

# SOVEREIGN: How does the multi-turn RL approach in LongNav-R1 compare to single-turn RL baselines on the RxR-CE benchmark

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 29, 2026

## Abstract

Recent advances in the areas of multimodal machine learning and artificial intelligence (AI) have led to the development of challenging tasks at the intersection of Computer Vision, Natural Language Processing, and Embodied AI. Whereas many approaches and previous survey pursuits have characterised one or two of these dimensions, there has not been a holistic analysis at the center of all three. Moreover, even when combinations of these topics are considered, more focus is placed on describing, e.g., current architectural methods, as opposed to also illustrating high-level challenges and oppor

## 1 Introduction

Analysis of: Core Challenges in Embodied Vision-Language Planning. Research goal: How does the multi-turn RL approach in LongNav-R1 compare to single-turn RL baselines on the RxR-CE benchmark when evaluated using normalized dynamic time warping (nDTW) for unseen environments?.

## 2 Methodology

Multi-query arXiv search (4 parallel queries, Relevance-sorted). TF-IDF cosine semantic verification (bigrams, threshold=0.15). NIM nv-embedqa-e5-v5 (dim=1024) for semantic indexing. Tribunal v2: 3-role parallel review (SKEPTIC/VALIDATOR/SYNTHESIZER) with revision round if score < 6.5.

## 3 Results

7 papers retrieved. 8 claims extracted, 7 verified. Tribunal: 8.2/10 \$\rightarrow\$ APPROVE (revision\_round=0). Policy: AUTO\_APPROVE.

## 4 Uncertainties

NIM free tier latency varies. TF-IDF verification is a weak signal. arXiv Relevance ranking is query-dependent. Tribunal consensus is LLM-based and prompt-sensitive.

## 5 Extracted Claims

| Claim   | Verified | Confidence |
|---|----------|------------|
| Recent advances in multimodal machine learning and AI have led to the development of challenging tasks at the intersection of     | ✓        | 0.38       |
| There has not been a holistic analysis at the center of Computer Vision, Natural Language Processing, and Embodied AI.            | ✓        | 0.31       |
| More focus is placed on describing current architectural methods rather than illustrating high-level challenges and opportunities | ✓        | 0.30       |
| This survey paper discusses Embodied Vision-Language Planning (EVLP) tasks, a family of prominent embodied navigation and         | ✓        | 0.43       |
| The paper proposes a taxonomy to unify EVLP tasks.  | ×        | 0.14       |
| The paper provides an in-depth analysis and comparison of new and current algorithmic approaches, metrics, simulated environments | ✓        | 0.34       |
| The paper presents core challenges that new EVLP works should seek to address.  | ✓        | 0.23       |
| The paper advocates for task construction that enables model generalizability and furthers real-world deployment.                 | ✓        | 0.24       |

## References

- <https://doi.org/10.48550/arxiv.2303.01396>
- <https://doi.org/10.1613/jair.1.13646>
- <https://openalex.org/W7117339700>