

Multimodal Segmentation Model Throughput Scaling with Input Resolution on GPUs

Assignee Research

May 31, 2026

Abstract

This report synthesises findings from 5 peer-reviewed papers addressing the following research question: How does the inference throughput of multimodal segmentation models scale with input resolution on GPU accelerators compared to standard CNN backbones. Multimodal referring segmentation aims to segment target objects in visual scenes, such as images, videos, and 3D scenes, based on referring expressions in text or audio format. This task plays a crucial role in practical applications requiring accurate object perception based. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Multimodal Referring Segmentation: A Survey. Research question: How does the inference throughput of multimodal segmentation models scale with input resolution on GPU accelerators compared to standard CNN backbones?.

2 Methodology

Systematic literature search across multiple databases yielded 5 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.3/10.

3 Results

5 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 6.3/10.

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.

References

- <http://arxiv.org/abs/1810.08137v3>
- <http://arxiv.org/abs/2508.00265v2>
- <http://arxiv.org/abs/2011.11052v1>