

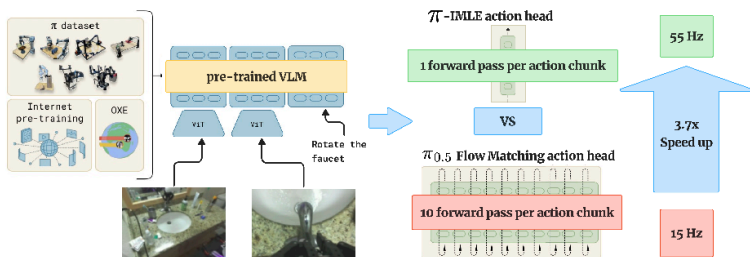


π -IMLE: Fast Single-Step Action Generation for Vision-Language-Action Models

Kian Hosseinkhani^{1*}, George Shramko¹, Mehran Aghabozorgi¹, Qinhe Peng², Jianing Qian², Tristan Engst¹, Alireza Moazeni¹, Dinesh Jayaraman², Ke Li¹

¹APEX Lab · School of Computing Science · Simon Fraser University ²University of Pennsylvania
{kian_hosseinkhani, george_shramko, maa143, tme3, sam62, kelij}@sfu.ca · {pengqh20, jianingq, dineshj}@seas.upenn.edu

Motivation

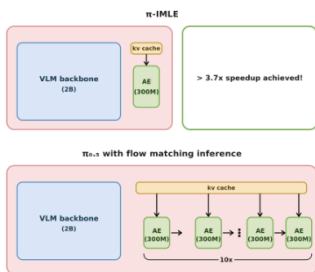


- Why VLAs are slow**
- VLAs → strong cross-task generalization → **slow** (built atop VLMs).
 - Stop-wait-execute → robot idles each prediction → **jerky motion**.
 - Wide perception-action loop → **policy acts on stale observations**.

- Our fix: π -IMLE**
- SOTA speed → single-step generation → **less jerky, higher reaction rate**.
 - Multimodal → many valid ways to do a task, cIMLE learns *all* modes.

55 Hz · 3.67x faster 98.3% LIBERO (best) Robust on LIBERO-plus Beats $\pi_0.5$ on every real task

Our Method: π -IMLE



- Frozen VLM backbone**
Replace only the action head — training cost is minimal.
- Single-step generator**
 G_θ maps (observation O_t , noise z) → action chunk \hat{A} in one forward pass.
- Trained via cIMLE**
Mode-covering by construction, no iterative steps needed.

Conditional IMLE: Formulation

- **Sample.** For each observation-action pair ($\alpha_t^{(i)}$, $A_{gt}^{(i)}$), draw m noise vectors $z_{i,1..m} \sim \mathcal{N}(0, I)$ and produce a set of m candidates:

$$\hat{A}_{i,j} = G_\theta(\alpha_t^{(i)}, z_{i,j}), \quad j = 1, \dots, m$$

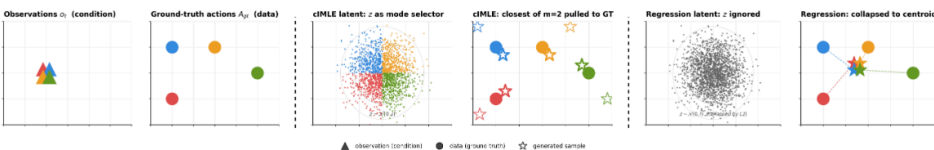
- **Assign.** Each $A_{gt}^{(i)}$ is assigned to its nearest candidate from $\alpha_t^{(i)}$'s candidate set:

$$j^*(i) = \arg \min_j \|\hat{A}_{i,j} - A_{gt}^{(i)}\|_2^2$$

- **Update.** Average L2 loss over every ground-truth action:

$$\mathcal{L}_{cIMLE} = \frac{1}{n} \sum_{i=1}^n \|\hat{A}_{i,j^*(i)} - A_{gt}^{(i)}\|_2^2$$

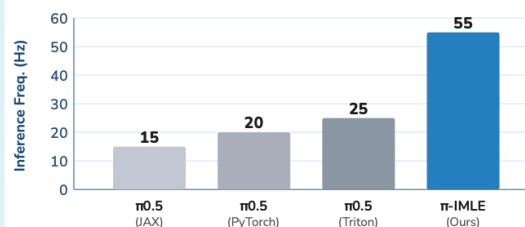
The loss penalizes every ground-truth action equally: any uncovered behavior mode raises the loss.



Single-Step Generation Lifts the Inference Ceiling

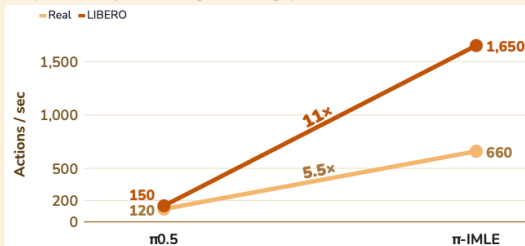
Inference Frequency (Hz)

π -IMLE attains 55 Hz, surpassing every system-level $\pi_{0.5}$ optimization capped by its 10-step flow-matching head.



Action Throughput (actions / sec)

Faster forward passes combined with a longer execution horizon compound to **up to 11x** higher throughput.



LIBERO Benchmark

40 tabletop tasks across 4 suites × 50 ep/task. π -IMLE is the only method that simultaneously leads in success rate AND inference frequency.

Model	H	Spat.	Obj.	Goal	Long	Avg ↑	Inf. ↑
$\pi_{0.5}$	10	97.2	99.0	97.8	96.0	97.5	1.0x
$\pi_{0.5}$	30	95.0	98.0	96.2	95.2	96.1	1.0x
OpenVLA-OFT	—	95.2	94.2	95.2	93.2	94.5	0.69x
Shallow- $\pi_{0.5}$ -L9	—	99.0	98.0	97.0	93.0	97.0	1.7x
Shallow- $\pi_{0.5}$ -L6	—	98.0	96.0	94.0	90.0	95.0	2.3x
π -IMLE	10	98.0	100.0	99.0	96.0	98.3	3.67x
π -IMLE	30	97.2	99.0	96.2	95.8	97.1	3.67x

VLA Wall-Clock per Episode (Simulation)

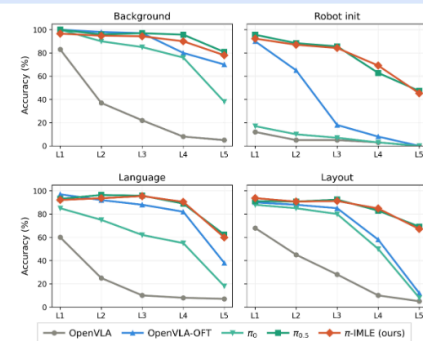
Def. Wall-clock on policy inference per episode (excluding robot execution/env overhead).

System	Fwd ↓	VLA-WC (s) ↓	Faster ↑
$\pi_{0.5}$ (JAX, H=10)	15.4	1.03	1.0x
π -IMLE (H=10)	15.5	0.28	3.6x
π -IMLE (H=30)	5.4	0.10	10.5x

Result. Longer horizon × faster forward pass → 10.5x less inference time.

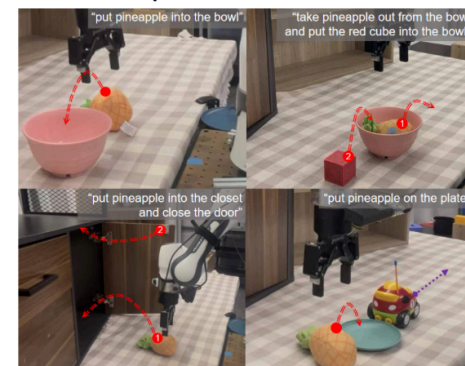
LIBERO-plus: Robustness Under Distribution Shift

Robustness under distribution shift. Systematic test-time perturbations across 4 axes (background, robot initial state, language, layout) × 5 severity levels (L1-L5).



π -IMLE retains $\pi_{0.5}$'s robustness across all axes and severities; baselines' success rates drop sharply as severity increases. OpenVLA-OFT's L1 regression collapses to the median mode under shift, while cIMLE enforces **mode coverage by construction**.

Real-World Experiments — Franka Panda



Zero-shot from DROID with two camera views (wrist + scene). $\pi_{0.5}$ at 15 Hz vs. π -IMLE at 55 Hz — same VLM, different head.

Task	Type	π -IMLE ↑	$\pi_{0.5}$ ↑	WC × ↓
Pineapple in bowl	single	10/10	9/10	6.6x
Swap pineapple & cube	multi	9/10	8/10	5.7x
Pineapple in cabinet	multi	8/10	7/10	3.9x
Moving plate	reactive	9/10	6/10	5.4x

9 / 10 vs **6 / 10** on the dynamic moving-plate task — faster inference closes the perception-action loop, $\pi_{0.5}$ chases stale targets.