

# Ask-to-Clarify 论文精读笔记

**论文标题:** Ask-to-Clarify: Resolving Instruction Ambiguity through Multi-turn Dialogue

**中文标题:** 先问再做：通过多轮对话消解指令歧义

**作者:** Xingyao Lin, Xinghao Zhu, Tianyi Lu, Sicheng Xie, Hui Zhang, Xipeng Qiu, Zuxuan Wu, Yu-Gang Jiang

**单位:** 复旦大学、上海创新研究院、UC Berkeley

**发表:** arXiv 2509.15061 (2025)

**阅读日期:** 2025年1月

**阅读目的:** 学习"模糊指令语义消歧"的核心方法论

## 📍 一句话总结

Ask-to-Clarify 是一个**协作式具身智能框架**，当收到模糊指令时，先通过多轮对话问清楚用户意图，再端到端生成低层动作执行任务——核心理念是\*\*"先问清楚，再动手"\*\*。

## ⌚ 核心问题与动机

### 现有VLA的致命缺陷

现状	问题
当前VLA是被动执行者(Executor)	收到指令就执行，不管指令是否清晰
单向模式(one-way mode)	没有反馈、没有澄清、没有交互
真实场景指令往往模糊	"把那个水果放盘子里"——哪个水果？

### 论文的核心洞察

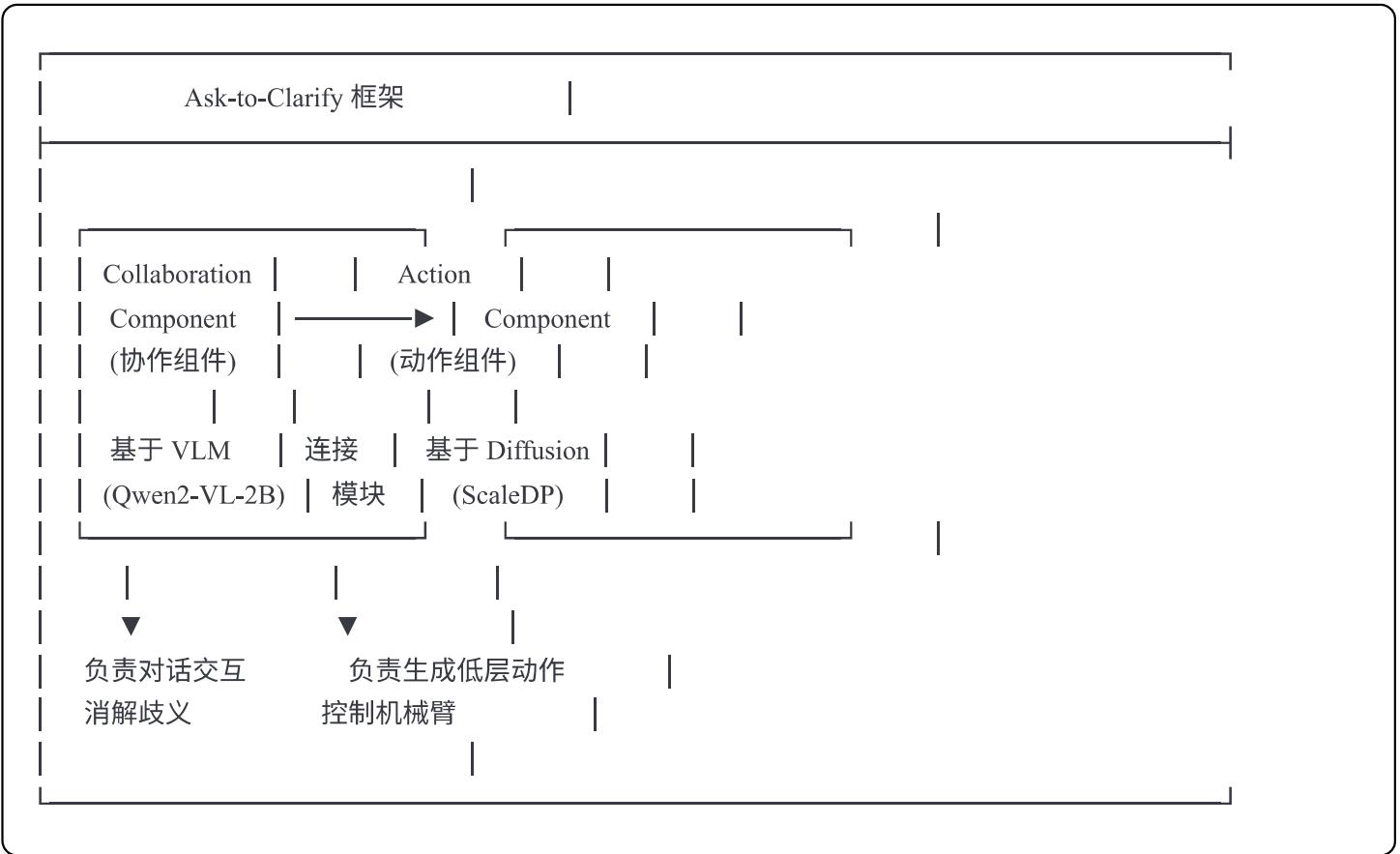
- ✖ Executor (执行者) : 指令 → 直接执行 → 可能失败
- ✓ Collaborator (协作者) : 指令 → 发现歧义 → 提问澄清 → 确认后执行 → 成功

### 生活类比

- **执行者:** 新来的实习生，师傅说"把那个拿过来"，随便拿一个，结果拿错了
- **协作者:** 聪明的实习生，师傅说"把那个拿过来"，会问"您说的是红色那个还是蓝色那个？"

## 🏗 框架架构（必须看懂！）

### 整体设计：双组件 + 连接模块



## 各组件职责

组件	技术实现	职责	类比
<b>Collaboration Component</b>	Qwen2-VL-2B (VLM)	理解指令、检测歧义、生成问题、推断正确指令	大脑：负责思考和沟通
<b>Action Component</b>	ScaleDP-Huge (Diffusion)	根据明确指令生成低层动作序列	手：负责执行
<b>Connection Module</b>	FiLM 调制	用指令信息调制视觉观测，为 Diffusion 提供条件	神经：连接大脑和手

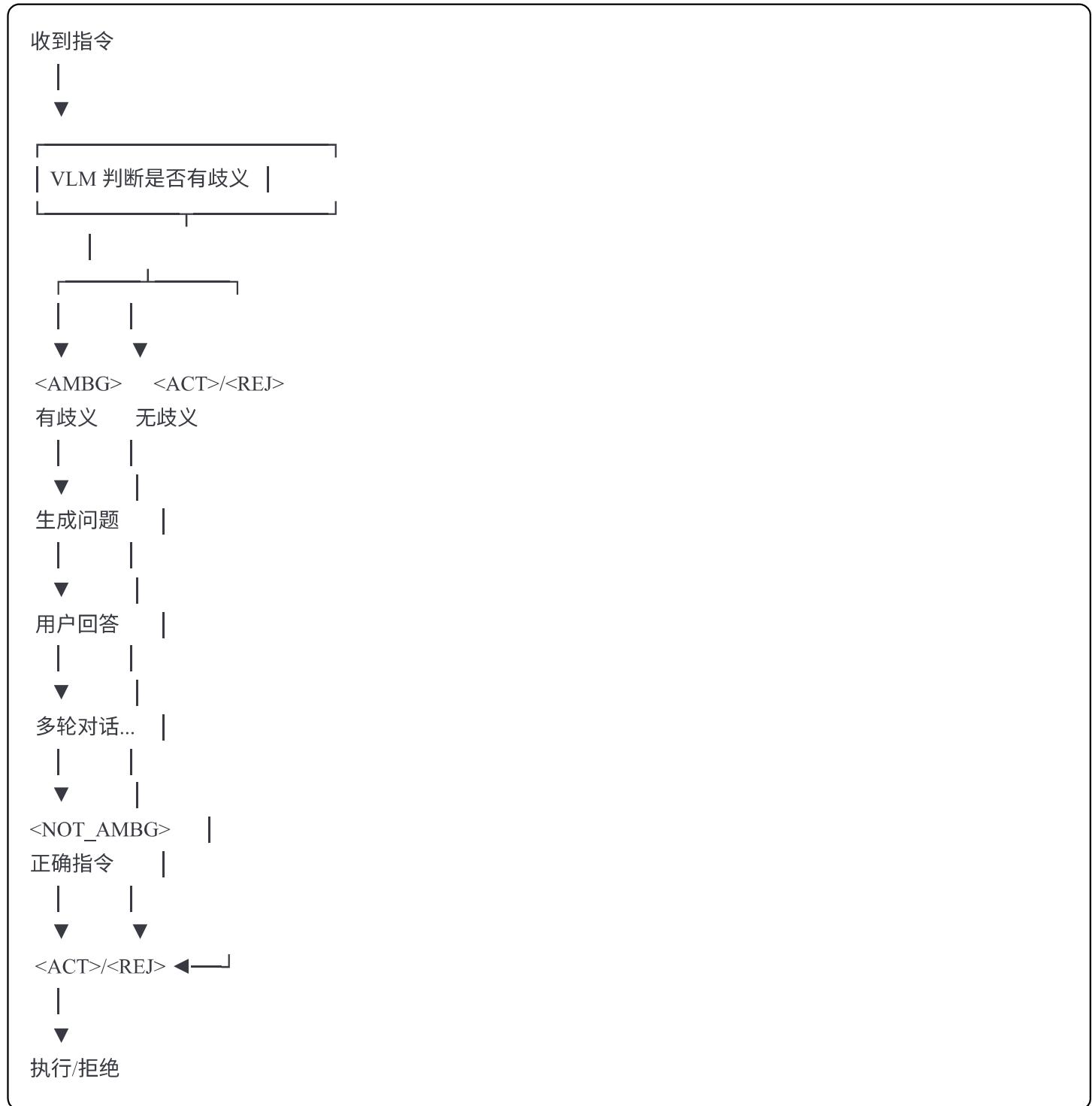
## 🔑 三大核心技术点

### 1 澄清触发条件：Signal Token 机制

论文引入了4个信号token来控制系统状态：

Token	中文	含义	触发时机
<AMBG>	有歧义	指令模糊，需要提问	检测到歧义时输出
<NOT_AMBG>	无歧义	已推断出正确指令	对话结束，歧义消除
<ACT>	执行	可以开始动作	目标物体在视野中
<REJ>	拒绝	不执行	目标物体不在视野中

### 状态转移图：



## 2 问句生成策略：数据驱动学习

不是手工设计规则，而是从数据中学习！

训练数据构建流程：

1. 收集包含多个相似物体的图片（如：两个只有颜色不同的方块）
2. 用 LLM (Qwen3-235B-A22B) 自动生成：
  - 模糊指令 (ambiguous instruction)
  - 问答对 (question-answer pairs)
  - 正确指令 (correct instruction)
3. 构建对话数据集用于训练 VLM

示例对话数据：

场景：桌上有苹果、桃子、橙子

[模糊指令] "把水果放到盘子里"

[对话轮次1]

机器人："您说的是哪个水果？是苹果、桃子还是橙子？" <AMBG>

用户："橙子"

[推断结果]

机器人："把橙子放到盘子里" <NOT\_AMBG>

[执行判断]

机器人:<ACT> (因为橙子在视野中)

## 3 澄清结果如何喂回 Action Generation

这是论文的**关键创新**：两阶段知识隔离训练策略 (Two-stage Knowledge-insulation Training)

Stage 1：学会问问题（训练协作能力）

训练目标：让 VLM 学会检测歧义、生成问题、推断正确指令

训练方式：

- 冻结 Vision Encoder
- 只微调 LLM 部分
- 使用对话数据

参数量：1.5B

学习率：1e-5

Epoch：50

Stage 2：学会执行动作（训练动作能力）

训练目标：让 Diffusion Expert 学会根据明确指令生成动作

训练方式：

- 【关键】冻结整个协作组件（知识隔离！）
- 只训练 Action Component + Connection Module
- 使用专家示教数据

参数量：978M

学习率：2e-5

Epoch：40

## 为什么要"知识隔离"？

问题：如果 Stage 2 继续训练 VLM，会发生什么？

答案：灾难性遗忘（Catastrophic Forgetting）！

VLM 会忘记 Stage 1 学到的对话能力。

解决：冻结 VLM，保护对话知识

用 Connection Module 补偿 VLM 和 Diffusion 之间的连接

## Connection Module 的作用

python

```
# 伪代码理解
def connection_module(instruction_tokens, observation_tokens):
    # 用指令信息去"调制"视觉观测
    # 相当于告诉Diffusion："重点关注指令提到的物体"
    modulated_obs = FiLM(observation_tokens, instruction_tokens)
    return modulated_obs # 作为Diffusion的条件输入
```

## 实验结果

### 主实验：8个真实世界任务

任务类型	具体任务	Ask-to-Clarify	$\pi_0$ (baseline)
放水果	Put Apple/Peach/Orange on plate	95.0%	91.7%
倒水	Pour from Red/Green/White cup	98.3%	93.3%
堆方块	Stack Blue/Yellow blocks	90.0%	57.5%

注意：Ask-to-Clarify 使用的是模糊指令，而 baseline 用的是明确指令！

## 消融实验：验证设计的必要性

Vision Encoder	LLM	Connection Module	Action Component	成功率
训练	训练	X	训练	0%
训练	训练	✓	训练	0%
冻结	冻结	✓	训练	90%
冻结	训练	X	训练	0%

**结论：**必须同时满足：

1.  冻结 VLM (知识隔离)
2.  使用 Connection Module

## 鲁棒性测试

条件	Ask-to-Clarify	$\pi_0$
正常光照	90.0%	57.5%
低光照 (关50%灯)	80.0%	22.5%
有干扰物 (苹果+石榴)	80.0%	65.0%

## 🔗 与你课题的直接联系

你的课题：工业机械臂的模糊指令语义消歧

Ask-to-Clarify 的设计	对你课题的启发
Signal Token 机制	可以用类似的状态机设计歧义检测触发器
数据驱动的问句生成	需要构建工业场景的对话数据集
两阶段知识隔离训练	如果你也用 VLM + Policy，可以借鉴
Connection Module (FiLM)	指令信息如何调制视觉特征

## 可以直接借鉴的点

1. 评估指标设计：

- 澄清后成功率
- 需要几轮对话才能消歧
- 拒绝执行的准确率（该拒绝时拒绝）

## 2. 任务设计：

- 先对话消歧 → 再执行任务
- 两步成功才算成功

## 3. 数据构建：

- 用 LLM 自动生成对话数据
- 收集包含相似物体的场景图片

---

## 📖 术语表

英文术语	中文翻译	解释
Ask-to-Clarify	先问再做	本文提出的框架名
Executor	执行者	被动执行指令的agent
Collaborator	协作者	主动交互的agent
Knowledge Insulation	知识隔离	冻结已训练模块防止遗忘
Signal Token	信号令牌	控制系统状态的特殊token
Connection Module	连接模块	连接VLM和Diffusion的桥梁
FiLM	特征线性调制	用语言调制视觉特征的技术
Catastrophic Forgetting	灾难性遗忘	学新知识时忘掉旧知识
End-to-end	端到端	从输入直接到输出，无需中间步骤

---

## gMaps 必看图表索引

图表	页码	内容	重要性
Figure 1	p.1	Executor vs Collaborator 对比	★★★
Figure 2	p.3	两阶段训练流程图	★★★

图表	页码	内容	重要性
Figure 3	p.4	Signal Detector 工作流程	★★★
Figure 4	p.4	实验任务示例图	★★
Table I	p.5	主实验结果对比	★★★
Table IV	p.5	训练策略消融实验	★★★
Table V	p.5	协作能力测试	★★

## ？我的疑问与思考

### 概念层面

1. **歧义检测的阈值**: 论文没有明确说 VLM 如何判断"有歧义" vs "无歧义", 是隐式学习的吗?
2. **问句的信息增益**: 论文的问句生成是数据驱动的, 没有显式计算信息增益, 这样够优吗?

### 技术层面

3. **Connection Module 细节**: FiLM 调制具体怎么做的? 需要看原始 FiLM 论文
4. **对话轮数限制**: 实验中平均几轮对话? 有没有设置最大轮数?

### 实验层面

5. **工业场景适用性**: 论文任务是家庭场景 (放水果、倒水), 工业场景 (螺丝、工具) 会更复杂
6. **数据集规模**: 每个任务只用 10 个专家示教, 够吗? 工业场景可能需要更多

### 可扩展方向

7. **主动学习**: 能否让机器人从错误中学习"什么时候该问"?
8. **多模态歧义**: 不只是物体歧义, 还有位置歧义、动作歧义 ("轻轻放"vs"用力按")

## ✓ Day 3 检验标准自查

能说出 3 种触发澄清的条件

- Signal Token `<AMBG>` 被 VLM 输出
- 场景中存在多个相似物体
- 指令包含指代词 ("那个"、"这个")

能设计一个简单的澄清问题

- "您指的是左边的红色螺丝还是右边的蓝色螺丝？"
  - 理解"信息增益"的直觉含义
  - 好问题 = 能最大程度减少候选数量的问题
- 

**笔记完成时间:** 2025年1月