

OpenVLA 论文精读笔记

论文标题: OpenVLA: An Open-Source Vision-Language-Action Model

作者: Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti 等

机构: Stanford, UC Berkeley, Toyota Research Institute, Google DeepMind

链接: <https://arxiv.org/abs/2406.09246>

代码: <https://github.com/openvla/openvla>

阅读日期: Day 4-5

阅读目的: 掌握一个"能跑的强基线", 为后续消歧实验做准备

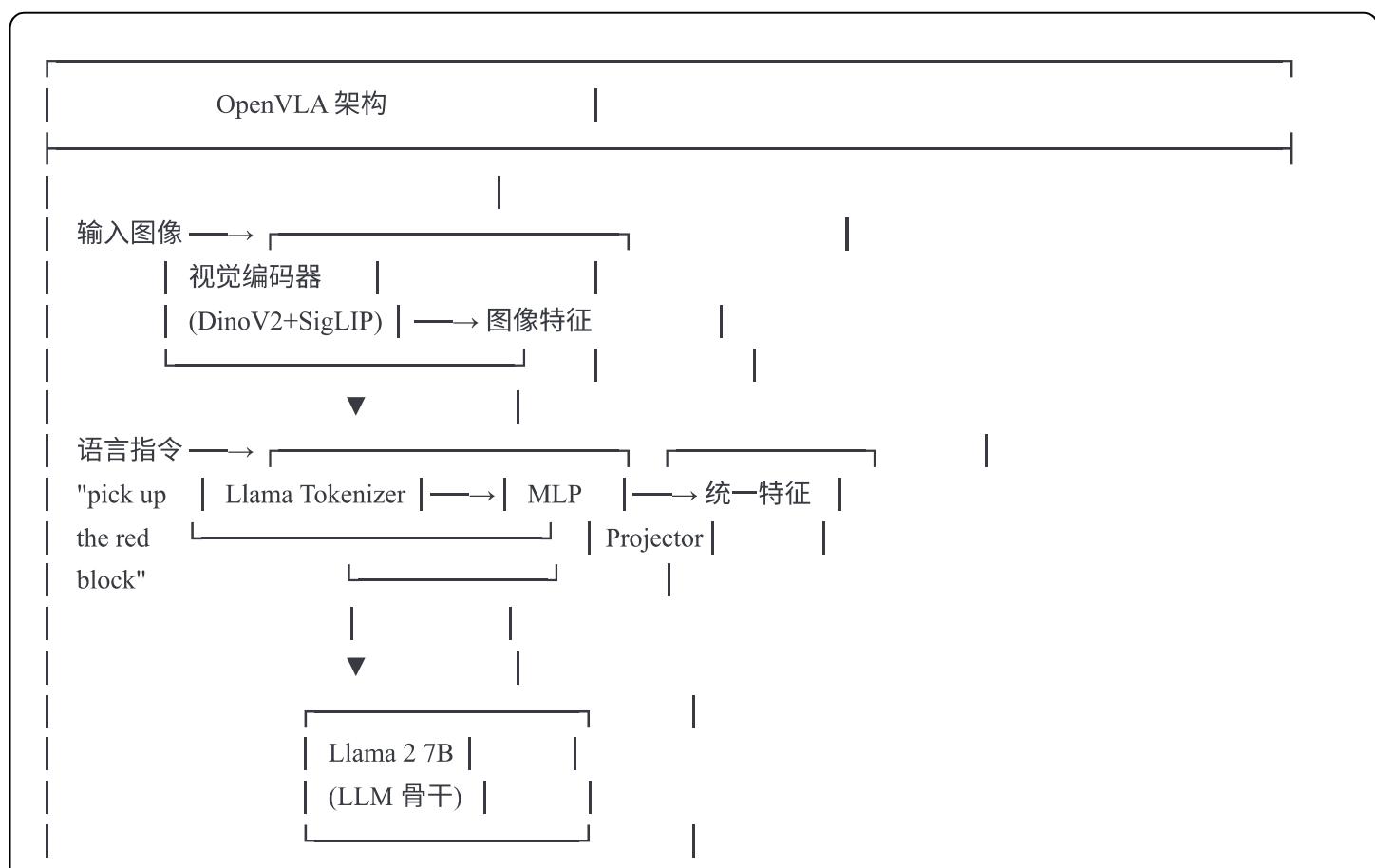
🎯 一句话总结

OpenVLA 是一个 7B 参数的开源 VLA 模型，把预训练的视觉-语言模型（VLM）直接微调成机器人控制策略，输入图像+语言指令，输出机器人动作。

核心卖点：

- **✓ 完全开源** (权重、代码、数据配方都有)
- **✓ 性能强** (比 55B 的 RT-2-X 还好 16.5%)
- **✓ 能微调** (LoRA 微调只需 1 张 A100, 10-15 小时)

📐 模型架构 (必看 Figure 2)



↓

动作 Token
 $[\Delta x, \Delta y, \Delta z, \Delta roll, \Delta pitch,$
 $\Delta yaw, \Delta gripper]$
 (7 维机器人动作)

三大核心组件

组件	具体模型	作用	中文解释
视觉编码器	DinoV2 + SigLIP (融合)	提取图像特征	DinoV2 擅长空间细节, SigLIP 擅长语义理解, 两个合起来更强
投影器	2层 MLP	对齐视觉和语言空间	把图像特征"翻译"成语言模型能理解的格式
语言模型骨干	Llama 2 7B	理解指令 + 生成动作	核心大脑, 既理解"拿红色方块", 又输出具体动作

🔑 关键设计：动作离散化（Action Tokenization）

问题：语言模型只能输出"词"（离散的 token），但机器人动作是连续的数值（比如 x 方向移动 0.03m）

解决方案：把连续动作"切成格子"

```
python

# 伪代码：动作离散化
# 假设 x 方向的动作范围是 [-0.1m, 0.1m]
# 切成 256 个格子 (bin)

bins = 256
action_range = (-0.1, 0.1) # 用数据的 1% 和 99% 分位数, 而不是 min/max (避免异常值)

# 连续动作 → 离散 token
continuous_action = 0.03 # 比如向右移动 3cm
bin_index = int((continuous_action - action_range[0]) / (action_range[1] - action_range[0]) * bins)
# bin_index ≈ 166

# 离散 token → 连续动作 (推理时反向操作)
```

为什么用 256 个 bin？

- 精度足够（每个 bin 约 0.8mm，机器人控制够用）
- 正好可以用一个字节表示（0-255）

📊 训练数据 (Section 3.3)

数据来源：Open X-Embodiment (OpenX)

指标	数值
总轨迹数	970k (97万条机器人演示)
数据集数量	70+ 个不同的机器人数据集
机器人类型	多种 (WidowX、Franka、Google Robot 等)

数据清洗 (很重要！)

论文发现原始数据有很多问题，做了以下清理：

原始 OpenX 数据 → 过滤 → 清洗后的训练数据

- 去除全零动作（机器人没动但记录了数据）
- 去除异常轨迹
- 统一动作空间格式

💡 对你课题的启发：

如果你以后要做工业机械臂的数据，也要注意数据清洗！比如过滤掉机械臂故障时的数据、人工标注错误的数据等。

👉 训练配方 (Section 3.4)

训练超参数

参数	值	备注
学习率	2e-5	固定，不用 warmup
训练轮数	27 epochs	比普通 LLM/VLM 多很多！
批大小	2048	
GPU	64 × A100	
训练时间	14 天	约 21,500 A100-hours

参数	值	备注
图像分辨率	224 × 224	384 没有提升，但慢 3 倍

🔥 关键发现

1. 必须微调视觉编码器

普通 VLM 训练：冻结视觉编码器 → 效果好

VLA 训练：冻结视觉编码器 → 效果差 ✗

微调视觉编码器 → 效果好 ✓

原因：机器人控制需要非常精细的空间信息（比如螺丝的精确位置），预训练的视觉特征不够“细腻”。

2. 需要更多训练轮数

模型类型	典型训练轮数
LLM	1-2 epochs
VLM	1-2 epochs
VLA	27 epochs

原因：动作预测比语言预测更难学，需要反复“看”数据。

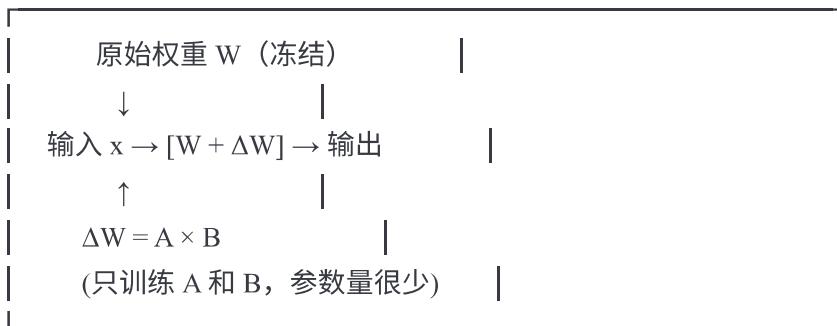
🔧 高效微调方法 (Section 5.3) —— 最实用的部分！

微调方法对比

方法	成功率	训练参数量	显存占用	推荐？
Full FT (全量微调)	69.7%	7188M (100%)	163 GB	效果最好但太贵
Last layer only	30.3%	465M	51 GB	✗ 效果很差
Frozen vision	47.0%	6760M	156 GB	✗ 效果差
Sandwich	62.1%	914M	64 GB	还行
LoRA (r=32)	68.2%	98M (1.4%)	60 GB	✓ 强烈推荐

LoRA 微调详解

LoRA (Low-Rank Adaptation, 低秩适应)：只训练一小部分“适配器”参数



LoRA 的核心思想：

- W：原始的 7B 参数，冻结不动
- A：降维矩阵， $\text{shape} = [\text{hidden}, r]$ ，r 很小（如 32）
- B：升维矩阵， $\text{shape} = [r, \text{hidden}]$
- 只训练 A 和 B，参数量 $= 2 \times \text{hidden} \times r \ll \text{原始参数量}$

⌚ 实用建议

python

```
# 推荐的 LoRA 微调配置
lora_config = {
    "rank": 32,           # r=32 足够, r=64 没有提升
    "target_modules": "all_linear", # 应用到所有线性层
    "learning_rate": 2e-5,
}

# 资源需求
# - 1 张 A100 GPU
# - 10-15 小时
# - 比全量微调节省 8 倍计算量!
```

⚡ 高效推理 (Section 5.4)

量化推理对比

精度	成功率	显存占用	推理速度
bfloat16 (默认)	71.3%	16.8 GB	基准
int8	58.1%	10.2 GB	慢 (不推荐)
int4	71.9%	7.0 GB	和 bf16 差不多

结论：用 int4 量化，显存砍半，性能不掉！

不同 GPU 的推理速度

GPU	推理频率
RTX 4090	~6 Hz
A100	~8 Hz
H100	~10 Hz

注意：OpenVLA 目前约 6Hz，对于需要 50Hz 的高频控制任务（如灵巧手）还不够快。

评估指标与实验结果

主要评估指标

指标	英文	含义	对消歧任务的意义
Success Rate	成功率	任务是否完成	最终指标
Generalization	泛化能力	换物体/换场景还能不能做	消歧后能否正确执行

泛化能力的四个维度

维度	英文	测试内容	例子
视觉泛化	Visual generalization	换光照、换背景	训练时白天，测试时晚上
运动泛化	Motion generalization	物体位置变了	杯子从左边挪到右边
物理泛化	Physical generalization	换物体形状/材质	训练用塑料杯，测试用玻璃杯
语义泛化	Semantic generalization	换物体类别	训练"拿苹果"，测试"拿橙子"

OpenVLA vs RT-2-X (核心对比)

指标	OpenVLA	RT-2-X
参数量	7B	55B
平均成功率	更高 16.5%	基准
语义泛化	稍弱	更强

指标	OpenVLA	RT-2-X
开源	<input checked="" type="checkbox"/> 是	<input type="checkbox"/> 否
可微调	<input checked="" type="checkbox"/> 是	<input type="checkbox"/> 否

为什么 OpenVLA 更小但更好？

- 训练数据更大 (970k vs 350k)
- 数据清洗更仔细
- 融合视觉编码器 (DinoV2 + SigLIP)

🔗 与我的课题的联系

课题：工业机械臂模糊指令语义消歧

OpenVLA 能给我什么？

用途	具体做法
强基线	用 OpenVLA 做"无消歧"的基线，对比你的"有消歧"方法
动作执行器	你的消歧模块输出明确指令后，交给 OpenVLA 执行
微调起点	在工业场景数据上 LoRA 微调 OpenVLA

可能的实验设计



OpenVLA 的局限（对消歧任务）

局限	影响	可能的解决方案
只支持单张图像	无法利用多视角消歧	等待后续版本 or 自己改
无主动询问能力	不会说“你指的是哪个？”	需要外接消歧模块
推理速度有限	6Hz 可能不够工业实时控制	用 int4 量化 + 更好的 GPU

✓ 检验标准 (Day 4-5)

学完这篇论文，你应该能回答：

- OpenVLA 的三大组件是什么？各自的作用？
- 动作是如何从连续值变成 token 的？
- 为什么 VLA 训练需要比 VLM 更多的 epoch？
- LoRA 微调的优势是什么？推荐的 rank 是多少？
- int4 量化对性能有影响吗？
- OpenVLA 和 RT-2-X 的主要区别？

📖 术语表

术语	中文	解释
VLA	视觉-语言-动作模型	输入图像+语言，输出动作
VLM	视觉-语言模型	输入图像+语言，输出语言（如 GPT-4V）
Action Tokenization	动作词元化	把连续动作变成离散 token
LoRA	低秩适应	一种高效微调方法，只训练少量参数
Quantization	量化	用更少的比特存储权重，换取更小显存和更快推理
OpenX	开放 X 体现	一个大规模机器人数据集联盟
DinoV2	—	Meta 的自监督视觉模型，擅长空间细节
SigLIP	—	Google 的视觉-语言对齐模型，擅长语义理解
Epoch	训练轮次	模型“看”完整个数据集一遍

⌚ 资源链接

- 论文: <https://arxiv.org/abs/2406.09246>
 - 代码: <https://github.com/openvla/openvla>
 - 项目主页: <https://openvla.github.io>
 - HuggingFace 模型: [openvla/openvla-7b](#)
-

笔记完成时间: Day 4-5

下一步: 尝试跑通 OpenVLA 推理 demo