

# Intermediate-Task Training on Multimodal Reasoning Datasets for Zero-Shot Cross-Lingual Visual Question Answering

Assignee Research

July 11, 2026

## Abstract

Visual question answering (VQA) is one of the crucial vision-and-language tasks. Yet, existing VQA research has mostly focused on the English language, due to a lack of suitable evaluation resources. Previous work on cross-lingual VQA has reported poor zero-shot transfer performance of current multilingual multimodal Transformers with large gaps to monolingual performance, without any deeper analysis. In this work, we delve deeper into the different aspects of cross-lingual VQA, aiming to understand the impact of 1) modeling methods and choices, including architecture, inductive bias, fine-tun

## 1 Introduction

This paper examines: Delving Deeper into Cross-lingual Visual Question Answering. Research question: To what extent does intermediate-task training on multimodal reasoning datasets improve zero-shot cross-lingual transfer on the visual question answering subset of XTREME-R relative to monolingual English intermediate training?.

## 2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.4/10.

## 3 Results

16 papers retrieved. 10 claims extracted; 9 independently verified. Quality review score: 7.4/10.

## 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 VQA task on the standard English GQA dataset requires predicting a correct answer from 1,853 possible classes.	✓	0.22
The xGQA dataset is a multilingual extension of the English GQA dataset covering 7 typologically diverse languages.	✓	0.18
The xGQA dataset relies on the same set of 1,853 answer classes as the English GQA dataset.	✓	0.20
The xGQA dataset is included in the multimodal multilingual evaluation benchmark IGLUE.	✓	0.16
The study empirically compares two pretrained multimodal multilingual Transformer architectures: M3P and UC2.	×	0.11
Standard approaches from text-only cross-lingual transfer scenarios do not leverage the full multilingual capabilities o	✓	0.20
A simple modified fine-tuning regime achieves gains of more than 10 absolute accuracy points over baselines in cross-lin	✓	0.25
The modified fine-tuning strategies have no substantial impact on model performance in the source language (English).	✓	0.18
The original work on xGQA evaluated only a simple 'shallow' linear classification head using the output [CLS] token.	✓	0.23
In the few-shot setup, the model is optimized on a handful of task-annotated examples in the target language after sourc	✓	0.20

## References

- <http://arxiv.org/abs/2005.13013v2>

- <http://arxiv.org/abs/2209.02982v2>
- <http://arxiv.org/abs/2202.07630v2>