

VIMA 论文精读笔记

论文标题: VIMA: General Robot Manipulation with Multimodal Prompts

作者: Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang 等 (Stanford, NVIDIA, Caltech, UT Austin)

发表: ICML 2023

阅读日期: 2025年1月

阅读目的: 理解多模态提示如何帮助机器人理解模糊指令，为"工业机械臂模糊指令语义消歧"课题提供技术参考

1 一句话总结

VIMA 是一个能够理解"图文混排"指令的机器人智能体——当纯文字描述不够精确时（比如"把那个零件放到那里"），可以在指令中直接嵌入图片来消除歧义，就像发微信时插入表情包一样自然。

2 为什么要读这篇论文？（与你课题的关系）

你的课题痛点

VIMA 提供的思路

"**把那个螺丝拧到那里**" 中的"那个"和"那里"怎么消歧？

用图片代替模糊的指代词：["把 [图片:螺丝A] 拧到 [图片:位置B]"]

工人可能不会精确描述物体特征

VIMA 的 Novel Concept Grounding 能让机器人"临时学习"新概念

SayCan 和 Inner Monologue 只处理纯文本指令

VIMA 扩展到**多模态**，可以结合图像、视频、文字

定位：如果说 SayCan 教会机器人"我能做什么"，Inner Monologue 教会机器人"边做边反思"，那么 VIMA 教会机器人"看图理解指令"。

3 核心概念解释

3.1 什么是 Multimodal Prompt（多模态提示）？

顾名思义：多模态 = 多种形式（文字 + 图片 + 视频），提示 = 给机器人的指令

生活比喻：

- 纯文字指令："把红色的杯子放到蓝色的盘子里" → 如果有多个红杯子呢？
- 多模态指令："把 [这个杯子] 放到 [这个盘子] 里" → 一目了然！

论文中的形式化定义：

$$\mathcal{P} := [x_1, x_2, \dots, x_l], \quad \text{其中 } x_i \in \{\text{文本, 图像}\}$$

也就是说，一个 prompt 就是**文本和图像交替排列的序列**。

3.2 VIMA 支持的 6 种任务类型

任务类型	英文名	举例	与消歧的关系
简单物体操作	Simple Object Manipulation	"把 [图:苹果] 放进 [图:碗]"	★★★ 直接消歧
视觉目标达成	Visual Goal Reaching	"把物体摆成 [图:目标场景]"	★★ 用图片指定目标
新概念落地	Novel Concept Grounding	"这是 blicket [图:A]，把 blicket 放进盒子"	★★★ 临时定义新词
单次视频模仿	One-shot Video Imitation	"像这样做 [视频:演示]"	★★ 动作消歧
视觉约束满足	Visual Constraint Satisfaction	"把物体扫进去，但不能碰到 [图:障碍物]"	★★ 约束消歧
视觉推理	Visual Reasoning	"把所有和 [图:这个] 相同纹理的物体放进 容器"	★ 属性推理

⌚ 对你最重要的是：Novel Concept Grounding（新概念落地）

这个任务测试的是：机器人能否“临时学会”一个新词汇的含义。比如：

| "This is a **blicket** [图片]. This is a **wug** [图片]. Put the wug into the blicket."

机器人从未见过 "blicket" 和 "wug" 这两个词，但通过图片，它能理解这是什么，并完成任务。

对你的启发：工业场景中，工人可能用方言或非标准名称描述零件（比如“那个长得像蘑菇的东西”），如果能结合图片，机器人就能“临时学习”这个非标准名称的含义。

3.3 Object-Centric Representation（物体中心表示）

顾名思义：不是把整张图片扔给模型，而是先“认出”图中有哪些物体，然后分别处理每个物体。

为什么这样做？

方法	输入	问题
原始像素	整张图片 → ViT → patches	模型需要自己学习"哪里有物体"，浪费算力
物体中心	检测框 + 裁剪图 → 物体 token	直接告诉模型"物体在这里"，学习更高效

技术细节：

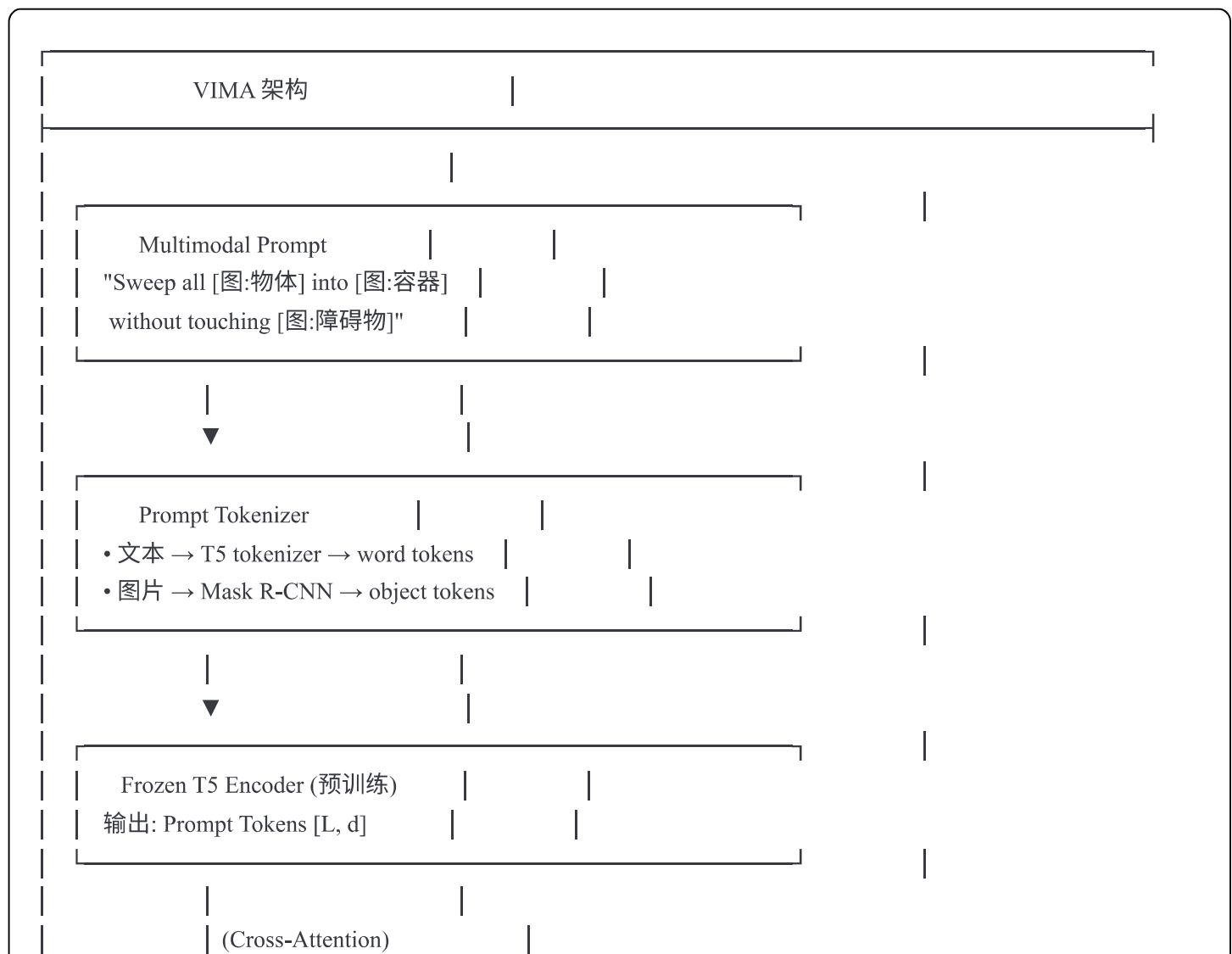
每个物体 token = Bounding Box 特征 + 裁剪图像的 ViT 特征

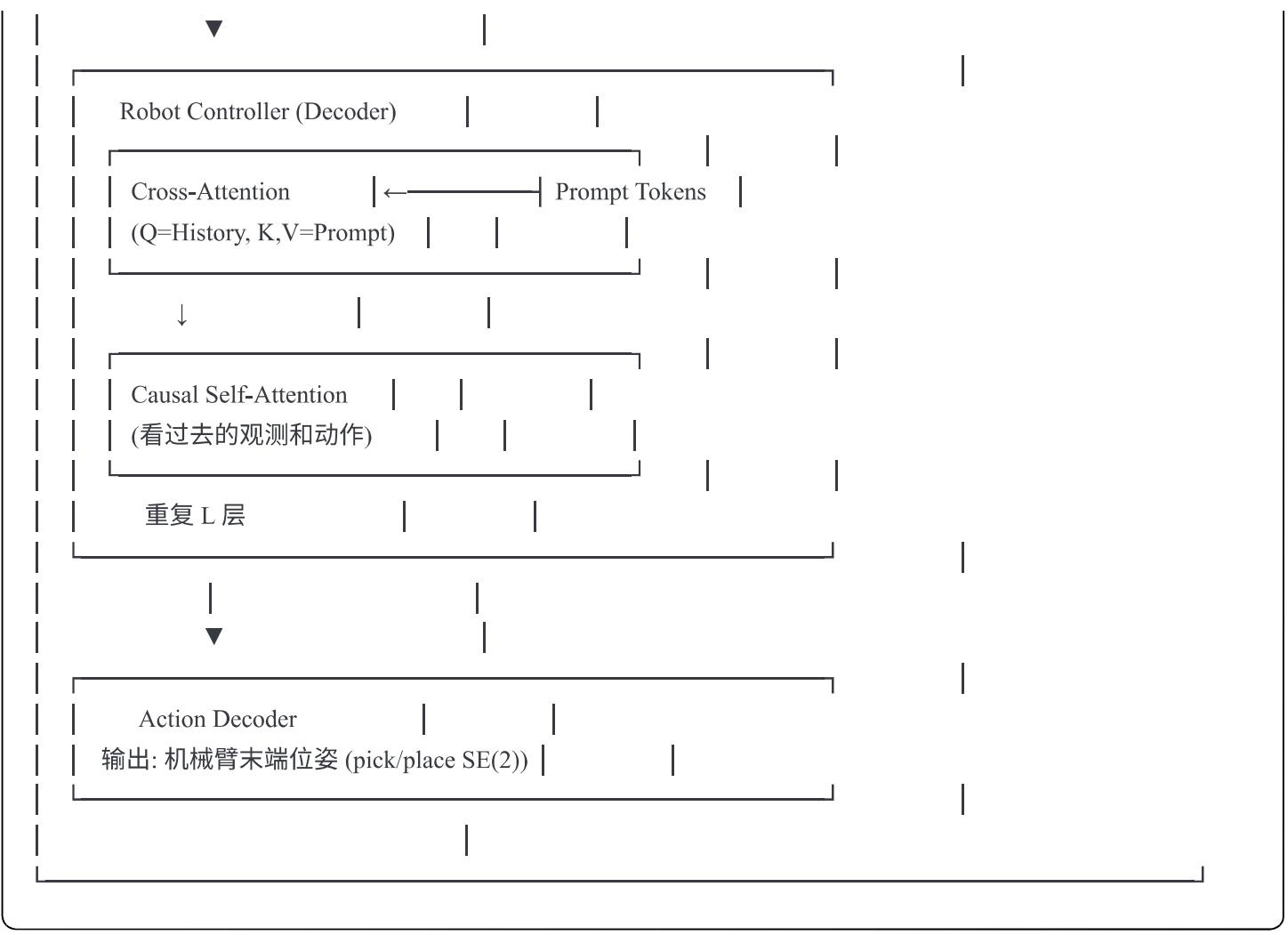
↓ ↓
[x_center, y_center, [视觉外观]
height, width]

实验结果：物体中心表示比 patch tokens 在所有泛化级别上都更好，尤其是在**数据量少**的时候优势更明显（10% 数据就能超过其他方法用 100% 数据的效果）。

4 模型架构详解

4.1 整体架构图（必须看懂！）





4.2 关键设计选择：Cross-Attention vs Direct Modeling

这是论文的核心贡献之一，必须理解！

方法	做法	问题
Direct Modeling (GPT 风格)	把 Prompt + History 拼成一个大序列，用 causal self-attention	序列太长，计算量大；Prompt 信息容易被“稀释”
Cross-Attention (VIMA 用的)	Prompt 单独编码，Controller 通过 cross-attention “查询”它	Prompt 信息保持完整；每一步都能“重新审视”指令

PyTorch 伪代码对比：

python

```

# ===== VIMA 的 Cross-Attention 方式 =====
def xattn_sequence_modeling(prompt_tokens, obs_tokens, act_tokens):
    # prompt_tokens: [L, d] - L 是 prompt 长度
    # obs_tokens: [T, d] - T 是时间步

    traj_tokens = interleave(obs_tokens, act_tokens) # [2T-1, d]
    x = traj_tokens

    for i in range(num_layers):
        # 关键: 用 trajectory 作为 Query, 向 prompt 提问
        x = x + cross_attn(q=x, kv=prompt_tokens) # "这个指令让我做什么? "
        x = x + ffn(x)
        x = x + causal_self_attn(q=x, kv=x)      # "之前我做了什么? "
        x = x + ffn(x)

    return x[-1] # 预测下一个动作

# ===== GPT 风格的 Direct Modeling =====
def causal_sequence_modeling(prompt_tokens, obs_tokens, act_tokens):
    # 直接拼接成一个超长序列
    x = concat([prompt_tokens, sep_token, interleave(obs_tokens, act_tokens)])

    for i in range(num_layers):
        x = x + causal_self_attn(q=x, kv=x) # 只有 self-attention
        x = x + ffn(x)

    return x[-1]

```

实验结论：Cross-Attention 在小模型和难泛化任务上优势明显。

4.3 动作空间 (Action Space)

VIMA 继承自 Ravens 仿真器的动作空间：

动作 = (pick_pose, place_pose)
= ((x1, y1, θ1), (x2, y2, θ2))

- **pick_pose**: 从哪里抓取
- **place_pose**: 放到哪里
- 每个 pose 是 SE(2): 平面位置 (x, y) + 旋转角度 θ

离散化：

- X 轴: 50 个 bins

- Y 轴：100 个 bins
 - 旋转：50 个 bins（四元数表示）
-

5 四级泛化评估协议（必须理解！）

这是 VIMA-Bench 的精髓，测试模型“举一反三”的能力：

Level	名称	训练时	测试时	难度
L1	Placement	见过相同的 prompt	只是物体摆放位置不同	★
L2	Combinatorial	见过所有物体和纹理	新的物体-纹理组合	★★
L3	Novel Object	见过一些物体	出现新物体、新纹理	★★★
L4	Novel Task	见过一些任务模板	全新的任务类型	★★★★

举例（以“把 X 放进 Y”为例）：

Level	训练	测试
L1	把红苹果放进蓝碗（苹果在左边）	把红苹果放进蓝碗（苹果在右边）
L2	红苹果+蓝碗, 绿香蕉+黄盘	红香蕉+黄碗（新组合）
L3	苹果、香蕉、碗、盘	草莓、杯子（新物体）
L4	“把 X 放进 Y”	“把所有和 X 相同颜色的物体放进 Y”（新任务）

VIMA 的优势：从 L1 到 L4 的性能下降幅度比其他方法小得多（只有其他方法的一半）。

6 实验结果速览

6.1 模型规模实验（Model Scaling）

模型大小	参数量	L1 成功率	L4 成功率
VIMA-2M	200万	~40%	~15%
VIMA-20M	2000万	~55%	~25%
VIMA-200M	2亿	~75%	~45%

结论：模型越大，泛化能力越强（符合 Scaling Law）。

6.2 数据效率实验 (Data Efficiency)

训练数据量	VIMA	VIMA-Gato	VIMA-Flamingo
100%	最好	次之	次之
10%	仍然很好	下降明显	下降明显
1%	超过其他方法100%数据的效果	很差	很差

结论：物体中心表示让 VIMA 非常数据高效。

6.3 消融实验核心发现

组件	选择	为什么
视觉编码器	Object Tokens > Image Patches > Single Image	物体级表示最好
Prompt 条件机制	Cross-Attention > Direct Modeling	小模型优势明显
语言编码器	T5 (预训练) > 从头训练	继承语义理解能力

7 关键公式

7.1 策略函数

$$\pi(a_t | \mathcal{P}, \mathcal{H})$$

- \mathcal{P} : 多模态 prompt (图文混排的指令)
- $\mathcal{H} = (o_1, a_1, o_2, a_2, \dots, o_t)$: 交互历史
- a_t : 当前要执行的动作

7.2 Cross-Attention 计算

$$\mathcal{H}' = \text{softmax} \left(\frac{Q_{\mathcal{H}} K_{\mathcal{P}}^{\top}}{\sqrt{d}} \right) V_{\mathcal{P}}$$

- $Q_{\mathcal{H}}$: 从交互历史计算的 Query ("我想问什么")
- $K_{\mathcal{P}}, V_{\mathcal{P}}$: 从 Prompt 计算的 Key/Value ("指令说了什么")
- d : 嵌入维度

直觉理解：每一步决策时，机器人都会“回头看看”指令，确保自己没有理解错。

7.3 训练目标（行为克隆）

$$\min_{\theta} \sum_{t=1}^T -\log \pi_{\theta}(a_t | \mathcal{P}, \mathcal{H})$$

就是标准的**负对数似然**，让模型预测的动作尽量接近专家演示的动作。

8 必看图表索引

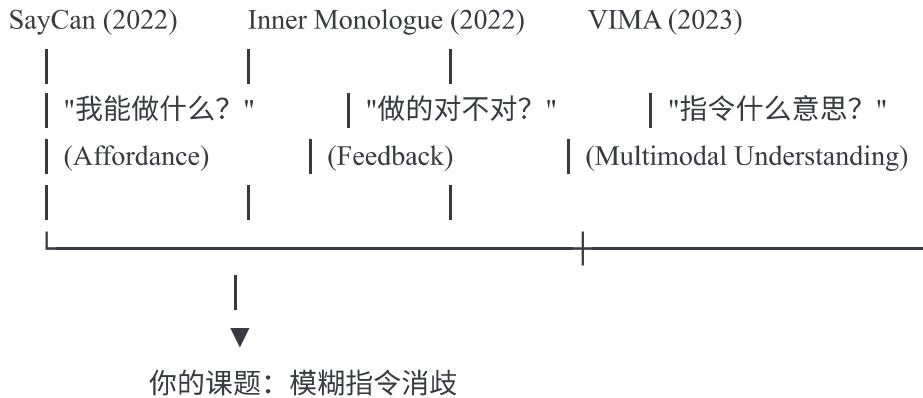
图表	页码	内容	重要性	为什么重要
Figure 1	p.2	多模态 Prompt 示例	★★★	理解什么是“图文混排指令”
Figure 2	p.3	四级泛化协议	★★★	理解评估方法
Figure 3	p.4	VIMA 架构图	★★★	核心架构必看
Figure 4	p.6	规模和数据实验	★★	了解 Scaling 性质
Figure 5	p.5	泛化性能下降对比	★★	VIMA 的鲁棒性
Figure 6	p.7	视觉编码器消融	★★	理解设计选择
Table 1	p.33	不同方法对比表	★★	快速了解 baseline

9 术语表

术语	中文翻译	解释
Multimodal Prompt	多模态提示	图文混排的指令
Object-Centric	物体中心的	以物体为单位处理，而不是像素
Cross-Attention	交叉注意力	用一个序列去“查询”另一个序列
Encoder-Decoder	编码器-解码器	先压缩再生成的架构
Behavior Cloning	行为克隆	直接模仿专家动作
Zero-shot Generalization	零样本泛化	不额外训练就能处理新情况

术语	中文翻译	解释
Novel Concept Grounding	新概念落地	临时学习新词汇的含义
SE(2)	平面刚体变换	2D 平移 + 旋转

10 与前序论文的联系



可能的融合方案：

1. 用 VIMA 的多模态 prompt 消解"那个/那里"的歧义
2. 用 SayCan 的 affordance 判断消歧后的指令能否执行
3. 用 Inner Monologue 的反馈机制确认理解是否正确

1 1 我的疑问（待后续解答）

概念层面

1. **多模态 vs 纯语言消歧：**如果用户不方便提供图片（比如语音指令），VIMA 的方法还适用吗？
2. **Novel Concept Grounding 的局限：**如果图片中有多个相似物体，只提供一张图片够吗？

技术层面

3. **Mask R-CNN 的依赖：**工业场景的零件可能和 COCO 数据集差异很大，检测器需要重新训练吗？
4. **Cross-Attention 的计算开销：**每一步都要和整个 Prompt 做 attention，长指令会不会很慢？

实验层面

5. **真实机器人实验：**论文主要在仿真器里做的，真实世界效果如何？
6. **与 RT-2、OpenVLA 的对比：**这些更新的大型 VLA 是否已经覆盖了 VIMA 的能力？

1 2 后续阅读计划

优先级	论文	原因
★★★	MOO (2023)	同样做多模态指令，但基于 RT-1
★★★	RT-2 (2023)	大型 VLA 的代表作
★★	OpenVLA (2024)	开源可实验
★★	Grounding DINO	更强的物体检测，可能替代 Mask R-CNN

1 3 代码资源

- **官方代码:** <https://github.com/vimalabs/VIMA>
- **Benchmark:** <https://github.com/vimalabs/VIMABench>
- **项目主页:** <https://vimalabs.github.io>

✓ 阅读检查清单

- 能用自己的话解释什么是"多模态提示"
- 能画出 VIMA 的大致架构图
- 理解 Cross-Attention 和 Direct Modeling 的区别
- 理解四级泛化评估协议
- 能说出 VIMA 对你课题的 3 个启发

笔记完成时间: 2025年1月