OpenVLA : An Open-Source Vision-Language-Action model

Team 5 Kim Jaemin Jung Geumyoung



Learning Trajectory Priors with Diffusion, Guided Sampling for Optimal Paths

Recap Problem:

- Optimization-based planners need good initialization
- Sampling-based planners can be inefficient

• MPD's Solution:

- Learn Prior : Diffusion models learn a generative model of expert trajectories
- **Guided Sampling :** Sample from the posterior by guiding the diffusion reverse process with gradients from motion planning costs





Contents

- 1. Introduction
- 2. OpenVLA
 - Architecture
 - \circ Training
- 3. Experiment
 - Out-of-box generalization
 - Fine-tuning
 - Quantization
- 4. Limitations & Summary
- 5. Quiz



Introduction | Motivation

How AI used to work



segmentation model





drilling platform

er ship tainer ship A lifeboat sh amphibian ou fireboat



captioning model

A group of people shopping at an outdoor market.



bananas

"Horrible services. The room was dirty and unpleasant. Not worth the money."



6

NEGATIVE

Processing Natural language processing (NLP) is a subfield of impurities, computer sources, and antificial intelligence concerned with the interactions between computers and human remotions between computers and human language, in particular how to original computers to process and analyze large amounts of natural language data. The result is a computer capatile of "understanding" the contents of documents, including the contents of documents, including the language. contentual muscless of the landuate within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themastures summarization model Natural Language Processing National language processing (NUP) a subfeet of imposites, composite torestore, and artificial intelligence compared with the interactions between compares and human language in particular from the program computers to process and anotyce large amounts of natural

Natural Language

How AI works now





Introduction | Motivation

How robotic learning works now



How robotic learning will work in the future





Introduction | Challenges

Existing works are :

- 1. Parameter-wise Heavy (~55B parameters)
- 2. Closed-Source
- 3. Lacking fine-tuning exploration

Octo: An Open-Source Generalist Robot Policy



RT-2: Vision-Language-Action Models

Transfer Web Knowledge to Robotic Control





OpenVLA | Robotic Vision - Language - Action model

consists of 7B parameter, fully open-source, support efficient fine-tuning





- 1. Outperform SOTA RT-2-X (55B) by 16.5% in absolute task success rate
 - work across 29 tasks, multiple robot embodiments, with fewer parameters(7B)
- 2. Demonstrate effectiveness of modern parameter-efficient fine-tuning and quantization
- 3. First open-source generalist VLA thus supports future research



Google Robot



Bridge V2 WidowX Robot

Put Eggplant into Pot



OpenVLA | Architecture





OpenVLA | Architecture





OpenVLA | Architecture





OpenVLA architecture consists of three key components :

- 1. Vision encoder that concatenates Dino V2 and SigLIP features
- 2. Projector that maps visual features to the language embedding space
- 3. Llama 2 7B-parameter large language model



OpenVLA architecture consists of three key components :

- 1. Vision encoder that concatenates Dino V2 and SigLIP features
- 2. Projector that maps visual features to the language embedding space
- 3. Llama 2 7B-parameter large language model



OpenVLA architecture consists of three key components :

- 1. Vision encoder that concatenates Dino V2 and SigLIP features
- 2. Projector that maps visual features to the language embedding space
- 3. Llama 2 7B-parameter large language model



Enabling VLM to predict robot actions : next-token prediction (CE loss)



OpenVLA I Training

Training Data : Open X-Embodiment dataset 70 robot dataset with 2M trajectories

Curated under below conditions, selected **970k robot episodes**

- 1. Single arm manipulations, with 3rd person view camera
- 2. Ensure a balanced mix of embodiments, tasks, scenes



OpenVLA Training Dataset Mixture			
Fractal [92]	12.7%		
Kuka [45]	12.7%		
Bridge[6, 47]	13.3%		
Taco Play [93, 94]	3.0%		
Jaco Play [95]	0.4%		
Berkeley Cable Routing [96]	0.2%		
Roboturk [97]	2.3%		
Viola [98]	0.9%		
Berkeley Autolab UR5 [99]	1.2%		
Toto [100]	2.0%		
Language Table [101]	4.4%		
Stanford Hydra Dataset [102]	4.4%		
Austin Buds Dataset [103]	0.2%		
NYU Franka Play Dataset [104]	0.8%		
Furniture Bench Dataset [105]	2.4%		
UCSD Kitchen Dataset [106]	<0.1%		
Austin Sailor Dataset [107]	2.2%		
Austin Sirius Dataset [108]	1.7%		
DLR EDAN Shared Control [109]	<0.1%		
IAMLab CMU Pickup Insert [110]	0.9%		
UTAustin Mutex [111]	2.2%		
Berkeley Fanuc Manipulation [112]	0.7%		
CMU Stretch [113]	0.2%		
BC-Z [55]	7.5%		
FMB Dataset [114]	7.1%		
DobbE [115]	1.4%		
DROID [11]	$10.0\%^{6}$		



OpenVLA I Training

Other considerations during training :

- OpenVLA model is trained on a cluster of 64 A100 GPUs for 14 days
 - Inference : 15GB of GPU memory when loaded bfloat16



Experiments

Experiments validate 3 key aspects of OpenVLA:

1. Performance as a Generalist Policy

Zero-shot generalization tests

2. Effectiveness of Fine-tuning

New robot setups & tasks

3. Performance with Limited Hardware

Parameter-efficient FT, quantization



Experiments

Experiments validate 3 key aspects of OpenVLA:

1. Performance as a Generalist Policy

Zero-shot generalization tests

2. Effectiveness of Fine-tuning

New robot setups & tasks

3. Performance with Limited Hardware

Parameter-efficient FT, quantization





Experiment | Generalization

Settings : Robot and tasks from pre-trained data

WidowX robot from BridgeData V2 evaluation

Mobile manipulation robot from RT-1 and RT-2 evaluation

Categorizing the term "Generalization" :

Visual : unseen backgrounds, distractor objects, appearances of objects

Motion : unseen object positions/orientations

Physical : unseen object sizes/shapes

Semantic : unseen target objects, instructions, concepts from the Internet









1. Performance as a Generalist Policy Zero-shot generalization tests

Experiment | Generalization

1. Performance as a Generalist Policy Zero-shot generalization tests



KAIST

Experiment

Experiments validate 3 key aspects of OpenVLA:

1. Performance as a Generalist Policy

Zero-shot generalization tests

2. Effectiveness of Fine-tuning

New robot setups & tasks

3. Performance with Limited Hardware

Parameter-efficient FT, quantization



"Adapting a pretrained VLA model to a new robot, new environment, or new task"

- Quick adaptation to a new setup, with a much smaller dataset
- Implicitly captures per-setup differences
 - Camera pose, embodiment, environment, reference frame, ...





"Is OpenVLA adaptable to new robot setups and tasks?"

- Prepared Franka robot arm setups not in pretraining data
- Models trained / fine-tuned on 7 tasks, 10-150 demonstrations each
 - Diffusion policies, Octo, OpenVLA (scratch and pretrained)





Experiments | Fine-Tuning

2. Effectiveness of Fine-tuning New robot setups & tasks

1. Diffusion Policy \rightarrow good for narrow, single-instruction tasks





Experiments | Fine-Tuning

- 1. Diffusion Policy \rightarrow good for narrow, single-instruction tasks
- 2. Fine-tuned VLAs \rightarrow better with multiple objects & language conditioning





Experiments | Fine-Tuning

- 1. Diffusion Policy \rightarrow good for narrow, single-instruction tasks
- 2. Fine-tuned VLAs \rightarrow better with multiple objects & language conditioning
- 3. Robot data pretraining \rightarrow significant performance boost





- 1. Diffusion Policy \rightarrow good for narrow, single-instruction tasks
- 2. Fine-tuned VLAs \rightarrow better with multiple objects & language conditioning
- 3. Robot data pretraining \rightarrow significant performance boost
- 4. OpenVLA shows strong performance across task types



Experiments

Experiments validate 3 key aspects of OpenVLA:

1. Performance as a Generalist Policy

Zero-shot generalization tests

2. Effectiveness of Fine-tuning

New robot setups & tasks

3. Performance with Limited Hardware

Parameter-efficient FT, quantization



3. Performance with Limited Hardware Parameter-efficient FT, quantization

Tested different fine-tuning strategies

- Strategies: Full FT, Last layer only, Frozen vision, Sandwich, LoRA
- Criteria: memory requirement, performance
- Remark: LoRA is efficient, with minimal performance degradation!
 Move <object> Put Carrot

	OpenVLA	Action De-Tokenizer	
	3	Llama 2 7B	$\Delta \theta$ $\Delta Grip$
Input Image			7D Robot Action
"Put eggplant in bowl"	DinoV2 SigLIP	Llama Tokenizer	

καιςτ





in Bowl

Strategy	Success Rate	Train Params ($\times 10^6$)	VRAM (batch 16)
Full FT	69.7 \pm 7.2 %	7,188.1	163.3 GB*
Last layer only	$30.3 \pm 6.1 ~\%$	465.1	51.4 GB
Frozen vision	$47.0\pm6.9~\%$	6,760.4	156.2 GB*
Sandwich	$62.1\pm7.9~\%$	914.2	64.0 GB
LoRA, rank=32	$\textbf{68.2} \pm \textbf{7.5\%}$	97.6	59.7 GB
rank=64	$\textbf{68.2} \pm \textbf{7.8\%}$	195.2	60.5 GB

LoRA: Low-Rank Adaptation of Large Language Models

- Vanilla fine-tuning: update (d x h) weights, expensive
- LoRA: only train B and A, where $BA = \Delta W$ (weight updates)
- Insight
 - Pretrained W is near a good solution
 - For fine-tuning, updating all weights is unnecessary!





Many foundation models require large GPU memory (even for inference)

- Idea: reduce the precision (# of bits used) of model weights
 - Less memory footprint
 - Extra quantization / dequantization overhead
 - Increased arithmetic error
- Memory ↔ Performance tradeoff





Experiments | Quantization

- Tested on BridgeData V2 tasks with different quantization levels
- Quantization overhead can reduce throughput
 - 8-bit: throughput too low for a 5Hz control loop
- Quantization can lower GPU memory transfer (improve throughput)
 - 4-bit: reduced memory transfer compensates overhead!





Figure 6: **OpenVLA inference speed for various GPUs.** Both bfloat16 and int4 quantization achieve high throughput, especially on GPUs with Ada Lovelace architecture (RTX 4090, H100). Further speed-ups are possible with modern LLM inference frameworks like TensorRT-LLM [89]. A: Model sharded across two GPUs to fit.

OpenVLA | Limitation

- Only support single-image observation
- Inference throughput (Hz)
- Room for performance improvement (90+%)



OpenVLA | Summary

- Open-source VLA model for manipulators
 - Image + language instruction \rightarrow robot action
- CV & NLP advancements on robotic applications
- Towards large, 'generalist', widely-deployable robot models
 - $\circ \quad \text{Robot data pretraining} \rightarrow \text{general performance}$
 - $\circ~$ Fine-tuning (LoRA) \rightarrow task-specific adaptation
 - $\circ \quad \text{Quantization} \rightarrow \text{memory-efficient inference}$





Thank you



- According to OpenVLA's fine-tuning experiments, the key benefit of LoRA over full fine-tuning is the drastic improvement in **performance** (task success rate).
 - a. True
 - b. False
- 2. Which is **NOT** true about OpenVLA?
 - a. OpenVLA can be fine-tuned on new robot tasks with 10-150 demonstrations.
 - b. OpenVLA strictly outperforms Diffusion Policy across all tasks.
 - c. OpenVLA supports low-memory deployment using quantization.

