

# Vision-to-Text: Benchmarking Multimodal LLMs on Extremely Low-Resource Languages



A visually grounded and metric-driven pipeline for building SilkRoad-VL



Shuoshuo Hou<sup>1234</sup>, Mieradilijiang Maimaiti<sup>1234\*</sup>, Zhexin Li<sup>1234</sup>, Jiaxin Wang<sup>1234</sup>, Ahmad Hassan<sup>1234</sup>, Kaishaer Jiapaer<sup>1</sup>, Nilufar Abdurakhmonova<sup>5</sup>, Roza Urinbayeva<sup>5</sup>, Madina Mansurova<sup>6</sup>, Shormakova Assem<sup>6</sup>, Gulnar Murat<sup>7</sup>, Le Wu<sup>8</sup>, Wushouer Silamu<sup>1234</sup>

<sup>1</sup> School of Computer Science and Technology, Xinjiang University;

<sup>2</sup> Xinjiang Laboratory of Multi-Language Information Technology;

<sup>3</sup> Xinjiang Multilingual Information Technology Research Center;

<sup>4</sup> Joint International Research Laboratory of Silk Road Multilingual Cognitive Computing, Urumqi, China

<sup>5</sup> National University of Uzbekistan;

<sup>6</sup> Al-Farabi Kazakh National University;

<sup>7</sup> Kapshagay Bidai Onimderi, Kapshagay, Kazakhstan

<sup>8</sup> Integrated Laboratory for Space, Air, and Ground Systems

\* Corresponding author: miradeljan51@xju.edu.cn



01

## Background and Motivation

Why it matters

02

## Research Goal and Contributions

What we build

03

## Methodology

How we build it

04

## Dataset and Evaluation

How we evaluate it

05

## Results and Evidence

What we find

06

## Conclusion and Q&A

Takeaways

## Research gap

High-resource dominance | Limited multimodal coverage | Low-resource underrepresentation

## Challenges

Translationese artifacts | Semantic drift | Weak visual grounding

## Our goal

Build SilkRoad-VL for six extremely low-resource languages: ug, kk, ky, tg, uz, and ur.

The key idea is to use visual grounding as a quality constraint during corpus construction.

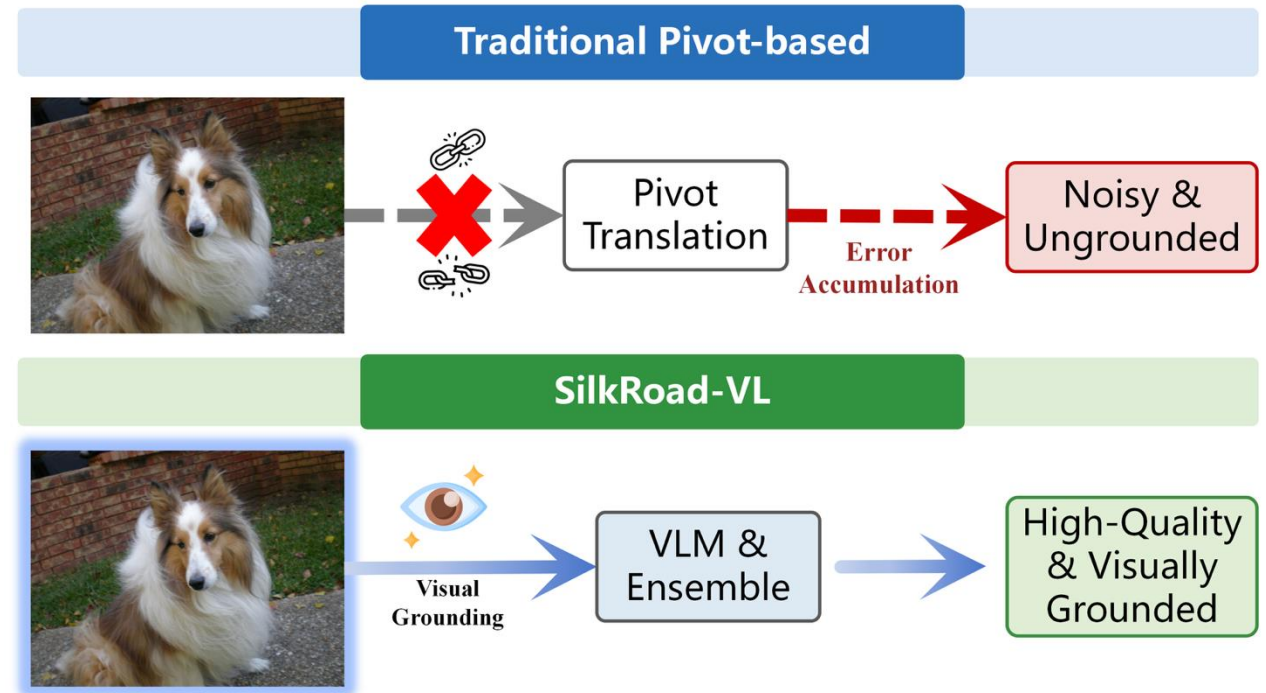


Fig. 1. Traditional pivot-based generation vs. SilkRoad-VL.

01 Background

**02 Research**

03 Methodology

04 Dataset

05 Results

06 Conclusion

## Goal

**Build a high-quality, visually grounded benchmark that improves multimodal translation for extremely low-resource languages.**

**01**

### Metric-driven captioning

Select grounded source captions with controllable short/long descriptions.

**02**

### Hybrid ensemble generation

Use complementary MT and LLM systems to create diverse candidates.

**03**

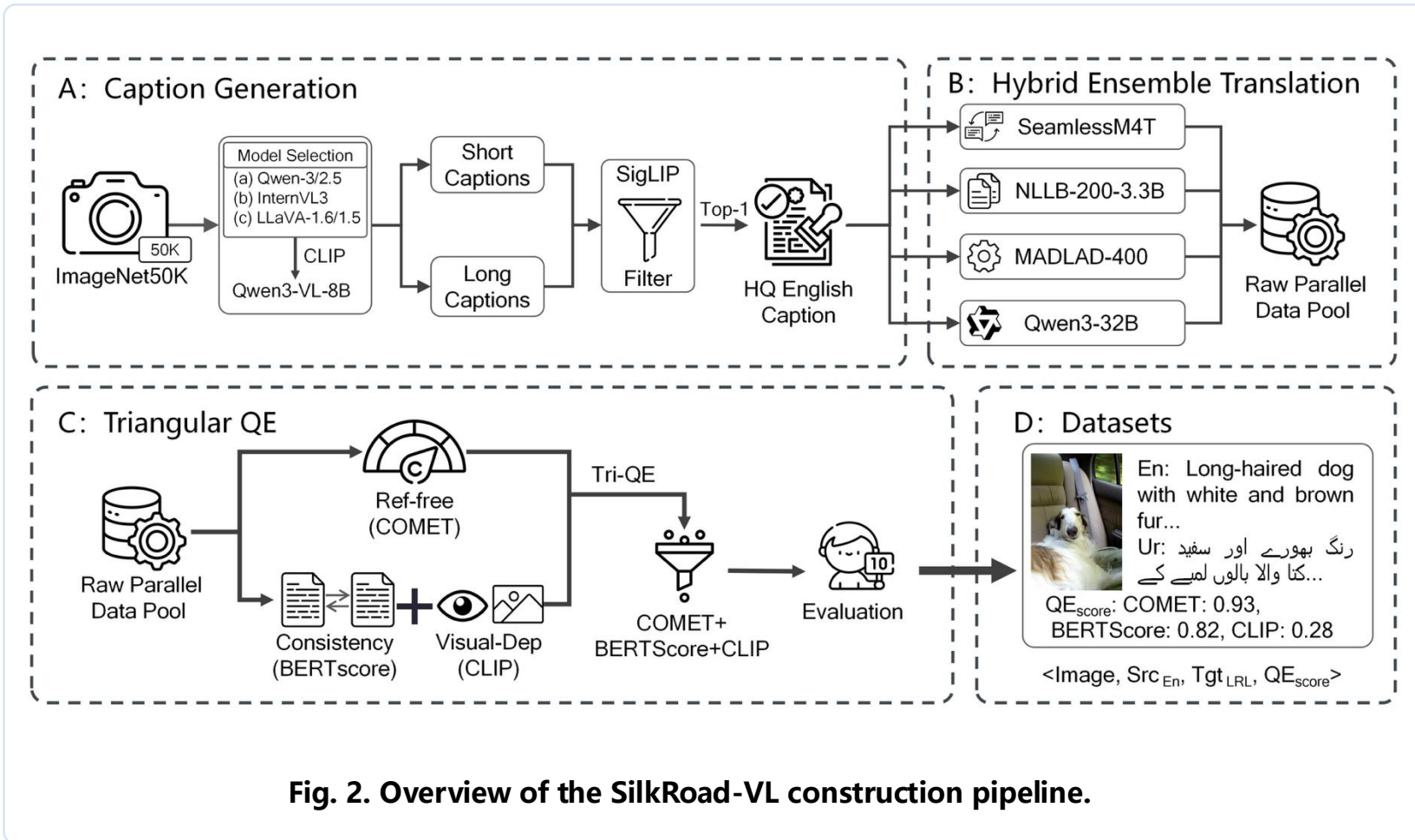
### Tri-QE filtering

Jointly enforce quality, semantic consistency, and visual grounding.

**04**

### Benchmark validation

Verify gains on SilkRoad-VL, Multi30K, and Visual Genome.



**Fig. 2. Overview of the SilkRoad-VL construction pipeline.**

## Three-stage workflow

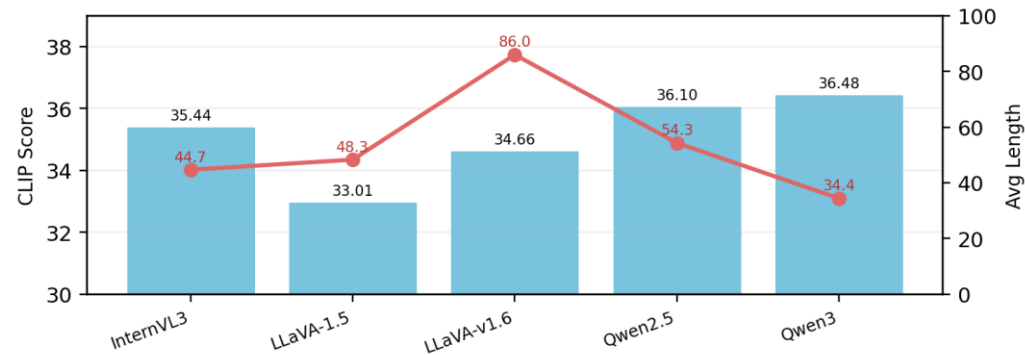
- 1 Caption generation**  
 Select strong visual captions as English anchors.
- 2 Hybrid ensemble**  
 Generate diverse translations from complementary models.
- 3 Tri-QE filtering**  
 Retain candidates that are high-quality and visually grounded.

## Metric-driven caption generation

VLM Selection: compare candidate VLMs and select Qwen3-VL-8B as the optimal caption backbone.

SigLIP Filtering: generate short and long captions, then retain visually faithful descriptions.

**VLM Backbone Selection**



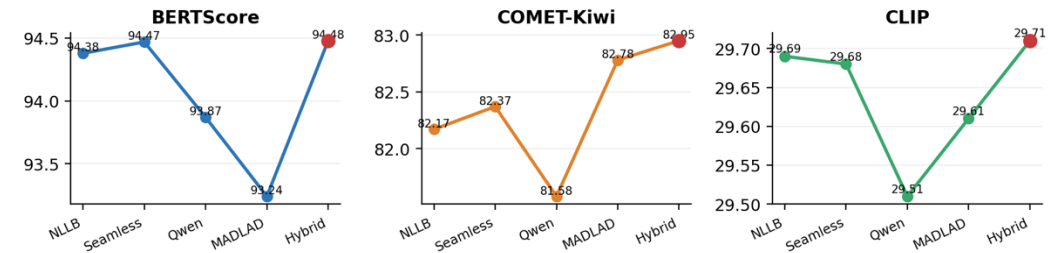
Qwen3-VL balances visual grounding and length controllability.

## Hybrid ensemble translation

Diverse Models: NLLB-200, MADLAD-400, SeamlessM4T-v2, and Qwen3-32B generate candidates.

Complementary Strengths: combine lexical fidelity, fluency, and multimodal alignment while reducing single-model bias.

**Hybrid Ensemble vs. Individual Models**



The ensemble-then-filter paradigm produces stronger supervision.

01 Background

02 Research

**03 Methodology**

04 Dataset

05 Results

06 Conclusion

## Triangular Quality Estimation

Quality: COMET-Kiwi

Semantic consistency: BERTScore

Visual grounding: CLIPScore

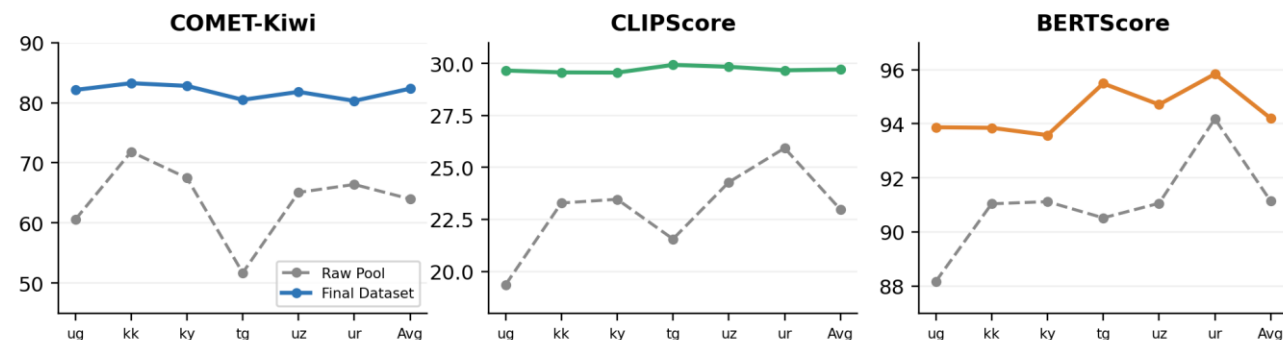
## Strict gate

Only candidates passing all three constraints are retained.

## Thresholds

**78.0 / 0.90 / 0.27**

## Tri-QE Filtering Improves Corpus Quality



Tri-QE removes noisy and hallucinated samples before downstream fine-tuning.

## SilkRoad-VL corpus

**84.9k pairs**

82.7k train · 0.4k dev · 1.8k test

### Target languages

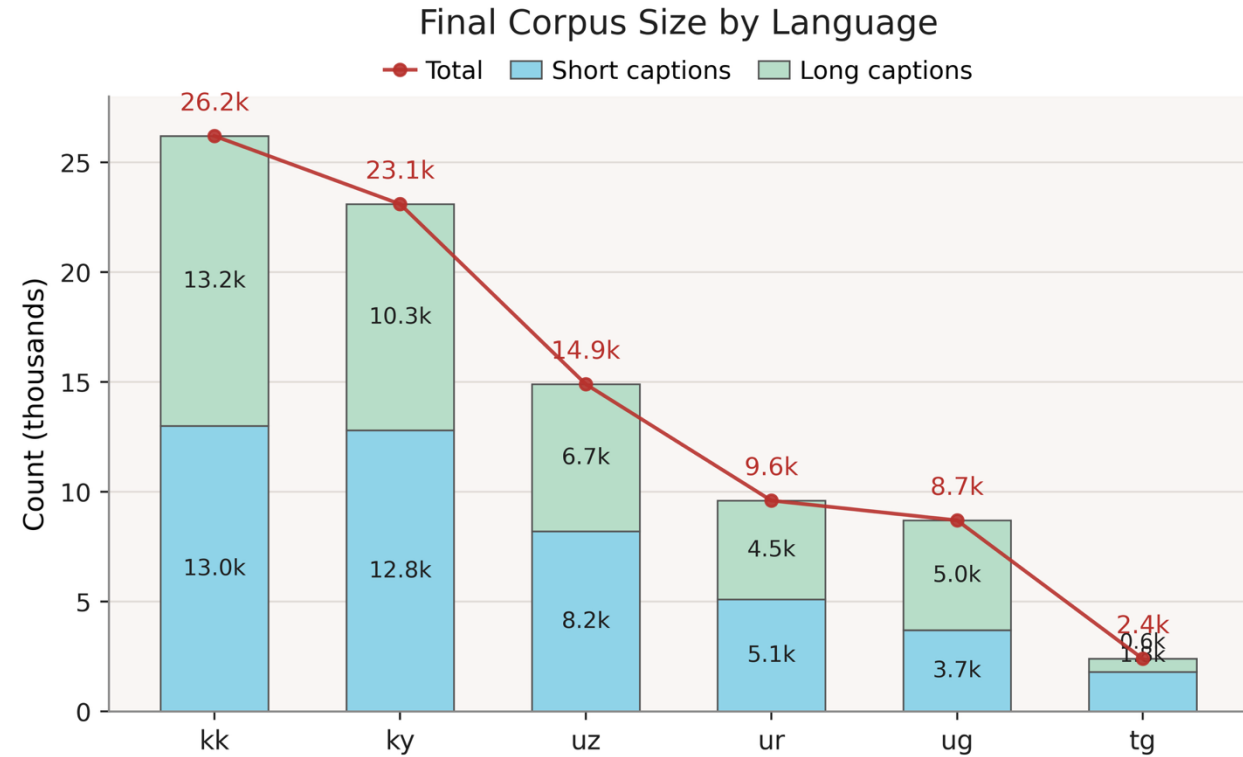
ug · kk · ky · tg · uz · ur

### Metrics

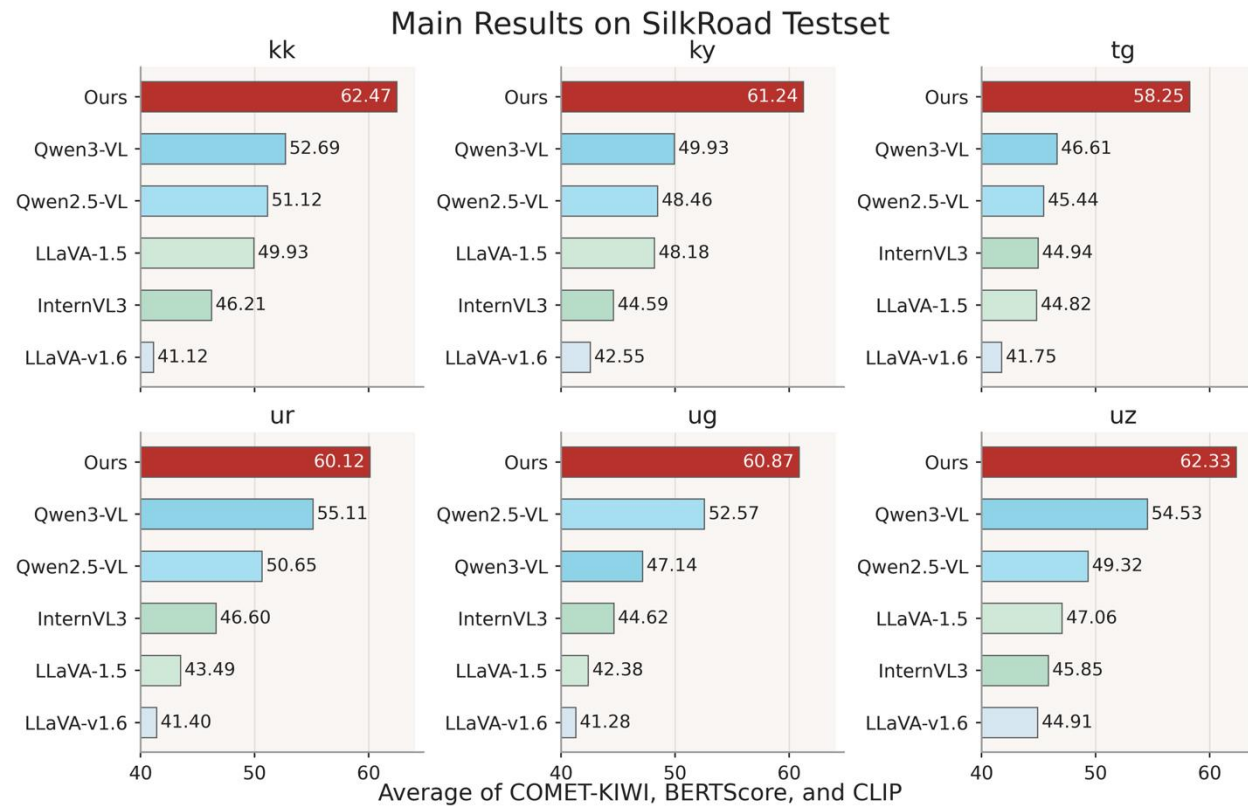
COMET-Kiwi: adequacy

BERTScore: semantic consistency

CLIPScore: visual grounding



Language-wise distribution of the final SilkRoad-VL corpus.



**Fig. 3. Main results on the SilkRoad-VL test set across six languages.**

## Key findings

Ours outperforms all strong baselines on all six languages.

Highest score: 62.47 in Kazakh; 62.33 in Uzbek; 61.24 in Kyrgyz.

Four languages exceed 60 points.

Tajik rises from 46.61 to 58.25 (+11.64).

## Takeaway

**SilkRoad-VL provides strong and consistent gains in extremely low-resource settings.**

01 Background

02 Research

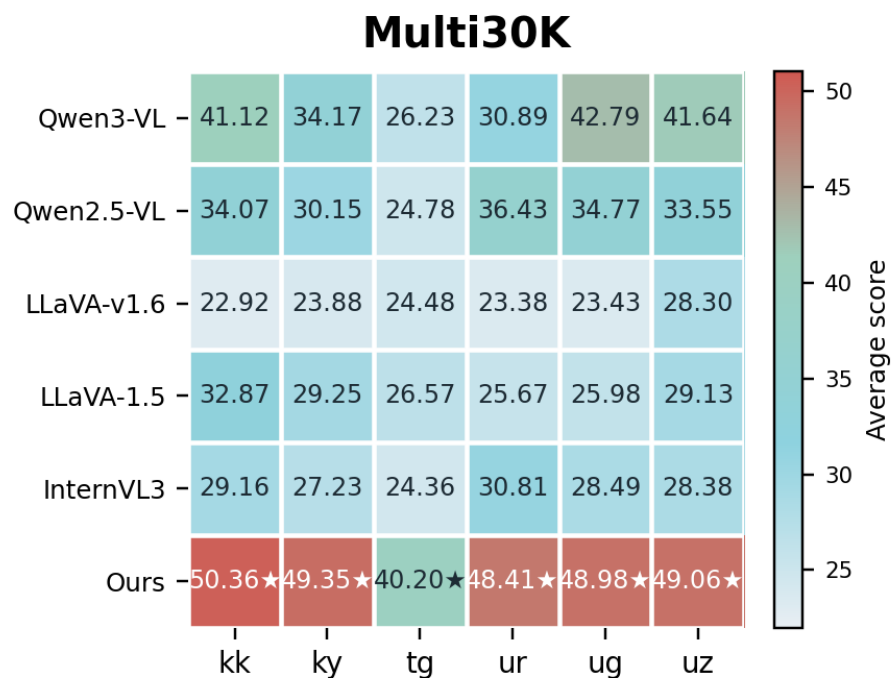
03 Methodology

04 Dataset

**05 Results**

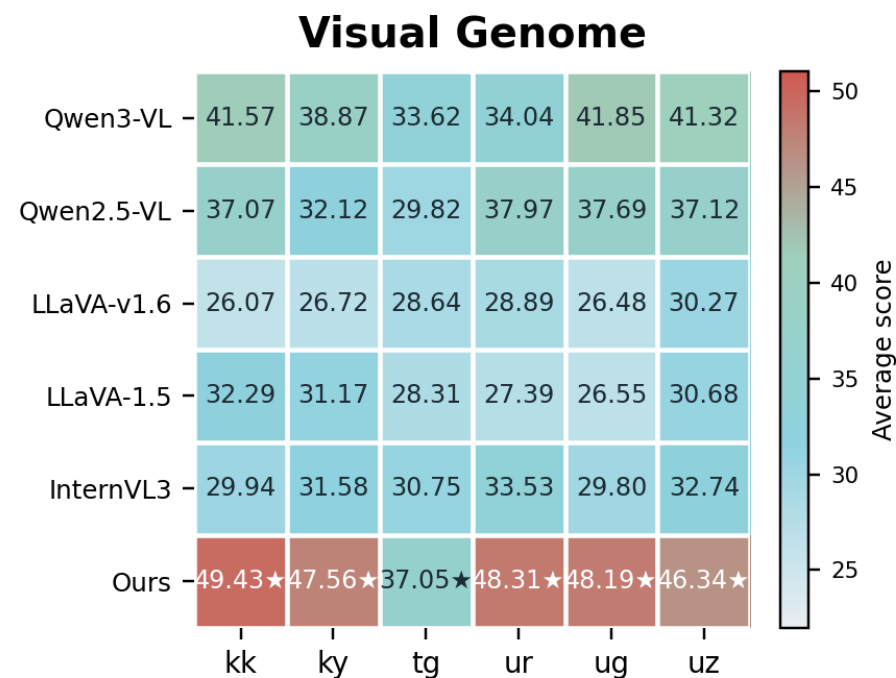
06 Conclusion

## Multi30K



Tajik: 26.57 → 40.20; Kazakh reaches 50.36.

## Visual Genome



Urdu: 37.97 → 48.31; Tajik reaches 37.05.



[01 Background](#)[02 Research](#)[03 Methodology](#)[04 Dataset](#)[05 Results](#)[06 Conclusion](#)

## Takeaways

We construct SilkRoad-VL, a high-quality multimodal benchmark for six extremely low-resource languages.

Metric-driven captioning, hybrid ensemble generation, and Tri-QE improve corpus quality.

Fine-tuned models generalize well to Multi30K and Visual Genome.

Future work will extend the pipeline to more low-resource languages and modalities.

## Thank you!

Questions are welcome.

## Code & Dataset



Corresponding author:  
miradeljan51@xju.edu.cn