

# The Effects of Pretraining on Video Guided Machine Translation

Ammon Shurtz, Lawry Sorenson, Stephen D. Richardson

LREC-COLING 2024

# RESEARCH QUESTIONS:

How can pretraining strategies improve Video-Guided Machine Translation (VMT)?

What modalities make the biggest impact for pretraining strategies in VMT?

# What is Video-Guided Machine Translation?

- Video-Guided Machine Translation provides video context to aid in translation.
- This can help the model disambiguate between potential meanings of the sentence.

# VaTeX Dataset

- Approx. 30,000 YouTube video clips labeled in English and Chinese by Amazon Mechanical Turk.
- 10 labels per video in both languages, 5 of which are translations of each other.
- Used for a shared task (VMT Challenge held at the ALVR Workshop) in 2020.

# VaTeX Example

A group of dogs of varying breeds pull a man in a sled across the finish line.



# VaTeX Example

## Original Sentence:

A group of dogs of varying breeds pull a man in a sled across the finish line.

## Masked input sentence:

A group \_\_\_\_\_ of varying \_\_\_\_\_ pull \_\_\_\_ man \_\_\_\_ a sled across the finish line.

## Produced translation:

一群不同品种的狗拉着一个乘坐雪橇的人穿过了终点线。

## Backtranslated by Google Translate:

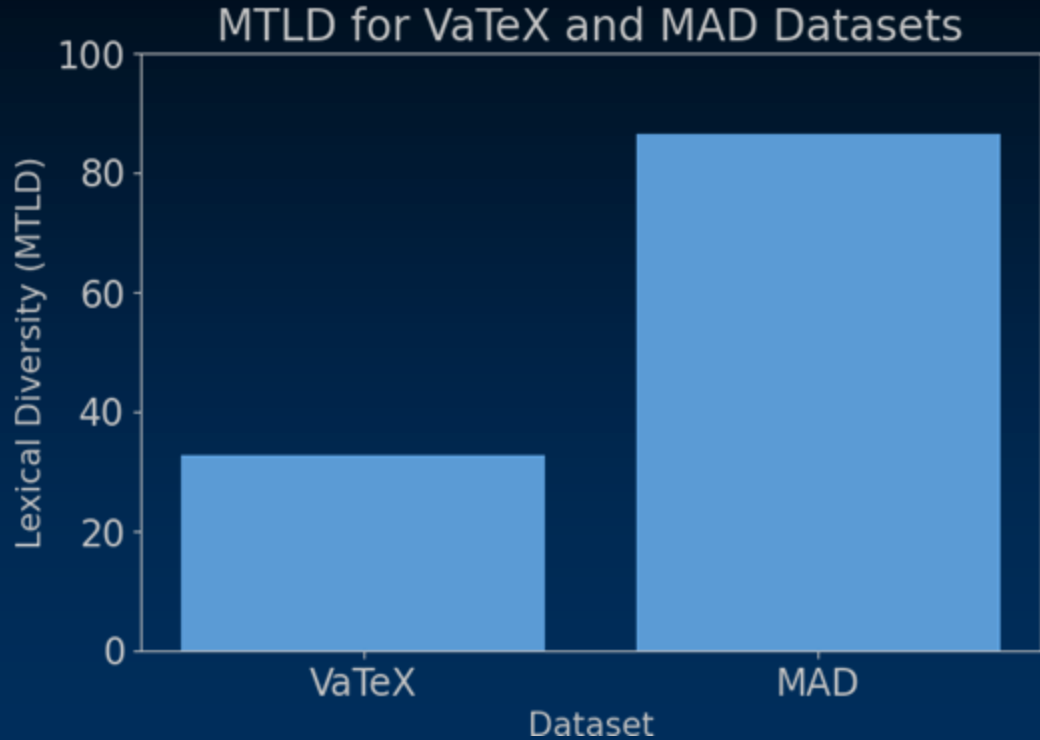
A pack of dogs of various breeds pulls a man on a sled across the finish line.

# Movie Audio Description (MAD) Dataset

- 600 Movies ~ 180k sentences
- Transcribed English Audio Descriptions (not dialogue)
- We created synthetic Chinese translations using Google Translate.
- Higher quality and more variety than the VaTeX dataset

# Lexical Richness

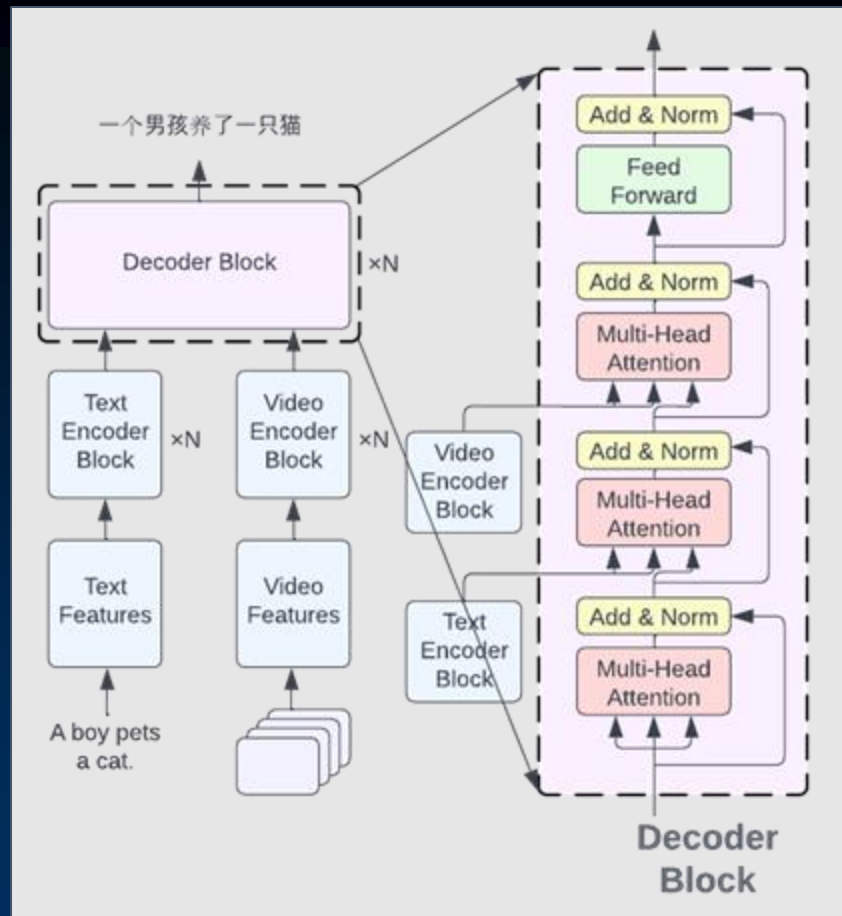
MTLD (Measure of Textual Lexical Diversity) reflects the amount of variety of words in text.





# Architecture

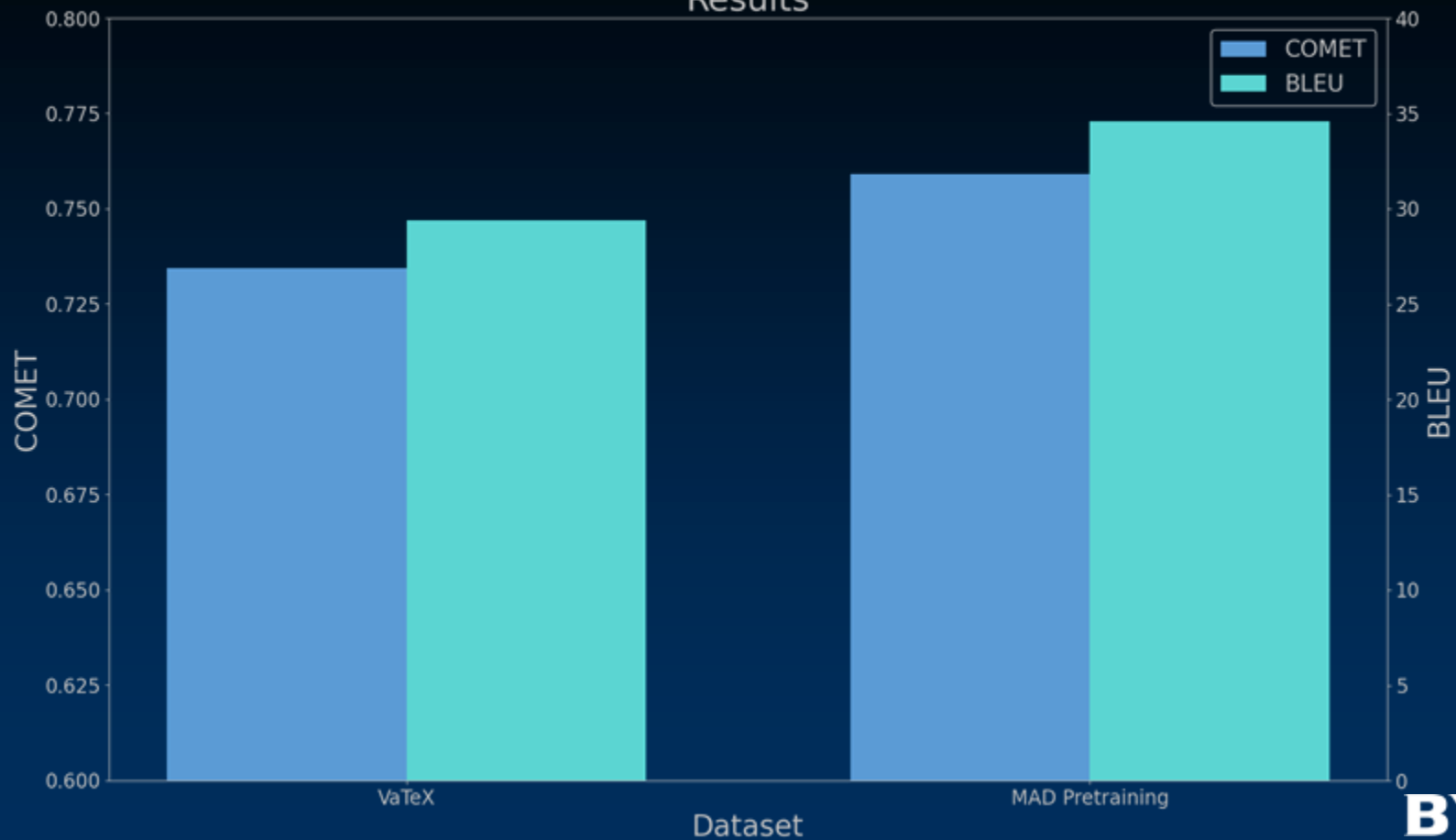
Transformer architecture with separate video and text encoders, decoder has dual attention.



# Experiments

- Try pretraining on the MAD dataset.
- Finetune on the VaTeX dataset.
- Test performance on the VaTeX dataset.

## Results

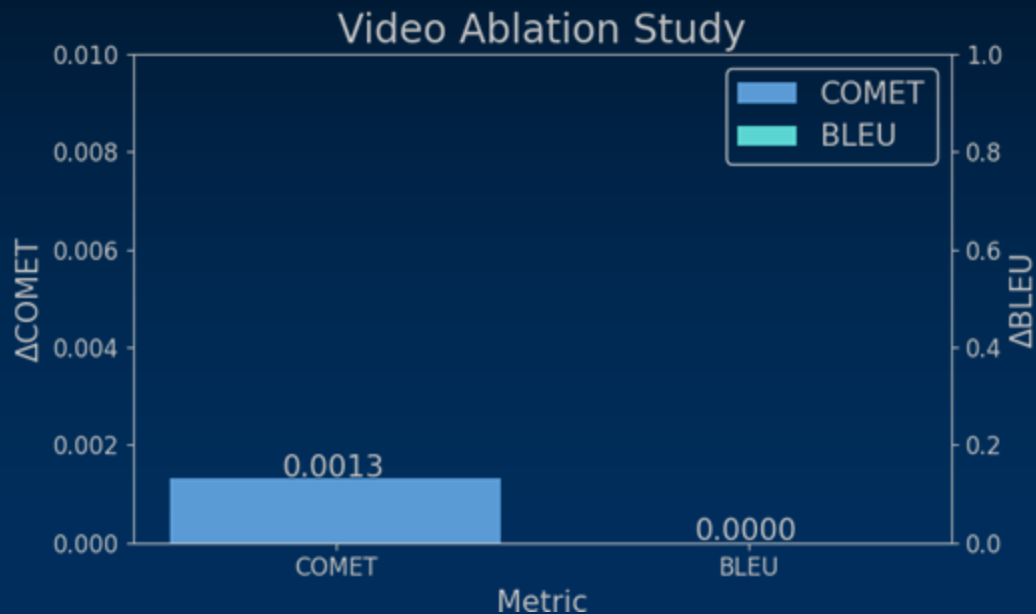


# Video Ablation Study

How much do the videos really matter compared to pretraining with more text?

# Video Ablation Study

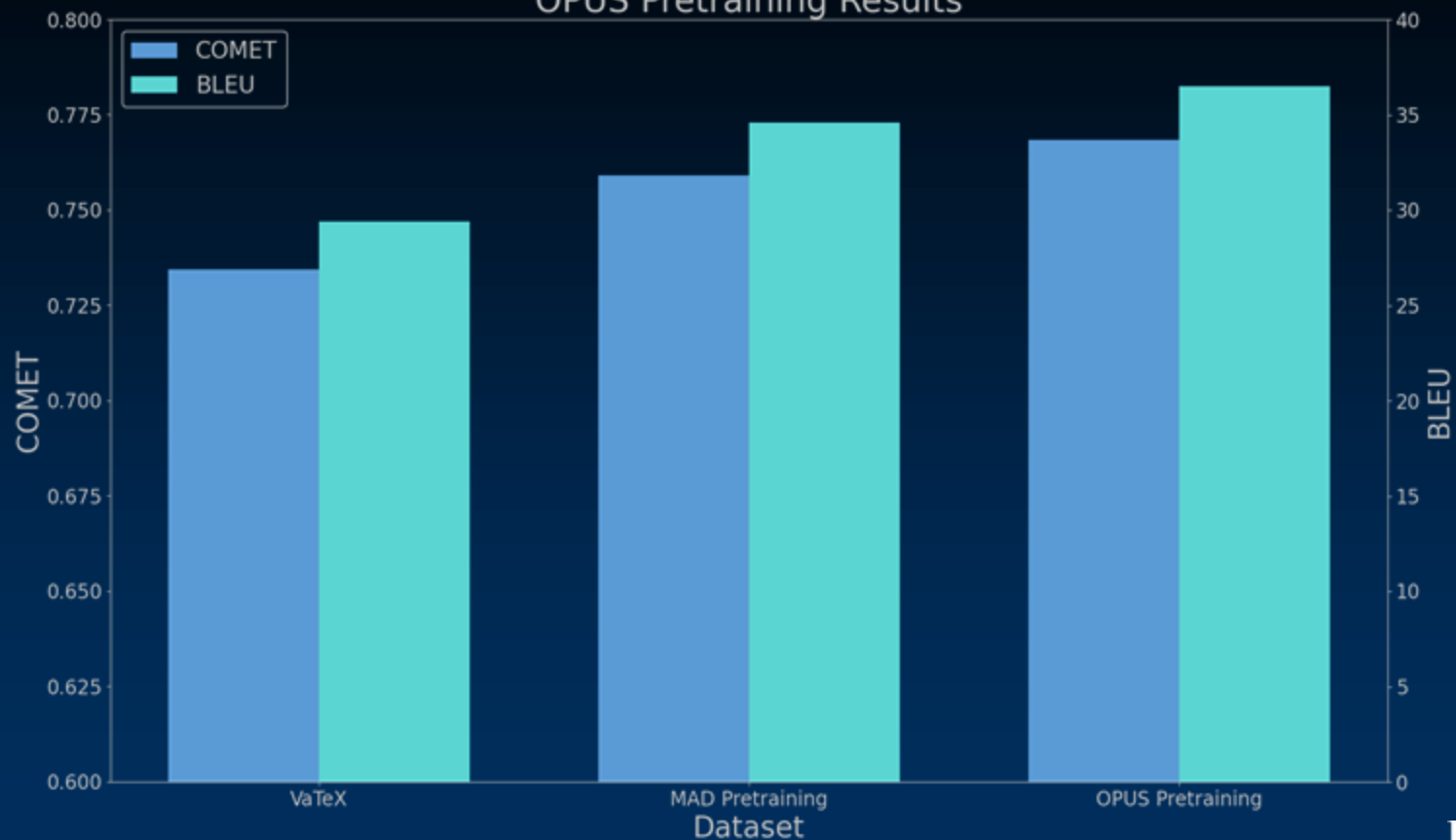
How much do the videos really matter compared to pretraining with more text?



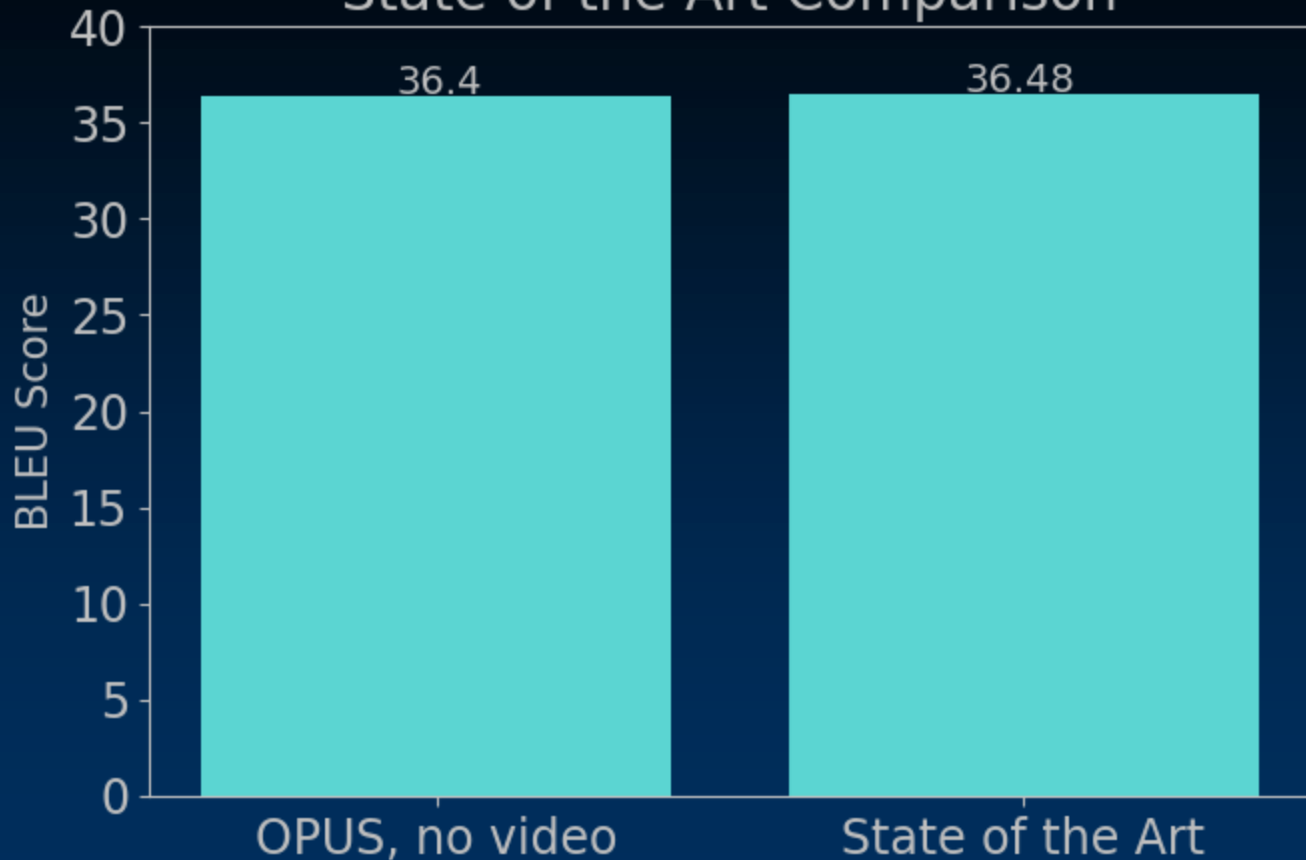
# Added Dataset: OPUS Subtitles

- Collection of 10M parallel subtitles in English and Chinese
- Similar domain to VaTeX and MAD, but only dialogue

## OPUS Pretraining Results



## State of the Art Comparison





# Conclusions

- Pretraining on diverse text significantly boosts performance.
- The VaTeX dataset sentences do not need video to be translated.

# Future Work

- BigVideo dataset
- MAD 2.0 dataset
- Translate from morphologically simple to more complex languages
- Unsupervised alignment between languages