

# Language-Action Alignment in LAP Enhances Cross-Embodiment Generalization on BridgeData V2

Assignee Research

June 7, 2026

## **Abstract**

This report synthesises findings from 9 peer-reviewed papers addressing the following research question: What is the impact of LAP's language-action alignment on cross-embodiment generalization accuracy when evaluated on the BridgeData V2 dataset without fine-tuning. 19 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## **1 Introduction**

This paper examines: LAP: Language-Action Pre-Training Enables Zero-shot Cross-Embodiment Transfer. Research question: What is the impact of LAP's language-action alignment on cross-embodiment generalization accuracy when evaluated on the BridgeData V2 dataset without fine-tuning?.

## **2 Methodology**

Systematic literature search across multiple databases yielded 9 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

## **3 Results**

9 papers retrieved. 19 claims extracted; 1 independently verified. Quality review score: 4.5/10.

## 4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.



## 5 Extracted Claims

Claim	Verified	Confidence
LAP-3B achieves performance comparable to $\pi$ 0.5-DROID on the seen embodiment (DROID).	×	0.05
Across three previously unseen embodiments and six real-world manipulation tasks, LAP-3B attains over 50% average zero-s	✓	0.25
LAP-3B delivers approximately a 2 $\times$ improvement in zero-shot success over the strongest baselines on unseen embodiments.	×	0.14
All evaluated open-sourced VLAs collapse to a zero success rate in zero-shot cross-embodiment transfer on unseen embodim	×	0.14
The evaluation of LAP covers four robot embodiments.	×	0.10
The evaluation of LAP covers ten real-world manipulation tasks.	×	0.09
The evaluation of LAP includes the LIBERO simulation benchmark.	×	0.09
$\pi$ 0.5-DROID is currently the strongest VLA on the DROID benchmark.	×	0.04
$\pi$ 0.5-Base is the strongest publicly available VLA base model known to the authors.	×	0.04
The FAST tokenizer used in the $\pi$ 0.5-replicated baseline is pre-trained on substantially more robot data than the LAP mod	×	0.07
LAP-3B adopts a Mixture-of-Transformers architecture combining a LAP-trained VLM backbone with a lightweight flow-matchi	×	0.05
Auto-regressive generation of language-actions at inference time is slow and impractical for real-time robot control.	×	0.06
LAP-3B differs from $\pi$ 0.5 only in the action representation used to supervise the VLM (language-actions versus FAST token	×	0.07
Under the LAP architecture, the VLM and action expert communicate solely through cross-attention.	×	0.04
LAP supports an inference frequency of 25Hz.	×	0.02
Grover et al. [71] VLA-0 has an inference frequency of approximately 0.5Hz.	×	0.02
VLM2VLA [23] has an inference frequency of 4Hz.	×	0.01
LAP enables cross-embodiment <sup>4</sup> transfer, whereas Grover et al. [71], VLA-0 [6], and VLM2VLA [23] do not.	×	0.10
LAP utilizes large-scale training, whereas VLA-0 [6] and VLM2VLA [23] do not.	×	0.07

## References

- <http://arxiv.org/abs/2603.16806v2>
- <http://arxiv.org/abs/2605.27759v1>
- <http://arxiv.org/abs/2602.10556v2>