

Model Scale and Hallucination Error Distribution in Federated Vision-Language Models on MS COCO

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: What is the correlation between model scale and the distribution of hallucination errors in federated vision-language models evaluated on MS COCO captioning. 13 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: COCO-Urdu: A Large-Scale Urdu Image-Caption Dataset with Multimodal Quality Estimation. Research question: What is the correlation between model scale and the distribution of hallucination errors in federated vision-language models evaluated on MS COCO captioning?.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.7/10.

3 Results

14 papers retrieved. 13 claims extracted; 2 independently verified. Quality review score: 4.7/10.

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
The COCO-Urdu dataset evaluation utilized an ensemble quality estimation (QE) pipeline integrating COMET-Kiwi, BERTScore	✓	0.19
Reference translations for COCO-Urdu were generated using the NLLB-3B model due to the unavailability of human reference	×	0.06
Zero-shot Urdu translations were obtained using the SeamlessM4T model.	×	0.06
The refined COCO-Urdu dataset contains 59,000 images and 319,000 captions.	✓	0.20
The refined COCO-Urdu dataset achieved a BLEU score of 0.53, a SacreBLEU score of 53, and a CHRF score of 74.	×	0.07
The zero-shot COCO-Urdu dataset achieved a BLEU score of 0.52, a SacreBLEU score of 52, and a CHRF score of 73.23.	×	0.06
The UICD dataset contains 31,000 images and 135,000 captions with a reported BLEU score of 0.86.	×	0.07
The Flickr8k Urdu dataset contains 700 images and 700 captions with a reported BLEU score of 0.83.	×	0.05
The translation efforts for COCO-Urdu were limited to a 50% subset of the MS COCO dataset due to computational constrain	×	0.10
A stratified sampling strategy based on the iterative stratification algorithm by Sechidis et al. was used to select the	×	0.05
The stratified sampling approach preserved label co-occurrence patterns and relative class distributions compared to the	×	0.05
SeamlessM4T v2 was used to translate all captions from the stratified MS COCO subset into Urdu in a zero-shot manner.	×	0.09
The zero-shot translation process did not require parallel Urdu training data.	×	0.03

References

- <http://arxiv.org/abs/2306.09265v1>

- <http://arxiv.org/abs/2504.09480v1>
- <http://arxiv.org/abs/2509.09014v1>