

问题	答案
CLIP 是什么？	一个把图像和文本映射到同一向量空间的预训练模型
怎么训练的？	对比学习：配对的拉近，不配对的推远
为什么强大？	4亿数据 + 开放词汇 = 超强泛化能力
对你有什么用？	理解模糊语言指令、定位目标物体的基础能力

CLIP是对比语言图像预训练，搭建了语言和这个图像识别的桥梁

CLIP 是让机器人“睁眼看世界”并“理解人话”的桥梁。

~~~~~

## 1 Success Detection (成功检测)

通俗理解：每个动作执行完，问一句“这步做成了吗？”

| 项目   | 说明                          |
|------|-----------------------------|
| 输入   | 执行前的图像 + 执行后的图像 + 技能名称      |
| 输出   | True (成功) 或 False (失败)      |
| 实现方式 | 仿真环境用规则判断；真实环境用 CLIP 微调的分类器 |

相当于是加了一个是否成功的判断

idea1

路线 B (工业向)

如果没有 RL Value Function，还能不能做 SayCan?

用几何可达性

用 grasp success 预测

用规则 / simulator

👉 这会直接变成你论文的 Method Section

idea2

你的课题核心是：“那个/那里”不清楚时怎么办。

Success Detection 给你一个特别实用的思路：

消歧不一定要在执行前一次性解决，可以先做一个“安全的试探动作”，然后用 success/failure 作为证据缩小歧义。

举个工业机械臂例子：

指令：“把那个螺丝放到那里”

你不确定“那个螺丝”是哪一个

你可以先 pick 候选 A

Success=False (没抓到/抓错/抓不起来) → 候选 A 被排除

Success=True → 继续确认“那里”的放置

这就是把“语言歧义”转成“可交互验证”的过程。

## 底层策略的位置

### 2) 它在 Inner Monologue 里处在什么位置? (层级分工)

你可以记这个最关键的分工:

- **LLM (InstructGPT)** : 决定 下一步做什么 (高层离散决策)
- **CLIPort (低层策略)** : 决定 这一步具体怎么做 (像素到动作/位姿)
- **Feedback (Object/Success/Scene)** : 告诉 LLM 刚才做得对不对、世界变成啥了

就像人类:

- 大脑: 决定“把黄色积木放进蓝碗”
- 手眼: 找准黄色积木和蓝碗的位置, 完成抓放
- 眼睛: 确认“放进去了吗?”

### MDETR vs CLIP (对比总结)

|    | CLIP              | MDETR                 |
|----|-------------------|-----------------------|
| 任务 | 图文匹配 (这张图和这句话配不配) | 图文定位 (这句话说的东西在图里哪个位置) |
| 输出 | 相似度分数             | 检测框 (bounding box)    |
| 粒度 | 整张图               | 图中的具体区域               |
| 类比 | "这张照片是不是在说猫? "    | "猫在照片的哪个位置? 框出来"      |

**MDETR** = Modulated Detection TransfoRmer

中文可以理解为: "被语言调控的检测器"

这个MDERT更像是专门用来处理这个位置信息的一种检测器

## LLM 在具身智能中的不同角色

根据你读的论文，我把 LLM 的角色分成 4 个层次：

### 角色 1：高层规划者（High-level Planner）

做什么：把复杂任务拆解成一步步的子任务

比喻：像一个项目经理，不亲自干活，但告诉团队“先做A，再做B，最后做C”

例子（SayCan）：

用户指令：“我把饮料洒了，能帮我清理一下吗？”

LLM 的规划：

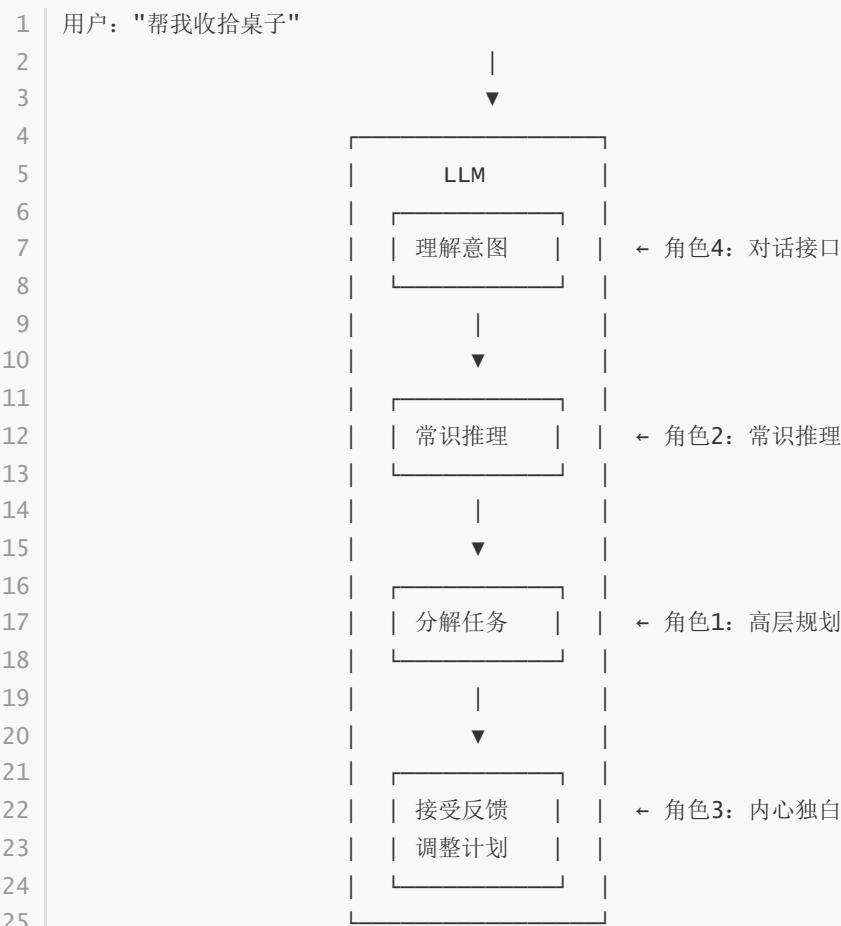
1. 找到海绵
2. 拿起海绵
3. 走到洒饮料的地方
4. 擦掉饮料
5. 把海绵放回原处



这个LLM像是一个这个CPU，相当于是一个大脑的决策层

在不同的层次里面有着不同的用法，你可以将他理解为这个大脑即可

这个就是LLM的功能



```

26 |
27 |
28 |    ▼    |    ▼    |    ▼
29   视觉系统      机械臂执行      成功检测
30   (MDETR)       (低层控制)     (反馈信号)
31

```

## 非常重要的三种反馈场景

### 2 Figure 2 (第4页): 三种反馈类型对比

这是论文最核心的概念图，必须看懂！

| 反馈类型                      | 中文     | 什么意思                 | 例子                                                                 |
|---------------------------|--------|----------------------|--------------------------------------------------------------------|
| Success Detection         | 成功检测   | 告诉 LLM "刚才那个动作做成了没"  | "Pick up coke" → Success: True <input checked="" type="checkbox"/> |
| Passive Scene Description | 被动场景描述 | 自动告诉 LLM "现在场景里有什么"  | "Scene: lime soda, coke, energy bar"                               |
| Active Scene Description  | 主动场景描述 | LLM 主动提问，人或 VQA 模型回答 | Robot Ask: "抽屉开了吗？" Human: "关着的"                                   |

## 局限性：

### 6 Section 5 (第8页): 局限性

这部分告诉你方法的边界在哪，非常重要！

| 局限性         | 说明                                  | 未来方向             |
|-------------|-------------------------------------|------------------|
| 低层策略是瓶颈     | LLM 再聪明，如果机械臂的基础技能不行，还是做不到了         | 需要更强的低层控制策略      |
| 依赖人类反馈      | 目前 Active Scene Description 需要人回答问题 | 未来用更好的 VQA 模型自动化 |
| 反馈可能不准      | 如果 MDETR 检测错了，LLM 会被误导              | 需要处理不确定性的机制      |
| LLM 有时会忽略反馈 | 偶尔 LLM 会"固执己见"，忽略环境反馈               | 需要更好的 prompt 设计  |
| 安全和伦理       | 没有考虑危险动作的检测                         | 需要增加安全模块         |

原来如此，这个prompt为什么要加上这个历史的记录，是为了抛弃已经尝试过的方法，更官方的去说，实际上他是保留了这个日志信息

### 先回答你的疑问

你说的“每次重新检测当前环境”，论文确实也做了！这就是【当前场景】那部分。

但问题是：只有当前场景是不够的。

让我用一个例子说明为什么。

### 🎯 核心问题：LLM 没有记忆！

LLM 就像一个每次对话都失忆的人。

每次你调用 LLM，对它来说都是“第一次见面”。它不知道：

- 之前发生了什么
- 为什么现在是这个状态
- 哪些事情已经尝试过了

## 这就是他的这个历史记录的意义

### 📊 对比总结：历史 vs 只有当前场景

| 信息类型      | 只有当前场景 | 有完整历史  |
|-----------|--------|--------|
| 知道现在的状态   | ✓ 知道   | ✓ 知道   |
| 知道哪些步骤完成了 | ✗ 不知道  | ✓ 知道   |
| 知道哪些尝试失败了 | ✗ 不知道  | ✓ 知道   |
| 知道失败的原因   | ✗ 不知道  | ✓ 可能知道 |
| 能避免重复错误   | ✗ 不能   | ✓ 能    |
| 能理解任务进度   | ✗ 困难   | ✓ 清晰   |

所以两者是互补的：

- **历史**：告诉 LLM “发生了什么”（因果关系）
- **当前场景**：告诉 LLM “现在是什么样”（最新状态）

可以，历史信息和当下信息的相辅相成

好一个文本区域对其，这个YOLO不能很好的根据这个语言找到这个对应的这个物体，但是这个MDETR可以，它对这个模糊语言拥有更加精准的识别能力

---

## MDETR 的默认工作方式：文本—区域对齐 (grounding)

MDETR是“你给一句话，它去找最符合那句话的框”。

它天然在训练时就被鼓励学：

- 文本里提到的属性、关系、位置词（纹着/贴着、左边那辆、红色的...）
- 和图像区域之间的对齐

所以它更像在回答：

“这句话描述的是图里哪一块？”

这就解释了你说的“泛化更强”的直觉：**它不是记住固定类别，而是在学“描述  $\leftrightarrow$  匹配”的能力。**

也就是我们所说的这个语言和对象的匹配能力

### 1) 你的结论怎么说更准确？

更准确的说法是：

- MDETR 是“文本条件检测/grounding”模型：

输入 = 图像 + 文本 query，输出 = 与文本语义最匹配的框。

所以它天生擅长做“语言  $\leftrightarrow$  区域”的对齐 (matching/grounding)。

- YOLO (典型版本) 是“类别条件检测 (closed-set detection)”：

输入 = 图像，输出 = 一组预定义类别的检测框。

它擅长“把图里有哪些已知类框出来”，但**不擅长直接根据自然语言描述去找那个物体**（除非你为它额外加模块或重新训练）。

一个是语言和图片的对齐，一个是这个图片和图片中检测物体的对齐，只能说是这个侧重点不一样