# BUPT-MCPRL at TRECVID 2023 Video to Text Description
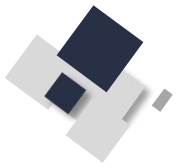
ZeLiang Ma , Shuai Jiang , Zhe Cui, Yanyun Zhao
Beijing University of Posts and Telecommunications
{mzl, js, cuizhe, zyy}@bupt.edu.cn

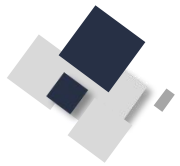Beijing University of Posts and Telecommunications

# CONTENTS
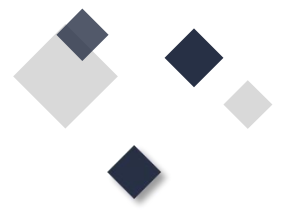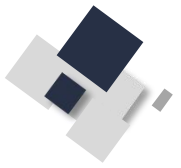
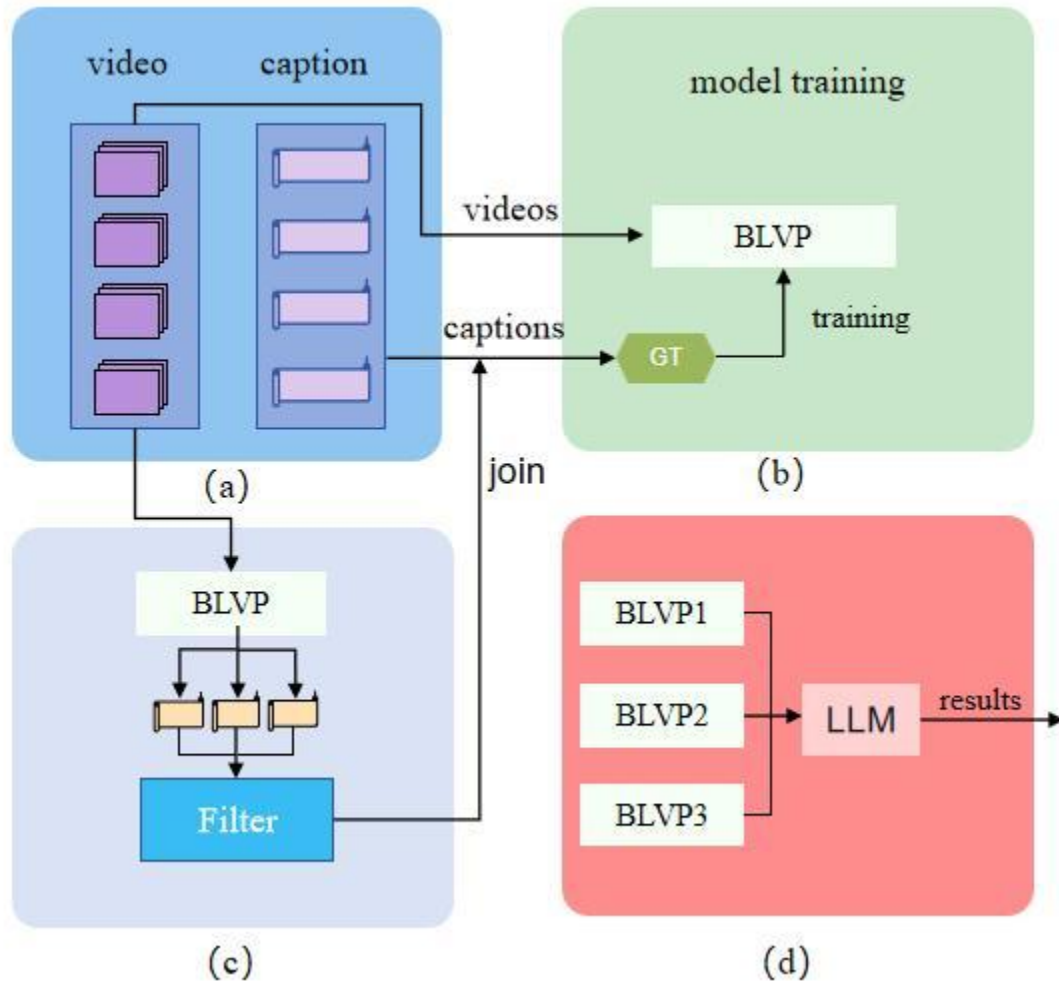Beijing University of Posts and Telecommunications

# 01 Challenges



where
when
who
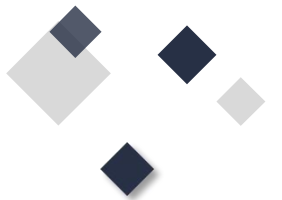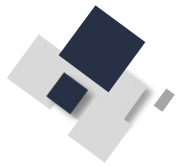what

# 02 Method

- Overall workflow of our system
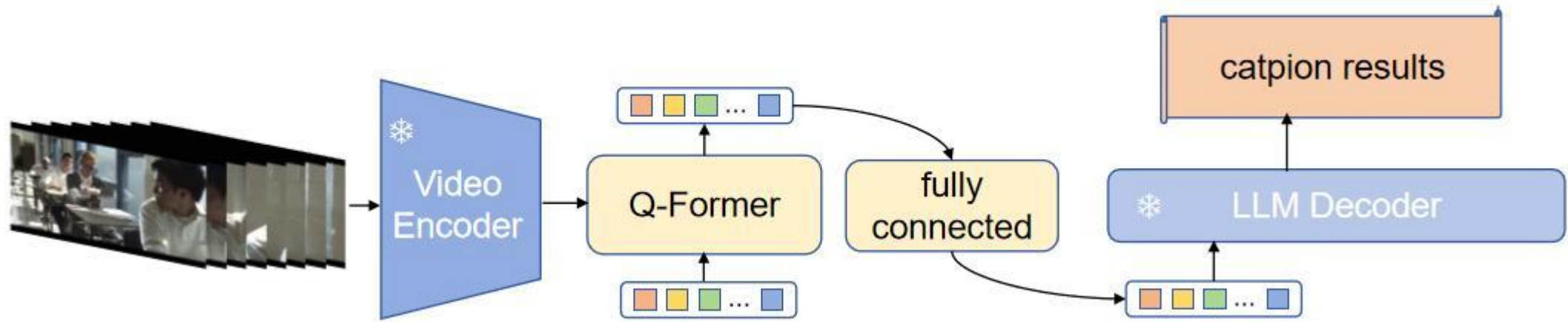


(a) the raw data

(b) the process of model training

(c) Process of New Data Generation and Data Filtering
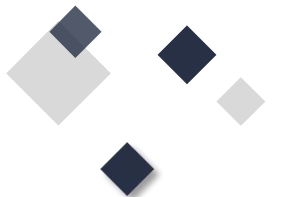
(d)the multi-model fusion of LLM

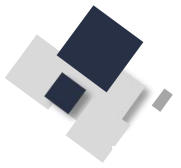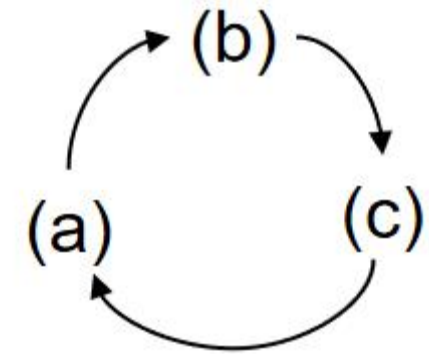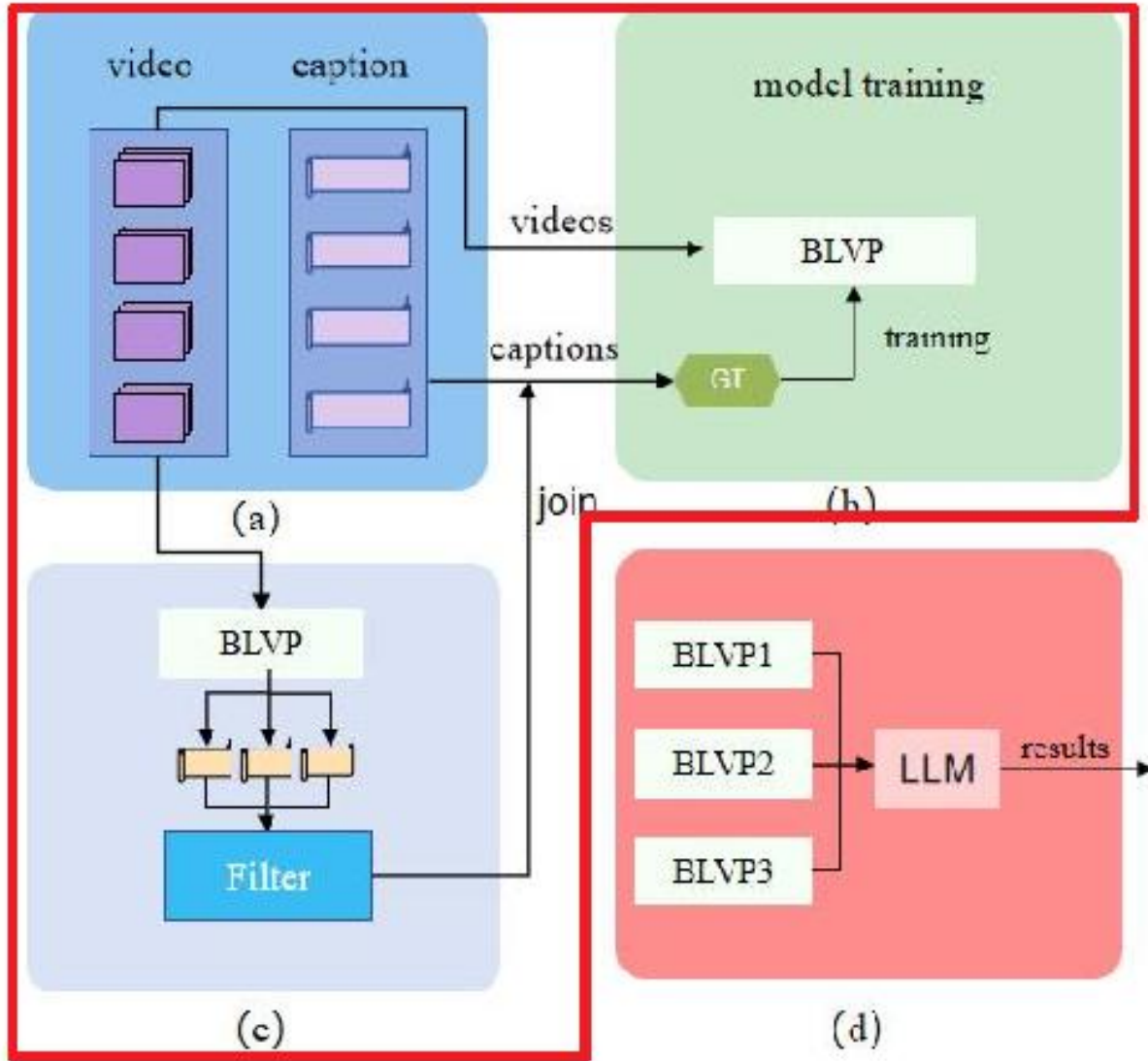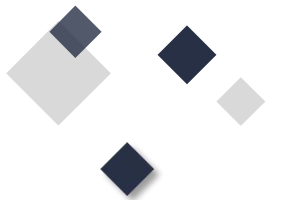# BLVP : BLIP2 for video to text



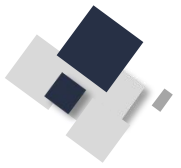Video Encoder: CLIP-G

LLM Decoder: OPT 2.7
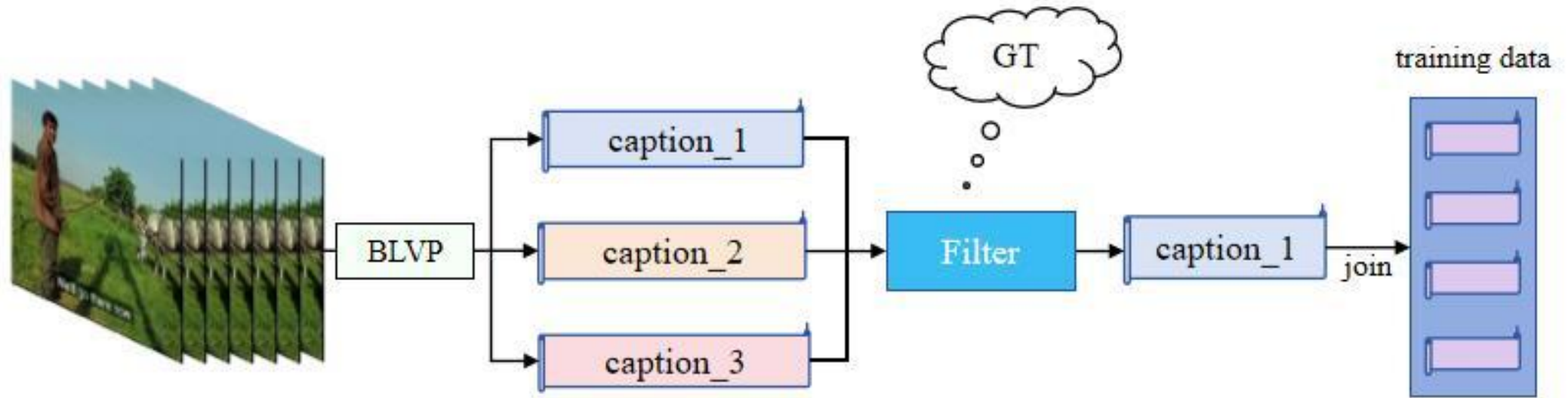
# Cyclic data Augmentation
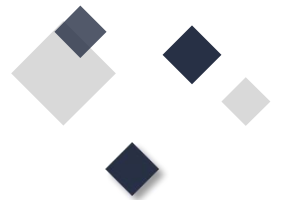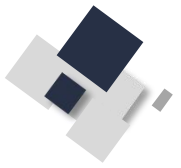


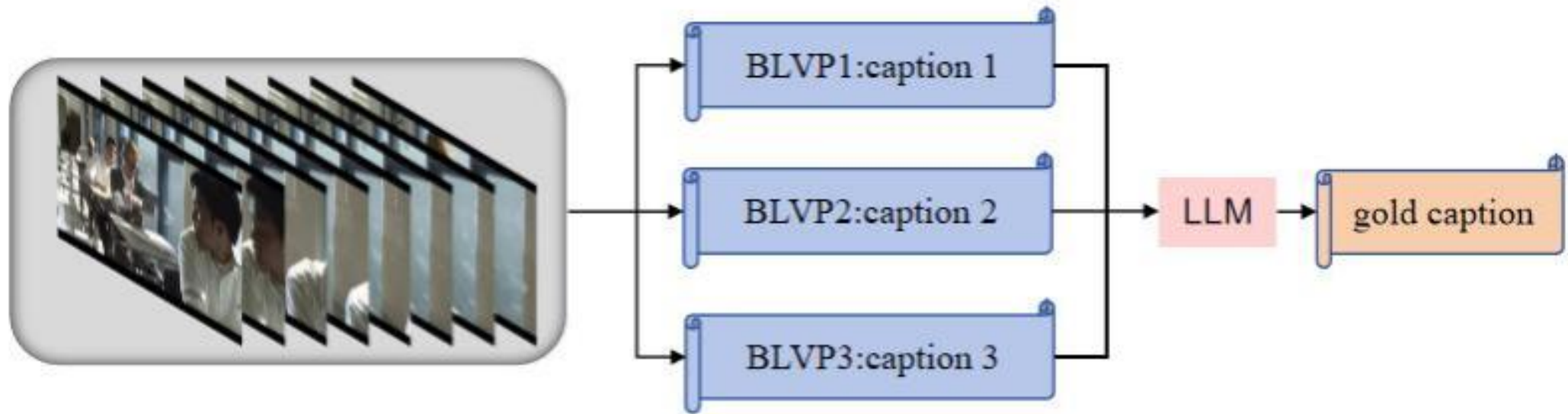The cyclic workflow.

# The process of data selection.



Filter rely on the CIDEr score calculated with respect to the Ground Truth.
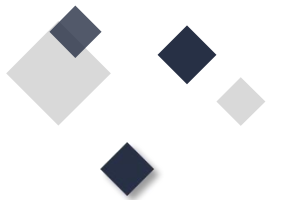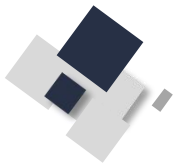
# LLM for Ensemble



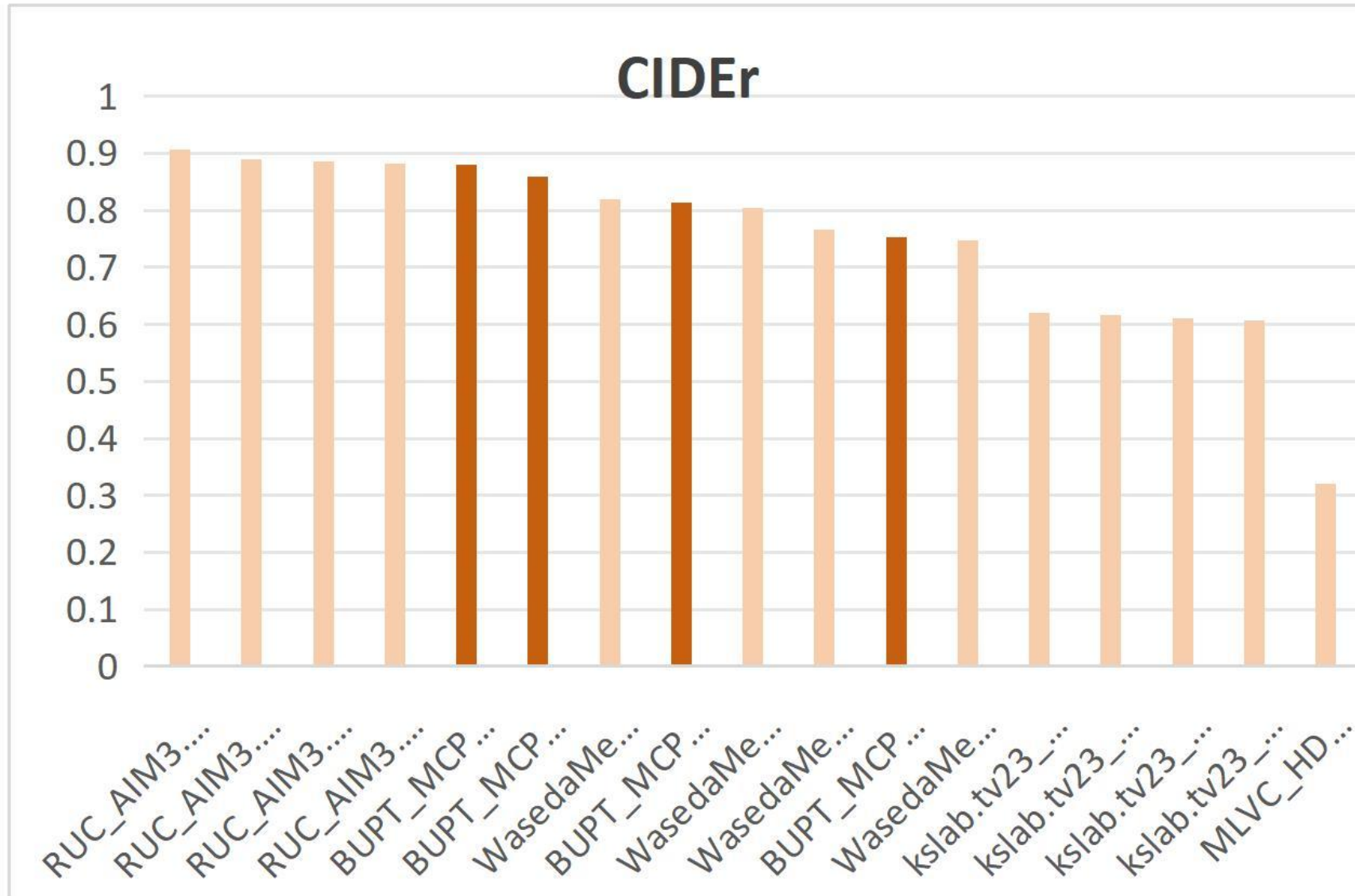LLM: Alpaca-Lora 7B

# ablation experiments

### data volume

| | |
|---|---|
| raw data | 54679 |
| the first round | 57217 |
| the second round | 60271 |
| the third round | 61620 |
| the fourth round | 64235 |

### CIDEr

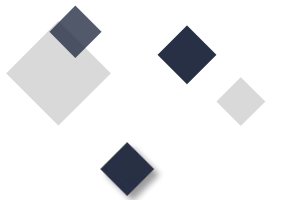| | CIDEr |
|---|---|
| base | 0.636213457 |
| round1 | 0.97475915 |
| round2 | 1.093636548 |
| round3 | 1.173122574 |
| round4 | 1.23467886 |

# 04 Conclusion

In this study, we transferred the BLIP2 model to the VTT task and proposed a cyclic data augmentation approach. Our experimental results on the TRECVID VTT dataset achieved a CIDEr score of 87.9, ranking second in the competition. Additionally, to extract the semantic information, we used LLM to integrate the results of multiple models, obtaining a state-of-the-art SPICE evaluation metric.

However, there is room for improvement in our temporal modeling approach, particularly in understanding complex motion behaviors. Additionally, due to time constraints, we were only able to perform cyclic data augmentation, for a limited five rounds.

THANKS

BUPT
MCPRL