# 提纲

**菜鸟**

# 咏卷心菜小鸡

卷心菜里藏玄机，
神鸡孕育其中栖。
眼似星辰闪橙光，
嘴如朝阳染霞霓。
天赐灵物人间现，
地育仙禽世间稀。
莫道此物无用处，
寻常百姓盘中餐。

叶随风舞舞翩翩，

林间鸟语声声甜。

奇峰异石入眼帘，

人间仙境在心田。

普通话 🔊　　上海话 🔊

广东话 🔊　　河南话 🔊

东北话 🔊　　陕西话 🔊

山东话 🔊　　四川话 🔊

香港话 🔊　　台湾话 🔊

https://www.text-to-speech.cn/

AI生成

https://png3d.com/                    https://pika.art/

The leaves dance with the wind, The birds in the forest

TEXT TO SONG

Stay With Me

by: Amy

https://tuna.voicemod.net/text-to-song

The leaves are dancing with the wind.

大家好我是沐沐

所以让我们可以这样

# TEXT-TO-VIDEO

PIKA

https://pika.art/

https://research.runwayml.com/gen2

上海大学首部AIGC宣传片：《智绘上大·未来可期》

广告天才专题
获奖作者：希希叔叔
《Kling汉堡》

可灵AI 视频生成

它把兔子叫到了自己家中

## 文心一言
https://yiyan.baidu.com/



## 通义千问
https://tongyi.aliyun.com/qianwen/



## 讯飞星火
https://xinghuo.xfyun.cn/



## Kimi
https://kimi.moonshot.cn/



## 天工
https://www.tiangong.cn/



## 豆包
https://www.doubao.com/



## 智谱
https://chatglm.cn/



## 可灵
https://klingai.kuaishou.com/

AI模型可大致分为决策式AI（Discriminant AI）和生成式AI（Generative AI）两类。

| 类型 | 决策式AI | 生成式AI |
|---|---|---|
| 技术路径 | 已知数据分别求解输出类别标签，区分不同类型数据，例如将图像区分为猫和狗 | 分析归纳已有数据后创作新的内容，例如生成逼真的猫或狗的图像 |
| 成熟程度 | 技术成熟，应用广泛，辅助提高非创造性工作效率 | 2014年开始快速发展，近期发展速度呈指数级爆发，部分领域应用落地 |
| 应用方向 | 推荐系统、风控系统、决策智能体等 | 内容创作、科研、人机交互以及多个工业领域 |
| 应用产品 | 人脸识别、精准广告推送、金融用户评级、智能辅助驾驶等 | 文案写作、文字转图片、视频智能配音、智能海报生成、视频智能特效、代码生成、语音人机交互、智能医疗诊断等 |

# AIGC（AI-Generated Content，人工智能生成内容）



内容数量

AIGC

AI 辅助用户创作
AIUGC

用户创作
UGC

专业制作
PGC

内容创作模式

图：内容创作模式的四个发展阶段

- 1957 年莱杰伦·希勒(Leiaren Hiller)和伦纳德·艾萨克森（Leon-ard Isaacson)完成了人类历史上第一支由计算机创作的音乐作品就可以看作是 AIGC 的开端。

- 2022 年才真正算是 AIGC 的爆发之年，人们看到了 AIGC 无限的创造潜力和未来应用可能性。

图：AIGC 技术累积融合 [02]

| 模型 | 提出时间 | 模型描述 |
|---|---|---|
| 变分自动编码 (Variational Autoencoders, VAE) | 2014年 | 基于变分下界约束得到的Encoder-Decoder模型对 |
| 生成对抗网络 (GAN) | 2014年 | 基于对抗的Generator-Discriminator模型对 |
| 基于流的生成模型 (Flow-based models) | 2015年 | 学习一个非线性双射转换 (bijective transformation)，其将训练数据映射到另一个空间，在该空间上分布是可以因子化的，整个模型架构依靠直接最大化log-likelihood来完成 |

2014 年，伊恩·古德费洛(Ian Goodfellow)提出的生成对抗网络(Generative Adversarial Network，GAN) 成为早期最为著名的生成模型。GAN 使用合作的零和博弈框架来学习，被广泛用于生成图像、视频、语音和三维物体模型等。GAN 也产生了许多流行的架构或变种，如 DCGAN，StyleGAN，BigGAN，StackGANPix2pix，Age-cGAN，CycleGAN、对抗自编码器(Adversarial Autoencoders，AAE )、对抗推断学习（Adversarially Learned Inference，ALI)等

# 生成

美国科罗拉多州上月举办艺术博览会，一幅名为《太空歌剧院》的画作最终获得数字艺术类别冠军。该作品先由AI制 图 工 具 Midjourney生成，再经Photoshop润色而來。

https://www.midjourney.com/

# 生成式AI

GAN网络应用：AI换脸 是指用另一个人脸来替换一张图片或视频中的一个人脸，合成新的媒体物，它是Deepfake技术最广为人知的一种应用形式。

TITANIC

Transformer、基于流的生成模型（Flow-based models）、扩散模型(Diffusion Model)等深度学习的生成算法相继涌现。

从最优化模型性能的角度出发，扩散模型相对GAN 来说具有更加灵活的模型架构和精确的对数似然计算，已经取代 GAN 成为最先进的图像生成器。2021 年6 月，OpenAI 发表论文已经明确了这个结论和发展趋势。

| 扩散模型 (Diffusion Model) | 2015年 | 扩散模型有两个过程，分别为扩散过程和逆扩散过程。在前向扩散阶段对图像逐步施加噪声，直至图像被破坏变成完全的高斯噪声，然后在逆向阶段学习从高斯噪声还原为原始图像的过程。经过训练，该模型可以应用这些去噪方法，从随机输入中合成新的"干净"数据。 |
| Transformer模型 | 2017年 | 一种基于自注意力机制的神经网络模型，最初用来完成不同语言之间的文本翻译任务，主体包含Encoder和Decoder部分，分别负责对源语言文本进行编码和将编码信息转换为目标语言文本 |

扩散模型(Diffusion Model)是受非平衡热力学的启发，定义一个扩散步骤的马尔可夫链，逐渐向数据添加随机噪声，然后学习逆扩散过程，从噪声中构建所需的数据样本。扩散模型最初设计用于去除图像中的噪声。随着降噪系统的训练时间越来越长并且越来越好，它们最终可以从纯噪声作为唯一输入生成逼真的图片。

# 提纲

角色动画生成

Yao, Heyuan, et al. "Controlvae: Model-based learning of generative controllers for physics-based characters." *ACM Transactions on Graphics (TOG)* 41.6 (2022): 1-16.

# ControlVAE: Model-Based Learning of Generative Controllers for Physics-Based Characters

Heyuan Yao[1,2]    Zhenhua Song[1,2]    Baoquan Chen[2,3]    Libin Liu[2,3]

[1]School of Computer Science, Peking University

[2]Key Laboratory of Machine Perception (MOE), Peking University

[3]School of Intelligence Science and Technology, Peking University

SIGGRAPH ASIA 2022 DAEGU

Yao, Heyuan, et al. "Controlvae: Model-based learning of generative controllers for physics-based characters." *ACM Transactions on Graphics (TOG)* 41.6 (2022): 1-16.

Starke, Sebastian, et al. "Categorical Codebook Matching for Embodied Character Controllers." *ACM Transactions on Graphics (TOG)* 43.4 (2024): 1-14.

Hassan, Mohamed, et al. "Synthesizing physical character-scene interactions." *ACM SIGGRAPH 2023 Conference Proceedings.* 2023.

Hassan, Mohamed, et al. "Synthesizing physical character-scene interactions." *ACM SIGGRAPH 2023 Conference Proceedings*. 2023.

Taming Diffusion Probabilistic Models for Character Control

Rui Chen*[1], Mingyi Shi*[2], Shaoli Huang[3], Ping Tan[1], Taku Komura[2], Xuelin Chen†[3]
* Joint First Author    † Corresponding author

This video contains voice

Style: FlickLegs

# 角色动画生成

https://anyskill.github.io/

# 提纲

## Robotics at Google

LM-Nav

SayCan

Inner Monologue

RT-1

PLAM

PLAM-E

Code as Policies

RT-2

*Florida is a lawless place.*

## Other Reseachers

VoxPoser - Fei-Fei Li

Diffusion Policy - Shuran Song

3D-LLM - Chuang Gan

etc.

# Timeline



**LM-Nav**

**Inner Monologue**

**Build Feedback System for SayCan**

**Navigation based LM**

**LLM Planner + BC-Z Skills**

**SayCan**

*cite from Robotics at Google

机器人动作生成

**Code as Policies**

**Prompt to generate codes**

Internet-Scale VQA + Robot Action Data

Q: What is happening in the image?
A: 311 423 170 55 244
A grey donkey walks down the street.

Q: Que puis-je faire avec ces objets?
A: 3455 1144 189 25673
Faire cuire un gâteau.

Q: What should the robot do to <task>?
A: 132 114 128 5 25 156
ΔTranslation = [0.1, -0.2, 0]
ΔRotation = [10°, 25°, -7°]

Vision-Language-Action Models for Robot Control

Q: What should the robot do to <task>? A: ...
RT-2
Large Language Model
ViT
A: 132 114 128 5 25 156
De-Tokenize
ΔT = [0.1, -0.2, 0]
ΔR = [10°, 25°, -7°]
Robot Action

Co-Fine-Tune          Deploy

Closed-Loop Robot Control
Put the strawberry into the correct bowl
Pick the nearly falling bag
Pick object that is different

User
Stack the blocks on the empty bowl.
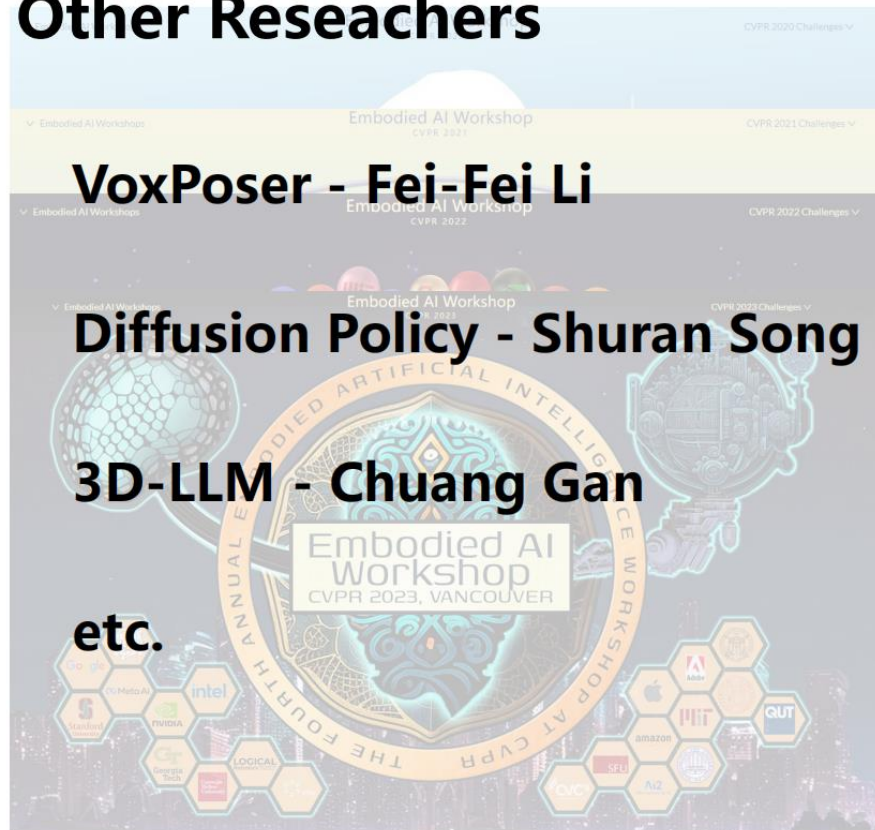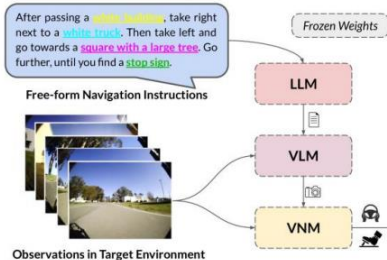
**L**arge **L**anguage **M**odel

**Policy Code**

Perception APIs
Control APIs

```
block_names = detect_objects("blocks")
bowl_names = detect_objects("bowls")
for bowl_name in bowl_names:
    if is_empty(bowl_name):
        empty_bowl = bowl_name
        break
objs_to_stack = [empty_bowl] + block_names
stack_objects(objs_to_stack)

        def is_empty(name):

def stack_objects(obj_names):
    n_objs = len(obj_names)
    for i in range(n_objs - 1):
        obj0 = obj_names[i + 1]
        obj1 = obj_names[i]
        pick_place(obj0, obj1)
```

**end-to-end low-level controller**

**RT-2**

*cite from Robotics at Google

## SayCan

**Paper**



**Do As I Can, Not As I Say:**
Grounding Language in Robotic Affordances

Michael Ahn*  Anthony Brohan*  Noah Brown*  Yevgen Chebotar*  Omar Cortes*  Byron David*  Chelsea Finn*
Chuyuan Fu*  Keerthana Gopalakrishnan*  Karol Hausman*  Alex Herzog*  Daniel Ho*  Jasmine Hsu*  Julian Ibarz*
Brian Ichter*  Alex Irpan*  Eric Jang*  Rosario Jauregui Ruano*  Kyle Jeffrey*  Sally Jesmonth*  Nikhil Joshi*
Ryan Julian*  Dmitry Kalashnikov*  Yuheng Kuang*  Kuang-Huei Lee*  Sergey Levine*  Yao Lu*  Linda Luu*  Carolina Parada*
Peter Pastor*  Jornell Quiambao*  Kanishka Rao*  Jarek Rettinghouse*  Diego Reyes*  Pierre Sermanet*  Nicolas Sievers*
Clayton Tan*  Alexander Toshev*  Vincent Vanhoucke*  Fei Xia*  Ted Xiao*  Peng Xu*  Sichun Xu*  Mengyuan Yan*  Andy Zeng*

Robotics at Google        Everyday Robots

**Overview**



**Demo**



*cite from Robotics at Google

## RT-1

### Paper

RT-1: ROBOTICS TRANSFORMER
FOR REAL-WORLD CONTROL AT SCALE

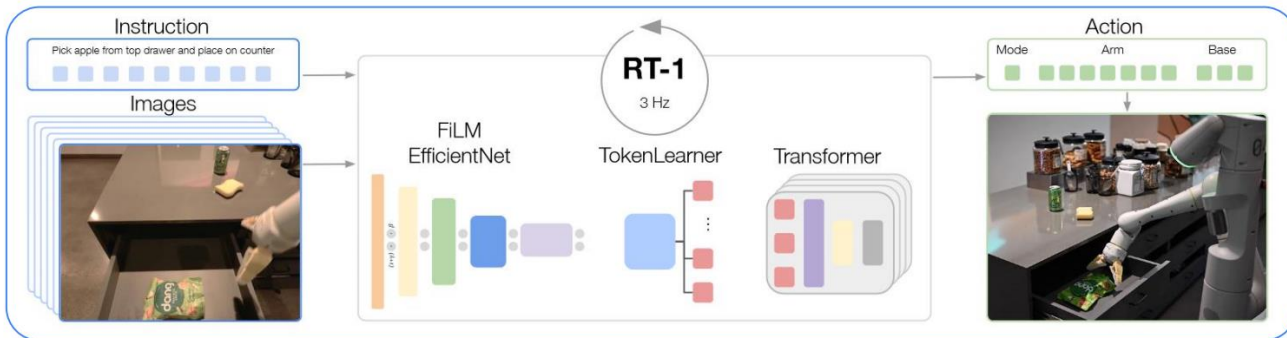Anthony Brohan*, Noah Brown*, Justice Carbajal*, Yevgen Chebotar*, Joseph Dabis*, Chelsea Finn*, Keerthana Gopalakrishnan*, Karol Hausman*, Alex Herzog†, Jasmine Hsu*, Julian Ibarz*, Brian Ichter*, Alex Irpan*, Tomas Jackson*, Sally Jesmonth*, Nikhil J Joshi*, Ryan Julian*, Dmitry Kalashnikov*, Yuheng Kuang*, Isabel Leal*, Kuang-Huei Lee‡, Sergey Levine*, Yao Lu*, Utsav Malla*, Deeksha Manjunath*, Igor Mordatch†, Ofir Nachum†, Carolina Parada*, Jodilyn Peralta*, Emily Perez*, Karl Pertsch*, Jornell Quiambao*, Kanishka Rao*, Michael Ryoo*, Grecia Salazar*, Pannag Sanketi*, Kevin Sayed*, Jaspiar Singh*, Sumedh Sontakke‡, Austin Stone*, Clayton Tan*, Huong Tran*, Vincent Vanhoucke*, Steve Vega*, Quan Vuong*, Fei Xia*, Ted Xiao*, Peng Xu*, Sichun Xu*, Tianhe Yu*, Brianna Zitkovich*

*Robotics at Google, †Everyday Robots, ‡Google Research, Brain Team

### Overview



### Details

#### 1. Action Spaces

arm: x, y, z, roll, pitch, yaw, opening of grasper

base: x, y, yaw

switch mode: {control the arm, the base, termination}

control frequency = 3Hz

**Action tokenization.** To tokenize actions, each action dimension in RT-1 is discretized into 256 bins. As mentioned previously, the action dimensions we consider include seven variables for the arm movement ($x$, $y$, $z$, roll, pitch, yaw, opening of the gripper), three variables for base movement ($x$, $y$, yaw) and a discrete variable to switch between three modes: controlling arm, base or terminating the episode. For each variable, we map the target to one of the 256 bins, where the bins are uniformly distributed within the bounds of each variable.

*cite from Robotics at Google

# 机器人动作生成

## PALM-E

**Paper**

### PaLM-E: An Embodied Multimodal Language Model

Danny Driess[1,2]    Fei Xia[1]    Mehdi S. M. Sajjadi[3]    Corey Lynch[1]    Aakanksha Chowdhery[3]
Brian Ichter[1]    Ayzaan Wahid[1]    Jonathan Tompson[1]    Quan Vuong[1]    Tianhe Yu[1]    Wenlong Huang[1]
Yevgen Chebotar[1]    Pierre Sermanet[1]    Daniel Duckworth[3]    Sergey Levine[1]    Vincent Vanhoucke[1]
Karol Hausman[1]    Marc Toussaint[2]    Klaus Greff[3]    Andy Zeng[1]    Igor Mordatch[3]    Pete Florence[1]

[1] Robotics at Google    [2] TECHNISCHE UNIVERSITÄT BERLIN    [3] Google Research

**Details**

**1. Multi-modal Imputs**

**1.1 State estimation**

**1.2 Images**

**1.3 Language**

**Overview**

106

# Code as Polices

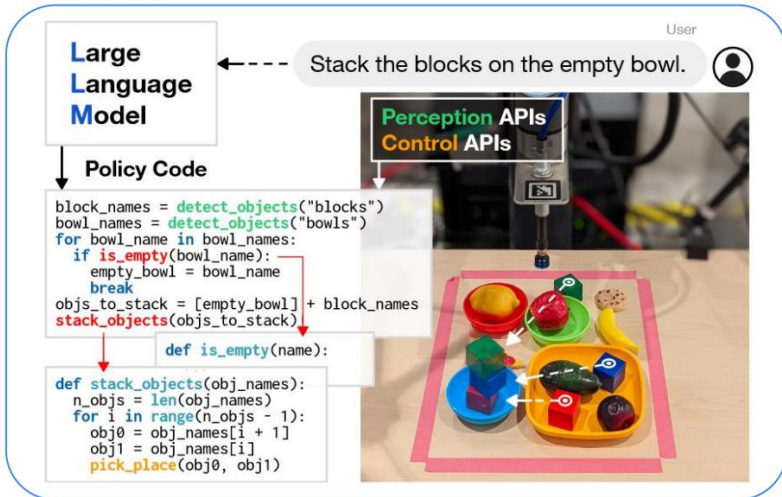## Paper

**Code as Policies:**
Language Model Programs for Embodied Control

Jacky Liang    Wenlong Huang    Fei Xia    Peng Xu    Karol Hausman    Brian Ichter    Pete Florence    Andy Zeng

Robotics at Google

## Overview



## Details

### 1. APIs

#### Perception APIs

of LMP-based policies. For example, in real-world experiments below, we use recently developed open-vocabulary object detection models like ViLD [3] and MDETR [2] off-the-shelf to obtain object positions and bounding boxes.

#### Control APIs

architect a dynamic codebase. We demonstrate across several robot systems that LLMs can autonomously interpret language commands to generate LMPs that represent reactive low-level policies (e.g., PD or impedance controllers), and waypoint-based policies (e.g., for vision-based pick and place, or trajectory-based control).

where `put_first_on_second` is an existing open vocabulary pick and place primitive (e.g., CLIPort [36]). For new embodiments, these active function calls can be replaced with available control APIs that represent the action space (e.g., `set_velocity`) of the agent. Hierarchical code-gen with verbose variable names

## RT-2

**Paper**



# RT-2: Vision-Language-Action Models
## Transfer Web Knowledge to Robotic Control

Anthony Brohan    Noah Brown    Justice Carbajal    Yevgen Chebotar    Xi Chen    Krzysztof Choromanski    Tianli Ding
Danny Driess    Avinava Dubey    Chelsea Finn    Pete Florence    Chuyuan Fu    Montse Gonzalez Arenas    Keerthana Gopalakrishnan
Kehang Han    Karol Hausman    Alex Herzog    Jasmine Hsu    Brian Ichter    Alex Irpan    Nikhil Joshi    Ryan Julian
Dmitry Kalashnikov    Yuheng Kuang    Isabel Leal    Lisa Lee    Tsang-Wei Edward Lee    Sergey Levine    Yao Lu    Henryk Michalewski
Igor Mordatch    Karl Pertsch    Kanishka Rao    Krista Reymann    Michael Ryoo    Grecia Salazar    Pannag Sanketi    Pierre Sermanet
Jaspiar Singh    Anikait Singh    Radu Soricut    Huong Tran    Vincent Vanhoucke    Quan Vuong    Ayzaan Wahid    Stefan Welker
Paul Wohlhart    Jialin Wu    Fei Xia    Ted Xiao    Peng Xu    Sichun Xu    Tianhe Yu    Brianna Zitkovich

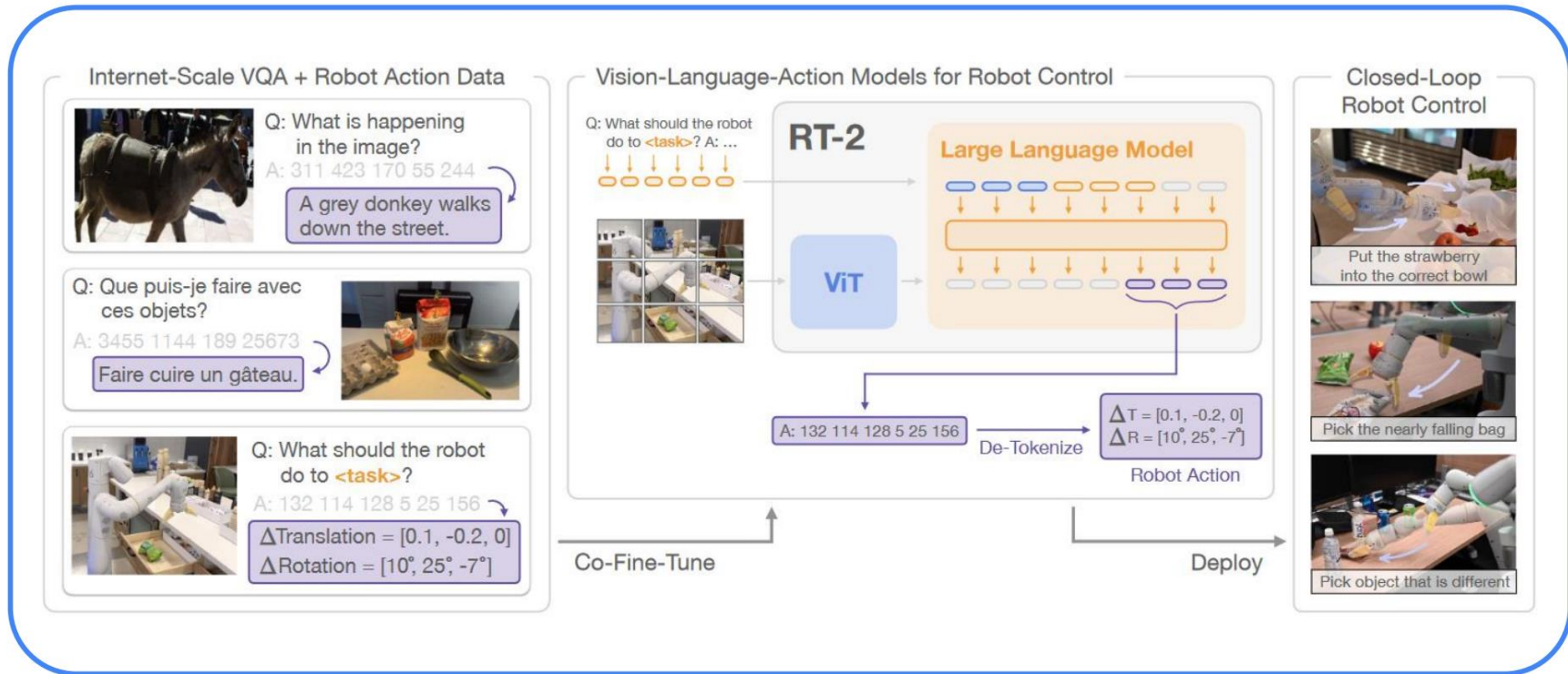*Authors listed in alphabetical order (see paper appendix for contribution statement).*

Google DeepMind

# RT-2

## Overview

# RT-2

**Demo**

**Figure 1:** VOXPOSER extracts language-conditioned **affordances** and **constraints** from LLMs and grounds them to the perceptual space using VLMs, using a code interface and without additional training to either component. The composed map is referred to as a 3D value map, which enables **zero-shot** synthesis of trajectories for large varieties of everyday manipulation tasks with an **open-set of instructions** and an **open-set of objects**.

# 机器人动作生成

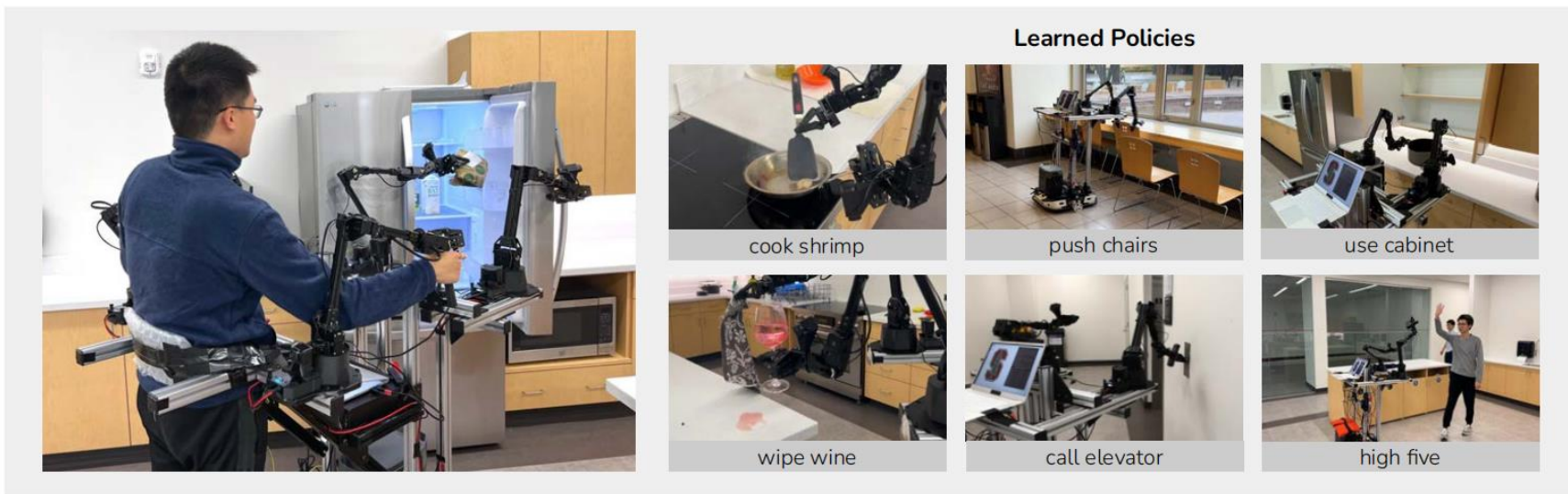斯坦福大学的科研团队近日开发了 Mobile ALOHA，可以执行打开厨房用具柜、洗锅、炸虾、做菜、打扫卫生、整理衣物、套被套等 50 多项家务。这款家用机器人成本仅 3.2 万美元。



**Figure 1:** *Mobile ALOHA* 🏃. We introduce a low-cost mobile manipulation system that is bimanual and supports whole-body teleoperation. The system costs $32k including onboard power and compute. *Left:* A user teleoperates to obtain food from the fridge. *Right: Mobile ALOHA* can perform complex long-horizon tasks with imitation learning.

其算法 Action Chunking with Transformers （ACT）采用了神经网络模型 Transformers，因此**具备模仿学习能力。只需要15分钟的演示，机械臂就可以学会一个动作**——直接从真实演示中执行端到端模仿学习，并通过自定义远程操作界面收集。
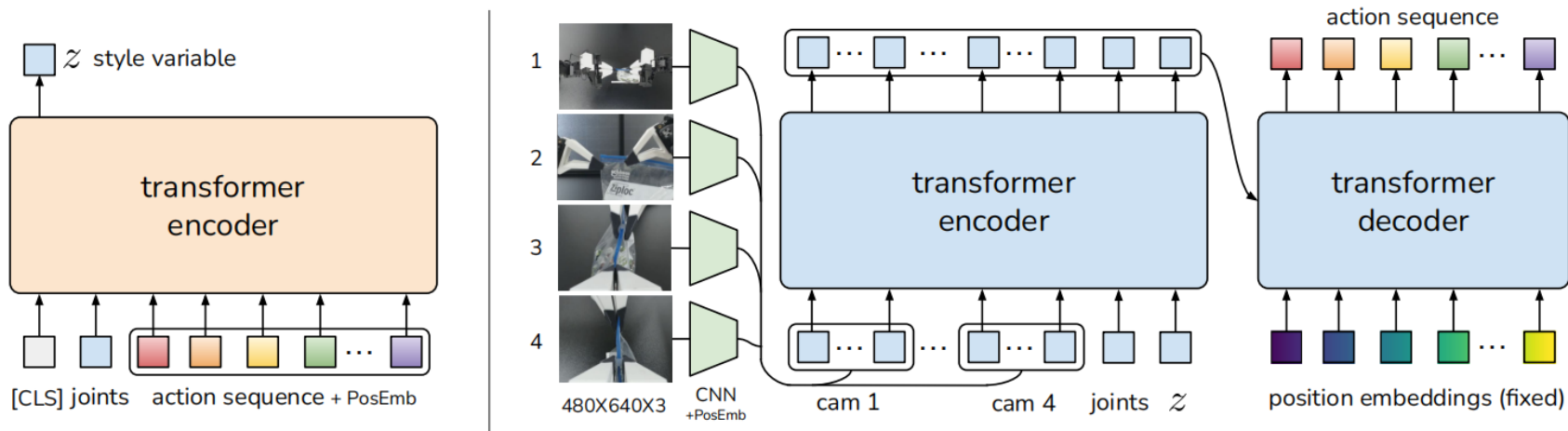


Fig. 4: *Architecture of Action Chunking with Transformers (ACT)*. We train ACT as a Conditional VAE (CVAE), which has an encoder and a decoder. *Left:* The encoder of the CVAE compresses action sequence and joint observation into $z$, the style variable. The encoder is discarded at test time. *Right:* The decoder or policy of ACT synthesizes images from multiple viewpoints, joint positions, and $z$ with a transformer encoder, and predicts a sequence of actions with a transformer decoder. $z$ is simply set to the mean of the prior (i.e. zero) at test time.

6x speed

Stanford
University

谢谢大家