeletons. That finding suggests that high local auxin signaling orients the whole tissue, coordinating growth directions. "It's really playing a fundamental role in determining how the plant tissue grows, and what shapes it takes," says Heisler. (*Current Biology*)

—Lesley Evans Ogden

**Directional Signals**

The periodic leaf patterns of plants can give rise to beautiful spirals and whorled patterns. Since the 1930s, it has been recognized that the plant hormone auxin stimulates leaf tissue formation, but how auxin is distributed within tissues to create such complex arrangements has remained

**Active Perception**

To enhance sensory perception, many animals—including humans, mice, and bats—in "active sensing" the environment use the body's motor systems to generate sensory information. In the case of hearing, dynamic movement can help to isolate important sounds—such as those made by prey species—from background noise. Echolocating bats are prime subjects in which to study this behavior, as they must coordinate their movements with feedback from the sonar vocalizations they use to locate prey.

Recently, neuroethologist Melvyn Wohlgemuth and his colleagues at Johns Hopkins University and the University of Maryland studied the mechanics of head and ear motion in big brown bats (*Eptesicus fuscus*). To overcome the difficulty of measuring these movements in flight, the team trained the bats to wait on a platform for their food—tethered mealworms—and used specialized motion-capture video to record the bats while they tracked prey.

The researchers found that as echolocating bats received auditory feedback, they made very quick adjustments, on a millisecond timescale, to their head and ear positions. When the target was far away, the bats held their ears upright. This positioning helps to accentuate sounds returning down the midline, explains Wohlgemuth.

Not only were head and ear movements coordinated with auditory feedback, the bats also adjusted the frequency and duration of their sonar calls depending on the number and closeness of prey. The researchers suggest that this coupling of vocalization adjustments and wagging sharply tunes bats' ability to hone in on their prey.

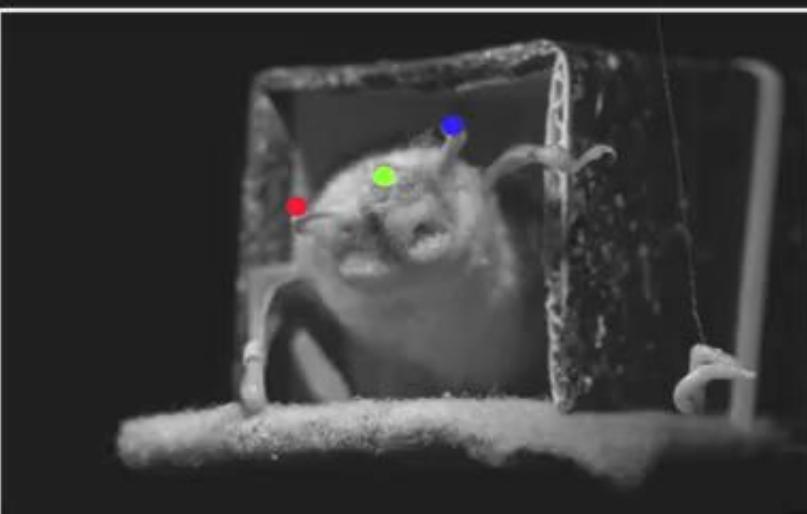
The findings clarify the remarkable hunting abilities of bats, as well as shed light on integrated auditory feedback from an engineering perspective, says Wohlgemuth. "These sorts of things could be very easily applied to artificial sensing technology to increase the resolution and acuity of those systems." (*PLOS Biology*)

Niki Wilson

6 | NATURAL HISTORY February 2017

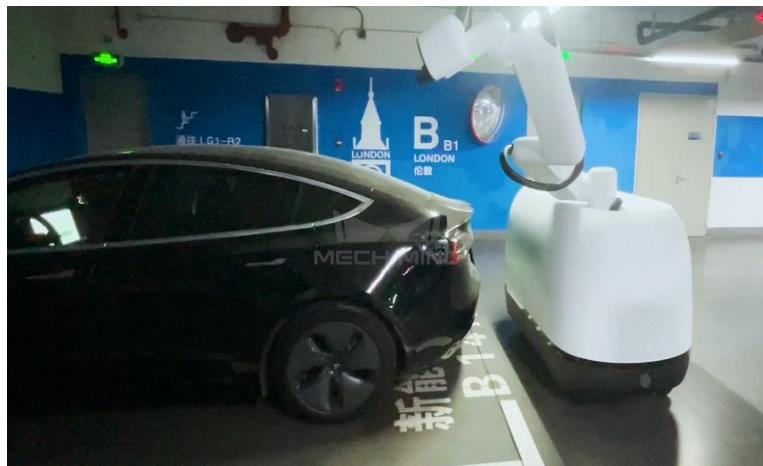
Not only were head and ear movements coordinated with auditory feedback, the bats also adjusted the frequency and duration of their sonar calls depending on the number and closeness of prey.

## Control of Head and Pinna Movements While Tracking a Moving Insect



# 1 主动感知

## ➤ 主动感知的应用



# 1 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.



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

# 1 主动感知

## ➤ 主动感知的优点

**The Past and Future of  
Robotics and Machine Intelligence  
Based on 250 Years of Research  
Experience**

Tues Sept 5 12-1pm  
Banatao Auditorium

Ruzena Bajcsy and colleagues:

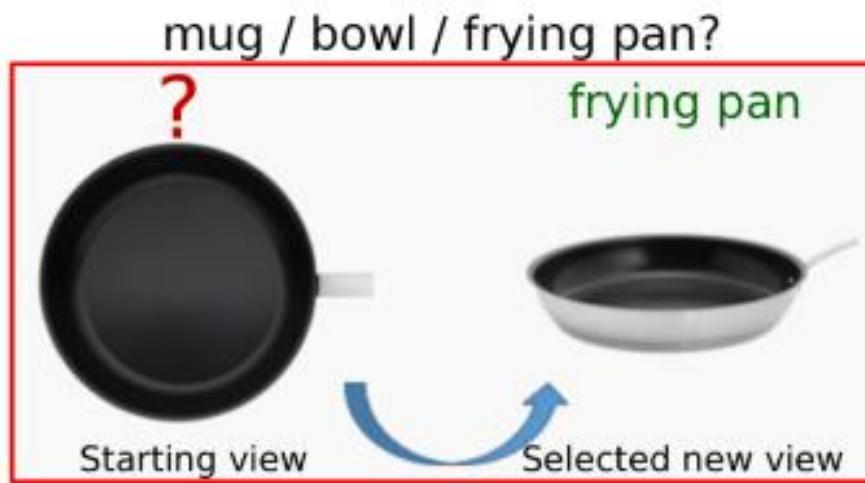
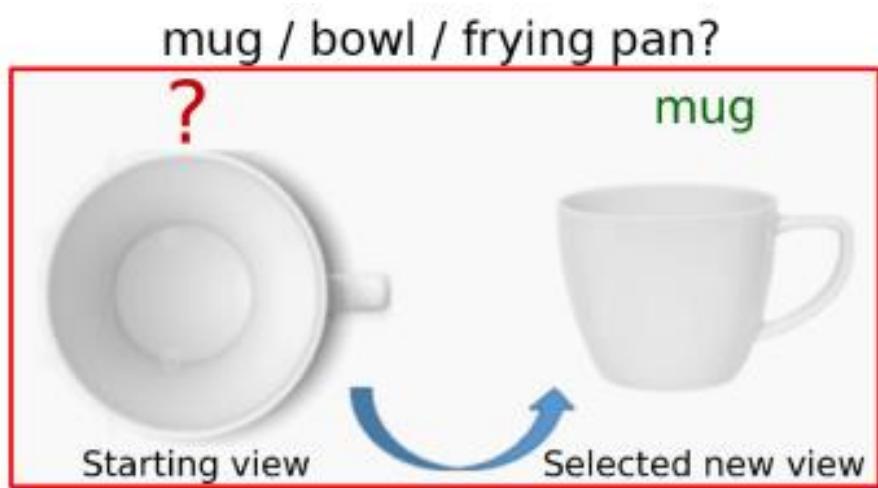
Rod Brooks, Ken Goldberg, Jitendra Malik  
Shankar Sastry, and Claire Tomlin



# 1 主动感知

## ➤ 主动感知的优点

很多困难的视觉感知问题，如果允许“动一动”，其难度会显著降低。

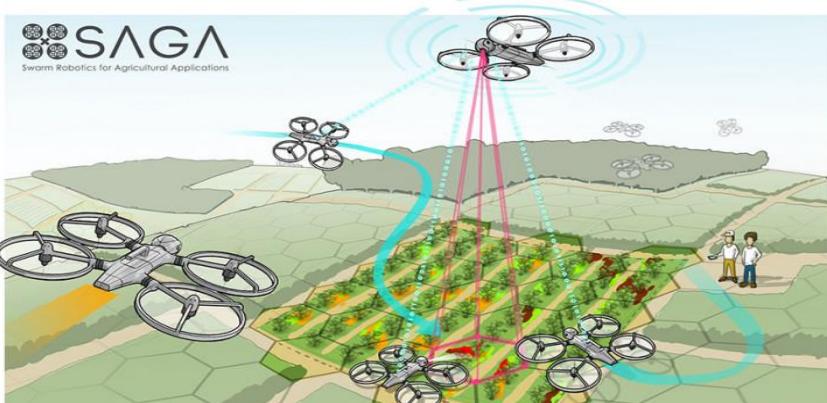


# 1 主动感知

## ➤ 主动感知的应用



A single view of a cluster can be deceiving (First three images from left). The first view suggests that there are three apples, the second one suggests five and finally from a bottom up view we see that there are actually six apples in the cluster. The rightmost image shows our platform for automating different precision agriculture and phenomics operation.



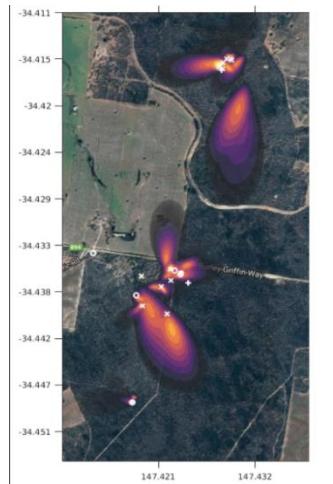
SAGA项目：农业机器人集群  
Swarm Robotics for Agricultural Applications



无人机跟踪高速无规律运动且运动范围大的小型动物（雨燕鹦鹉）。



(a) Experimental Setup (b) Planned Views. None of them shows all three apples clearly.



# 1 主动感知

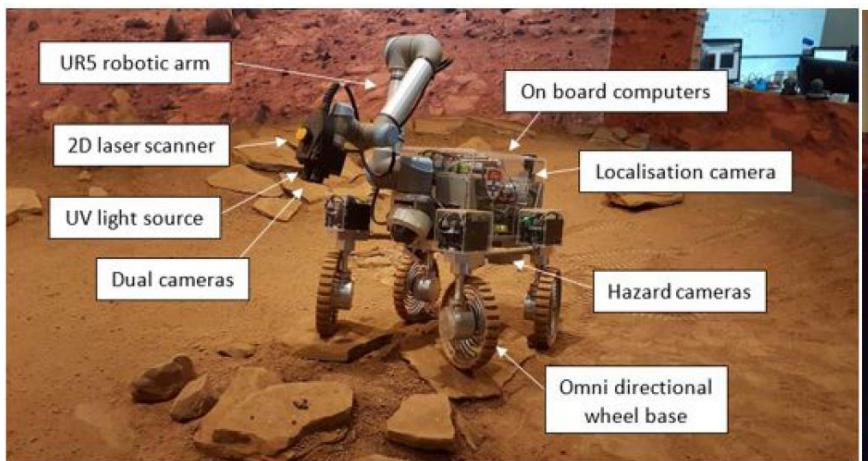
## ➤ 主动感知的应用



两种模态

Camera

Spectrometer



# 1 主动感知

---

## ➤ 主动感知的关键

- Single view上的感知
  - 常规认为，应该在single view上做到尽可能好。这种想法其实割裂了感知与动作的关系。
  - 其实未必需要这样。有时候在single view上做到最好，需要很大的代价（例如大数据的深度学习训练）。
  - 总结两种模式：“弱感知” + “强控制”；“强感知” + “弱控制”
- 多个view上的融合（证据积累）
  - 常规方法解决的比较多的是同构的，例如，摄像机转来转去，其实获得的都是图像
  - 实际中应该更多地考虑异构的多模态情形  
与常规多模态融合不同的是，主动感知中的融合一般是序贯融合
- 动作生成
  - 基于学习的方法，应该考虑训练时候的动作集与测试时的动作集有差异（比如训练时没有动作约束，但测试时有约束）的问题。这方面的工作在特征表示方面开展的多，但在动作方面很少

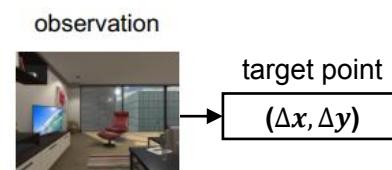
- 
- 从视觉导航到主动感知
  - 视觉语义导航
  - 具身场景描述
  - 具身语言问答

## 2 视觉语义导航

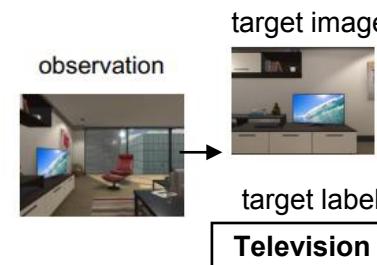
### ➤ 导航分类

- 视觉环境中常见的三类导航任务：
  - 1) PointGoal: 给定目标坐标点进行导航；
  - 2) ObjectGoal: 给定目标物体的类别，使智能体自主导航至特定物体旁；
  - 3) AreaGoal: 给定目标区域的类别如厨房，使智能体自主导航至特定场景内。

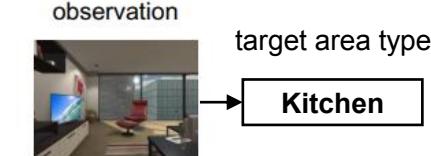
视觉语义导航对应其中的ObjectGoal问题，在视觉语义导航中不考虑PointGoal任务和AreaGoal任务。



PointGoal



ObjectGoal

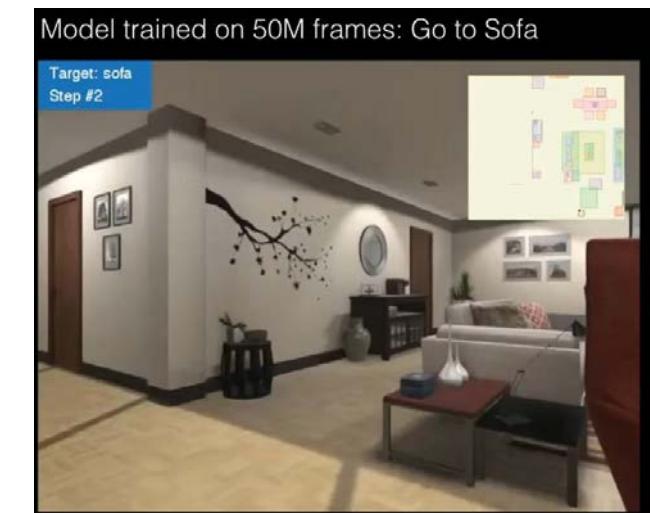
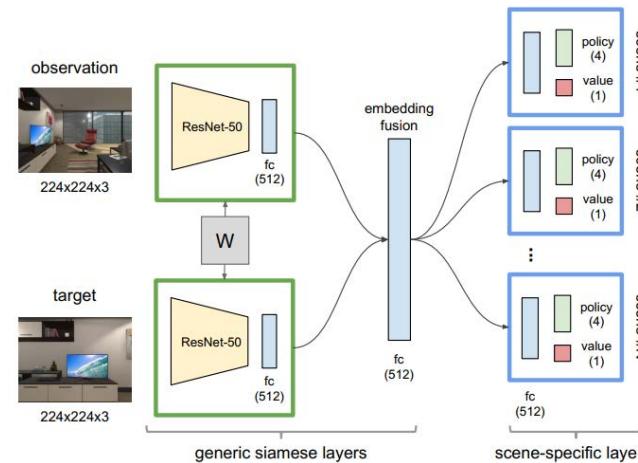
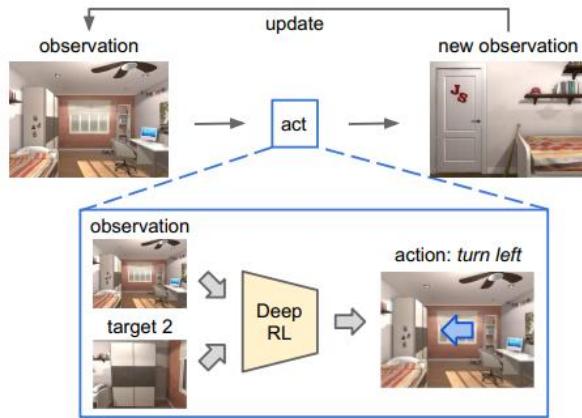


AreaGoal

## 2 视觉语义导航

### ➤ 目标驱动的导航问题的提出

- 提出室内场景下目标驱动 (target-driven) 的导航问题，即已知目标的图像信息，在室内场景中智能体仅使用第一视角图像信息通过与环境的交互完成对目标的导航任务。
- 提出一种孪生网络的基本结构，是最早的端到端解决视觉导航问题的模型，成为目标驱动导航问题模型的基本结构。



## 2 视觉语义导航

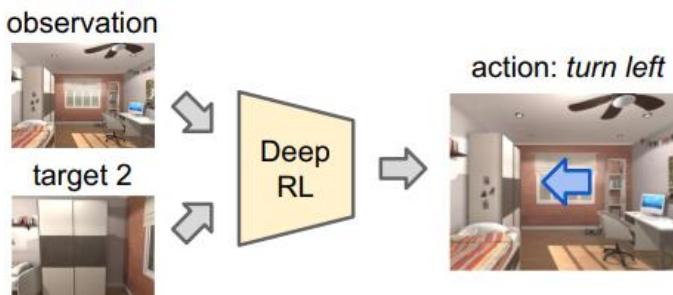
### ➤ 目标驱动的导航问题的提出

按其已知的输入信息将其分为两类：

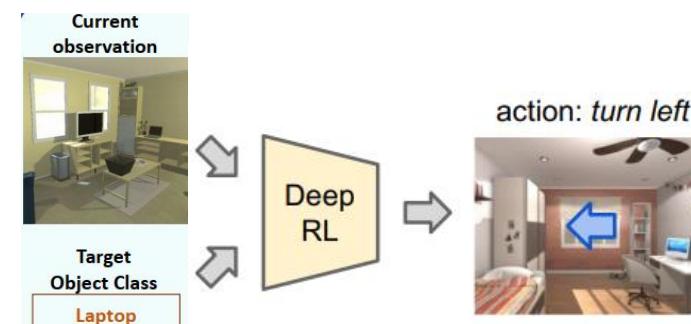
- 1) 已知目标物体的**图像**；
- 2) 已知目标物体的**语义或类别标签**

一般所指的标准的视觉语义导航是第二类问题，即仅已知目标物体的类别，如书、苹果等。

**知识图谱、深度图、地图重建**等方法被广泛应用于视觉语义导航问题。



目标物体图像



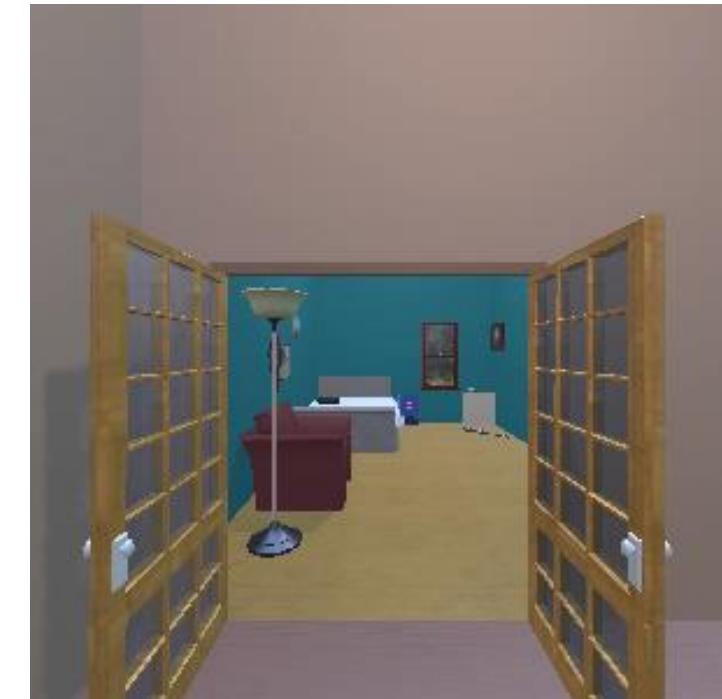
目标物体语义标签

## 2 视觉语义导航

---

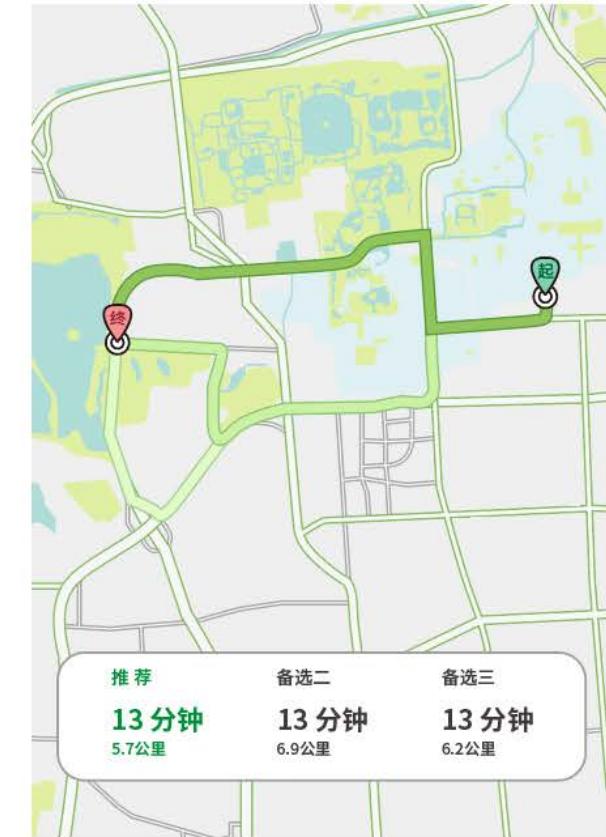
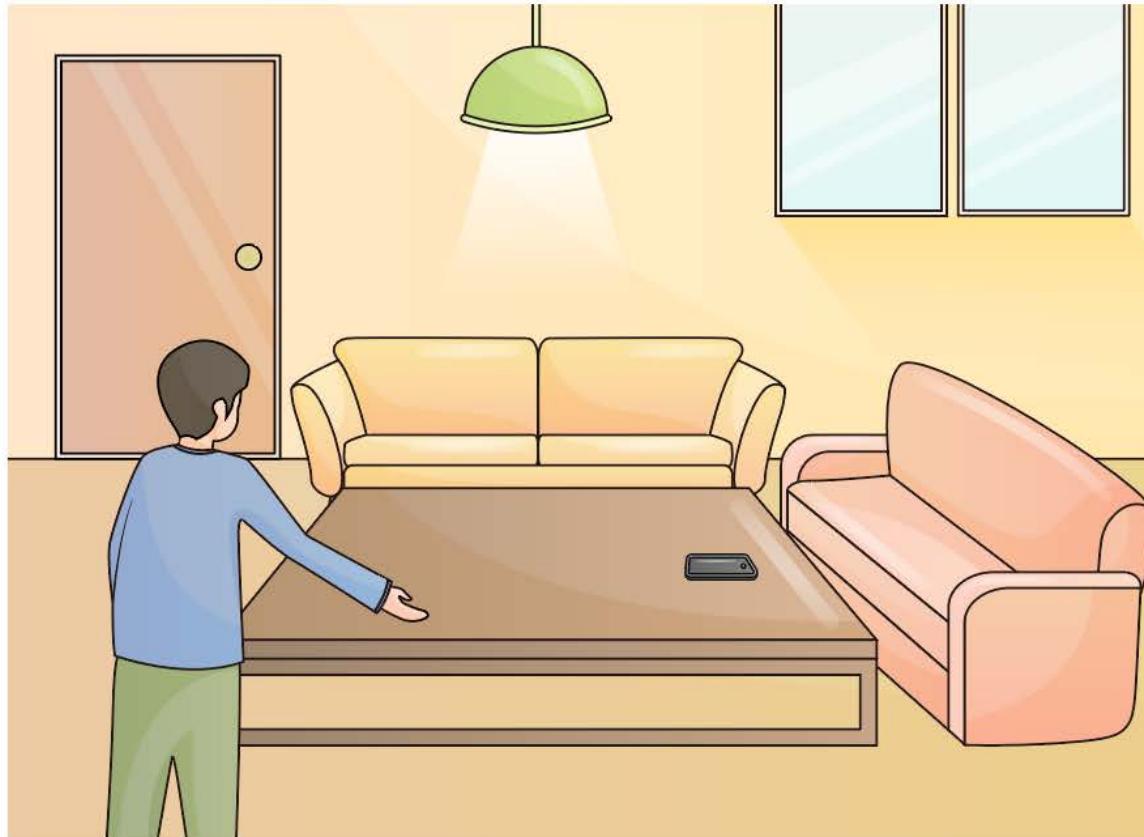
### ➤ 分类

- 视觉语义导航：给定目标的语义标签，仅使用视觉输入在环境中导航至特定目标。在视觉语义导航中，不考虑已知场景地图或者已知目标物体位置信息的情况。
- 任务类型：
  - 已见过场景&已见过目标
  - 已见过场景&未见过目标
  - 未见过场景&已见过目标
  - 未见过场景&未见过目标
- 常见方法：
  - 基于学习的导航方法



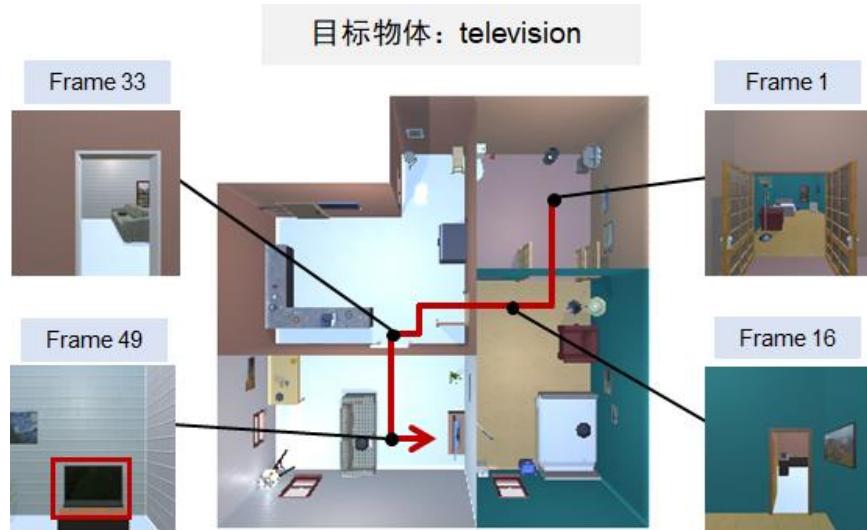
## 2 视觉语义导航

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



## 2 视觉语义导航

### ➤ VSN问题描述



$$s_t = \varphi(I_t)$$

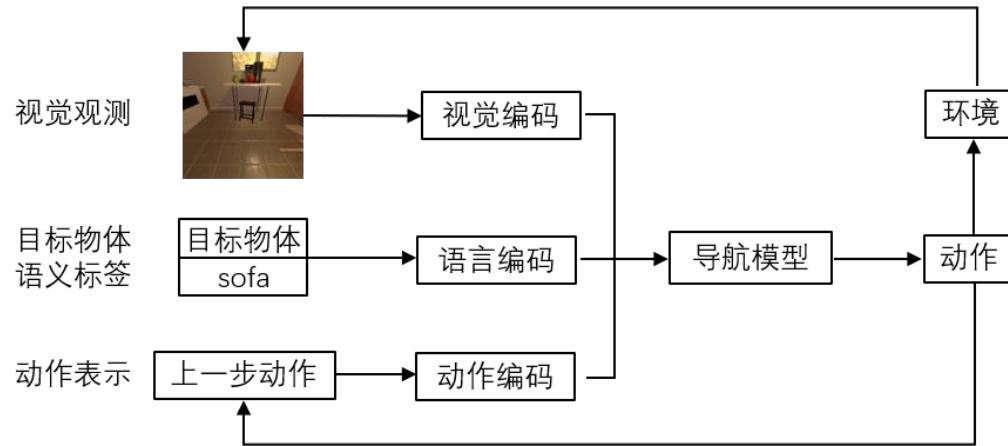
$$o = \phi(O)$$

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

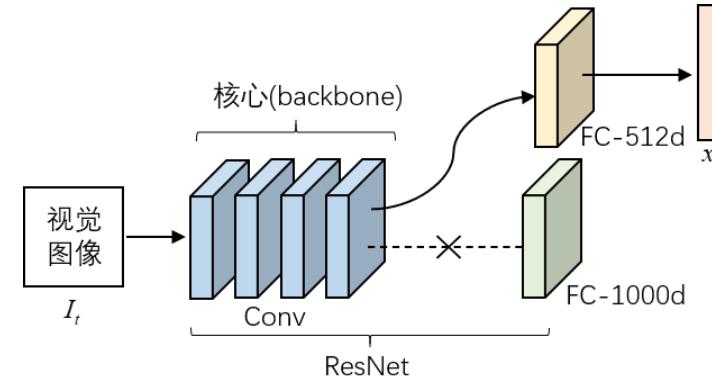
$$s_{t+1} = \text{Env}(s_t, a_t)$$

# 2 视觉语义导航

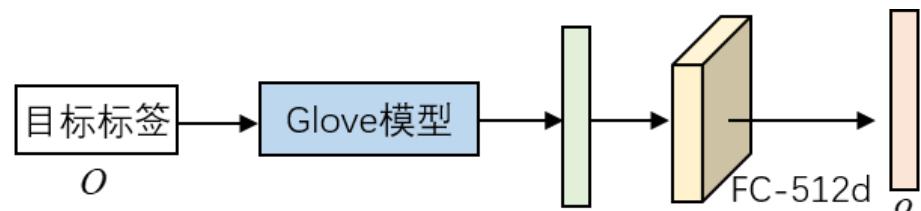
## ➤ VSN基本方法



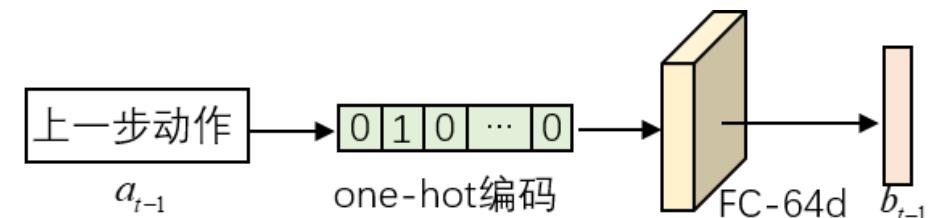
视觉图像特征提取



语义特征提取



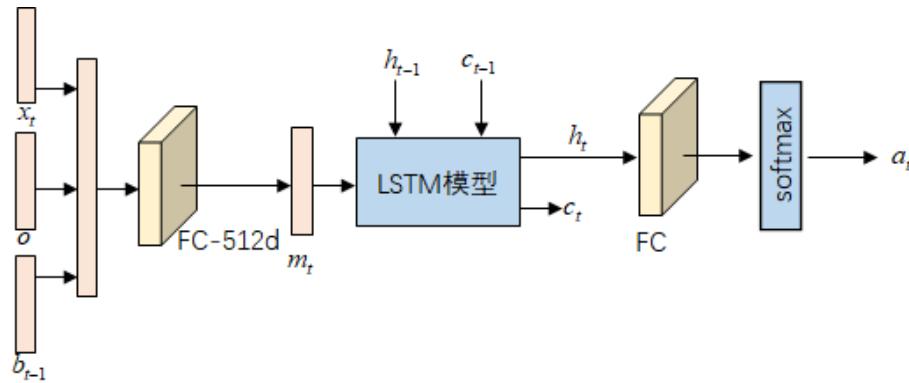
动作特征提取



## 2 视觉语义导航

### ➤ VSN基本方法

导航模型



- 强化学习

$$r_t = \begin{cases} 5 & \text{成功找到目标物体} \\ -0.01 & \text{其他} \end{cases}$$

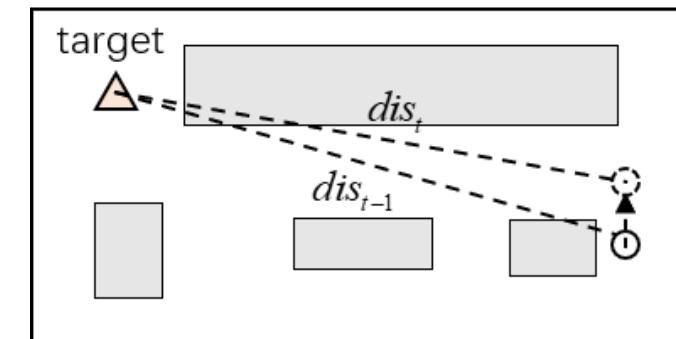
$$r_t = \begin{cases} 5 & \text{成功找到目标物体} \\ dis_t - dis_{t-1} + penalty & \text{其他} \end{cases}$$

- 模仿学习



Target: Pillow

$$L_\theta = \sum_{t=1}^T -\log \pi_\theta(a_t | s_1, a_1, \dots, s_{t-1}, a_{t-1})$$



# 2 视觉语义导航

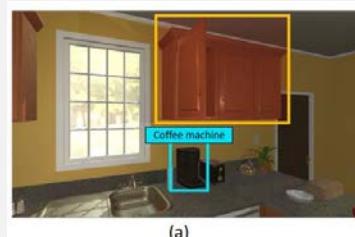
## ➤ 融入场景图谱的VSN

### 动机



### 启发

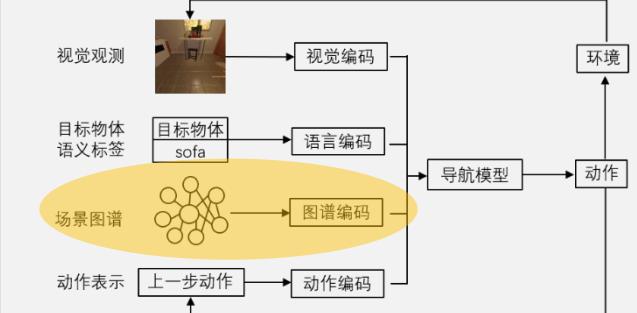
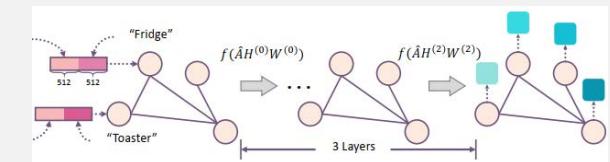
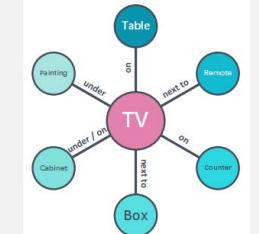
- 寻找马克杯时，当前视角中没有马克杯，但根据以往的经验，马克杯的位置很有可能是在咖啡机旁边的橱柜里。
- 寻找芒果时，智能体之前没有见过芒果，但根据经验，芒果和苹果是类似的物体都是水果，而智能体之前在冰箱中见过苹果，于是推断出芒果有可能的位置是在冰箱里。



#### 引入场景图谱的优势：

- 有效编码不同类物体之间的空间关系
- 对于智能体未见过的物体类别，场景图谱能够提供其与已见过物体之间的空间和视觉相对关系。

### 方法



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

# 2 视觉语义导航

## ➤ 融入场景图谱的VSN

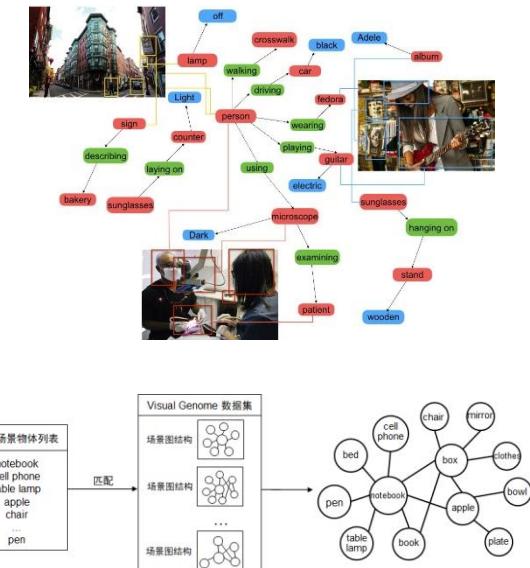
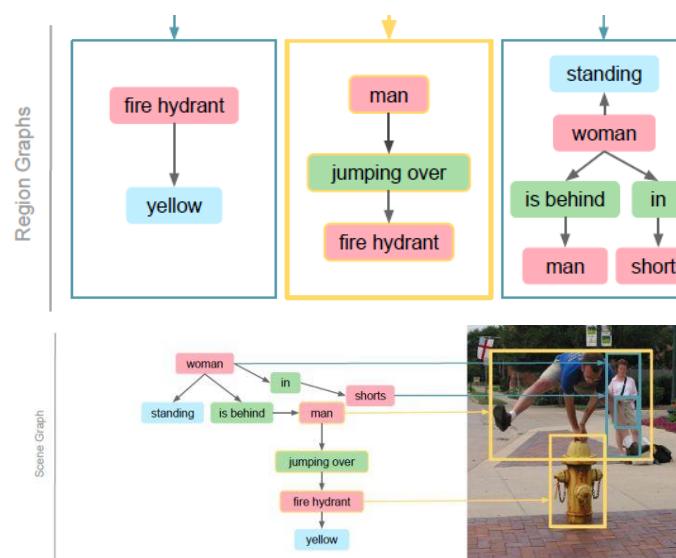
## ➤ 场景图谱的构建

将场景图谱记为  $G = (V, E)$ , 其中  $V, E$  分别代表图中的节点和边。每个节点  $v \in V$  代表一个物体类别, 每条边  $e \in E$  代表一种相对位置关系。使用 Visual Genome 数据集构建场景图谱, 该数据集包含超过十万张图片, 每张图片都被标注了对应的物体和物体间的相对关系, 其包含相当多的物体种类。



### • 关于Visual Genome

- Visual Genome (VG) 是李飞飞组于2016年发布的规模标注的图片语义理解数据集。
- VG数据集中的每张图片, 包含Region Description, Region Graph, Scene Graph
  - Region Description: 图片被分成若干region, 每个region有与其对应的一句自然语言描述
  - Region Graph: 提取每个region中的目标、属性和关系构成局部 Graph
  - Scene Graph: 把图片中所有Region Graph合并成一个全局 Scene Graph



## 2 视觉语义导航

### ➤ 融入场景图谱的VSN

### ➤ 场景图谱的处理

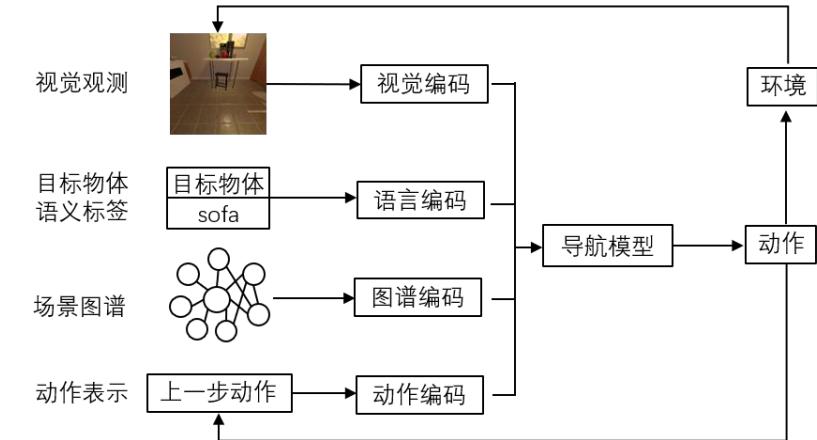
- 图卷积网络 (GCN)

GCN是CNN在图结构下的变种，用于从图结构 $G = (V, E)$ 中学习到一定的函数表达式。每个节点 $v$ 的输入是一个特征向量 $x_v$ 。所有节点的输入可表示为一个矩阵 $X = [x_1, \dots, x_{|V|}] \in \mathbb{R}^{|V| \times D}$ ，其中 $D$ 是输入向量的维度。整个图结构可以表示为一个二进制邻接矩阵 $A$ ，对其进行正交化得到 $\hat{A}$ 。

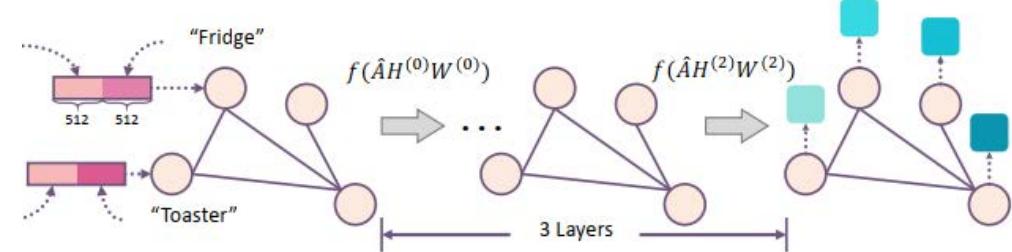
GCN网络输出一个节点级的向量表示 $Z = [z_1, \dots, z_{|V|}] \in \mathbb{R}^{|V| \times F}$ 。ReLU激活函数记为 $f(\cdot)$ ，对于GCN网络有

$$H^{(l+1)} = f(\hat{A}H^{(l)}W^{(l)})$$

其中 $H^{(0)} = X$ ， $H^{(L)} = Z$ ， $W^{(l)}$ 是GCN第 $l$ 层的权重矩阵， $L$ 是GCN网络的层数。



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



## 2 视觉语义导航

- 融入场景图谱的VSN
- 例子



将智能体所处环境进行网格化。动作空间由平移运动和旋转运动构成。

平移运动：包括停止，和向前、左、右四个方向的平移运动。每一次平移运动的移动距离固定为网格的边长。

旋转运动：包括左转、右转两个方向的旋转运动。每一次旋转运动的旋转角度为 $90^{\circ}$ 。

智能体通过执行动作空间中的动作，不断对环境进行探索，当智能体执行停止动作时，结束对环境的探索。

- 随机探索模型：智能体每一步从其动作空间中随机选择探索动作。
- 深度优先探索模型：智能体每一步优先沿着当前可行动作方向向前探索，并优先向右移动。
- IL：利用模仿学习训练动作生成模型。
- IL+A3C：利用模仿学习与强化学习相结合的方式训练动作生成模型。
- IL+A3C+场景图谱：引入场景图谱，采用IL+A3C的方法训练模型。

## 2 视觉语义导航

### ➤ 融入场景图谱的VSN



见过目标物体导航定量结果



未见过目标物体导航定量结果

$$d \geq 5 \quad d \geq 5$$

实验设置	模型	准确率(%)	SPL (%)	准确率(%)	SPL(%)
已见过 目标物体	随机	8.3	3.66	5.2	2.33
	深度优先探索	15.2	10.12	10.3	8.55
	IL	48.3	18.67	37.1	13.12
	IL+A3C	52.8	20.64	41.3	15.28
	IL+A3C+场景图谱	55.1	21.38	45.7	17.32

实验设置	模型	准确率(%)	SPL (%)	准确率 $\geq 5$ (%)	SPL $\geq 5$ (%)
未见过 目标物体	随机	8.2	3.58	5.2	2.32
	深度优先探索	15.0	10.11	10.2	8.53
	IL	37.1	13.65	24.3	8.78
	IL+A3C	43.7	15.23	28.1	10.05
	IL+A3C+场景图谱	45.5	17.13	30.3	12.65

## 2 视觉语义导航

### ➤ 融入场景图谱的VSN

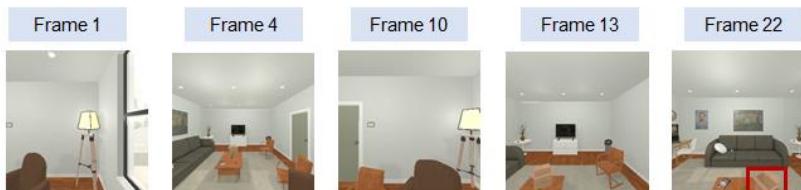
#### AI2-THOR



目标物体  
television



目标物体  
box



目标物体  
laptop



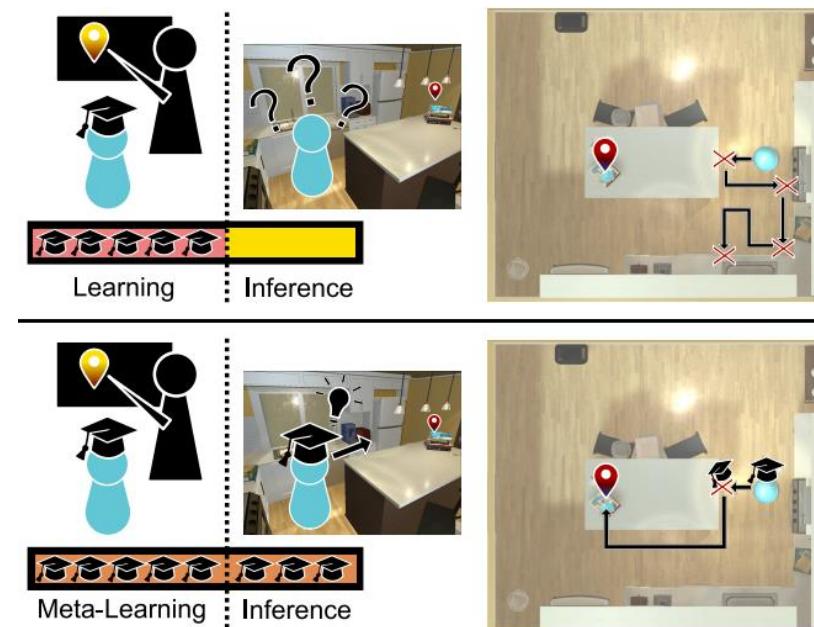
目标物体  
pillow



## 2 视觉语义导航

### ➤ 前沿：基于元学习的自适应视觉语义导航

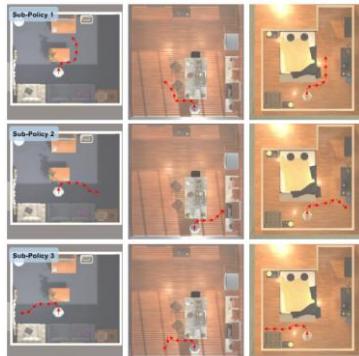
- 传统的强化学习或者深度学习方法，训练结束后执行某项任务时，网络结构及参数固定不变。
- 视觉语义导航任务的难点是智能体对在训练过程中未见过的新场景的泛化能力，因为不同场景中房间结构和物品摆放有较大的区别。为此，提出一种基于元学习的自适应方法，在训练和执行的过程中都在不断调整模型参数，使之能适应新场景。



## 2 视觉语义导航

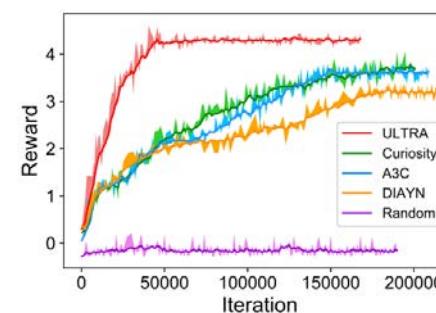
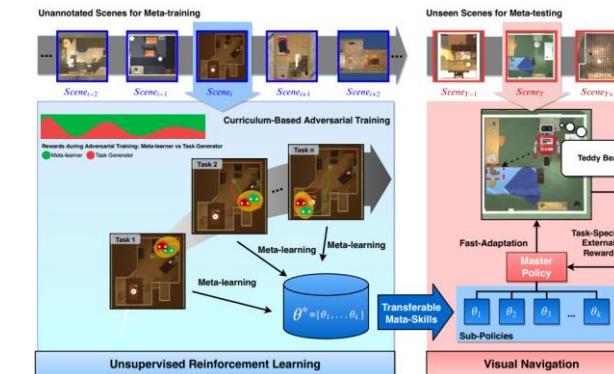
### ➤ 前沿：基于无监督学习的视觉语义导航

在无监督信号的情况下从环境中学习可迁移的元技能（meta-skills）和子策略（sub-policy）如绕过障碍物等。当新的环境提供相应的外部奖励时，智能体可以通过学习一个高级的master策略来组合之前学习到的元技能，从而快速适应视觉导航任务。



	All		$L \geq 5$	
	Success	SPL	Success	SPL
Random	8.21	3.74	0.24	0.09
A3C (learn from scratch)	19.20	7.48	9.43	4.13
DIAYN	17.23	6.30	8.72	3.79
Curiosity	21.07	8.51	10.31	4.37
<b>Ours</b>				
ULTRA	<b>27.74</b>	<b>11.47</b>	<b>20.57</b>	<b>8.04</b>
- hierarchical policy	24.27	10.54	14.13	5.61
- adversarial training	20.23	8.35	10.04	4.33

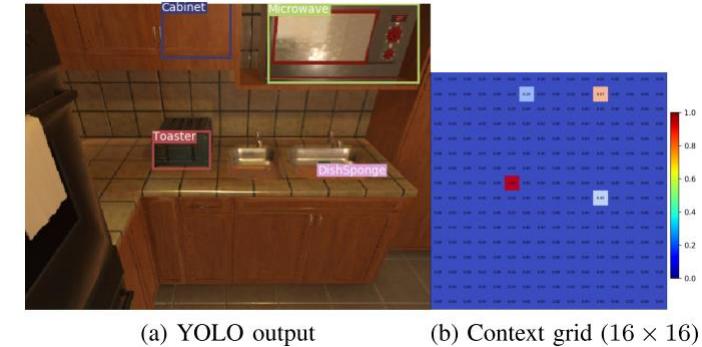
模型结构：蓝色部分是对抗训练过程，任务生成器自动生成越来越难的任务，而元学习模块用于学习如何完成这些任务。从这些任务中，元学习模块可以学习一组可迁移的子策略。在右侧，给定任务特定的外部奖励函数，元学习模块基于已学习到的可迁移策略通过学习新的master策略可以快速适应视觉导航任务。



## 2 视觉语义导航

### ➤ 前沿：基于上下文信息的视觉语义导航

- 人们寻找目标时不会将所有场景都看一遍，而是寻找最相关的物体。例如，智能体当前在位置(a)寻找海绵(b)，当前位置看不到海绵。利用空间的上下文信息（context）作为线索，如海绵应该在水槽附近，而水槽很容易被检测到，这使得智能体从看不见目标的位置开始也可以导航到目标。
- 提出了一种基于上下文信息的导航方式，使用context grid的表示方式来表示当前场景中的物体与目标物体的相似度，从而帮助机器人找到更小的或隐藏的目标。该方法具有更好的泛化性。

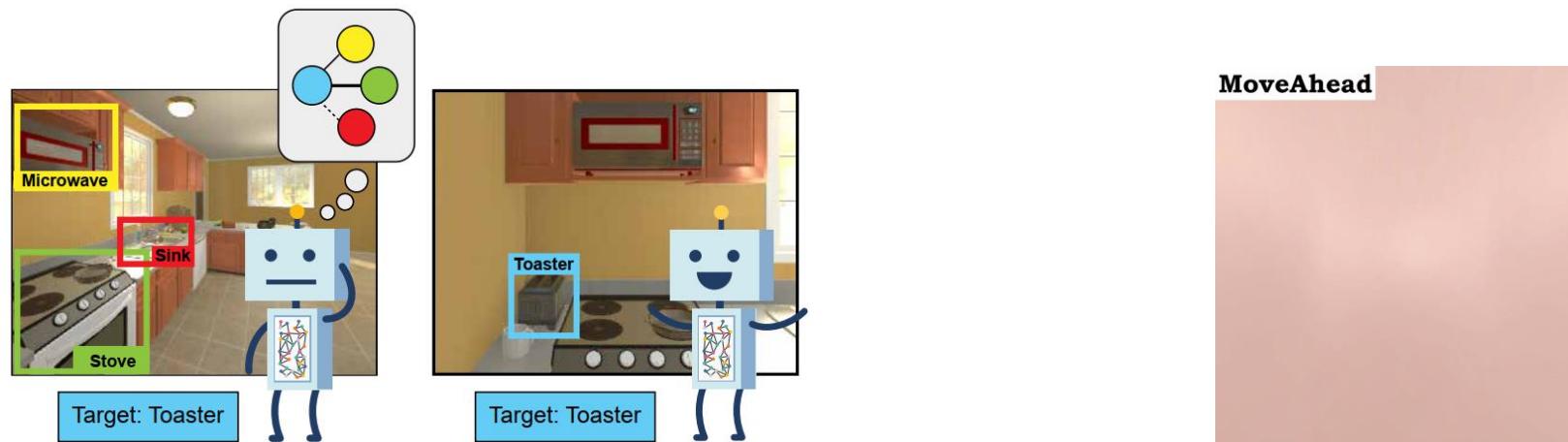


## 2 视觉语义导航

### ➤ 前沿：基于层次关系学习的视觉语义导航

人们在寻找物体时会建立不同物体间的内在层次关系。例如，在厨房中找烤面包机时，人会倾向于先找可能与之相邻的更大的候选物体如微波炉或者炉灶，然后再在这些物体周围寻找烤面包机。

将与目标物体有空间或者语义关系的目标称为parent-object，提出了一种层次物体关系学习方法，在图卷积网络中使用Context Vector作为结点向量表示。通过一种新的reward shaping机制引入层次物体关系。



# 2 视觉语义导航

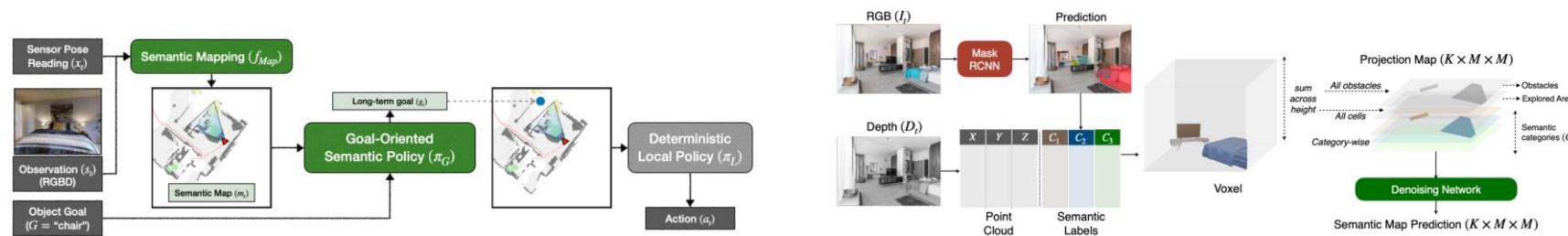
## ➤ 前沿：基于建图的视觉语义导航

- 在未知场景中进行目标导航时，智能体不仅需要识别目标，还需要知道在哪里更容易找到目标物体。
- 基于建图的导航方法，将第一视角通过Mask-RCNN获取语义信息，映射为Top-down视角的二维地图，利用该地图训练一种导航策略。



### 模型结构

模型包含两个模块，语义地图构建模块，和基于语义地图的语义决策模块。使用预训练的Mask-RCNN检测，映射到点云中生成带有语义的像素映射，将不同类别的物体映射到不同的地图层中生成语义地图。若当前状态发现目标，语义决策模块将目标作为导航点；若当前没有发现目标，则根据目标类别和当前语义地图，推测最有可能找到的位置作为导航点。



- Object Goal Navigation using Goal-Oriented Semantic Exploration (CVPR 2020) Devendra Singh Chaplot, Dhiraj Gandhi et al.

# 2 视觉语义导航

## ➤ 前沿：基于建图的视觉语义导航

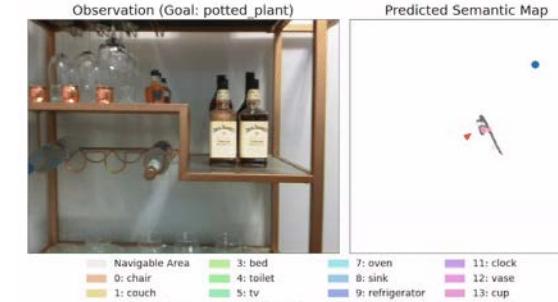
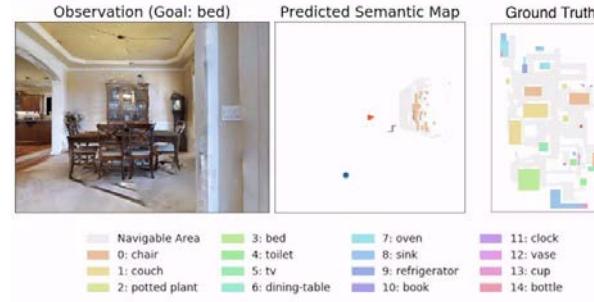
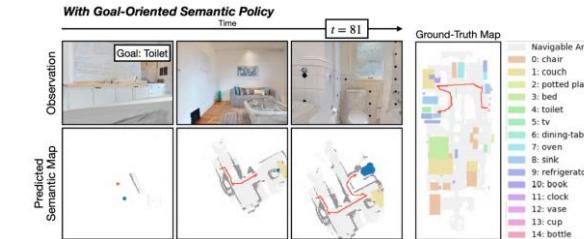
### • 实验结果

实验平台：Habitat

结果：采用语义地图的方法比单纯使用强化学习进行端到端训练的效果在成功率，SPL和DTS（智能体停止时与成功位置的距离）上都好很多。

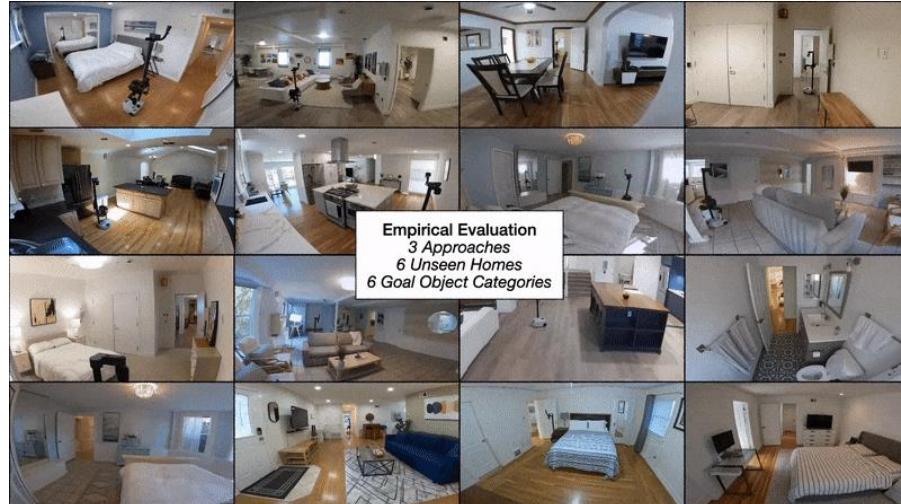
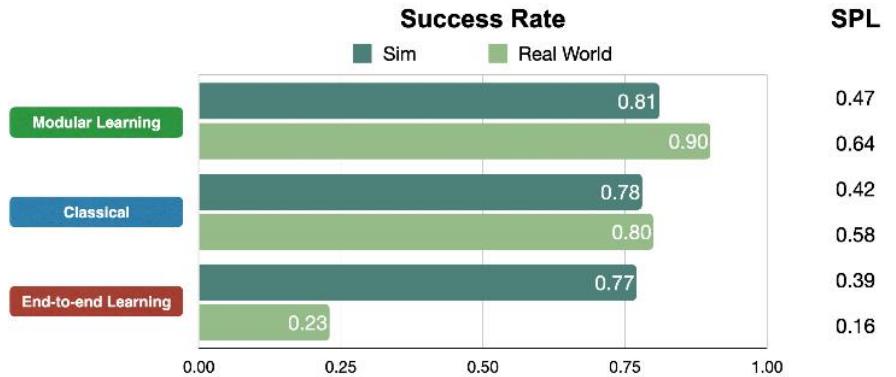
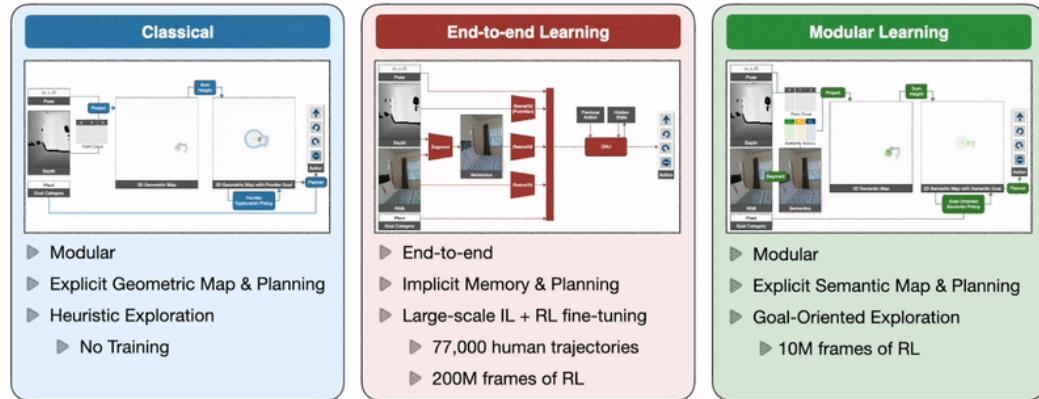
Table 1: Results. Performance of SemExp as compared to the baselines on the Gibson and MP3D datasets.

Method	Gibson			MP3D		
	SPL	Success	DTS (m)	SPL	Success	DTS (m)
Random	0.004	0.004	3.893	0.005	0.005	8.048
RGBD + RL [38]	0.027	0.082	3.310	0.017	0.037	7.654
RGBD + Semantics + RL [31]	0.049	0.159	3.203	0.015	0.031	7.612
Classical Map + FBE [46]	0.124	0.403	2.432	0.117	0.311	7.102
Active Neural SLAM [9]	0.145	0.446	2.275	0.119	0.321	7.056
SemExp	<b>0.199</b>	<b>0.544</b>	<b>1.723</b>	<b>0.144</b>	<b>0.360</b>	<b>6.733</b>



## 2 视觉语义导航

### ➤ 前沿：基于建图的视觉语义导航



## 2 视觉语义导航

### ➤ 前沿

#### Rearrangement

Current state      Goal state

- Weihs, Luca, et al. , Visual room rearrangement. (CVPR 2021).

#### Audio-Visual Navigation

In this space, we're creating smart agents that can respond to real-life situations like the fire alarm going off during a piano lesson.

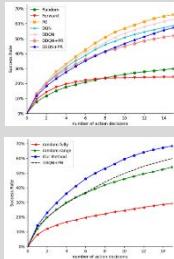
SoundSpaces: 一个声音模拟平台，用于实现两个视觉逼真的3D环境（Replica和Matterport3D）的视听导航

- Chen Changan, et al. , Soundspaces: Audio-visual navigation in 3d environments. ECCV 2020

# 2 视觉语义导航

## ➤ 前沿

### 主动目标发现

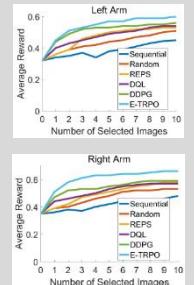


- Active object detection with multi-step action prediction using deep Q-network, IEEE Transactions on Industrial Informatics, 2019

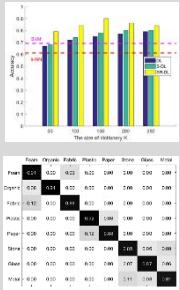
### 主动目标识别



- Extreme trust region policy optimization for active object recognition, IEEE Transactions on Neural Networks and Learning Systems, 2018

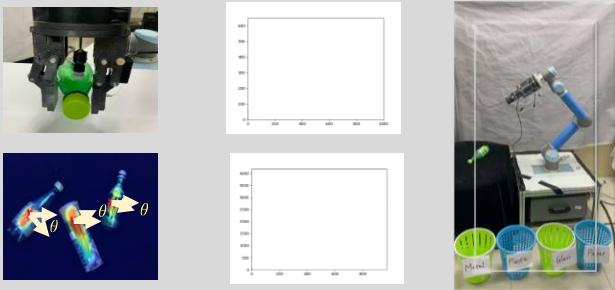


### 主动触觉识别



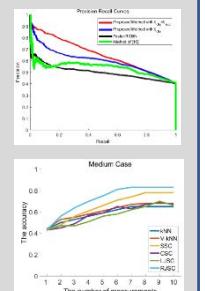
- Material identification using tactile perception: A semantics-regularized dictionary learning method, IEEE/ASME Transactions on Mechatronics, 2018

### 视觉引导的主动触觉感知



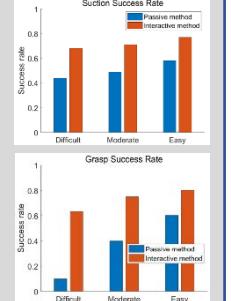
- Visual affordance guided tactile material recognition for waste recycling, IEEE Transactions on Automation Science and Engineering, 2022

### 主动目标定位



- Active object discovery and localization using sound-induced attention, IEEE Transactions on Industrial Informatics, 2021

### 主动交互感知

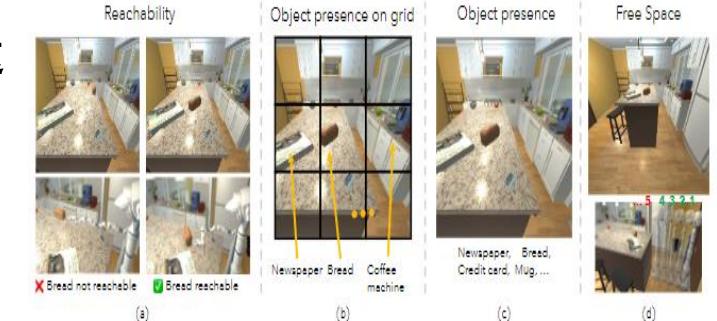


- An interactive perception method for warehouse automation in smart cities, IEEE Transactions on Industrial Informatics, 2021

## 2 视觉语义导航

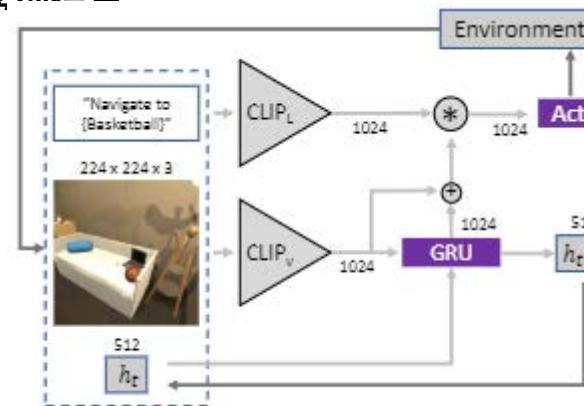
### ➤ EmbClip

- 在未知场景中进行目标导航时，良好的视觉表示可以实现对导航信息有效的编码，这对于导航非常重要
- 基于CLIP的视觉表示在物体存在、物体存在位置、物体可达性和自由空间等语义上实现了良好的编码



#### 模型结构

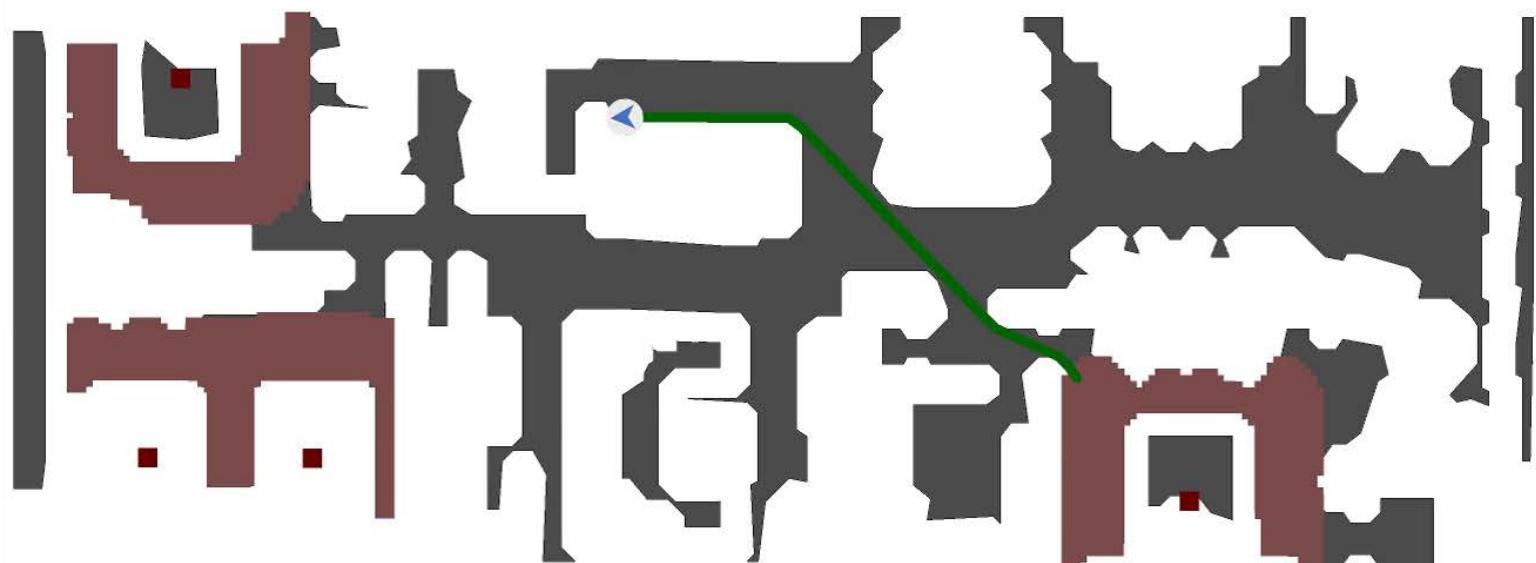
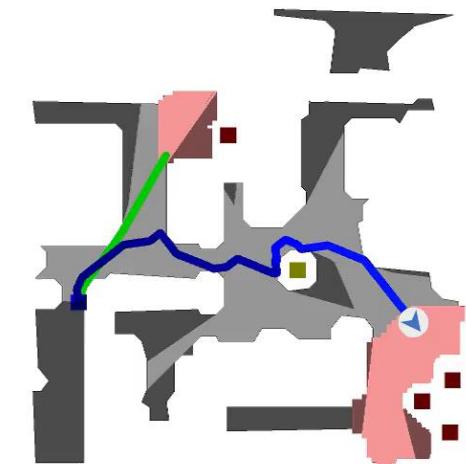
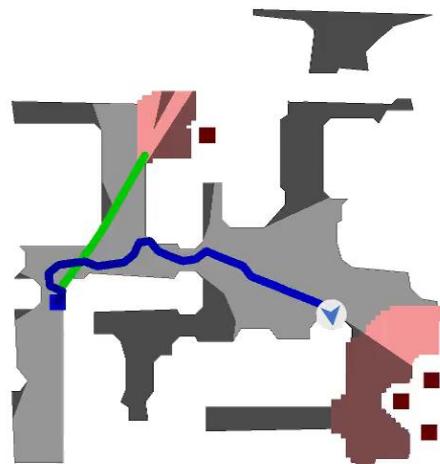
模型将当前观测图像和导航语义分别经过CLIP Vison编码器和 CLIP Language编码器获取特征，特征融合后放入动作选择模型（Act）获取导航动作，另外模型使用GRU处理序列信息



- Simple but effective: Clip embeddings for embodied ai (CVPR 2022) Khandelwal A, Weihs L, Mottaghi R, et al.

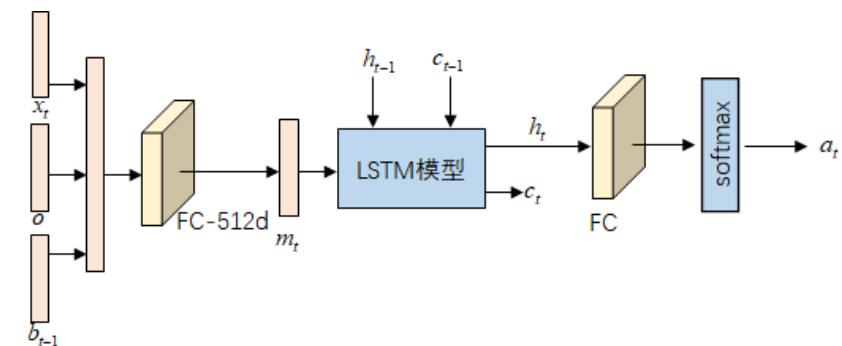
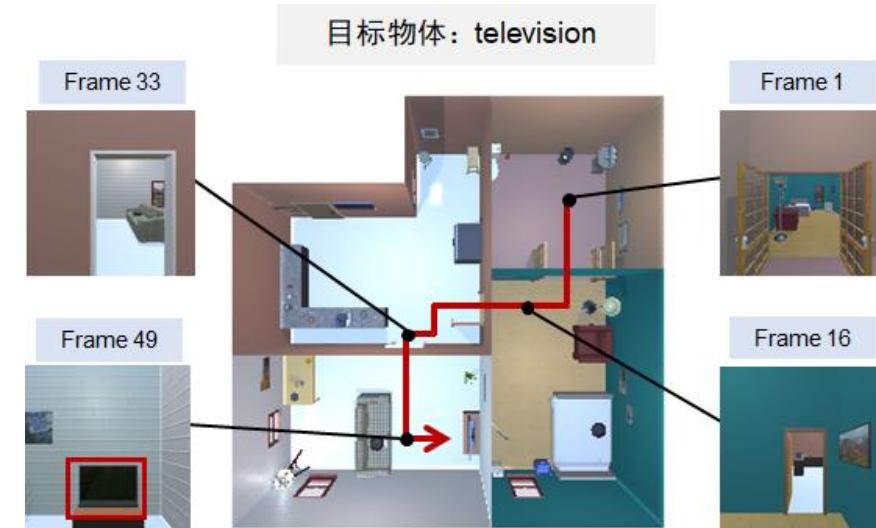
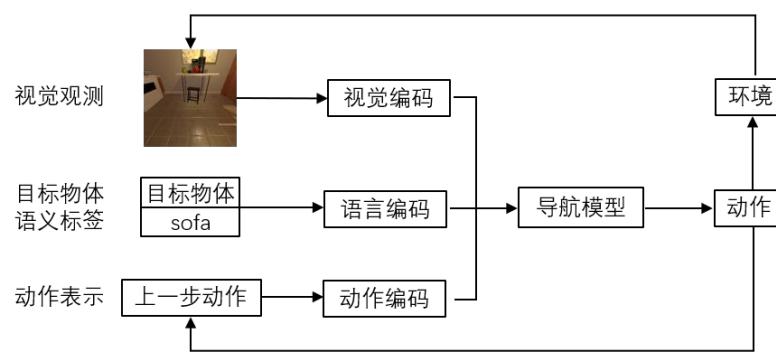
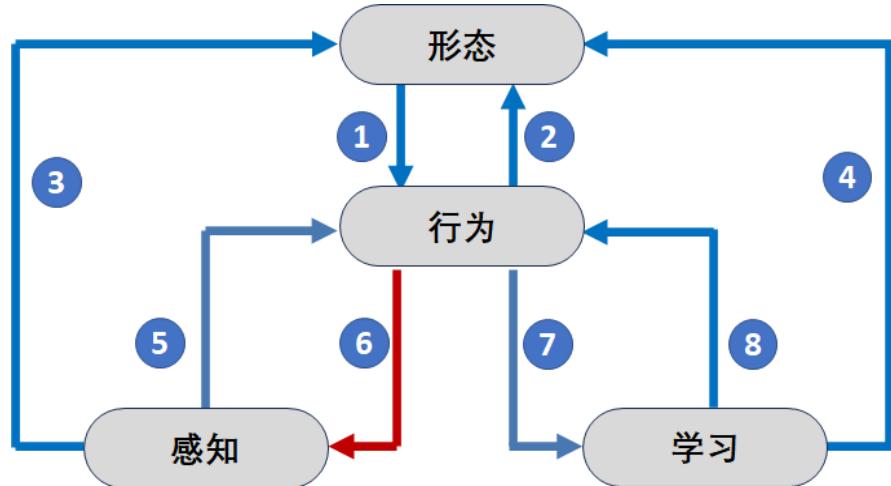
### 3 具身场景描述

#### ➤ 脆弱性



## 2 视觉语义导航

### ➤ 小结



- 场景图片、知识图谱, ... ...

- 
- 从视觉导航到主动感知
  - 视觉语义导航
  - 具身场景描述
  - 具身语言问答

# 3 具身场景描述

## ➤ 图像语言描述

Image  
Captioning



A living room with a couch and a tv.

Dense  
Captioning



a picture on the wall. the bed is white. the lights are on. pink top on the cake. a small plant in a vase. the chairs are on the table. a white shade. a white chair. the lights are on. a white door.

Image  
Paragraphing



A bedroom has a large bed inside it and a large plastic bed on which. a large high window is hanging over the bed bed. the bed is round and has a flower on top of it. there is a small lamp lamp on the table that has a picture hanging on it

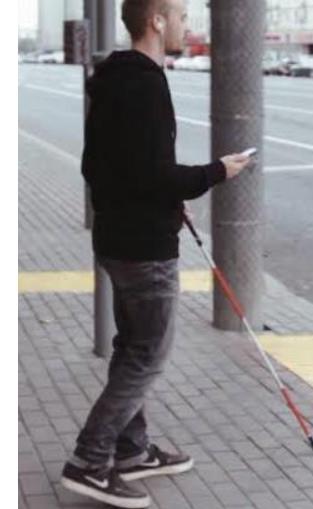
Video  
Captioning



A cat is playing with a dog.

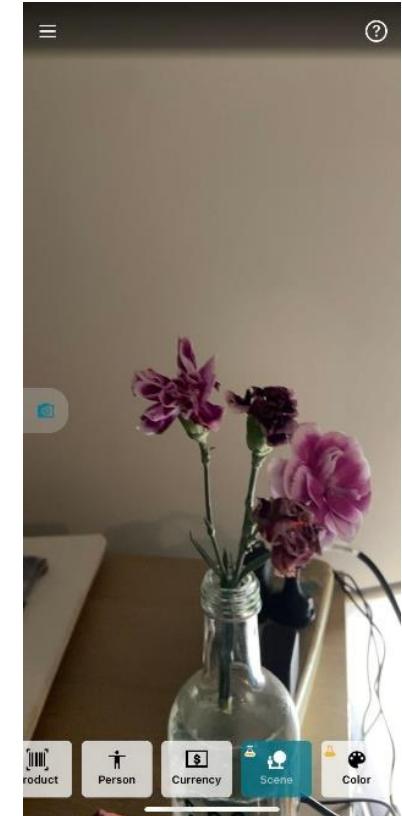
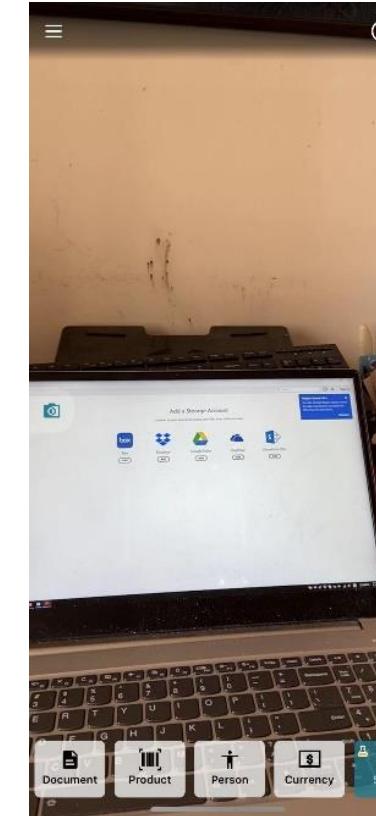
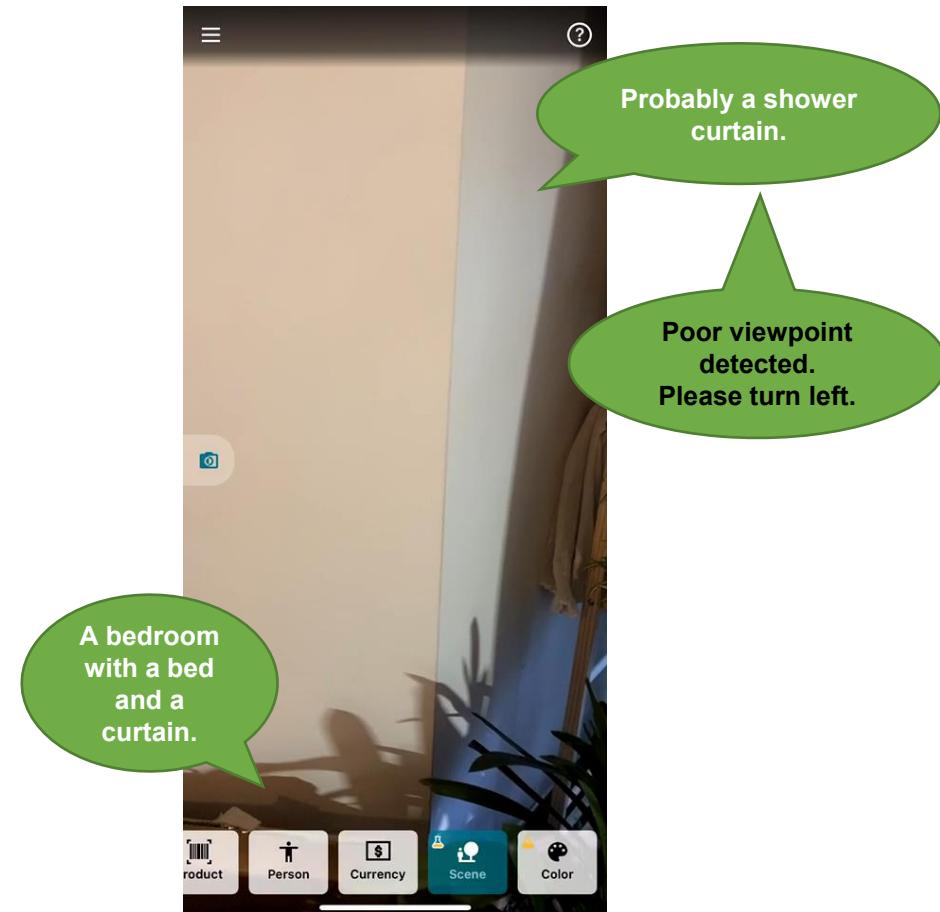
### 3 具身场景描述

#### ➤ 图像语言描述



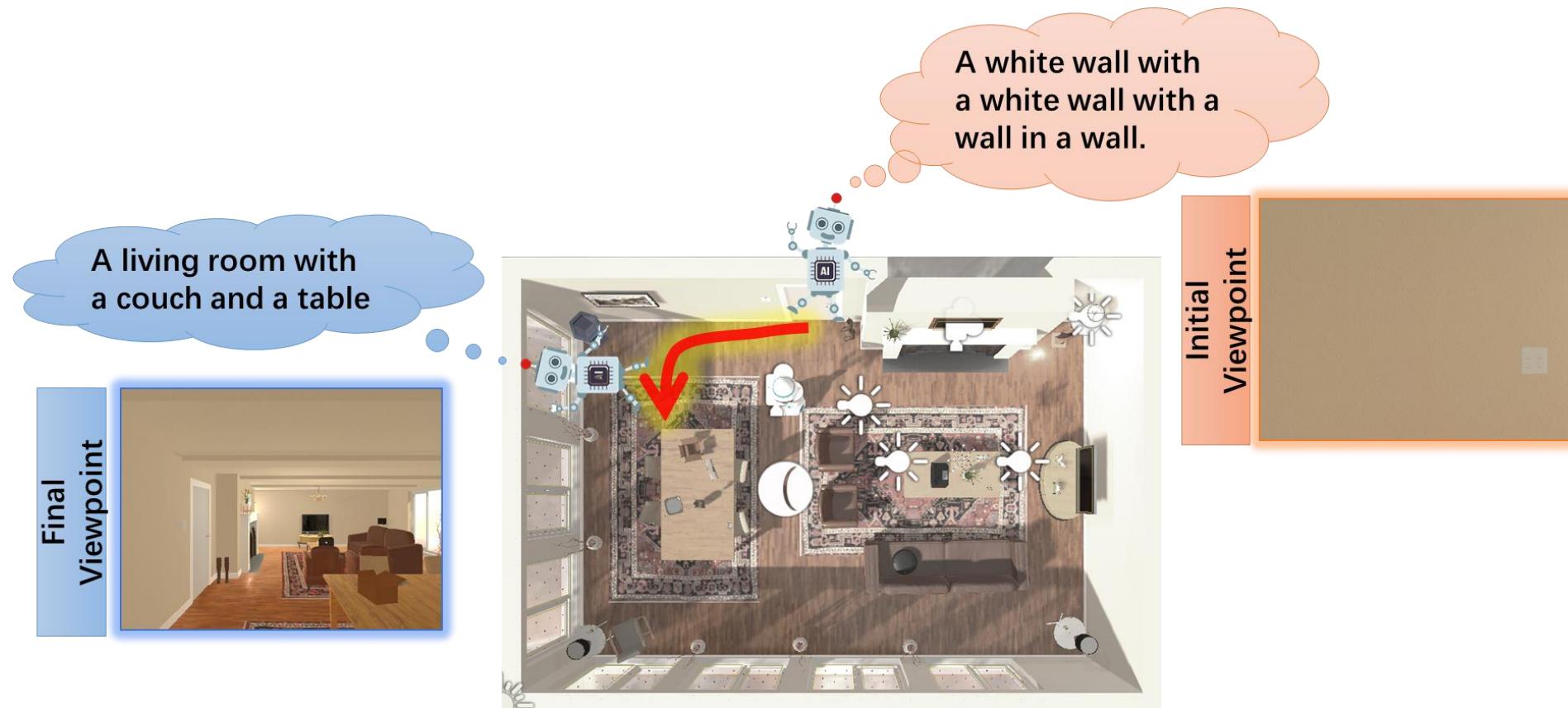
# 3 具身场景描述

## ➤ 图像语言描述



# 3 具身场景描述

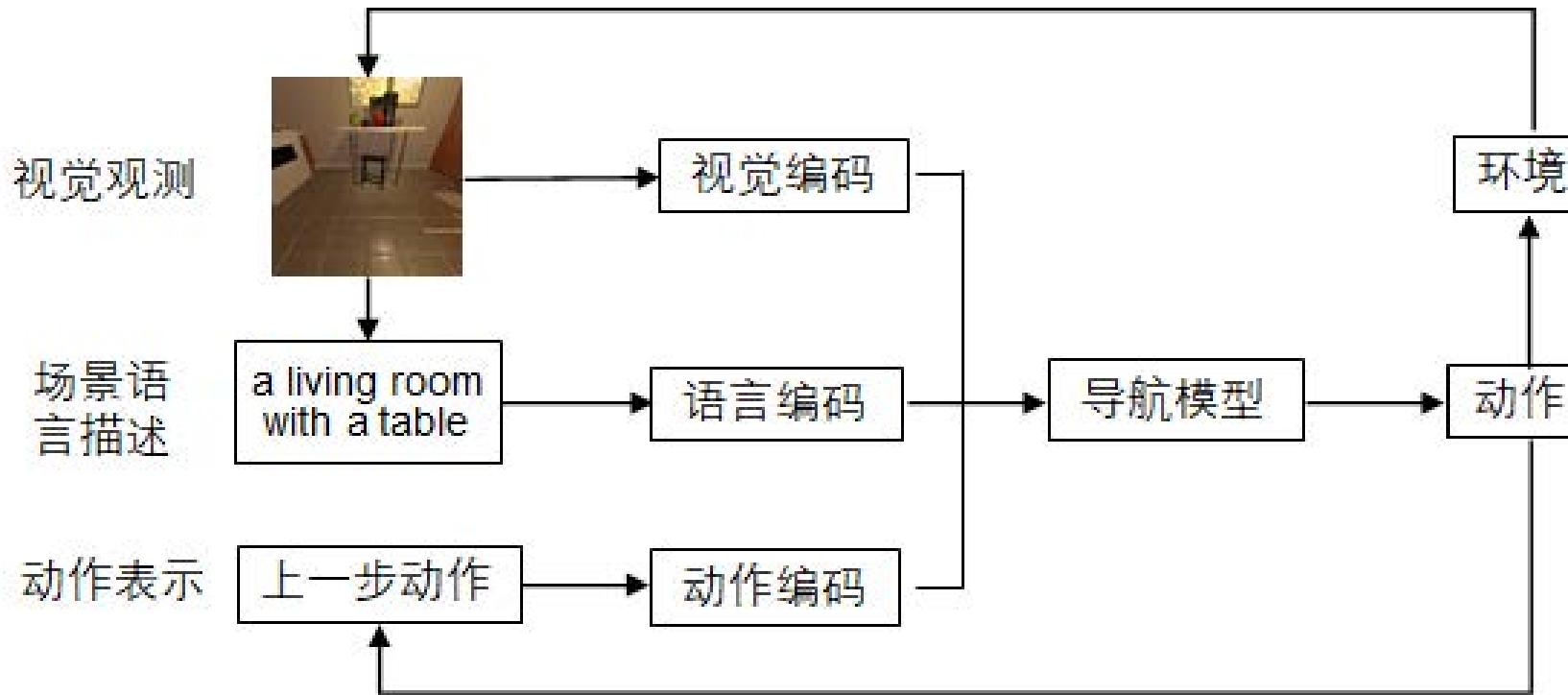
## ➤ 图像语言描述



- Embodied scene description, RSS, 2020

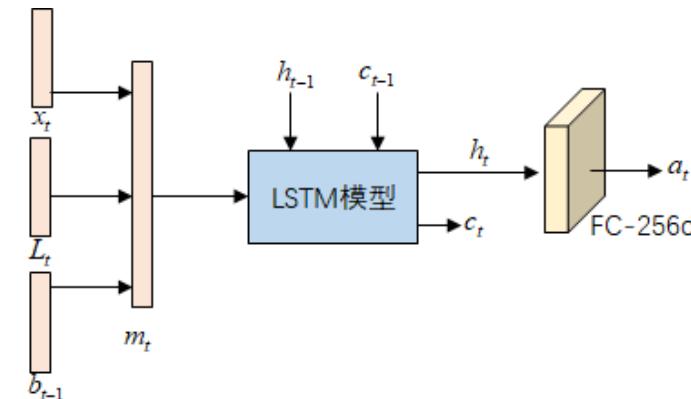
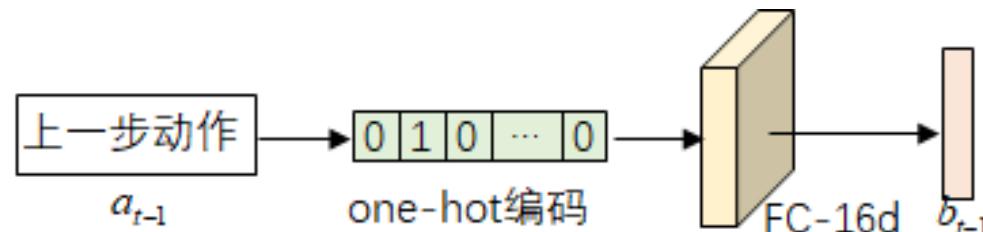
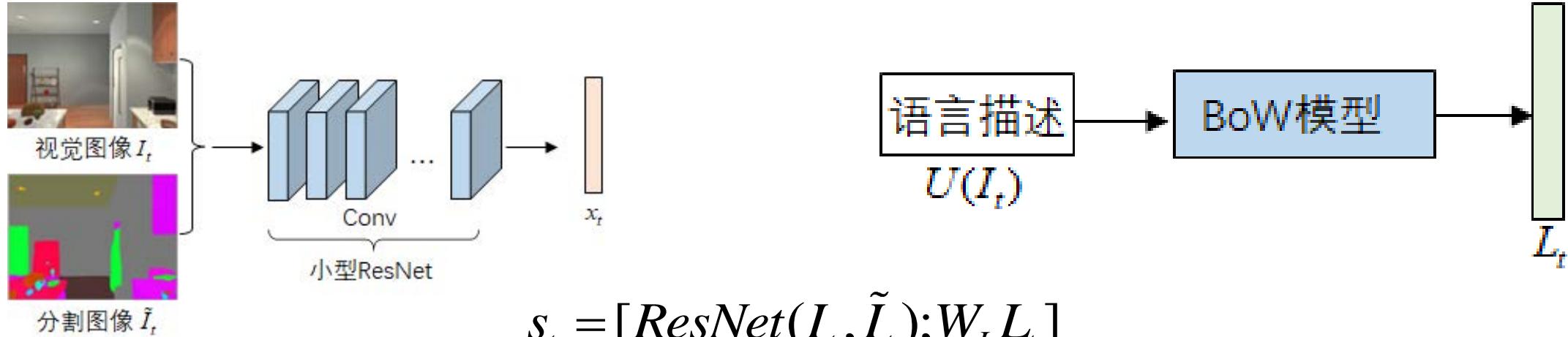
# 3 具身场景描述

## ➤ 算法框架



# 3 具身场景描述

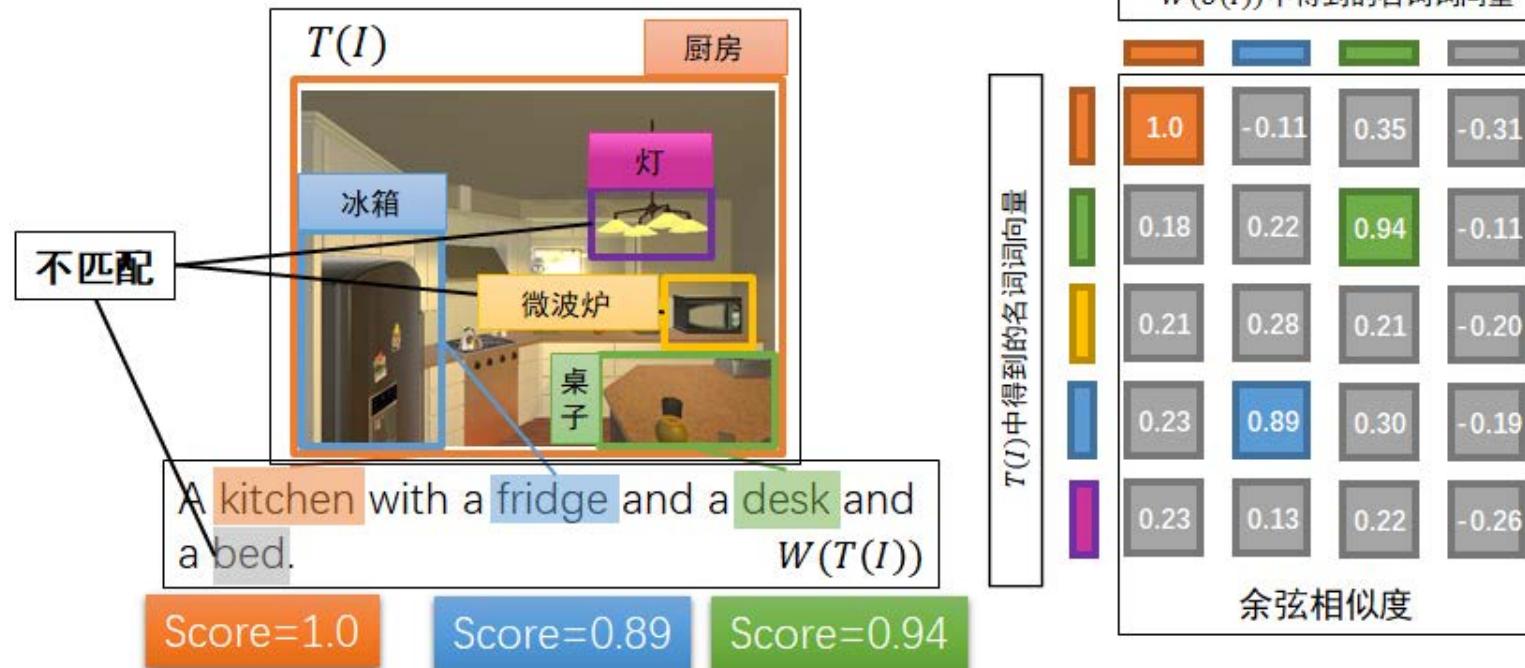
## ➤ 编码



# 3 具身场景描述

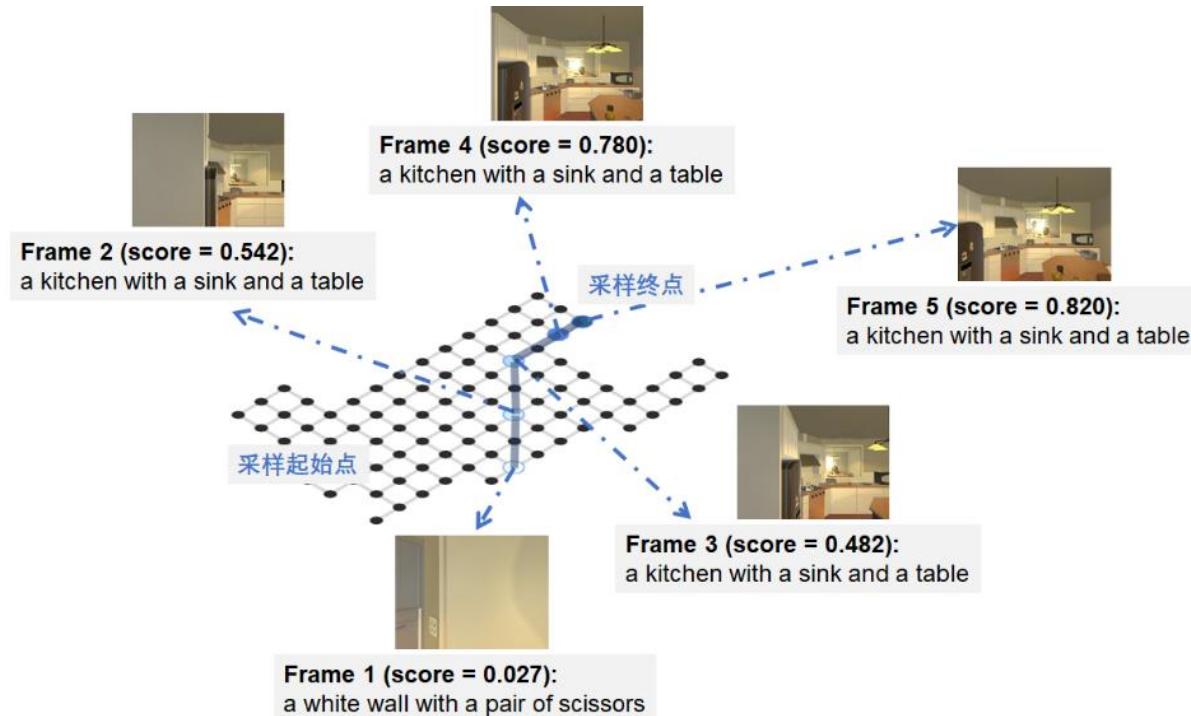
## ➤ 语言描述的评价

$$score(I) = sim(I, U(I)) + \lambda \frac{|T(I)|}{N}$$



### 3 具身场景描述

#### ➤ 模型训练：模仿学习+强化学习



$$L_{\theta} = \sum_{t=1}^T -\log \pi_{\theta}(a_t | s_1, a_1, \dots, s_{t-1}, a_{t-1})$$

# 3 具身场景描述

## ➤ 实验验证：数据集



将智能体所处环境进行网格化。动作空间由平移运动和旋转运动构成。

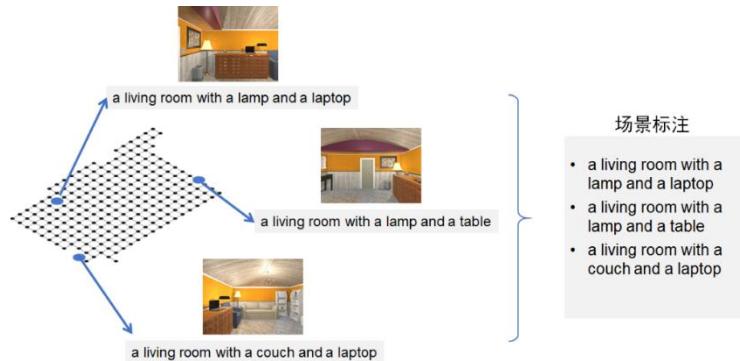
平移运动：包括前、后、左、右、左前、左后、右前、右后，以及停止。每一次平移运动的移动距离固定为网格的边长。

旋转运动：以等间隔进行旋转运动，每一次旋转角度为 $45^\circ$ 。

# 3 具身场景描述

## ➤ 实验验证：数据集

- Number of Steps (NoS): 智能体停止之前走的步数。
- Score of the Last Image (SoL): 触发 Stop 动作的位置的得分。
- Natural Language Metrics: 利用为每个场景自动生成的自然语言标注，自然语言任务的指标就可以用来评估所提出的任务。我们选择了几个指标，包括BLEU (Bilingual Evaluation understudy) , Meteor, 以及ROUGE\_L。其中，BLEU和ROUGE\_L指标都是基于n-gram判断生成语句与参考语句相似度的指标，BLEU指标侧重于准确率，ROUGE\_L指标侧重于召回率，而 Meteor指标同时考虑了准确率与召回率。



$$\text{bleu}_n = \frac{\sum_{c \in \text{candidates}} \sum_{n\text{-gram} \in c} \text{Count}_{clip}(n\text{-gram})}{\sum_{c' \in \text{candidates}} \sum_{n\text{-gram}' \in c'} \text{Count}(n\text{-gram}')}$$

candidate: the cat sat on the mat

reference: the cat is on the mat

$$\text{bleu}_1 = \frac{5}{6} = 0.83$$

$$\text{bleu}_2 = \frac{3}{5} = 0.60$$

$$\text{bleu}_3 = \frac{1}{4} = 0.25$$

$$\text{bleu}_4 = \frac{0}{3} = 0.00$$

- Embodied scene description, RSS, 2020

# 3 具身场景描述

## ➤ 实验验证：方法比较

- 随机探索模型：智能体每一步从其动作空间中随机选择探索动作。
- IL (RGB)：利用模仿学习和 RGB 图像训练动作生成模型。
- IL (RGB+Segm.)：利用模仿学习同时使用 RGB 图像和语义分割图训练动作生成模型。
- IL+RL (RGB+Segm.)：同时以RGB图像和分割图像作为输入，使用模仿学习训练和强化学习微调模型的范式训练动作生成模型。

	NoS( $\downarrow$ )	SoL	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Meteor	ROUGE_L
Random	26.97	0.44	0.57	0.46	0.36	0.26	0.29	0.61
IL(RGB)	38.75	0.79	0.75	0.66	0.56	0.42	0.46	0.79
IL(RGB+Segm)	16.99	0.81	0.78	0.69	0.59	0.43	0.48	0.82
IL+RL(RGB+Segm)	15.27	0.81	0.77	0.69	0.58	0.43	0.48	0.82

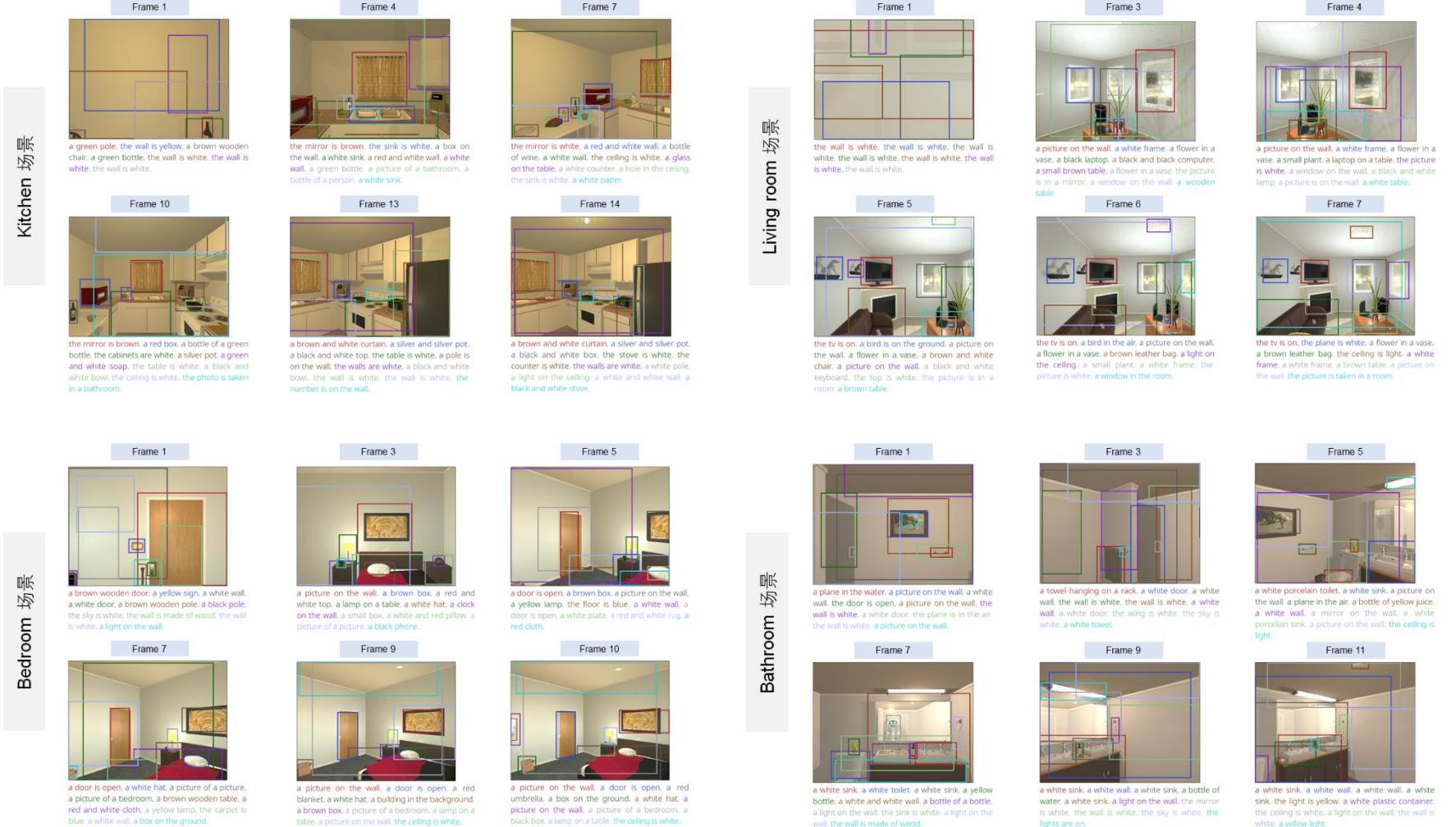
# 3 具身场景描述

## 实验验证



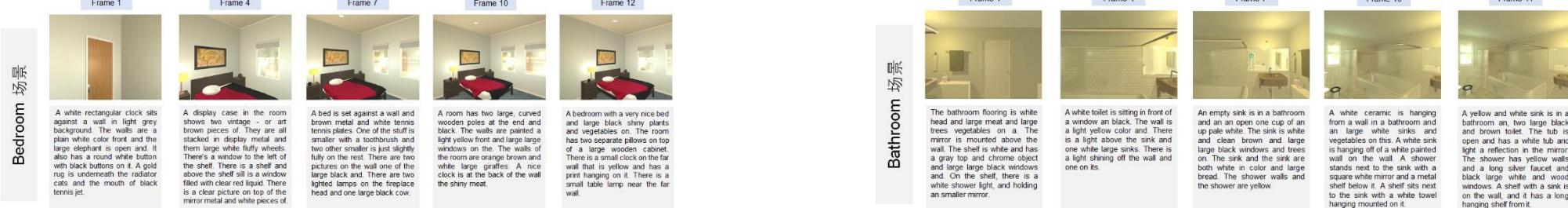
# 3 具身场景描述

## ➤ 实验验证：Dense Captioning



# 3 具身场景描述

## 实验验证: Image Paragraphing



	Nos	SoL	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Meteor	ROUGE_L
Dense captioning	14.41	0.706 3	0.7637	0.6315	0.5201	0.4275	0.2875	0.4815
Image paragraghing	13.28	0.693 6	0.5600	0.3568	0.2224	0.1510	0.1706	0.3601

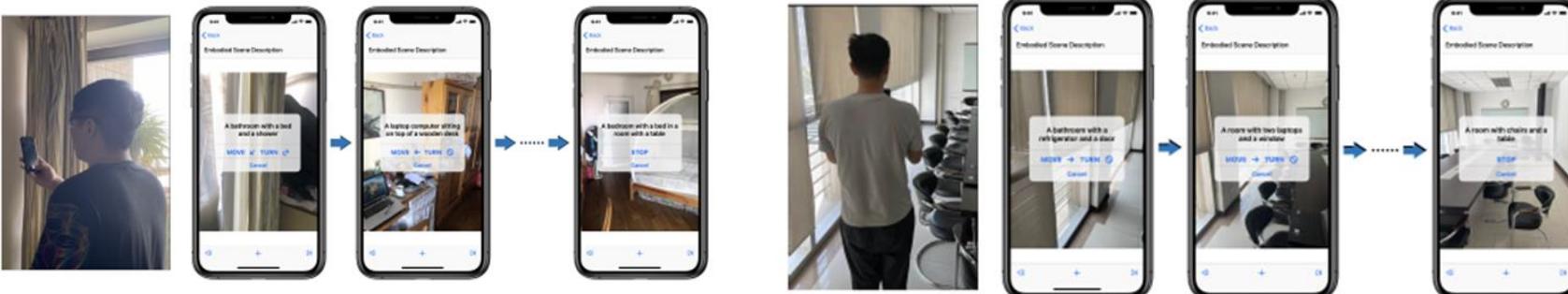
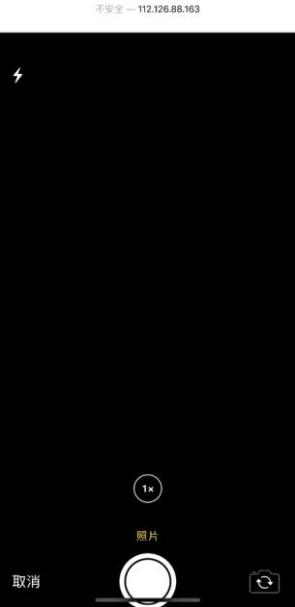
# 3 具身场景描述

## ➤ 物理验证



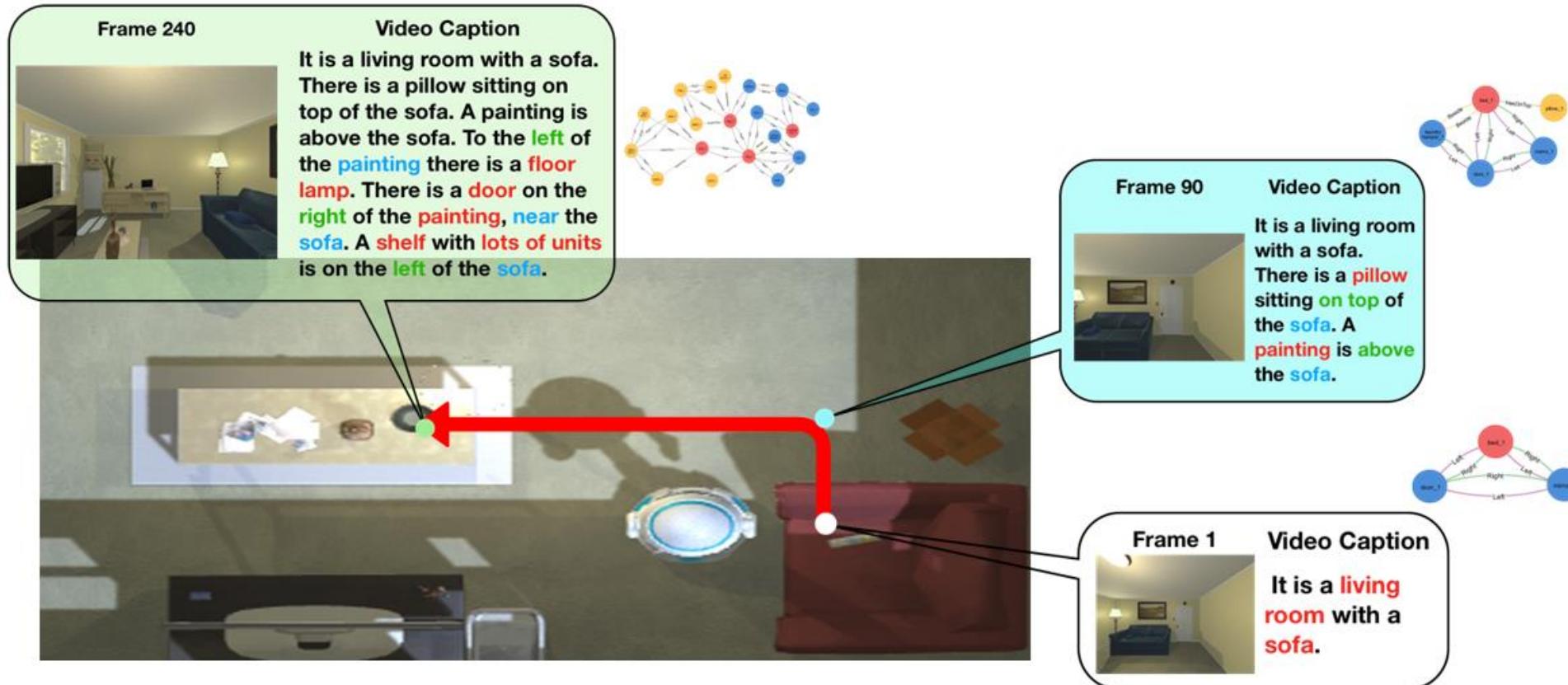
# 3 具身场景描述

## ➤ 交互验证



# 3 具身场景描述

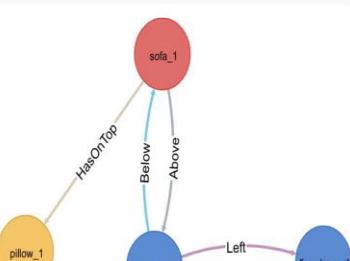
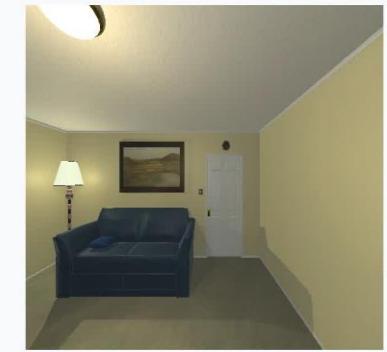
## ➤ 拓展：视频场景描述



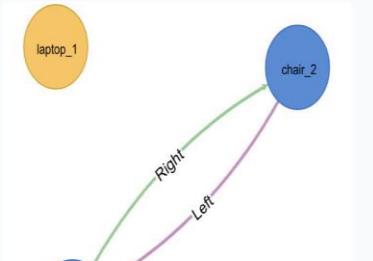
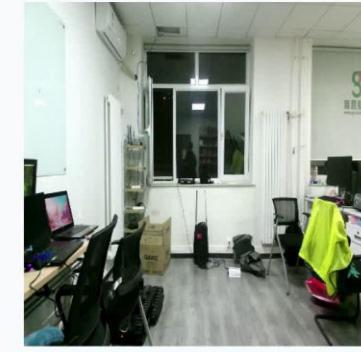
- Robotic indoor scene captioning from streaming video, ICRA, 2021

# 3 具身场景描述

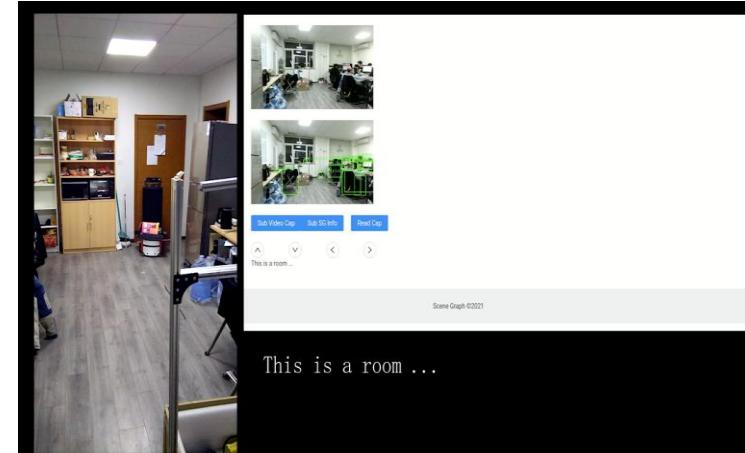
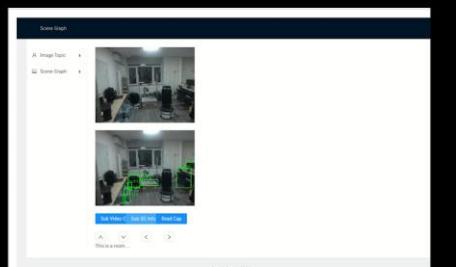
## ➤ 拓展：视频场景描述



It is a living room with a sofa. There is a small table in front of the sofa. There is a pillow sitting on top of the sofa. There is a painting above the sofa. To the left of the painting there is a floor lamp .



It is a living room with two chairs. There is also a laptop in the room .



### 3 具身场景描述

#### ➤ 拓展：红外场景描述



A blurry photo of a person in a rain.



A city street with cars and a traffic light.



A man riding a snowboard down a snow covered slope.

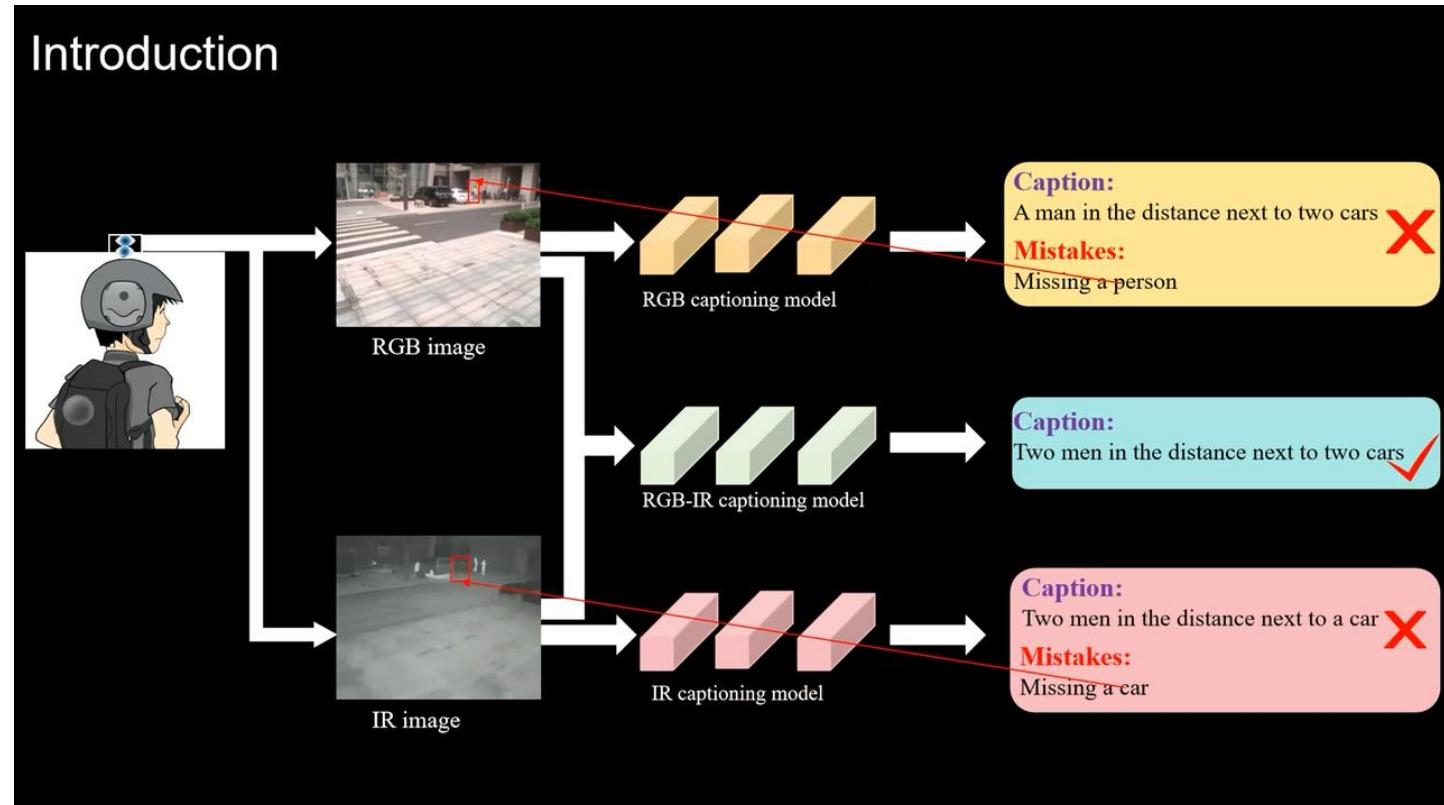
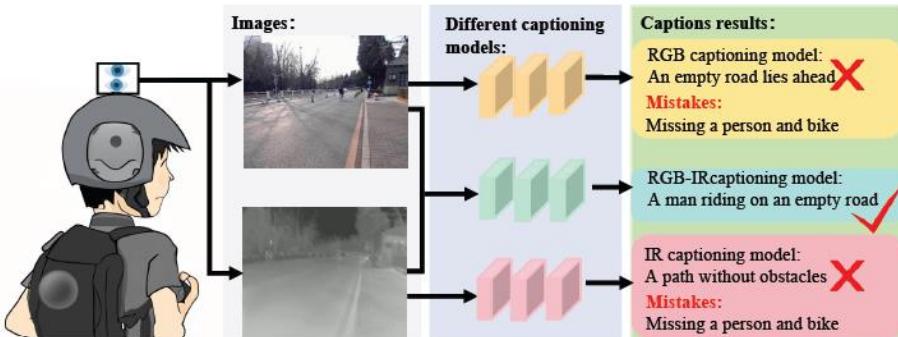


A man and a woman are sitting on a bench.



# 3 具身场景描述

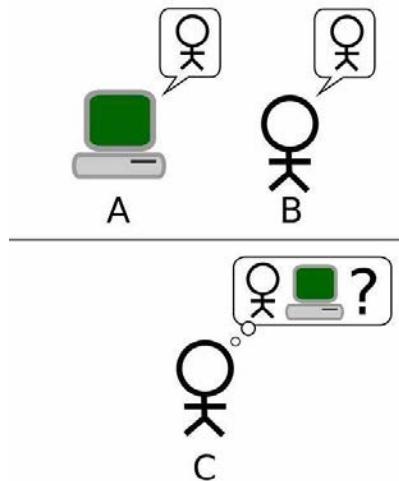
## ➤ 多模态



- 
- 从视觉导航到主动感知
  - 视觉语义导航
  - 具身场景描述
  - 具身语言问答

# 4 具身语言问答

## ➤ 图灵测试



# 4 具身语言问答

## ➤ Visual Question Answering

Visual Question Answer (VQA) 是对视觉图像的自然语言问答，作为视觉理解 (Visual Understanding) 的一个研究方向，连接着视觉和语言，模型需要在理解图像的基础上，根据具体的问题然后做出回答。



Q: How many white objects in this picture ?

A: 9



Q: What color is the chair in front of the wall on the left side of the stacked chairs ?

A: blue



Q: What is the largest white object on the left side of the picture ?

A: printer

## 4 具身语言问答

### ➤ Visual Question Answering



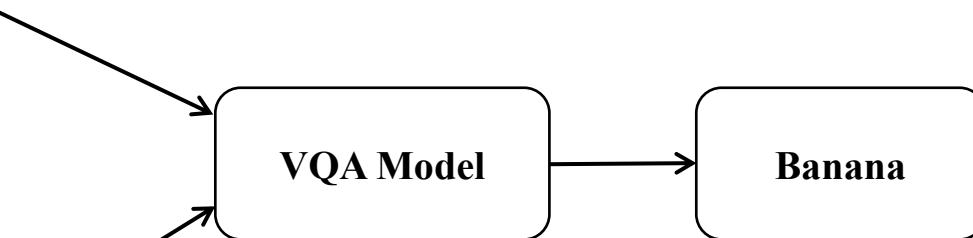
<https://visualqa.org/>

Input Image



Input Question

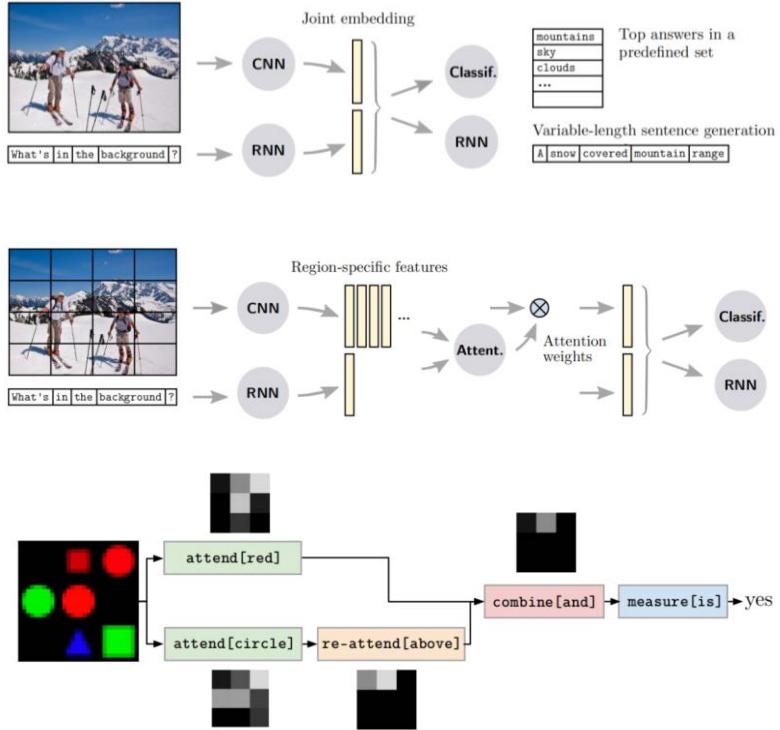
What is the mustache made  
of?



一个VQA系统以一张图片和一个关于这张图片形式自由、开放式的自然语言问题作为输入，以生成的自然语言答案作为输出。

# 4 具身语言问答

## ➤ Visual Question Answering



- Malinowski等人作为提出“开放世界”视觉问答的第一个尝试人之一，在[1]中描述了一种在贝叶斯公式中结合语义文本解析和图像分割的方法
- Tu等人在[2]中对VQA的另一个早期尝试是基于文本和视频的联合解析图。
- 在[3]中，German等人提出了一种自动的“查询生成器”，它对带注释的图像进行训练，然后从任何给定的测试图像中生成一系列二值问题。

主要方法：

- 联合嵌入方法: **Joint embedding approaches**
- 注意力机制: **Attention mechanisms**
- 组合模型: **Compositional Models**
- 知识增强方法: **Models using external knowledge base**

- [1] M. Malinowski and M. Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In Proc. Advances in Neural Inf. Process. Syst., pages 1682–1690, 2014.
- [2] K. Tu, M. Meng, M. W. Lee, T. E. Choe, and S.-C. Zhu. Joint video and text parsing for understanding events and answering queries. IEEE Trans. Multimedia, 21(2):42–70, 2014.
- [3] D. German, S. German, N. Hallonquist, and L. Younes. Visual Turing test for computer vision systems. Proceedings of the National Academy of Sciences, 112(12):3618–3623, 2015.

## 4 具身语言问答

---

### ➤ QA(ChatGPT)+Robot



# 4 具身语言问答

## ➤ 问题描述

Embodied question answering (EQA) 在3D环境中，智能体出现在随机位置，然后问智能体一个问题（汽车是什么颜色的？）为了找到答案，智能体首先探索环境，以第一人称视角收集视觉信息，然后回答问题（橙色）

**Input:**

**First-person view visual Scene**

(no GPS, no map)

**question**

What color is the car?

**Output:**

**Actions: forward, forward,**

**turn left ..... stop**

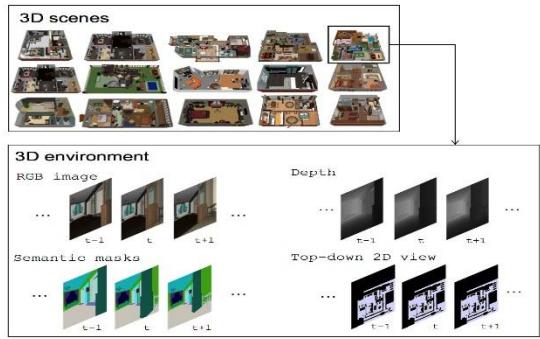
**Answer : orange**



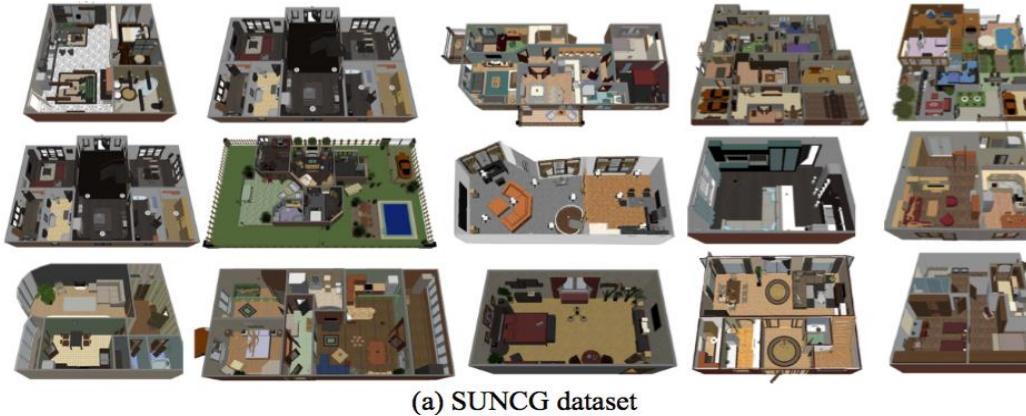
- Embodied question answering. CVPR, 2018

# 4 具身语言问答

## ➤ 数据集



(b)House3D  
environment



(a) SUNCG dataset

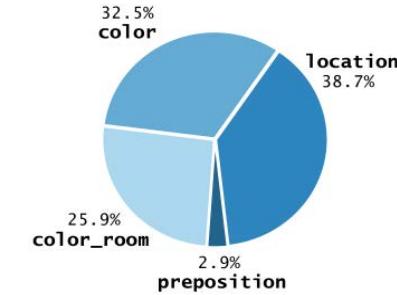
- ◆ containing more than 40K manually created indoor environments.
- ◆ All these scenes are composed of individually labeled 3D objects.
- ◆ allowing us to compute full volumetric ground truth labels

# 4 具身语言问答

## ➤ 数据集

EQAv1 dataset:

- consists of 4 question types—location, color, color\_room, preposition.
- consists of over 5000 question across over 750 environments
- 7 unique room types, 45 unique objects
- Approximately 6 questions are asked per environment on average, 22 at most, and 1 at fewest.



上图显示了数据集分割和问题类型分布

EQAv1	location:	'What room is the <OBJ> located in?'
	color:	'What color is the <OBJ>?'
	color_room:	'What color is the <OBJ> in the <ROOM>?'
	preposition:	'What is <on/above/below/next-to> the <OBJ> in the <ROOM>?'
	existence:	'Is there a <OBJ> in the <ROOM>?'
	logical:	'Is there a(n) <OBJ1> and a(n) <OBJ2> in the <ROOM>?'
	count:	'How many <OBJS> in the <ROOM>?'
	room_count:	'How many <ROOMs> in the house?'
	distance:	'Is the <OBJ1> closer to the <OBJ2> than to the <OBJ3> in the <ROOM>?'

gym	dining room
patio	living room
office	bathroom
lobby	bedroom
garage	elevator
kitchen	balcony

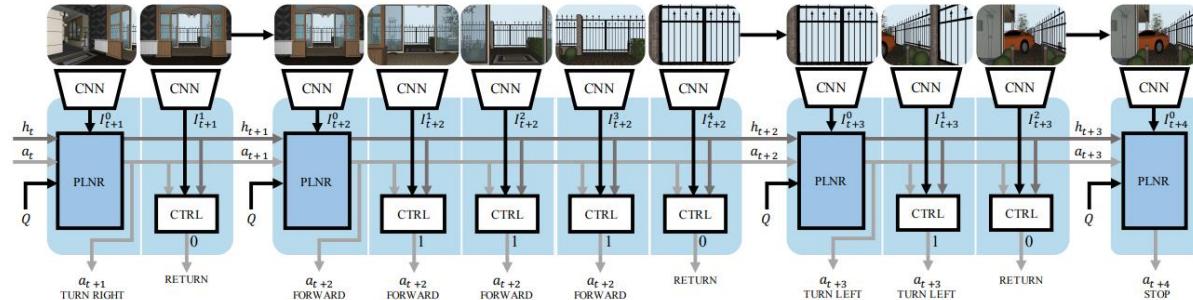
(b) Queryable Rooms

rug	piano	dryer	computer	fireplace	whiteboard	bookshelf	wardrobe	cabinet
pan	toilet	plates	ottoman	fish tank	dishwasher	microwave	water dispenser	
bed	table	mirror	tv stand	stereo set	chessboard	playstation	vacuum cleaner	
cup	xbox	heater	bathtub	shoe rack	range oven	refrigerator	coffee machine	
sink	sofa	kettle	dresser	knife rack	towel rack	loudspeaker	utensil holder	
desk	vase	shower	washer	fruit bowl	television	dressing table	cutting board	

(c) Queryable Objects

# 4 具身语言问答

## ➤ 基本方法

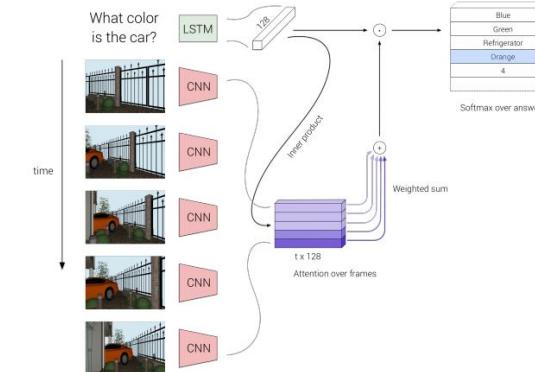


导航模块PACMAN导航器将导航分解为规划器和控制器。

计划者选择动作，控制器执行这些动作的次数不定。这使规划者可以在更短的时间范围内进行操作，从而增强了梯度流。使用第一人称时间的RGB和depth图像作为输入，输出一系列的动作：

$$< a_{t+1}, a_{t+2}, \dots, \text{stop} >$$

问题回答模块以导航帧和问题为条件，计算最后五帧的点积注意力，并将图像特征的注意力加权组合与问题编码相结合来预测答案



# 4 具身语言问答

## ➤ 评价标准

$T$ : 步长, 即初始位置到目标位置的步数 (向前一步最多0.25米)

$d_T$  在 $T$ 的情况下导航终点到目标物之间的距离

$d_\Delta$  在 $T$ 的情况下导航从起始到终止位置朝目标方向移动的距离

$d_{\min}$  在 $T$ 的情况下结束位置到目标的最小距离

为了评估不同难度的导航性能与问答性能指标, 原作者利用不同距离目标位置的步长来划分不同的难度等级即 (10步, 30步, 50步三个难度等级)

MR: 准确答案在智能体回答的基于置信度的答案列表中的排名



Evaluate Navigation  
in Embodied QA



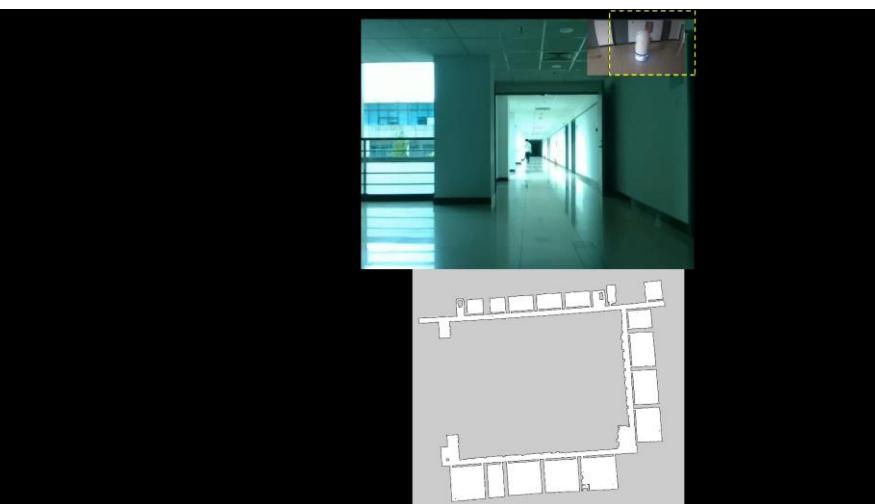
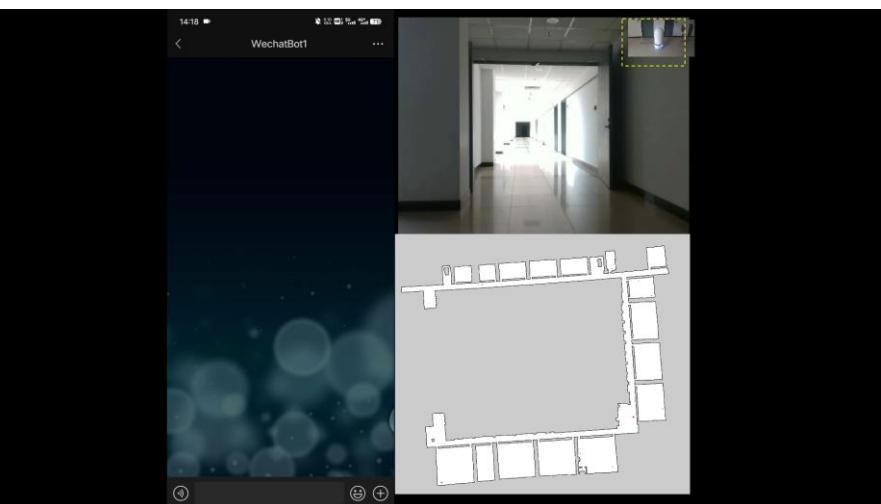
Evaluate QA in Embodied QA

	Navigation																		QA			
	d <sub>T</sub>			d <sub>Δ</sub>			d <sub>min</sub>			%r <sub>T</sub>			%r <sub>Δ</sub>			%stop			MR			
	T <sub>-10</sub>	T <sub>-30</sub>	T <sub>-50</sub>	T <sub>-10</sub>	T <sub>-30</sub>	T <sub>-50</sub>	T <sub>-10</sub>	T <sub>-30</sub>	T <sub>-50</sub>	T <sub>-10</sub>	T <sub>-30</sub>	T <sub>-50</sub>	T <sub>-10</sub>	T <sub>-30</sub>	T <sub>-50</sub>	T <sub>-10</sub>	T <sub>-30</sub>	T <sub>-50</sub>	T <sub>-10</sub>	T <sub>-30</sub>	T <sub>-50</sub>	
Baselines	Reactive	2.09	2.72	3.14	-1.44	-1.09	-0.31	0.29	1.01	1.82	50%	49%	<b>47%</b>	52%	53%	48%	-	-	-	3.18	3.56	3.31
	LSTM	1.75	2.37	2.90	-1.10	-0.74	-0.07	0.34	1.06	2.05	55%	53%	44%	59%	57%	50%	80%	75%	80%	3.35	3.07	3.55
	Reactive+Q	1.58	2.27	2.89	-0.94	-0.63	-0.06	0.31	1.09	1.96	52%	51%	45%	55%	57%	<b>54%</b>	-	-	-	3.17	3.54	3.37
	LSTM+Q	1.13	2.23	2.89	-0.48	-0.59	-0.06	0.28	0.97	1.91	<b>63%</b>	53%	45%	64%	59%	<b>54%</b>	80%	71%	68%	3.11	3.39	3.31
Us	ACT+Q	<b>0.46</b>	<b>1.50</b>	<b>2.74</b>	<b>0.16</b>	<b>0.15</b>	<b>0.12</b>	0.42	1.42	2.63	58%	54%	45%	60%	56%	46%	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>3.09</b>	3.13	3.25
	ACT+Q-RL	1.67	2.19	2.86	-1.05	-0.52	0.01	<b>0.24</b>	<b>0.93</b>	<b>1.94</b>	57%	<b>56%</b>	45%	<b>65%</b>	<b>62%</b>	52%	32%	32%	24%	3.13	<b>2.99</b>	<b>3.22</b>
Oracle	HumanNav*	0.81	0.81	0.81	0.44	1.62	2.85	0.33	0.33	0.33	86%	86%	86%	87%	89%	89%	-	-	-	-	-	-
	ShortestPath+VQA	-	-	-	0.85	2.78	4.86	-	-	-	-	-	-	-	-	-	-	-	-	3.21	3.21	3.21

- Embodied question answering. CVPR, 2018

# 4 具身语言问答

## ➤ 评价标准



# 4 具身语言问答

## ➤ 前沿：层次化导航

### Neural Modular Control for Embodied Question Answering

引入了一个用于具体问题回答的分层策略。给定一个问题(“客厅里的沙发是什么颜色的?”)和观察，我们的主策略预测一个子目标序列——出房间，查找房间[生活]，查找对象[沙发]，回答——然后由专门的子策略执行，导航到目标对象并回答问题(“灰色” )。



给定一个问题，智能体的目标是预测一系列导航子目标，并执行它们，最终找到目标对象并给出正确答案。因此提出NMC作为为每个子目标生成子目标和子策略的主策略。每个子目标被分解为一个任务和一个参数 $\langle g_{subgoal}, g_{arguments} \rangle$ 。有4个可能的任务-出口房间，寻找房间，寻找对象和回答。

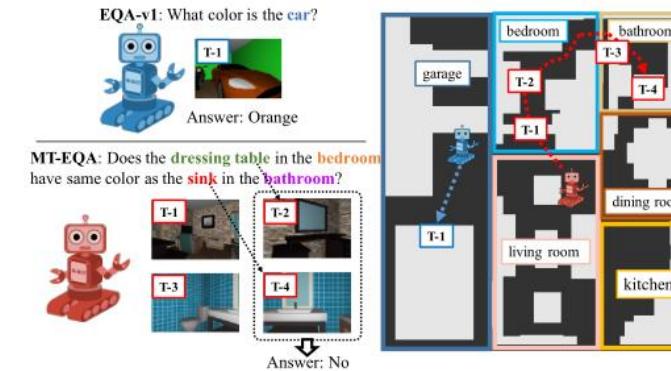
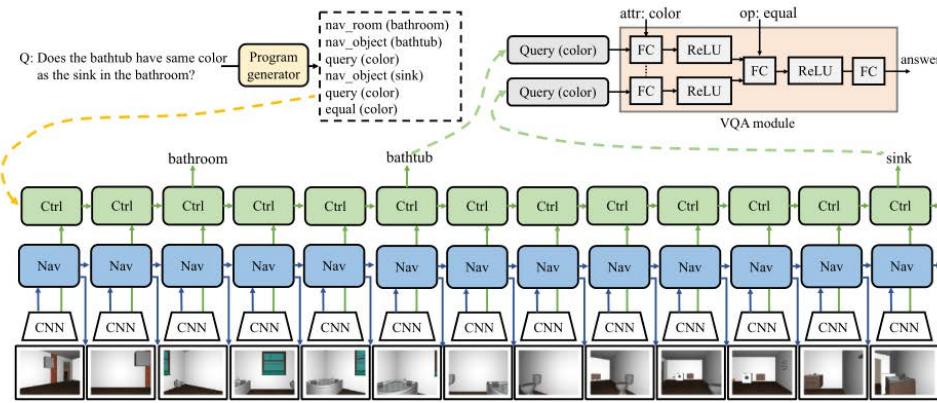
Subgoal	Argument(s)	Description	Success
Exit-room	None	When there is only 1 door in spawn room, or 1 door other than door entered through in an intermediate room; agent is forced to use the remaining door.	Stopping after exiting through the correct door.
Find-room	Room name (gym, kitchen, ...)	When there are multiple doors and the agent has to search and pick the door to the target room.	Stopping after entering target room.
Find-object	Object name (oven, sofa, ...)	When the agent has to search for a specific object in room.	Stopping within 0.75m of the target object.
Answer	None	When the agent has to provide an answer from the answer space.	Generating the correct answer to the question.

# 4 具身语言问答

## ➤ 前沿：多目标问答

Multi-Target Embodied Question Answering

MT-EQA的问题涉及到需要导航的多个目标(如卧室、梳妆台、浴室、水槽)，以及多个目标(如梳妆台、水槽)之间的属性比较。将原来的EQA问题从有限的单目标设置扩展到更具挑战性的多目标设置，要求智能体在回答问题之前进行比较推理



- Multi-target embodied question answering. CVPR, 2019.

# 4 具身语言问答

## ➤ 前沿：引入点云

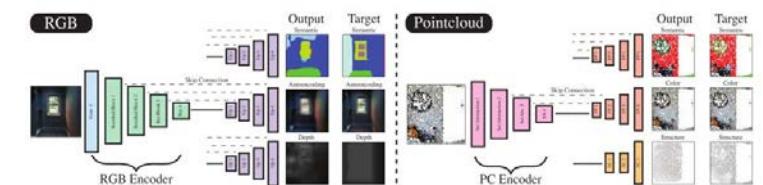
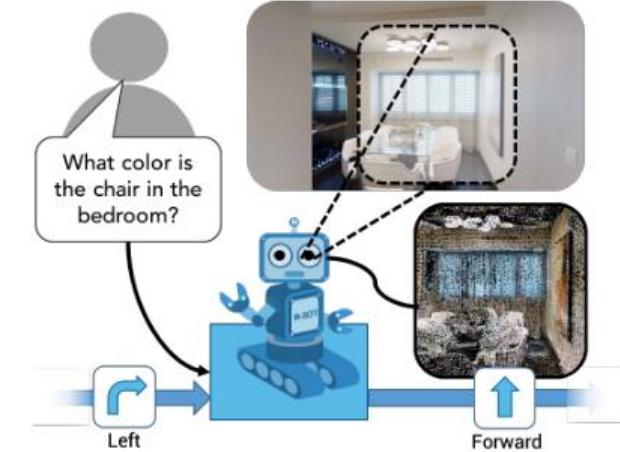
Embodied Question Answering in Photorealistic Environments with Point Cloud Perception

仿真环境：Matterport 3D

数据集：MP3D-EQA dataset

发现：点云为学习避障提供了比RGB图像更丰富的信息

网络：开发了一个端到端的可训练导航模型，能够直接从3D点云学习目标驾驶导航策略。设计了两种网络，一个简单的前馈网络，它将最近的5个视觉观察结果的嵌入串联起来作为输入来预测一个动作。第二个是一个两层的GRU+RNN，将当前观察和之前动作的编码作为输入来预测当前动作。通过网络与输入组合进行对比（网络的输入为RGB, Pointcloud, RGB+Pointcloud），证明点云信息对导航更有效



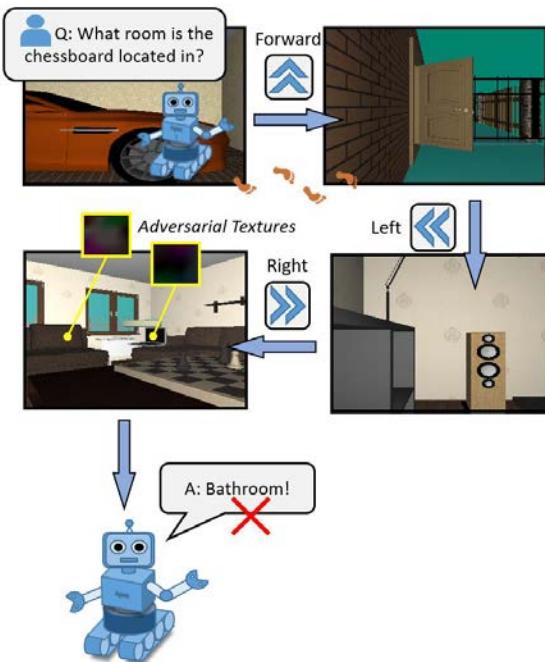
- Embodied question answering in photorealistic environments with point cloud perception. CVPR, 2019.

# 4 具身语言问答

## ➤ 前沿：对抗

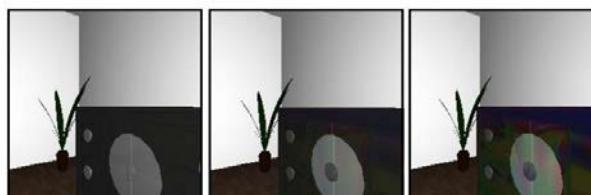
### • 问题背景

- I. 大量的对抗攻击研究都关注于静态场景，而对于具身机器人的鲁棒性仍未有相关研究；
- II. 对抗攻击有助于我们提升对具身机器人的理解并进一步提升其鲁棒性。

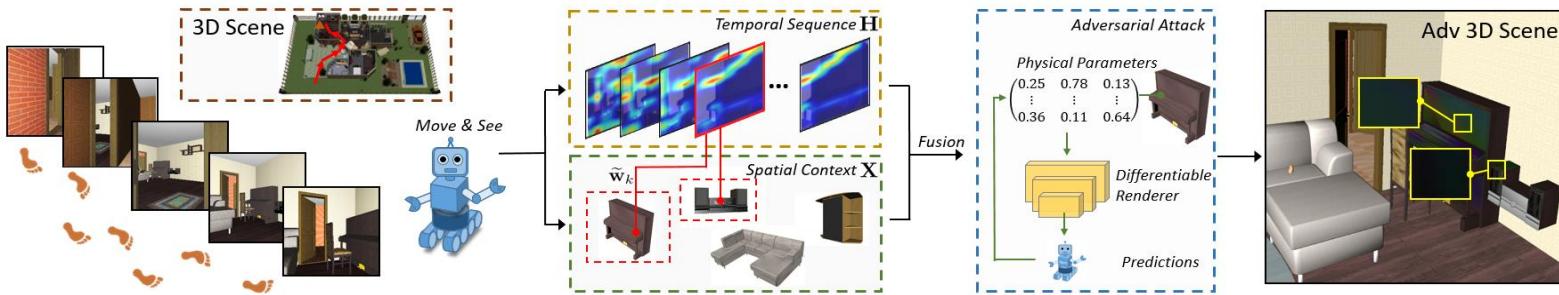


### • 方法与效果

- I. 提出了时空融合的对抗攻击算法，生成人眼不可区分的对抗纹理，攻击具身机器人；
- II. 在EQA和EVR等任务上成功攻击具身机器人。



新浪网：“忽悠”智能机器人，竟然改改物品纹理就成功了  
<http://tech.sina.com.cn/roll/2020-07-14/doc-iivhuipn2860336.shtml>



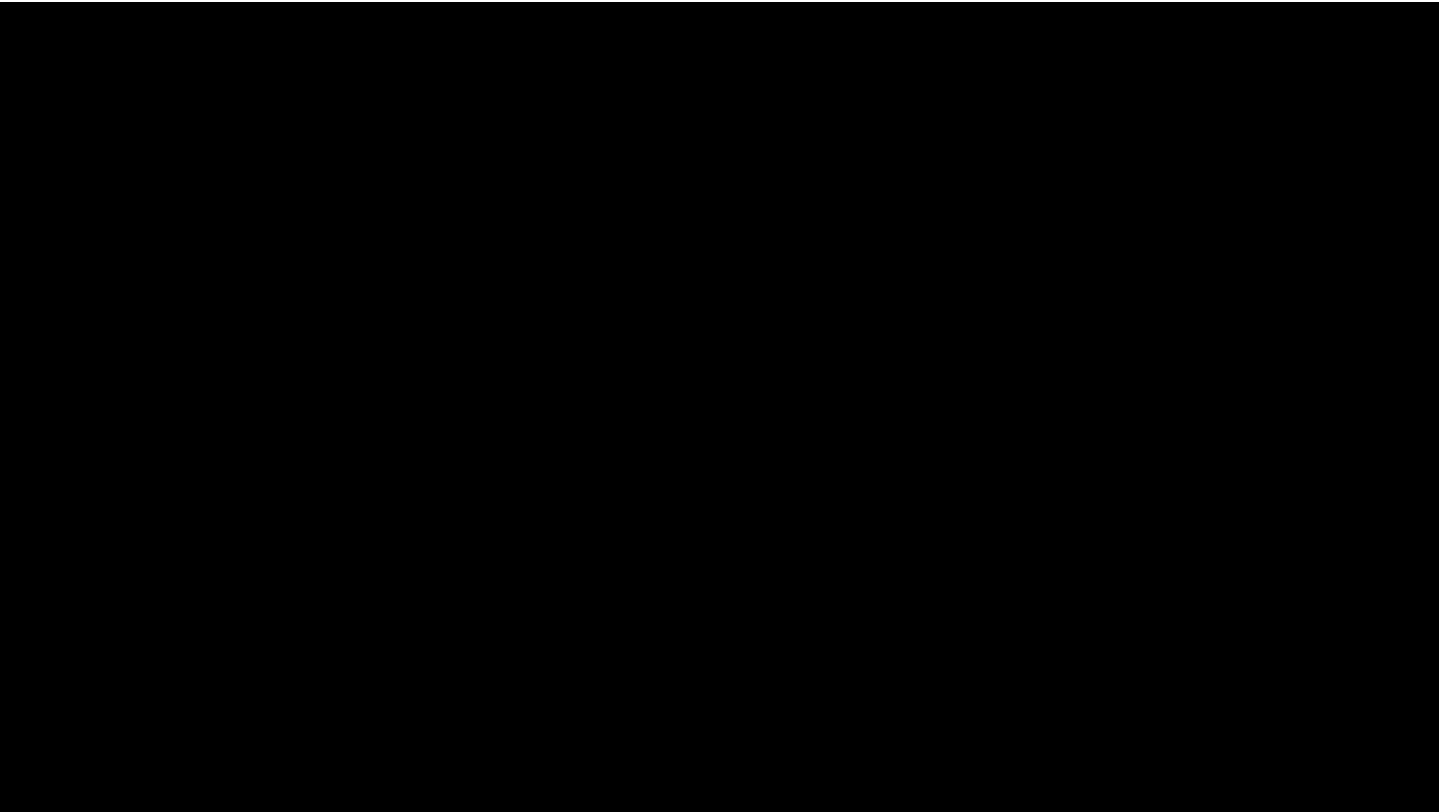
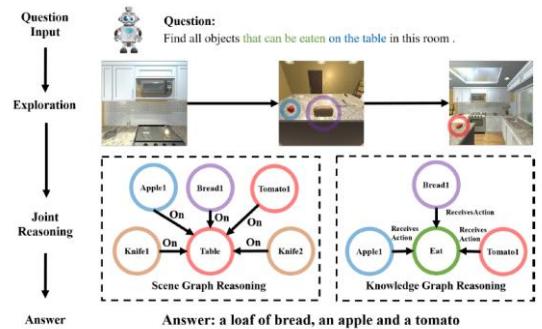
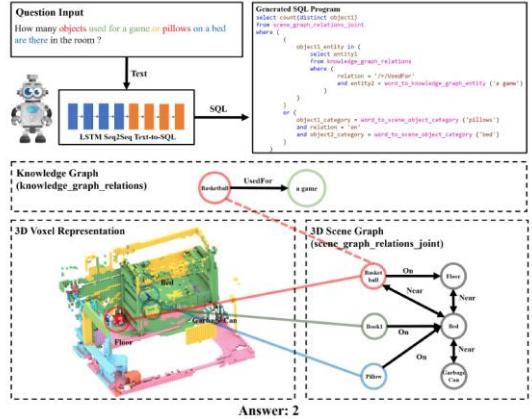
### • 研究意义

- I. 提升对深度学习模型的理解 (texture vs shape)；
- II. 通过对抗训练可以增强模型在高斯噪音等环境中的泛化能力。

	QA		Navigation	
	Adv	Gaussian	Adv	Gaussian
Vanilla	5.67%	22.14%	1.39	1.20
AT	23.56%	38.87%	1.17	1.01
GT	8.49%	32.90%	1.32	1.09

# 4 具身语言问答

## ➤ 前沿：引入知识



- Knowledge-based Embodied Question Answering, T-PAMI, 2023

# 4 具身语言问答

## ➤ 前沿：交互问答

- 交互式问题回答(IQA)，即回答问题的任务，需要一个智能体与动态视觉环境交互。IQA向智能体提出了一个场景和一个问题，比如：“冰箱里有苹果吗？”智能体必须在场景中导航，获得对场景元素的视觉理解，与对象互动(例如打开冰箱)，并根据问题计划一系列行动。

- 仿真环境：AI2thor
- 数据集：IQUAD V1

- 与Embodied QA区别：

- Embodied QA要求智能体导航到目标附近，可以“看到”目标，通过“看到的”像素进行回答
- IQA要求智能体导航到目标附近，并且对场景中的对象进行交互，最终“看到”目标进行回答

### IQA: Visual Question Answering in Interactive Environments

Daniel Gordon<sup>2</sup> Aniruddha Kembhavi<sup>1</sup> Mohammad Rastegari<sup>1</sup>  
Joseph Redmon<sup>2</sup> Dieter Fox<sup>2</sup> Ali Farhadi<sup>1,2</sup>

<sup>1</sup>Allen Institute for Artificial Intelligence

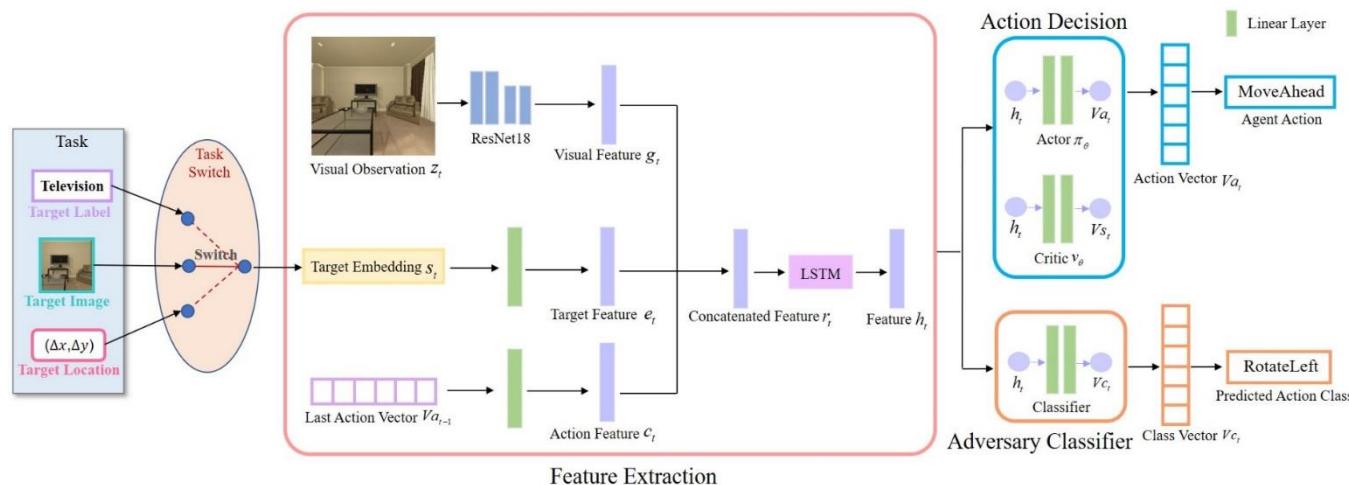
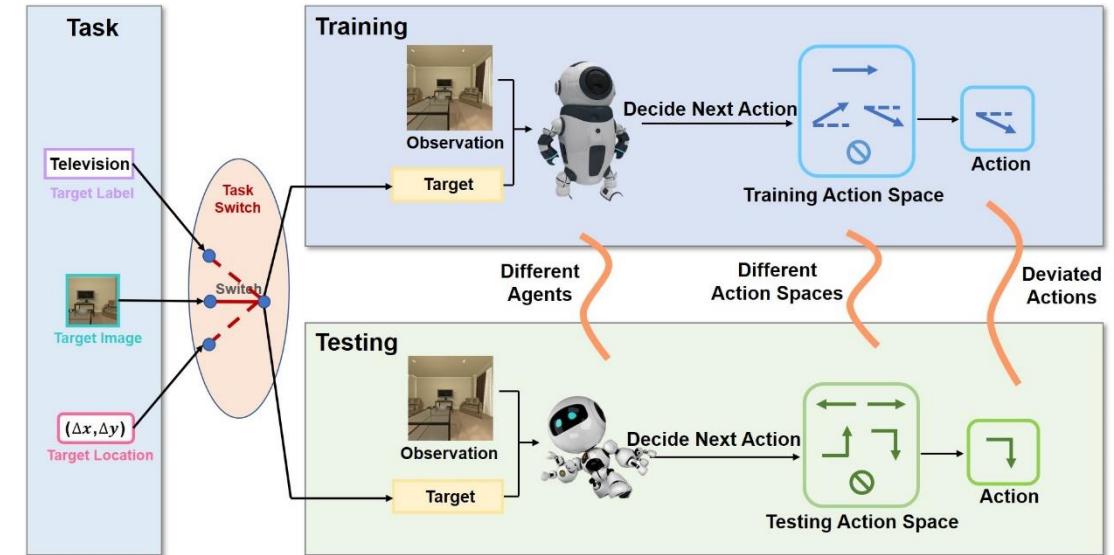
<sup>2</sup>Paul G. Allen School of Computer Science, University of Washington



# 4 具身语言问答

## ➤ 前沿：具身失配

- 对抗训练 —— 学习鲁棒的状态特征表示
  - 特征提取：提取视觉、目标以及历史动作特征，融合后得到状态特征
  - 动作决策：Actor-Critic结构生成下一步动作
  - 对抗学习：预测产生当前状态的前一步动作，希望对抗学习层无法分辨前一步动作类型
- 适应训练 —— 策略迁移至新的动作空间
  - 更新Actor网络维数，使用较少的样本训练适应于新测试动作空间的策略
- 该模块可插入不同的导航决策模型中



# 4 具身语言问答

## ➤ 前沿：具身失配

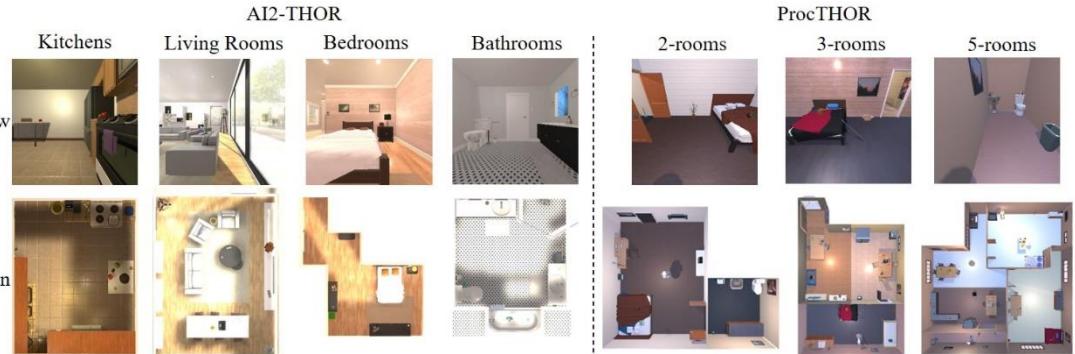
- 实验数据集
  - AI2-THOR
  - Sub-Proc (ProcTHOR-10K子集)
- 测试动作空间 (5种设置)

$$A_{test} = A_{train} \quad A_{test} \subseteq A_{train} \quad A_{test} \supseteq A_{train} \quad A_{test} \cap A_{train} \neq \emptyset \quad A_{test} \setminus A_{train} = \emptyset$$

- 评价指标

- 成功率  $SR = \frac{1}{N} \sum_{i=1}^N R_i$

- 路径加权成功率  $SPL = \frac{1}{N} \sum_{i=1}^N R_i \frac{L_i}{\max(L_i, G_i)}$



QUANTITATIVE RESULTS OF THE PROPOSED MODEL MODEL IN TASK VSN-Label in AI2-THOR AND SUB-PROC											
Dataset	Methods	$A_{test} = A_{train}$ SR(%) SPL(%)	$A_{test} \subseteq A_{train}$ SR(%) SPL(%)	$A_{train} \subseteq A_{test}$ SR(%) SPL(%)	$A_{train} \cap A_{test} \neq \emptyset$ SR(%) SPL(%)	$A_{train} \cap A_{test} = \emptyset$ SR(%) SPL(%)					
AI2-THOR	Random	6.70	3.21	6.60	3.25	6.30	3.11	7.00	3.41	6.50	2.91
	A3C	39.82	12.76	29.60	9.99	9.16	28.30	9.55	28.00	9.67	—
	ScenePriors	44.40	12.75	31.00	12.73	33.00	14.01	33.70	13.10	34.00	11.70
	SAVN	47.50	13.62	39.20	10.63	38.80	11.50	37.60	10.56	37.60	12.80
	SemMap	53.20	18.65	45.60	14.23	44.10	13.69	43.30	13.37	43.50	13.06
	Noise-A3C	40.90	13.95	31.70	11.69	32.40	11.23	30.90	10.88	30.60	10.75
	Noise-ScenePriors	47.30	19.42	37.40	14.45	37.10	14.87	36.70	13.33	36.40	13.52
	Noise-SAVN	50.00	20.12	41.30	15.78	41.30	15.26	40.90	15.19	42.20	14.63
	Noise-SemMap	56.50	20.12	48.10	15.78	47.80	15.26	47.60	15.19	47.20	14.63
	EMAL-A3C	49.00	15.56	46.20	13.51	46.40	11.55	44.90	13.06	44.90	11.31
Sub-Proc	EMAL-ScenePriors	57.70	31.83	55.90	15.91	56.20	15.03	55.30	18.69	54.90	17.55
	EMAL-SAVN	60.30	16.33	58.20	15.72	58.40	13.93	58.10	13.61	57.60	12.88
	EMAL-SemMap	<b>66.70</b>	25.39	<b>64.30</b>	<b>20.08</b>	<b>64.10</b>	<b>18.29</b>	<b>63.90</b>	18.05	<b>63.60</b>	<b>17.63</b>

## Qualitative Performance

Scene: Kitchen Object: CoffeeMachine

Scene: Living room Object: Television

$$A_{test} = A_{train}$$

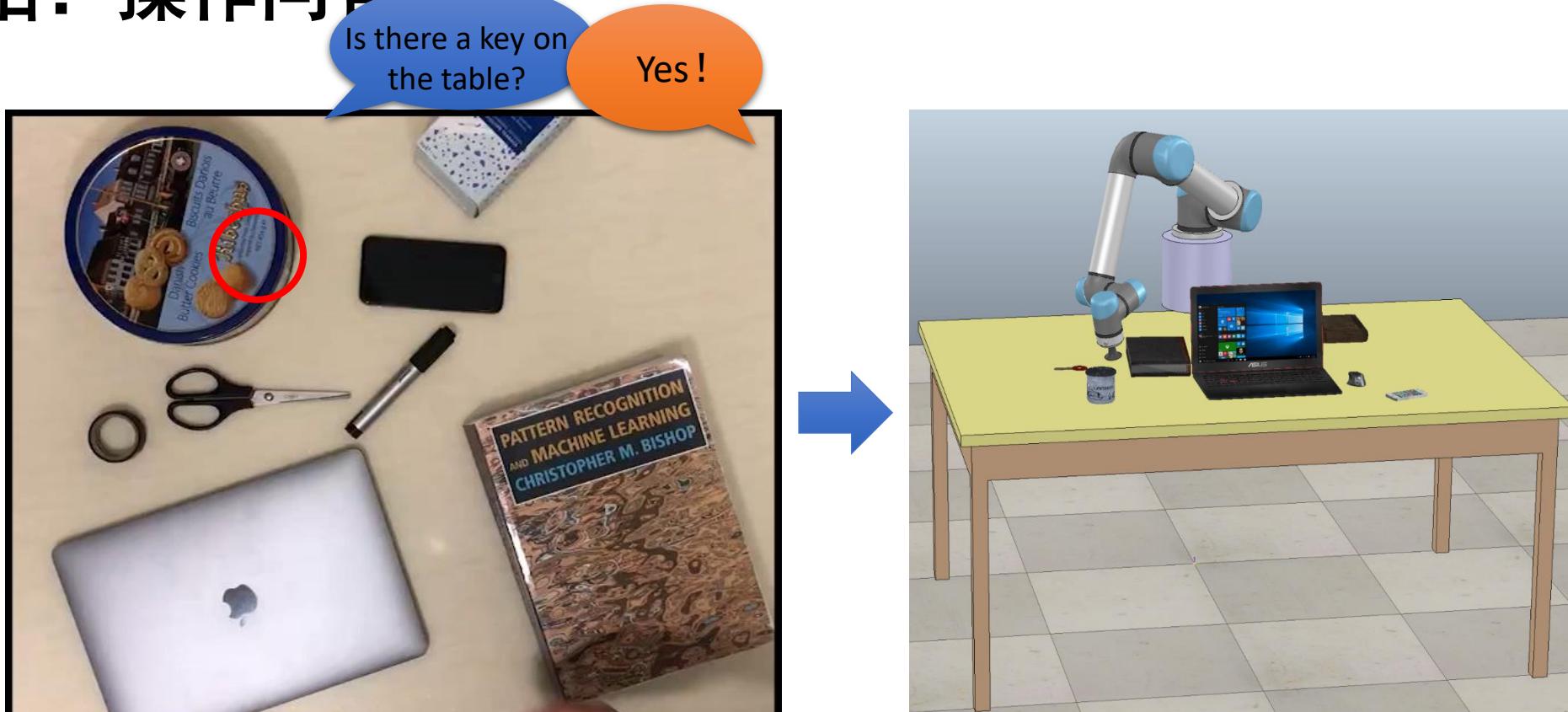
Scene: Bedroom Object: Lamp

Scene: Bathroom Object: Sink



# 4 具身语言问答

## ➤ 前沿：操作问答



How about robots?  
Can robots deal with this kind of QA like humans?

- MQA: Answering the question via robotic manipulation, RSS, 2021

# 4 具身语言问答

## ➤ 前沿：操作问答



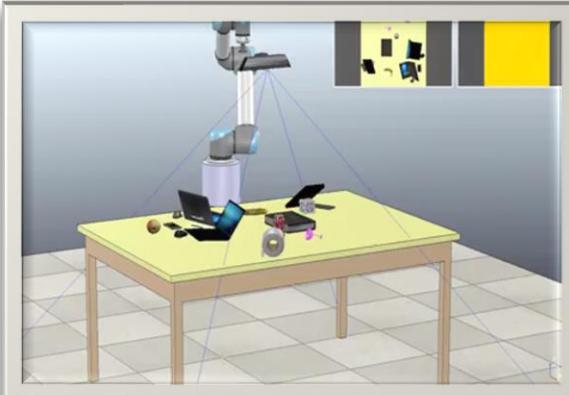
Visual Question Answering (VQA)



Embodied Question Answering (EQA)



Interactive Question Answering (IQA)



Manipulation Question Answering (MQA)

	VQA	EQA	IQA	MQA
Understanding	✓	✓	✓	✓
Exploration	—	✓	✓	✓
Interaction	—	—	✓	✓
Manipulation	—	—	—	✓

A comprehensive comparison of VQA, EQA, IQA and the proposed MQA tasks is illustrated in the above table.

- MQA: Answering the question via robotic manipulation, RSS, 2021

# 4 具身语言问答

## ➤ 前沿：操作问答

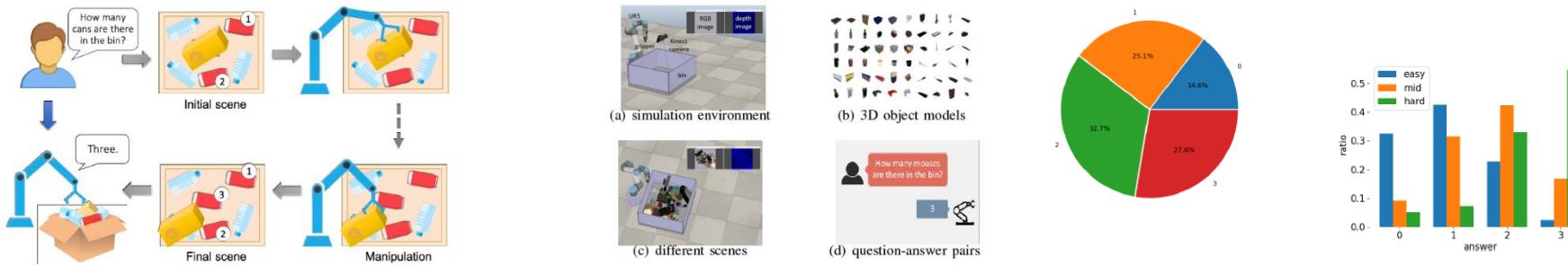


Fig. 2. The overview of the MQA dataset.

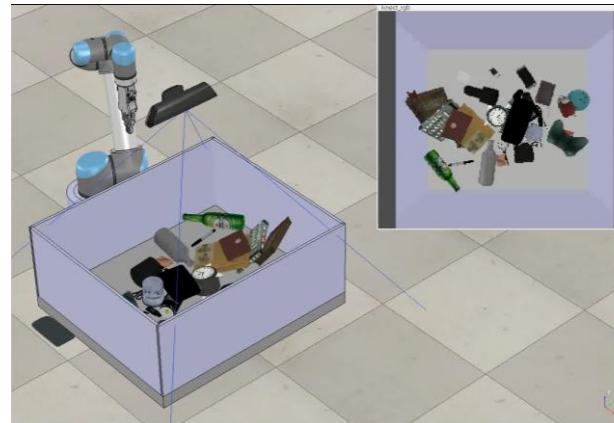


- MQA: Answering the question via robotic manipulation, RSS, 2021

# 4 具身语言问答

## ➤ 前沿：操作问答

question: Is there a key in the bin?



state 0



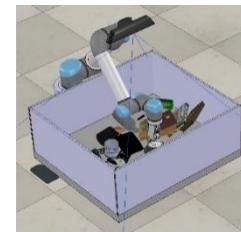
manipulation



state 1



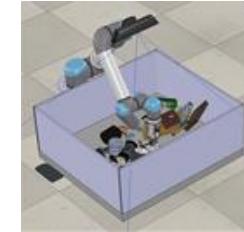
manipulation



State 2



manipulation



state 3



- MQA: Answering the question via robotic manipulation, RSS, 2021

# 4 具身语言问答

## ➤ 前沿：操作问答

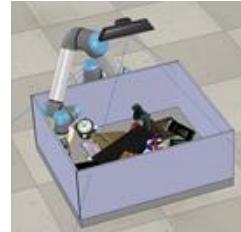
question: How many keyboards are there in the bin?



state 0



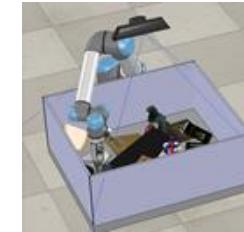
manipulation



state 1



manipulation



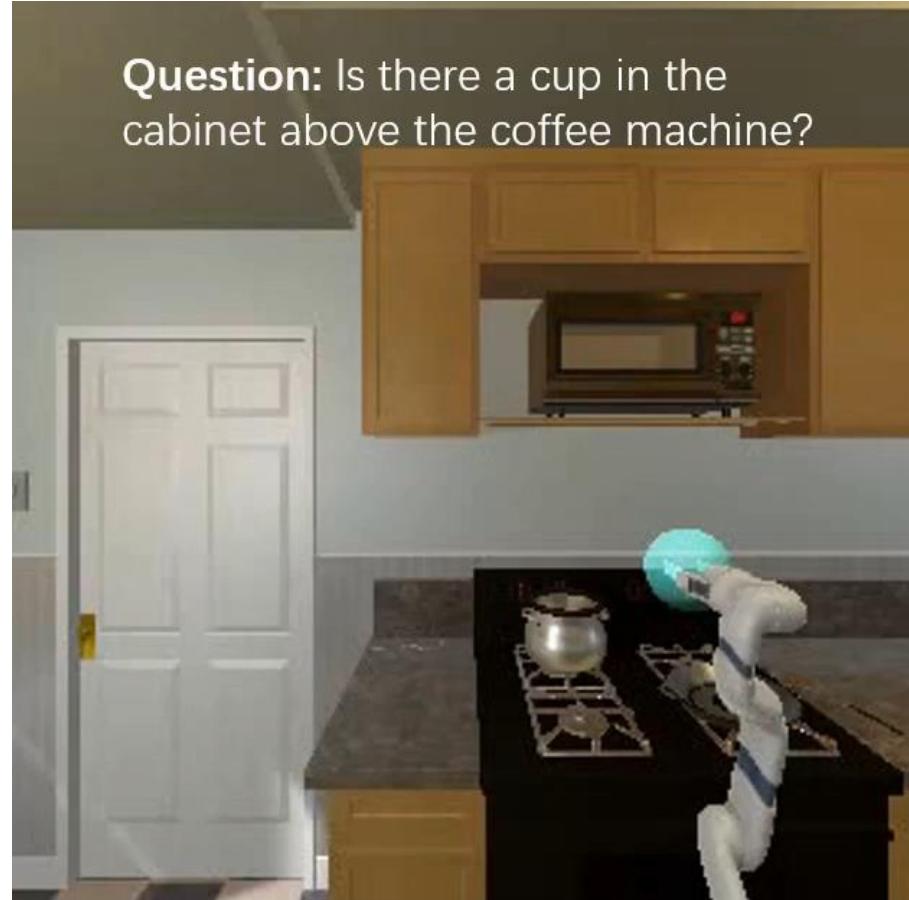
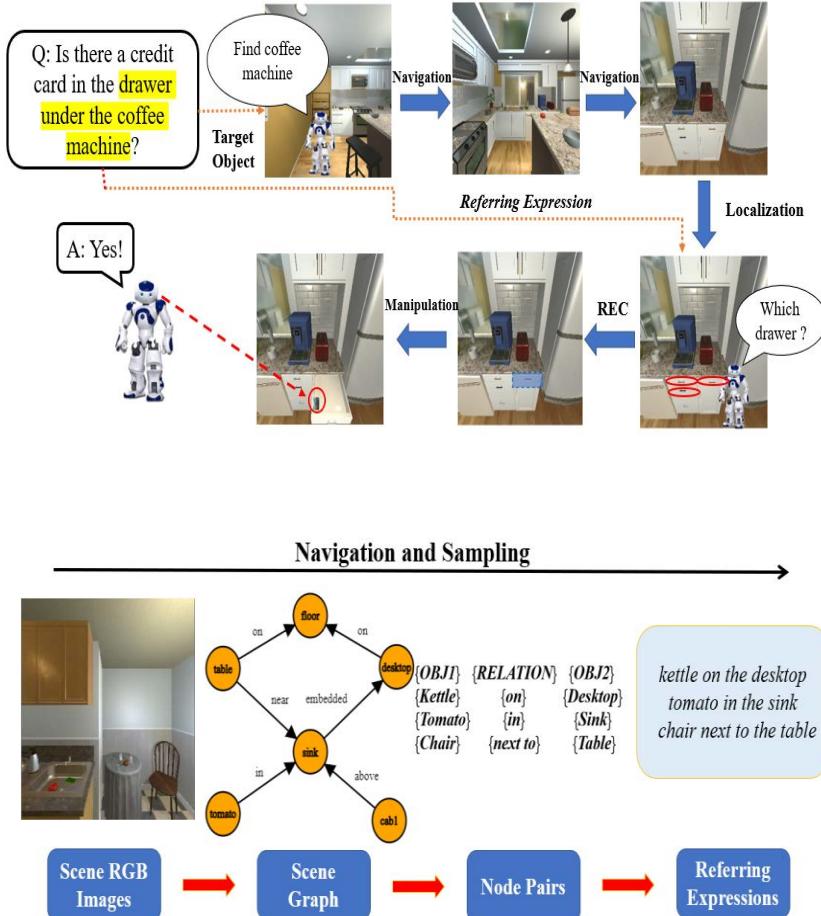
state 2



- MQA: Answering the question via robotic manipulation, RSS, 2021

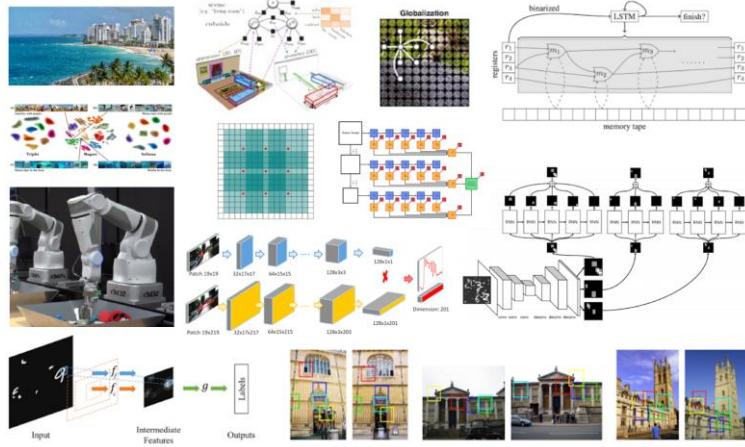
# 4 具身语言问答

## ➤ 前沿：操作问答

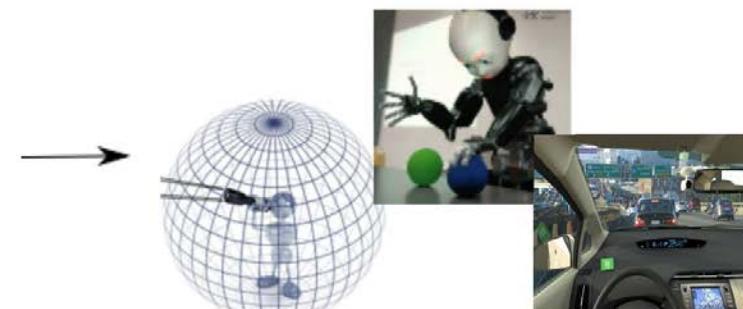
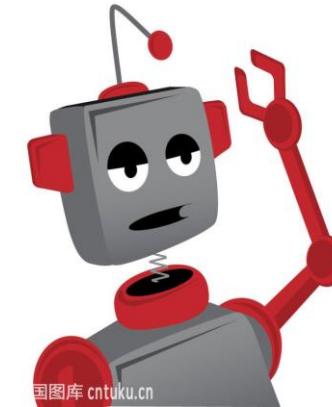


- Embodied Referring Expression for Manipulation Question Answering in Interactive Environment, ICRA, 2023

# 总结与展望



主动感知



# 总结与展望

