

# VLA综述论文精读笔记

论文标题: A Survey on Vision-Language-Action Models for Embodied AI

作者: Yueen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, Irwin King

阅读日期: 2025年1月

阅读目的: 为"工业机械臂模糊指令语义消歧"课题建立理论基础

## 1 一句话总结

VLA是一类能够接收视觉和语言输入、输出机器人动作的多模态模型，其核心架构分为高层任务规划器（理解指令、分解任务）和低层控制策略（执行具体动作），而语言指令的理解与落地（Grounding）是连接两者的关键桥梁。

## 2 核心分类体系

### VLA的三大研究主线

```
1 VLA Models
2   └── Components (组件研究)
3     ├── 强化学习: DQN → PPO → Decision Transformer
4     ├── 预训练视觉表征: CLIP, R3M, MVP, VC-1
5     ├── 世界模型: Dreamer, Genie, TWM
6     └── 推理能力: Embodied Chain-of-Thought
7
8   └── Control Policy (低层控制策略)
9     ├── 非Transformer: CLIPort, BC-Z, MCIL
10    ├── Transformer-based: RT-1, Gato, RoboCat
11    ├── 多模态指令: VIMA, MOO ← 【技术相关】
12    ├── 3D视觉: PerAct, Act3D, RVT
13    ├── 扩散策略: Diffusion Policy, DP3
14    └── 大型VLA: RT-2, OpenVLA
15
16   └── Task Planner (高层任务规划)
17     ├── 端到端: PaLM-E
18     └── 模块化
19       ├── 语言驱动: SayCan, Inner Monologue ← 【核心相关】
20       ├── 代码驱动: ProgPrompt, ChatGPT for Robotics
21       └── Grounded (落地) ← 【直接相关】
22         ├── Task Grounding (任务落地)
23         └── World Grounding (世界落地/Affordance)
```

## 🎯 我的课题在哪里?

"工业机械臂模糊指令语义消歧" 位于:

- Task Planner → Modular → Language-based 分支
- 核心挑战: 当用户说"把那个螺丝拧到那里"时, 如何消解"那个"和"那里"的歧义

### 3 与我课题的关系

#### 直接相关的概念

概念	论文中的位置	与我课题的联系
Grounding	§IV-A3, Figure 6	把模糊的语言落地到具体物体/位置
Language-conditioned Policy	§III-B, Table III	指令如何影响动作生成
Task Decomposition	§IV-B1	模糊指令需要先分解再执行
Affordance	§III-A, SayCan	判断"能不能做"来消歧
Multimodal Prompts	§III-B3, VIMA	用多模态信息辅助消歧

#### 关键发现

##### 1. SayCan的"task-grounding + world-grounding"框架

- Task grounding: LLM说"应该做什么"
- World grounding: Affordance判断"能不能做"
- 启发: 模糊指令消歧可以用类似的双重验证机制

##### 2. Inner Monologue的闭环反馈

- 执行过程中不断获取反馈，动态调整理解
- 启发: 消歧不一定是一次性的，可以在执行中逐步确认

##### 3. VIMA的多模态指令

- 除了语言，还可以用图像、演示等方式指定目标
- 启发: 当语言模糊时，可以请求用户提供更多模态的信息

### 4 待读论文清单

#### ★★★ 必读（与课题直接相关）

#	论文	原因	预计耗时
1	SayCan (2022)	Grounding LLM的开创性工作	3小时
2	Inner Monologue (2022)	语言反馈闭环控制	2小时
3	CLIP (2021)	视觉-语言对齐的基础	3小时

## ★★★ 推荐阅读 (技术基础)

#	论文	原因	预计耗时
4	RT-2 (2023)	理解Large VLA	2小时
5	VIMA (2022)	多模态指令处理	2小时
6	OpenVLA (2024)	开源可实验	2小时

## ★★ 扩展阅读 (深入理解)

#	论文	原因
7	Grounding DINO	Visual Grounding的SOTA
8	LLM-Planner	LLM做任务规划
9	ProgPrompt	代码驱动的任务规划

## 5 我的疑问

### 概念层面

1. **Grounding vs Disambiguation:** 这两个概念的边界在哪里？Grounding是消歧的一种手段，还是消歧是Grounding的前置步骤？
2. **Affordance的计算:** SayCan用value function作为affordance，但在工业场景中，如何定义和计算affordance？
3. **FiLM vs Cross-Attention:** 两种语言-视觉融合方式各有什么优缺点？在消歧任务中哪个更合适？
4. **Action Chunking:** ACT论文提出的action chunking对连续动作很有效，但对于需要精确消歧的任务，会不会因为"打包预测"而丢失细节？

### 技术层面

5. **数据集缺口:** 现有数据集（RefCOCO, CALVIN等）的指令都比较清晰，有没有专门针对**模糊指令**的benchmark？如果没有，我是否需要自己构建？
6. **评估指标:** 如何评估“消歧”的效果？除了最终任务成功率，有没有更细粒度的指标？

## 6 关键图表索引

图表	页码	内容	重要性
Figure 2a	p.2	VLA概念韦恩图	★★★

图表	页码	内容	重要性
Figure 2b	p.2	发展时间线	★★
Figure 3	p.3	完整分类体系	★★★
Figure 4	p.7	层级架构示意	★★★
Figure 5	p.9	5种VLA架构	★★★
Figure 6	p.13	模块化Task Planner	★★★
Table III	p.8	控制策略对比表	★★
Table V	p.14	数据集汇总	★★

## 7 术语表 (方便复习)

术语	中文	解释
VLA	视觉-语言-动作模型	接收视觉和语言，输出动作
Grounding	落地/接地	把抽象语言对应到具体物体/位置
Affordance	可供性	物体"能被怎么用"的属性
FiLM	特征调制层	用语言调制视觉特征的方法
BC	行为克隆	直接模仿专家动作的学习方式
DDPM	去噪扩散概率模型	生成式模型，用于动作生成
Task Planner	任务规划器	把长期目标分解为子任务
Control Policy	控制策略	执行具体动作的策略网络

笔记完成时间: 2025年1月