

RT-2 论文精读笔记

论文全称： RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control

中文直译： RT-2：视觉-语言-动作模型将互联网知识迁移到机器人控制

作者： Anthony Brohan 等 (Google DeepMind)

发表时间： 2023 年 7 月

项目主页： <https://robotics-transformer2.github.io>

学习计划位置： Day 9 下午 · 论文精读 (1)

一、一句话总结

RT-2 的核心思想：把机器人的动作（移动多少、旋转多少）当成一种“语言”，直接塞进已经在互联网上学了海量知识的视觉-语言大模型里一起训练，这样机器人就能“继承”大模型的常识和推理能力。

 **比喻理解：** 想象你雇了一个精通多国语言、博览群书的翻译官 (VLM)，现在你教他一门新“语言”——机器人动作语。因为他已经很聪明了，学这门新语言很快，而且他能把以前读过的百科知识用在新语言上。比如他以前读过“苹果是红色的”，现在你让他“把苹果放到颜色相同的碗里”，他就能做到——尽管你从来没在机器人训练数据里教过“颜色匹配”这件事。

二、论文要解决什么问题？

2.1 背景矛盾

存在一个核心矛盾：

	互联网大模型 (VLM)	机器人
优势	海量常识（“苹果是红色的”）、推理能力	能实际操作物理世界
短板	只会输出文字，不会控制机械臂	训练数据太少，缺乏常识

之前的方案（如 SayCan）是**两层分离的**：让大模型做“高层规划”（说话），让底层控制器做“动手”。这就像你请了一个军师 (LLM) 和一个武夫 (底层控制器)，军师出谋划策，武夫只管砍。但问题是：**武夫本身不懂策略，军师也不懂打架**，两者之间的协调很生硬。

2.2 RT-2 的核心提问

能不能直接让大模型“学会动手”？

即：让一个单一的端到端 (end-to-end，从输入到输出一步到位) 模型，既理解语言和图像，又能直接输出机械臂的具体动作？

三、核心方法 (Section 2 & 3) 必读

3.1 核心 idea：动作也是一种"语言" —— 动作分词 (Action Tokenization)

这是全文最关键的思想，也是 Figure 1 的精髓。

问题：VLM（视觉-语言模型）只会输出文字 token（词元），机器人需要的是连续的数字（比如"末端执行器向前移动 0.1m"）。怎么让一个"说话的模型"去"动手"？

解决方案：把动作离散化成整数，当作"特殊单词"。

具体来说，机器人的一个动作包含 8 个维度：

[终止信号, Δx , Δy , Δz , $\Delta roll$, $\Delta pitch$, Δyaw , 夹爪开合度]
是否完成 位置偏移量(3个) 旋转偏移量(3个) 抓紧/松开

每个连续维度被均匀切成 256 个格子（bins），也就是把一个连续的浮点数"四舍五入"到最近的整数编号（0~255）。

比喻：就像温度计上只有 256 个刻度。真实温度是 23.7°C，你取最近的刻度 24。精度有损失，但足够用了。

最终一个动作变成一串数字，比如："1 128 91 241 5 101 127"

这串数字对 VLM 来说和 "一只灰色的驴在街上走" 没有本质区别——都是 token 序列。

3.2 训练格式：伪装成 VQA（视觉问答）

RT-2 把机器人控制任务"伪装"成一个视觉问答任务：

输入（问题）：[机器人摄像头图片] + "Q: what action should the robot take to pick the apple? A:"
输出（回答）："1 128 91 241 5 101 127"

对于 VLM 来说，这就是一道看图答题。只不过"答案"不再是一句话，而是一串代表动作的数字。

3.3 关键训练技巧：共同微调 (Co-Fine-Tuning)

这是论文中非常重要的工程细节。

训练策略	做法	效果
从头训练 (from scratch)	不用预训练权重，直接用机器人数据训练	✗ 非常差（5B 模型几乎不 work）
纯微调 (fine-tuning)	只用机器人数据微调预训练模型	⚠ 还行，但会遗忘互联网知识
共同微调 (co-fine-tuning) <input checked="" type="checkbox"/>	同时用互联网数据 + 机器人数据微调	✓ 最好！既保留常识又学会动手

比喻：这就像一个英语老师去学日语。如果他完全放下英语只学日语（纯微调），时间久了英语会退步。但如果他每天一半时间教英语、一半时间学日语（共同微调），两门语言都能保持好。这在深度学习里叫灾难性遗忘 (Catastrophic Forgetting)，共同微调就是对抗遗忘的策略。

3.4 输出约束 (Output Constraint)

一个实际的工程问题：VLM 可以输出任意 token，但机器人只接受合法的动作 token。

解决方法：在推理（inference）时，如果模型被问的是机器人动作任务，就限制解码只从 256 个动作 token 中采样；如果是普通 VQA 任务，就正常输出所有 token。

| 这就像考试时，选择题只让你填 ABCD，不许你写散文。

3.5 实时推理 (Real-Time Inference)

55B 参数的模型太大了，机器人旁边的电脑跑不动。

解决方案：把模型部署到云端 TPU 集群，机器人通过网络请求动作。

模型	推理频率
RT-2-PaLI-X-55B	1-3 Hz (每秒 1~3 个动作)
RT-2-PaLI-X-5B	~5 Hz (每秒约 5 个动作)

| ! 这也是一个重要局限：延迟太高，无法做高频控制（如需要 100Hz 的灵巧操作）。

四、两个模型实例

RT-2 不是一个固定模型，而是一种方法论（recipe）。论文实例化了两个版本：

	RT-2-PaLI-X	RT-2-PaLM-E
底座 VLM	PaLI-X (视觉编码器 ViT + 语言模型 UL2)	PaLM-E (decoder-only LLM)
参数量	5B / 55B	12B
动作 token 方案	整数 0~255 直接有对应 token	覆盖最不常用的 256 个 token
特点	更灵活，可选不同大小	与 PaLM 生态兼容

五、涌现能力 (Emergent Capabilities) ★ Section 4.3 必读

这是全文最震撼的部分。所谓涌现能力，就是模型在机器人训练数据里从未见过，但因为继承了互联网知识而“自动获得”的能力。

5.1 三大涌现类别

① 符号理解 (Symbol Understanding)

- 指令: "move apple to 3" (把苹果移到数字 3 旁边)
- 机器人数据里从没出现过"数字"或"图标", 但 VLM 在互联网图片上见过
- RT-2-PaLI-X-55B 成功率 **82%**, 而 RT-1 只有 **16%**

② 推理 (Reasoning)

- 指令: "move banana near the sum of two plus one" (把香蕉移到 $2+1=3$ 旁边)
- 需要数学计算 + 物理操作
- 指令: "pick a healthy drink" (拿一个健康的饮料) → 模型选了水而不是可乐
- RT-2 平均 **46%**, RT-1 只有 **16%**

③ 人脸识别 (Person Recognition)

- 指令: "move coke can to Taylor Swift"
- RT-2 平均 **53%**, RT-1 只有 **20%**

💡 与你的课题的关系: 这直接说明 VLM 的**常识推理**可以帮助机器人理解**模糊指令**。比如"把那个东西放到对的地方"——什么是"对的地方"? 需要常识推理。RT-2 证明了这条路走得通。

5.2 整体泛化结果

在未见过的物体、背景、环境上, RT-2 比 RT-1 提升了约 **2 倍**, 比其他基线提升了约 **6 倍**。

六、思维链推理 (Chain-of-Thought Reasoning) · Section 4.4

RT-2 还可以像 ChatGPT 一样"先想再做":

指令: "I need to hammer a nail, what object from the scene might be useful?"
(我需要锤一颗钉子, 场景里什么东西可以当锤子用?)

模型输出:

Plan: Rocks. ← 先用自然语言"想"
Action: 1 129 138 122 132 135 106 127 ← 再输出动作去拿石头

实现方法: 在训练数据里加一个 "Plan:" 字段, 让模型先生成计划文字, 再生成动作 token。只需要额外微调几百步。

💡 哲学思考: 这体现了"知行合一"的思想——模型不是"会想不会做"或"会做不会想", 而是在同一个框架里把**思考 (Plan) 和行动 (Action) **统一起来。思考指导行动, 行动验证思考。

七、消融实验关键结论 (Section 4.3)

发现	启示
从头训练 5B 模型几乎完全失败	预训练至关重要 , 没有互联网知识的大模型是空壳
Co-fine-tuning > Fine-tuning	保留原始数据防止遗忘, "边学新边复习旧"
55B > 5B	模型越大, 泛化越好 (Scale 定律在机器人领域也成立)

八、局限性 (Section 5)

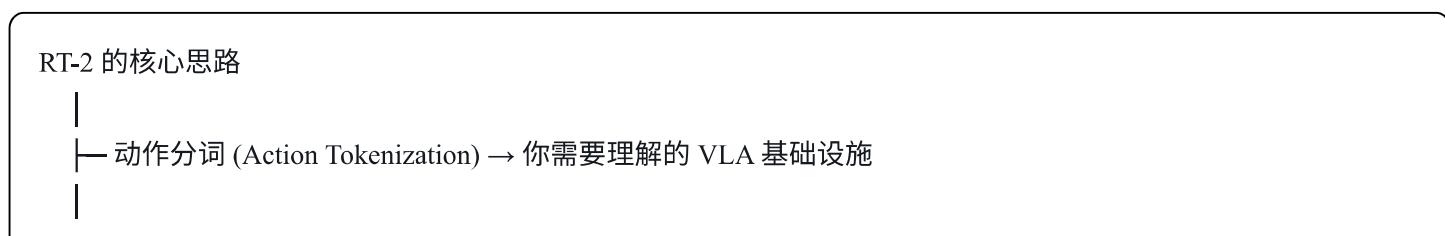
局限	说明	对你课题的影响
不能学新动作	只能用训练数据里已有的动作技能 (抓、放、推), 不能凭空学会翻跟头	你的机械臂如果需要特殊技能, 仍需采集对应数据
计算开销巨大	55B 模型需要云端 TPU, 1-3Hz 控制频率	工业场景需要考虑延迟
闭源模型受限	基于 PaLI-X/PaLM-E, 无法公开复现	你实际做实验建议用 OpenVLA (开源替代品)

九、与你的课题的深层联系

9.1 RT-2 为你的"模糊指令消歧"课题证明了什么?

1. **VLM 的常识可以直接用于底层控制:** 不需要复杂的中间层, 一个模型就能把"拿那个快掉下去的包"翻译成具体动作
2. **模糊指令的消歧依赖常识推理:** 比如"pick a healthy drink"就是一个模糊指令——什么算"healthy"? 这需要营养学常识
3. **共同微调是保留常识的关键:** 你未来做工业机械臂指令消歧模型, 也需要类似的策略来同时保留语言理解和动作执行能力

9.2 RT-2 → 你的课题的知识链



- 涌现能力 (Emergent Capabilities) → 证明常识可以消解指令歧义
- 思维链 (CoT) → 你可以让机器人"先想后做", 输出消歧理由
- Co-Fine-Tuning → 你训练自己模型时的核心策略

十、必看图表清单

图表	内容	为什么必看
Figure 1 ★★★	整个 RT-2 的一图流总览	理解"动作即语言"的核心 idea
Figure 2 ★★	涌现能力的定性展示	直观感受 VLM 常识迁移的效果
Figure 4 ★★	泛化性能柱状图	RT-2 vs 所有基线, 一目了然
Figure 6b ★★	消融实验: 大小 × 训练策略	理解 co-fine-tuning 和 scale 的重要性
Figure 7 ★★	思维链推理的实际执行	看"先想后做"的实际效果
Table 5 ★	涌现能力的定量评估	数字说话, 记住 RT-2 比 RT-1 好多少

十一、可以跳过的部分

- Appendix D: PaLI-X 和 PaLM-E 的详细架构描述 (除非你要自己实现)
- Appendix E: 超参数细节 (学习率、batch size 等, 知道大概即可)
- 离散化动作的具体工程实现细节 (知道"256 bins"的 idea 就够了)

十二、关键术语速查表

英文术语	中文翻译	通俗解释
VLA (Vision-Language-Action)	视觉-语言-动作模型	能看、能听、能动的统一模型
VLM (Vision-Language Model)	视觉-语言模型	能看图说话的模型 (如 GPT-4V)
Action Tokenization	动作分词/动作离散化	把连续动作切成整数编号
Co-Fine-Tuning	共同微调	用原始数据+新数据一起微调, 防遗忘
Emergent Capabilities	涌现能力	训练数据里没有, 但大模型"自动学会"的能力

英文术语	中文翻译	通俗解释
Chain-of-Thought (CoT)	思维链	先用自然语言推理，再输出动作
End-Effecter	末端执行器	机械臂最末端与物体接触的部分（如夹爪）
6-DoF	六自由度	3个平移方向 + 3个旋转方向
Closed-Loop Control	闭环控制	每一步都看当前状态再决定下一步
Catastrophic Forgetting	灾难性遗忘	学新知识时把旧知识忘了的现象
Discretization / Binning	离散化/分桶	把连续值切分到有限个格子里
Symbol Tuning	符号调优	用新含义覆写模型已有的 token

十三、RT 系列发展脉络

RT-1 (2022)

35M 参数，专门为机器人设计的 Transformer

能力：跟着指令做简单抓放

局限：不理解常识，泛化差



RT-2 (2023)  本文

5B~55B 参数，基于互联网 VLM

突破：动作即语言，常识迁移，涌现能力

局限：闭源，计算贵，不能学新动作



RT-X / Open X-Embodiment (2023)

跨机器人的统一数据集 + 统一模型



OpenVLA (2024)

开源的 VLA，你可以实际跑的模型

 **下一步行动：** Day 10 上午继续精读 PaLM-E 论文，它是 RT-2 的底座之一，理解它可以更深入地明白多模态融合的机制。