

Gemini 1.5 Pro Long-Term Dependency Modeling in Video-Language Understanding

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does Gemini 1.5 Pro handle long-term dependency modeling in video-language understanding tasks compared to prior models, and what metrics (e.g., F1 score, latency) best capture this performance. 17 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Enhancing Visual Question Answering through Ranking-Based Hybrid Training and Multimodal Fusion. Research question: How does Gemini 1.5 Pro handle long-term dependency modeling in video-language understanding tasks compared to prior models, and what metrics (e.g., F1 score, latency) best capture this performance delta?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.3/10.

3 Results

13 papers retrieved. 17 claims extracted; 1 independently verified. Quality review score: 4.3/10.

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.

5 Extracted Claims

Claim	Verified	Confidence
The RankVQA model was evaluated using the VQA v2.0 and COCO-QA datasets.	×	0.13
VQA v2.0 contains over 200,000 images and 600,000 questions.	×	0.07
COCO-QA comprises 123,287 images and over 117,000 questions.	×	0.05
The experimental environment used NVIDIA Tesla V100 (32GB) GPUs.	×	0.01
The experimental environment used Intel Xeon E5-2698 v4 CPUs.	×	0.04
The experimental environment had 256GB DDR4 memory.	×	0.07
The experimental environment had 2TB SSD storage.	×	0.02
The experimental environment used Ubuntu 20.04 LTS as the operating system.	×	0.01
The experimental environment used PyTorch 1.10.0 as the deep learning framework.	×	0.03
The experimental environment used CUDA Version 11.2.	×	0.04
The experimental environment used cuDNN Version 8.1.	×	0.02
The experimental environment used Python Version 3.8.10.	×	0.03
Images in the datasets were resized to a uniform size of 224x224 pixels.	×	0.02
The RankVQA model employs the Faster R-CNN model to extract visual features from images.	×	0.10
The RankVQA model utilizes a pre-trained BERT model to extract text features.	×	0.13
The RankVQA model uses a multi-head self-attention mechanism for multimodal fusion.	×	0.15
The RankVQA model includes a ranking learning module to optimize the relative ranking of answers.	✓	0.18

References

- <http://arxiv.org/abs/2412.16117v1>

- <http://arxiv.org/abs/2408.07303v2>
- <http://arxiv.org/abs/2403.05530v5>