

复习后下一步：少而精的“补齐知识点 + 精读论文”清单（面向工业机械臂模糊指令语义消歧）

你的已读基础：SayCan (affordance/world grounding)、Inner Monologue (闭环语言反馈)、VIMA (多模态提示与泛化)、VLA Survey (全景地图)。

目标：用最少的知识点 + 最少的论文，把你带到 2024-2025 年 VLA/具身智能的主流技术栈，并且直接服务你的课题：模糊指令语义消歧 + 工业机械臂执行。

A. 需要补齐的知识点（少而精，5 个就够）

1) Offline Imitation Learning (离线模仿学习) + BC (行为克隆)

中文意思：用离线专家演示数据直接学策略 (policy)。

为什么你必须补：OpenVLA / RT-X / VIMA-BENCH 这条线基本都以离线数据为核心；你要做工业机械臂消歧，后续很可能也要用离线日志/示教。

你只需要掌握这几件事：

- BC 的目标：最大化动作似然（等价于最小化 NLL）
- 分布偏移 (distribution shift) 与闭环误差累积（为什么需要闭环反馈/重规划）
- 评估：success rate、子目标成功率、交互轮次（消歧任务很关键）
练习 (2小时)：用一个小数据集（哪怕是 toy）实现 BC：`obs -> action`，跑通训练/验证/测试的完整 pipeline。

2) Affordance/Value Function (可供性/价值函数) 再“补一层”

你已经从 SayCan 理解了“能不能做”。这里再补：可供性 ≈ 在当前状态下某技能成功概率的估计。

为什么你必须补：工业场景的“能不能做”常常更复杂：夹具限制、工具约束、安全区、力控等。

你只需要掌握：

- MDP/价值函数基本定义（你已会）
- 离线/在线获取成功标签（真实工业里常用：规则 + 少量人工校验）
- 价值函数如何变成“安全约束/先验”去做消歧（例如：两个候选目标都符合语言，但只有一个可抓取）

3) Open-vocabulary Perception & Visual Grounding (开放词汇感知 + 视觉指代落地)

中文意思：“用自然语言去找物体/区域/属性”，包括 open-set detection (开放集检测) 与 referring expression (指代表达理解)。

为什么你必须补：你的课题里“那个/那里/左边那个红的”本质就是 grounding + disambiguation。

你只需要掌握：

- “检测” vs “指代落地 (REC) ”：前者找类别，后者找文本描述对应的那个实例
- 一个强工具链：GroundingDINO (找候选) → (可选) SAM 分割 → 作为后续规划/控制的输入

4) Diffusion Policy (扩散策略) : 生成式动作建模的主流路线

中文意思: 把动作序列当作“要生成的样本”，用扩散模型逐步去噪生成动作。

为什么你必须补: 2023 后操作任务里扩散策略成为强基线，很多新 VLA/控制方法会把 diffusion 当作底座或重要对比。

你只需要掌握:

- 扩散的直觉：从噪声动作一步步“修正”到可执行动作
- 为什么它对机器人好：能表示多峰动作分布（同一句话可能有多种可行抓取/路径）
- receding horizon (滚动时域) 怎么做闭环执行（避免一次生成太长导致漂移）

5) Ambiguity-aware Interaction (面向歧义的交互) : 何时该“问一句”

中文意思: 当指令不够明确时，系统需要识别不确定性，并发起澄清对话 (clarifying question)。

为什么你必须补: 你的课题就是“模糊指令消歧”；真正的工业系统往往比“强行猜”更安全的是“问一句”。

你只需要掌握:

- 歧义检测：什么时候系统该承认“不确定”
- 选择问题：问哪一句最省交互成本却最大幅度减少不确定性（信息增益直觉）
- 评估指标：平均澄清轮数、澄清后成功率、误澄清率（问了反而更糟）

B. 少而精的论文清单 (建议 7 篇，按“对你课题的直接收益”排序)

每篇都给你：读什么、跳什么、与你课题怎么连。

1) OpenVLA (2024) — *An Open-Source Vision-Language-Action Model*

你读什么:

- 训练数据/训练范式（为什么它能在多机器人上工作）
- fine-tuning (参数高效微调) 的实践建议
- 任务与评估（你要学它怎么“定义成功”）

你可以先跳过: 模型细节里很底层的工程实现

与你课题的连接:

- 你后续做实验最缺的是“能跑的强基线”，OpenVLA 给你一个开源起点（训练/微调代码也有）
链接：arXiv 2406.09246；代码：openvla/openvla

2) Open X-Embodiment & RT-X (2023) — *Robotic Learning Datasets and RT-X Models*

你读什么:

- 数据集怎么统一不同机器人的 demonstration 格式
- “cross-robot transfer”为何成立（你做工业臂也会遇到跨工位/跨工具迁移）

你可以跳过: 某些模型细节与 ablation

与你课题的连接:

- 消歧任务很可能需要你自己构建数据，RT-X 给你“数据组织与标准化”的参考范式

链接: arXiv 2310.08864; 项目页: robotics-transformer-x

3) Diffusion Policy (2023) — *Visuomotor Policy Learning via Action Diffusion*

你读什么:

- 扩散策略如何做控制 (receding horizon)
- 与传统 BC/Transformer policy 的对比优势

你可以跳过: 扩散推导里过细的数学符号 (先抓住直觉和训练/推理流程)

与你课题的连接:

- “语言模糊 → 多解”, 扩散天然适合表示多模态动作; 你做消歧时可以把它当强执行器底座
链接: arXiv 2303.04137; 项目页: diffusion-policy.cs.columbia.edu

4) Grounding DINO (2023) + Grounding DINO 1.5 (2024) — *Open-set detection / grounding*

你读什么:

- 它怎么把“语言”融合进检测器, 实现 open-set + referring
- 在 RefCOCO 等指代表达数据上的表现 (与你的“那个/那里”最像)

你可以跳过: 大规模训练细节 (先会用、会评估)

与你课题的连接:

- 给你的消歧系统提供“候选实体集合”: 先找出 3 个可能的“那个”, 再让 planner/交互去消歧
链接: arXiv 2303.05499; arXiv 2405.10300; 代码: IDEA-Research/GroundingDINO

5) RT-2 (2023) — *Vision-Language-Action Models Transfer Web Knowledge to Robotic Control*

你读什么:

- “VLM + 机器人动作 token”的端到端范式
- 它为什么能把 web 的视觉概念迁移到机器人动作上

你可以跳过: 一些离散化动作细节的工程实现

与你课题的连接:

- 你要跟上时代潮流, 必须理解“把 VLM 直接变成 policy”的主流范式; 这篇是里程碑
链接: arXiv 2307.15818; DeepMind 博客: RT-2

6) PaLM-E (2023) — *An Embodied Multimodal Language Model*

你读什么:

- 把“连续传感/状态”塞进语言模型的思路 (工业臂也需要状态: 夹爪开合、力矩、工位编号)
- 多任务联合训练与正迁移 (positive transfer)

你可以跳过: 超大模型规模相关的细节 (先看范式)

与你课题的连接:

- “消歧”不仅靠视觉语言, 还要靠**机器人自身状态**; PaLM-E 给你一个把状态纳入 LLM 的范式
链接: arXiv 2303.03378; 项目页: palm-e.github.io

7) Ask-to-Clarify (2025) — *Resolving Instruction Ambiguity through Multi-turn Dialogue*

你读什么：

- 什么时候触发澄清、澄清问句怎么生成
- 澄清后怎么把结果喂回 action generation (尤其是端到端 VLA + 对话)

你可以跳过：如果有很复杂的模块堆叠，先抓“框架与指标”

与你课题的连接：

- 这就是你的核心问题“模糊指令语义消歧”的近期代表路线：**先问清楚再动手**
链接：arXiv 2509.15061

可选加读（只加 1 篇，别贪多）：

NeurIPS 2024 Aligning Robots' Uncertainty with Inherent Task Ambiguity (把“不确定”建模得更像真正的歧义)

C. 推荐阅读顺序（10 天内快速跟上）

- Day 1–2: Grounding DINO (你要先能“找出那个”)
- Day 3: Ask-to-Clarify (你要先能“问一句”)
- Day 4–5: OpenVLA (你要先能“跑一个强基线”)
- Day 6: RT-X 数据范式 (你要知道未来数据怎么组织)
- Day 7–8: Diffusion Policy (你要知道当前强控制器怎么做)
- Day 9–10: RT-2 / PaLM-E (二选一先读；另一个当扩展)

D. 立刻能做的 3 个“小而硬”实验（直接服务你的课题）

1. **候选生成实验（Grounding）**：对工业场景图片 + 模糊指令（“把那个螺丝拧到那里”）输出 Top-K 候选实体（螺丝/孔位/工具）。
2. **澄清问句实验（Disambiguation）**：当 Top-K 置信度接近时，自动生成一个最省轮次的澄清问题（例如“你指的是左边还是右边那颗？”）。
3. **执行对比实验（Policy）**：用 BC (你自己写) vs OpenVLA (微调/提示) vs Diffusion Policy (如能跑) 比较：成功率、澄清轮数、误执行率。

Reference Links (建议直接收藏)

- OpenVLA: <https://arxiv.org/abs/2406.09246> | <https://github.com/openvla/openvla>
- RT-X / Open X-Embodiment: <https://arxiv.org/abs/2310.08864> | <https://robotics-transformer-x.github.io/>
- Diffusion Policy: <https://arxiv.org/abs/2303.04137> | <https://diffusion-policy.cs.columbia.edu/>
- GroundingDINO: <https://arxiv.org/abs/2303.05499> | <https://github.com/IDEA-Research/GroundingDINO>
- GroundingDINO 1.5: <https://arxiv.org/abs/2405.10300>
- RT-2: <https://arxiv.org/abs/2307.15818> | <https://deepmind.google/blog/rt-2-new-model-translates-vision>

[-and-language-into-action/](#)

- PaLM-E: <https://arxiv.org/abs/2303.03378> | <https://palm-e.github.io/>
- Ask-to-Clarify: <https://arxiv.org/abs/2509.15061>