# School of Computer Science and Technology
## (School of Cyberspace Security)

## LLM for Machine Translation in Low-Resource Languages: From Text-Only to Multimodal Learning

Mieradilijiang Maimaiti
https://www.miradeljan.com

ICNLP20260, Xi'An, China
20260321

# Outline

- **LLM Everywhere**

- Machine Translation with LLMs

- Challenges of LLM-based MT on LRLs

- Strategies for Text-only LLM-based MT

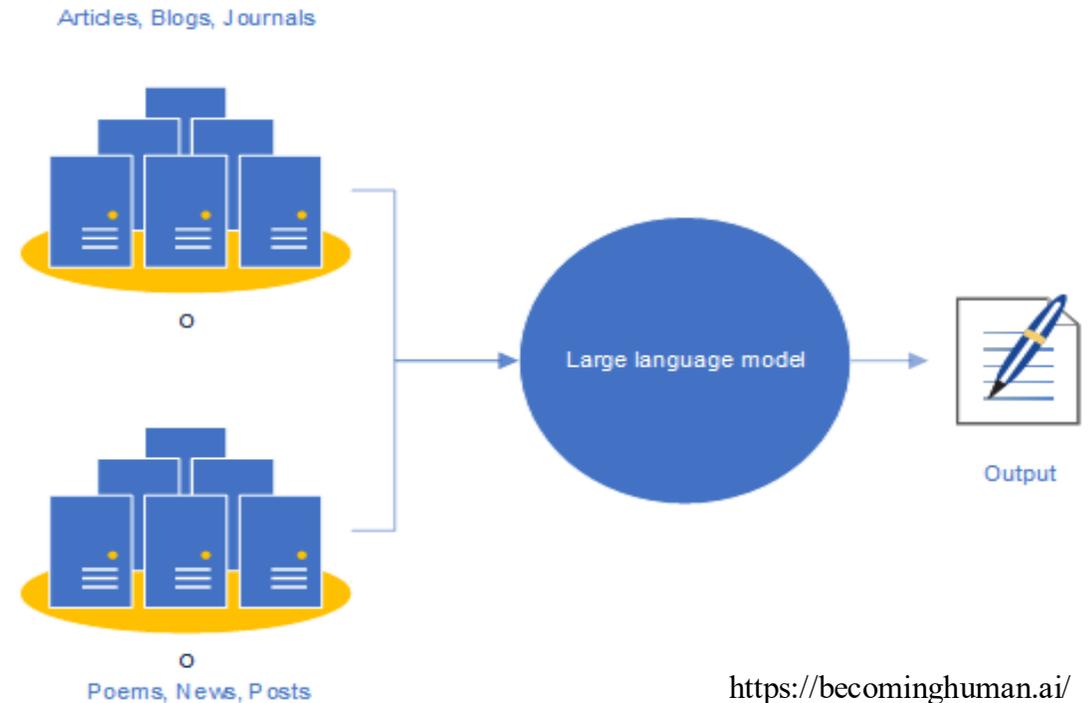- Findings on Multi-modal LLM-based MT

- Conclusion

# LLM Everywhere

## What Are Large Language Models?

Large Language Models (LLMs) are <span style="color:red">Transformer-based architectures</span> trained on massive unlabeled text via **self-supervised** and **semi-supervised learning**, enabling them to generate natural, context-aware text across diverse formats and styles.
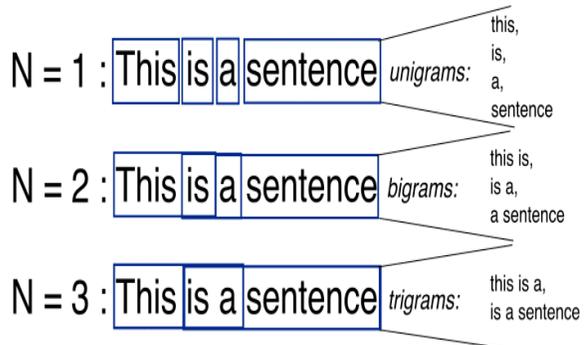
**LLM**



**Large-scale Pre-trained Language Model**

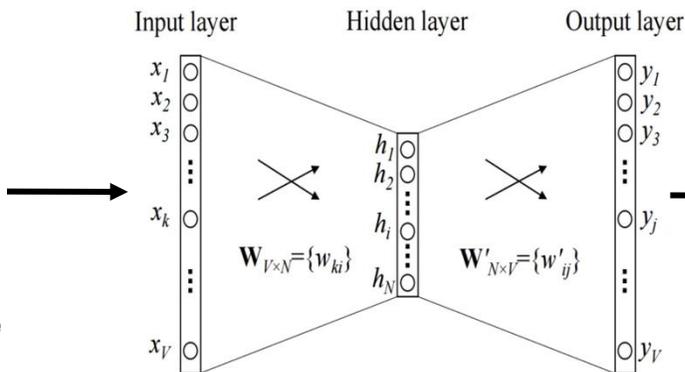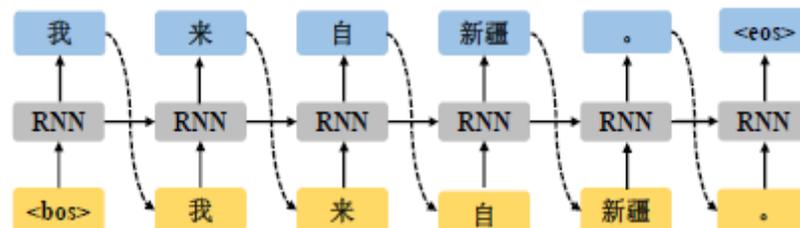Articles, Blogs, Journals

Large language model

Output

Poems, News, Posts

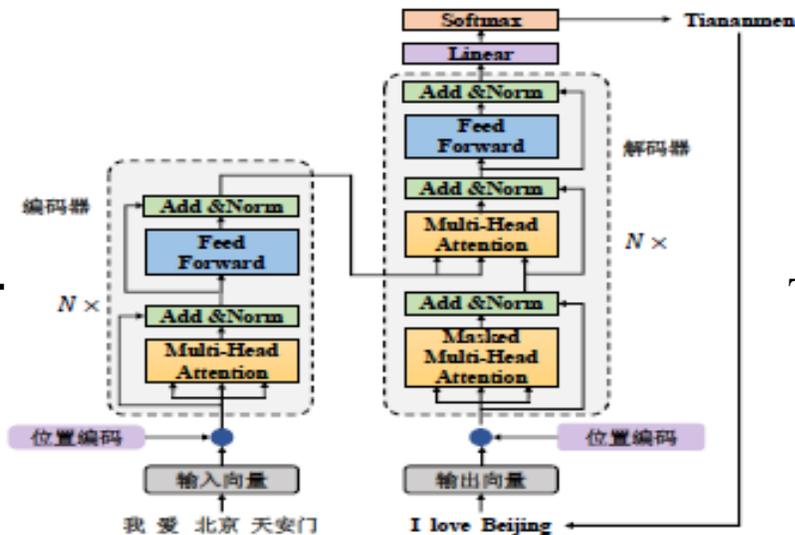https://becominghuman.ai/

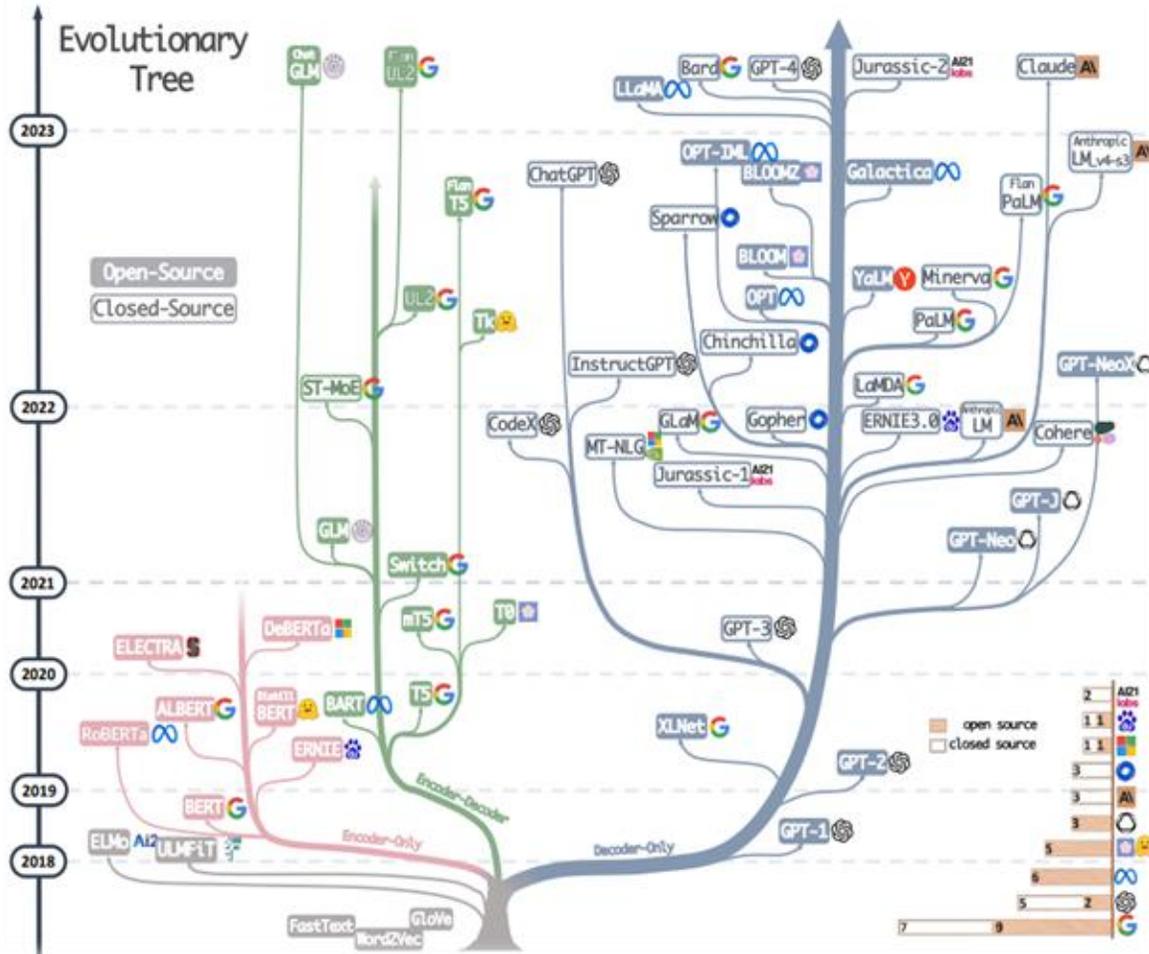# Language Models



N-gram LM

Word2Vec

RNNLM

Transformer

PLM

LLM

## The Evolutionary of LLM



(Yang et al., 2023)

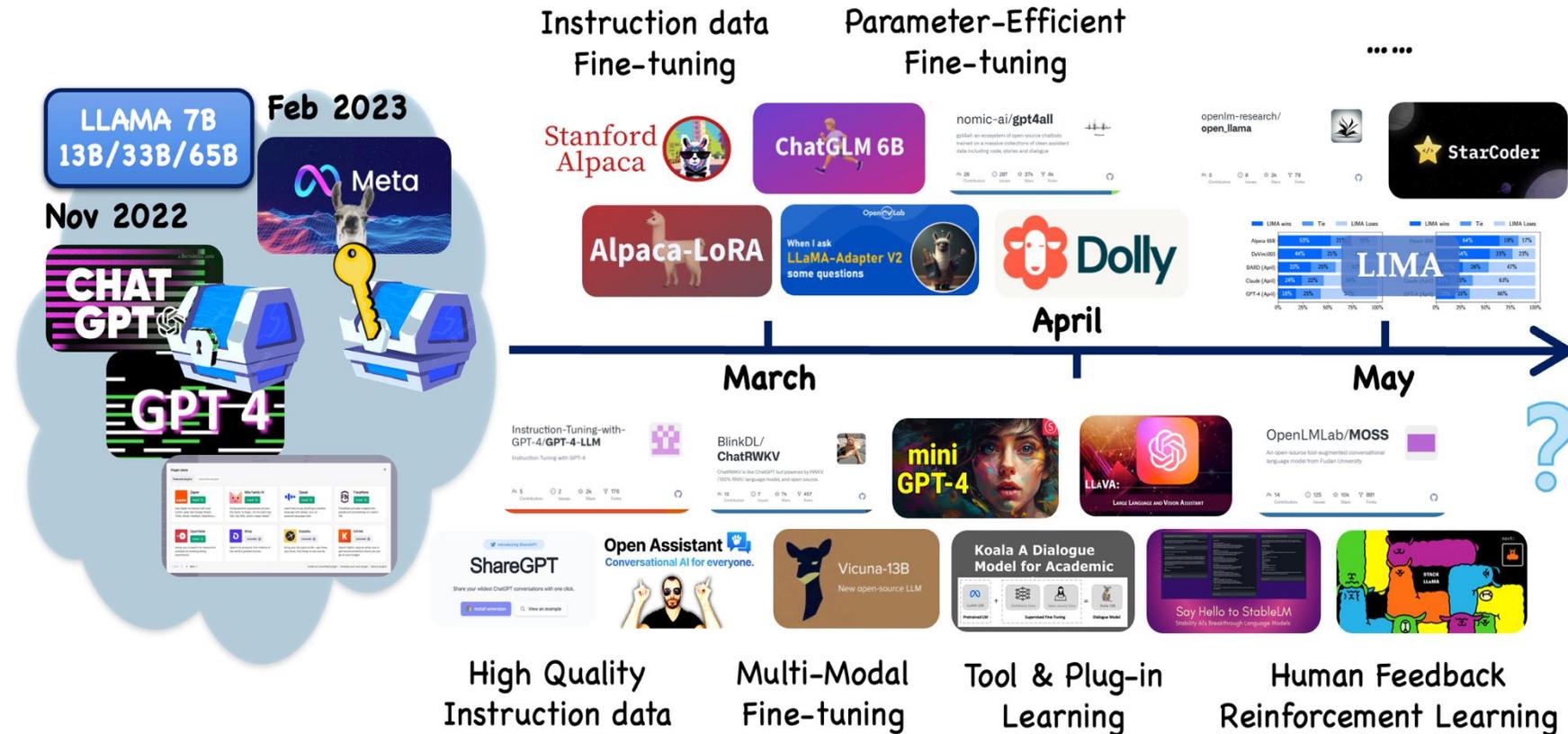•**GPT**Improving Language Understanding by Generative Pre-Training. 2018
•**GPT-2**Language Models are Unsupervised Multitask Learners. 2018.
•**GPT-3**"Language Models are Few-Shot Learners". NeurIPS 2020.
•**OPT**"OPT: Open Pre-trained Transformer Language Models". 2022.
•**PaLM**"PaLM: Scaling Language Modeling with Pathways". Aakanksha Chowdhery et al arXiv 2022.
•**BLOOM**"BLOOM: A 176B-Parameter Open-Access Multilingual Language Model".2022.
•**MT-NLG**"Using DeepSpeed and Megatron to Train Megatron-Turing NLG 530B, A Large-Scale GenerativeLanguage Model". 2021.
•**Gopher**"Scaling Language Models: Methods, Analysis & Insights from Training Gopher", 2021.
•**chinchilla** "Training Compute-Optimal Large Language Models". 2022.
•**LaMDA**"LaMDA: Language Models for Dialog Applications".2021.
•**LLaMA**"LLaMA: Open and Efficient Foundation Language Models". 2023.
•**GPT-4** "GPT-4 Technical Report".2023.
•**BloombergGPT**BloombergGPT: A Large Language Model for Finance, 2023,
•**GPT-NeoX-20B**: "GPT-NeoX-20B: An Open-Source Autoregressive Language Model".2022.

• • • • • •
**2026**

# Basic Fact - foundational models is rapidly evolving
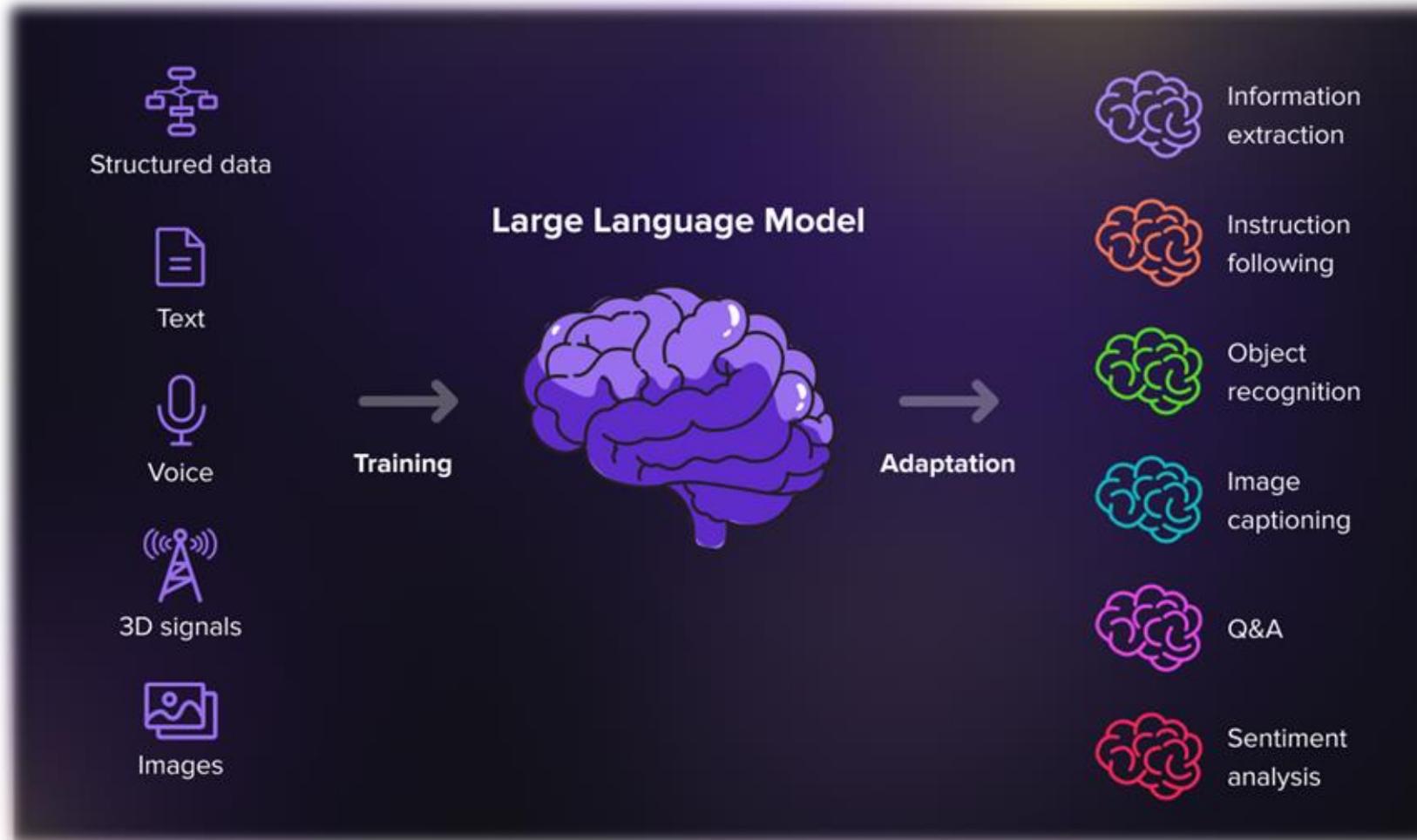
A growing number of instruction-tuned models have emerged



(Ding et al., CCL2023)

## LLM performance on NLP down-stream tasks



https://blog.gopenai.com/

## LLM Everywhere --- Conversation, QA, Outline, Recommend, MT, Coding

**You**
麻烦给我写一个python 脚本，需要把tsv 格式的文件转换成json 格式的文件。

**ChatGPT**
当然，我可以为您提供一个Python脚本来将TSV格式的文件转换成JSON格式的文件。以下是脚本的一个基本示例：

```python
python                                    Copy code

import csv
import json

# 将TSV文件转换为JSON文件的函数
def tsv_to_json(tsv_file_path, json_file_path):
    # 读取TSV文件
    with open(tsv_file_path, 'r', encoding='utf-8') as tsvfile:
        reader = csv.DictReader(tsvfile, delimiter='\t')

        # 创建一个列表来保存所有行的数据
        data = []
        for row in reader:
            data.append(row)

    # 写入JSON文件
    with open(json_file_path, 'w', encoding='utf-8') as jsonfile:
        json.dump(data, jsonfile, indent=4, ensure_ascii=False)

# 使用示例
tsv_file = 'example.tsv'  # TSV文件的路径
json_file = 'output.json' # 输出JSON文件的路径
tsv_to_json(tsv_file, json_file)
```

**You**
我不喜欢用with open 格式来文件操作，重写一下！

**ChatGPT**
当然，我可以按照您的要求修改脚本，不使用`with open`来进行文件操作。请注意，虽然`with open`是推荐的方式，因为它可以自动处理文件的关闭，但是如果您不想使用它，可以直接使用`open`和`close`来手动管理文件的打开和关闭。下面是修改后的脚本：

```python
python                                    Copy code

import csv
import json

def tsv_to_json(tsv_file_path, json_file_path):
    # 打开TSV文件
    tsvfile = open(tsv_file_path, 'r', encoding='utf-8')
    reader = csv.DictReader(tsvfile, delimiter='\t')

    # 创建一个列表来保存所有行的数据
    data = []
    for row in reader:
        data.append(row)

    # 关闭TSV文件
    tsvfile.close()

    # 打开JSON文件
    jsonfile = open(json_file_path, 'w', encoding='utf-8')
    json.dump(data, jsonfile, indent=4, ensure_ascii=False)

    # 关闭JSON文件
    jsonfile.close()

# 使用示例
tsv_file = 'example.tsv'   # TSV文件的路径
json_file = 'output.json' # 输出JSON文件的路径
tsv_to_json(tsv_file, json_file)
```
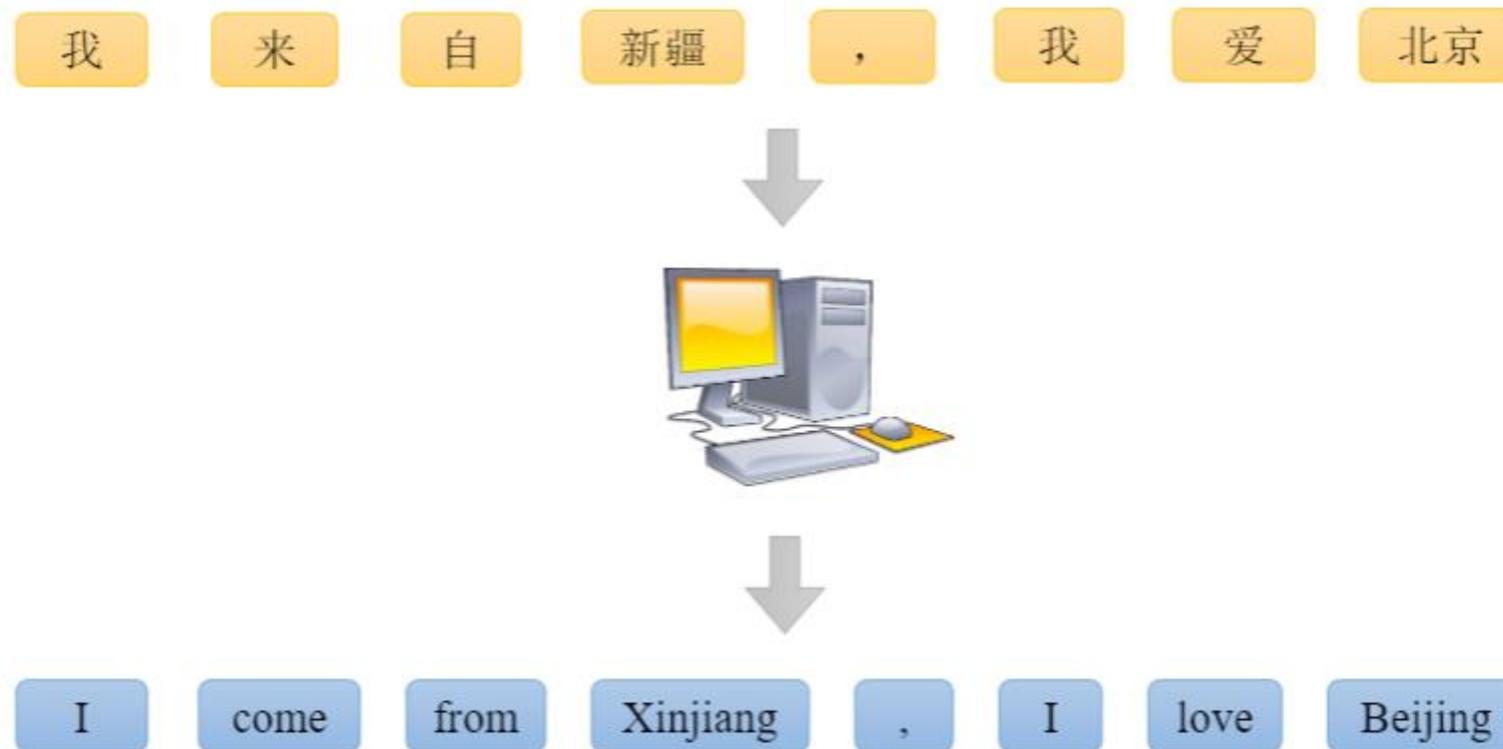
# Outline

- **LLM Everywhere**

- **Machine Translation with LLMs**

- Challenges of LLM-based MT on LRLs

- Strategies for Text-only LLM-based MT

- Findings on Multi-modal LLM-based MT

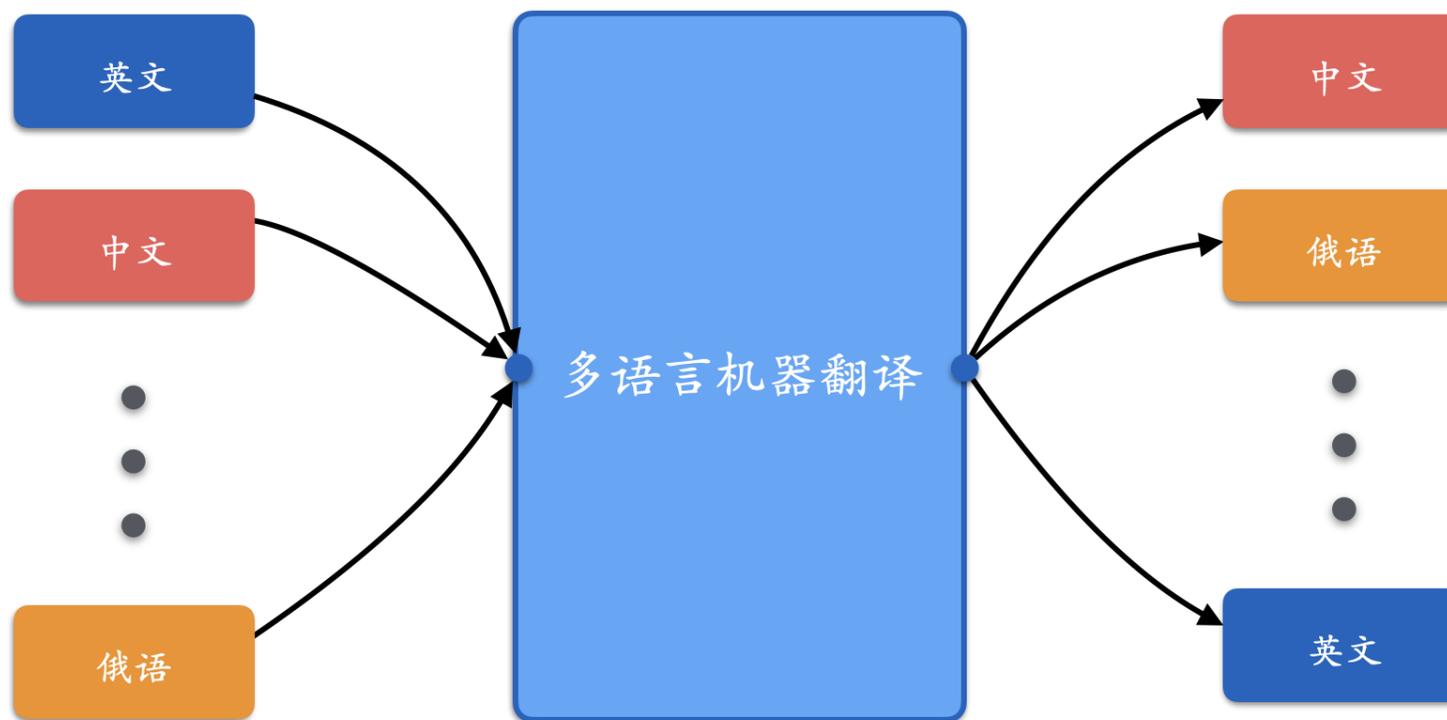- Conclusion

# Machine Translation with LLMs

## Machine Translation
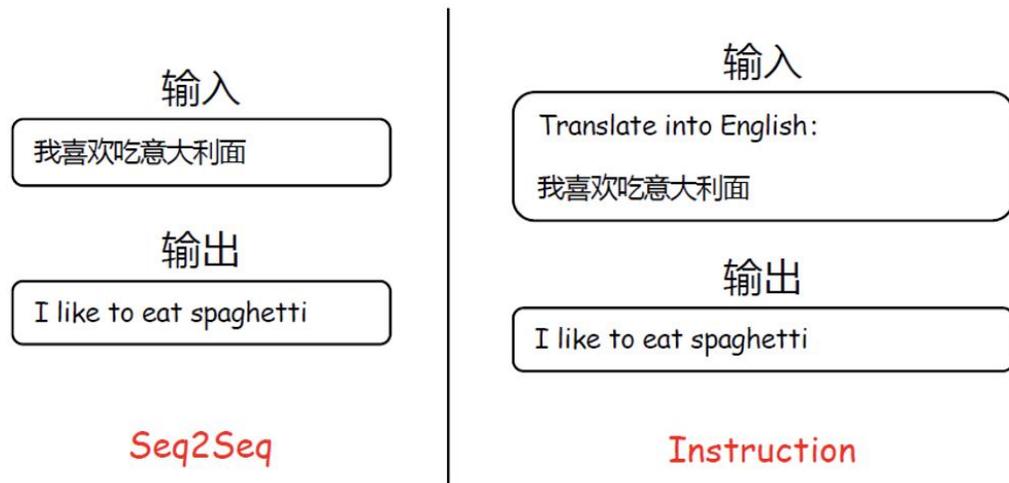
- Use computers to transform one sequence into another.

- Multilingual machine translation aims to enable translation between multiple languages using a single unified model, thereby significantly improving translation performance for low-resource languages while greatly reducing deployment costs.

英文

中文

俄语

多语言机器翻译

中文

俄语

英文

Huang, CCMT2023

## Machine Translation with LLM

- General Model v.s. Specific Task（Translation）

  - Specifying model behavior through instruction.
  - In-context Learning (ICL)

输入
Translate into English:

Chinese: 你喜欢吃什么?

English: What do you like to eat?

Chinese: 我喜欢吃意大利面

English:

输出

I like to eat spaghetti

**In-context Learning**

输入

我喜欢吃意大利面

输出

I like to eat spaghetti

**Seq2Seq**

输入

Translate into English:

我喜欢吃意大利面

输出

I like to eat spaghetti

**Instruction**

Acquiring the ability to directly **follow instructions** through instruction learning.

A **small number of examples** obtained during **pre-training**.

Prof. Huang CCL 2023

The paradigm of learning is shifting

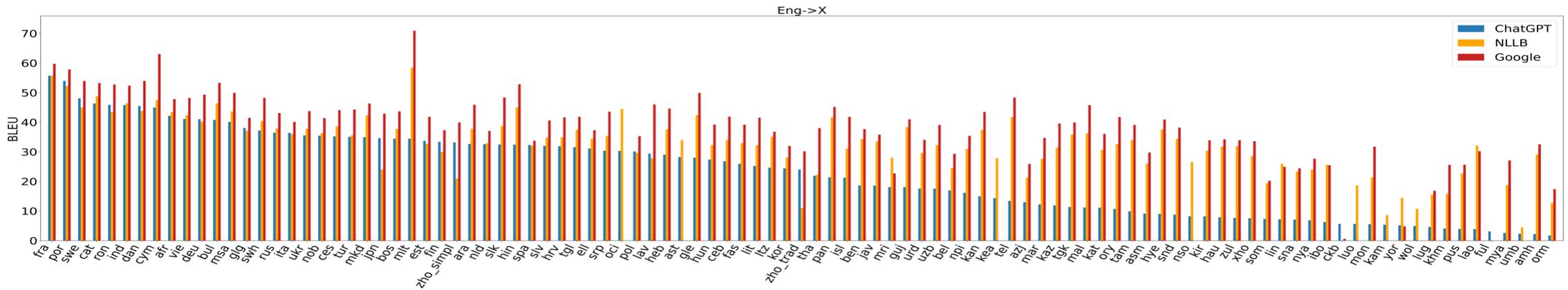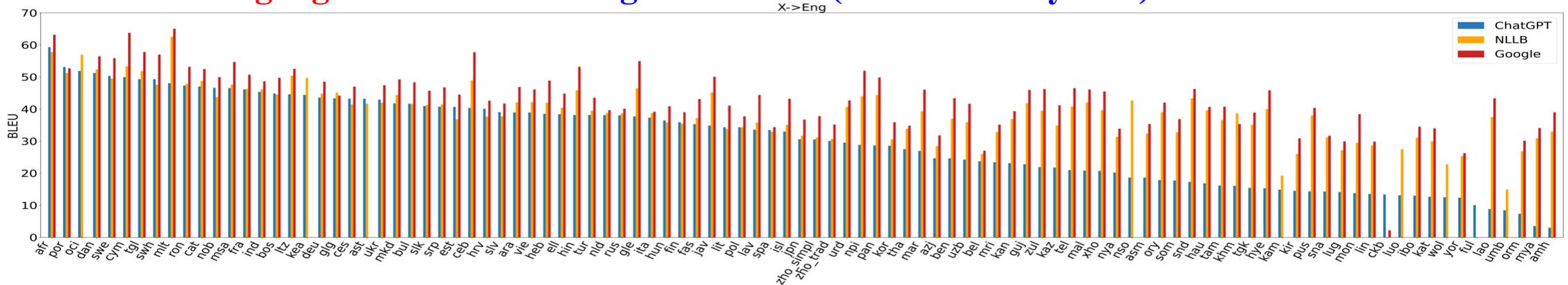- **NMT:** Translation knowledge is primarily learned from parallel data.



- **LLM:** General knowledge (including translation knowledge) is primarily learned from monolingual data.



In-context Learning

Prof. Huang CCL 2023

- **LLMs still lag behind strong supervised baselines：** ChatGPT (16.67%) v.s. NLLB (83.33%)
- **Low-resource languages: far behind Google Translate (commercial system)**



Zhu W, et al. Multilingual Machine Translation with Large Language Models: Empirical Results and Analysis, arXiv:2304.04675

Prof. Huang CCL 2023

# Pros of LLM-based MT

1. Context Understanding
→ Accuracy & fluency

2. Interactive Translation
→ Real-time user feedback / guidance

3. Versatile Text Processing
→ Formal, spoken, slang

4. Document-Level Translation
→ Consistency & context in long/complex texts

5. Multilingual Capability
→ Simultaneous multi-language translation

**Context**



Prof. Liu CCMT 2023

# Pros of LLM-based MT

**Interactive Translation**

Please translate it into English:
大型语言模型已被证明在使用少量样本学习的情况下，能在各种自然语言任务中表现出色。为了进一步了解规模对少量样本学习的影响，我们训练了一个具有5400亿参数、密集激活的Transformer语言模型，我们称之为Pathways语言模型。

The translation of the provided text into English is:

"Large language models have been shown to perform well in various natural language tasks with limited sample learning. To further understand the impact of scale on limited sample learning, we trained a Transformer language model with 540 billion parameters and dense activations, which we call the Pathways language model."

**Versatile Text Processing**

Please translate "远水救不了近火" into English.

The Chinese phrase "远水救不了近火" literally translates to "distant water cannot put out a nearby fire." This idiom is used to express the idea that help that is too far away is of no use in an emergency or urgent situation. In other words, when you're facing an immediate problem, solutions that are too far off in time or space won't be of any help.

Prof. Liu CCMT 2023

## Document-Level Translation



Wang et al. Document-Level Machine Translation with Large Language Models. EMNLP 2023.

Prof. Liu CCMT 2023

# Cons of LLM-based MT

1. Domain Expertise & Knowledge Transfer
   - Limited domain-specific terminology capture
   - Catastrophic forgetting during adaptation

2. Low-Resource Languages
   - Rare language pairs: room for improvement

3. Evaluation Bottlenecks
   - Traditional metrics fall short for high-quality translations

4. Computational Overhead
   - High cost for fine-tuning & adaptation

# Outline

- LLM Everywhere

- Machine Translation with LLMs

- Challenges of LLM-based MT on LRLs

- Strategies for Text-only LLM-based MT

- Findings on Multi-modal LLM-based MT

- Conclusion

# Challenges of LLM-based MT on LRLs

Original Input: " جۇڭگو پۇتبولچىلار جەمئىيتى بۇ قېتىمقى ۋەقەگە ۋ قەقەگە يۇكسەك ئەھمىيەت بېرىپ، قىسمەن ئەسەبىي تور ئەھلىنىڭ توردا پۇتبول كەسپكارلىرىغا بوھتان چاپلىغانلىقى، تۆھمەت چاپلىغانلىقى ۋە پىتنە-ئىغۋا توقۇپ ئىش تېرىغانلىقىنى كۈچلۈك ئەيىبلىدى. "



Ug→Zh translation sample

Original Input: "جۇڭگو پۇتبولچىلار جەمئىيتى بۇ قېتىمقى ۋەقەگە يۇكسەك ئەھمىيەت بېرىپ، قىسمەن ئەسەبىي تور ئەھلىننىڭ توردا پۇتبول كەسپكارلىرىغا بوھتان چاپلىغانلىقى، تۆھمەت چاپلىغانلىقى ۋە پىتنە-ئىغۋا توقۇپ ئىش تېرىغانلىقىنى كۈچلۈك ئەيىبلىدى."



Ug→Zh translation sample

Original Input: " جۇۈگگو پۇ‌تبولچىلار جەمئىيتى بۇ قېتىمقى ۋەقەگە يۇ‌كسەك ئەھمىيەت بېرىپ، قىسمەن ئەسەببي تور ئەھلىننىڭ توردا پۇ‌تبول كەسپكارلرىغا بوھتان چاپلىغانلىقى، تۆھمەت چاپلىغانلىقى ۋە پىتنە-ئىغۋا توقۇپ ئىش تېرىغانلىقنى كۇ‌چلۇك ئەيىبلىدى. "



Ug→Zh translation sample

Original Input: "جۇگۇگو پۇتبولچىلار جەمئىيتى بۇ قېتىمقى ۋەقەگە يۇكسەك ئەھمىيەت بېرىپ، قىسمەن ئەسەببى تور ئەھلىننىڭ توردا پۇتبول كەسپكارلىرىغا بوھتان چاپلىغانلىقى، تۆھمەت چاپلىغانلىقى ۋە پىتنە-ئىغۋا توقۇپ ئىش تېرىغانلىقنى كۇچلۇك ئەيىبلىدى."



Ug→Zh translation sample

Original Input："我来自我国最西部地理位置比较偏僻、经济不发达、教育水平和条件很低的以及医疗和医学落后的新疆，但是我一直以来坚持不懈的努力而向往美好的未来。"



Zh →Ug translation sample

Original Input："我来自我国最西部地理位置比较偏僻、经济不发达、教育水平和条件很低的以及医疗和医学落后的新疆，但是我一直以来坚持不懈的努力而向往美好的未来。"



Zh →Ug translation sample

Original Input："我来自我国最西部地理位置比较偏僻、经济不发达、教育水平和条件很低的以及医疗和医学落后的新疆，但是我一直以来坚持不懈的努力而向往美好的未来。"



Zh →Ug translation sample

Original Input：　"我来自我国最西部地理位置比较偏僻、经济不发达、教育水平和条件很低的以及医疗和医学落后的新疆，但是我一直以来坚持不懈的努力而向往美好的未来。"



Zh →Ug translation sample

# MT Performance Comparison

Zh →Ug translation sample   >     >     >  

Ug→Zh translation sample   >     >     >  

Original Input: " جۇڭگو پۇتبولچىلار جەمئىيىتى بۇ قېتىمقى ۋەقەگە يۇكسەك ئەھمىيەت بېرىپ، قىسمەن ئەسەبىي تور ئەھلىنىڭ توردا پۇتبول كەسپكارلىرىغا بوھتان "



Ug→Zh translation sample

# Open model (Qwen): MT

Original Input: " جۇڭگو پۇتبولچىلار جەمئىيىتى بۇ قېتىمقى ۋەقەگە يۇكسەك ئەھمىيەت بېرىپ، قىسمەن ئەسەببى تور ئەھلىنىڭ توردا پۇتبول كەسپكارلىرىغا بوھتان چاپلىغانلىقى، تۆھمەت چاپلىغانلىقى ۋە پىتنە-ئىغۋا توقۇپ ئىش تېرىغانلىقنى كۇچلۇك ئەيىبلىدى. "



Ug→Zh translation sample

# Open model (Yuanbao): MT

Original Input: " جۇڭگو پۇتبولچىلار جەمئىيتى بۇ قېتىمقى ۋەقەگە يۇكسەك ئەھمىيەت بېرىپ، قىسمەن ئەسەببي تور ئەھلىننىڭ توردا پۇتبول كەسپكارلىرىغا بوھتان چاپلىغانلىقى، تۆھمەت چاپلىغانلىقى ۋە پىتنە-ئىغۋا توقۇپ ئىش تېرىغانلىقنى كۈچلۈك ئەيىبلىدى. "



Ug→Zh translation sample

Original Input："我来自我国最西部地理位置比较偏僻、经济不发达、教育水平和条件很低的以及医疗和医学落后的新疆，但是我一直以来坚持不懈的努力而向往美好的未来。"



Zh →Ug translation sample

Original Input："我来自我国最西部地理位置比较偏僻、经济不发达、教育水平和条件很低的以及医疗和医学落后的新疆，但是我一直以来坚持不懈的努力而向往美好的未来。"



Zh →Ug translation sample

# Open model (Yuanbao): MT

**Original Input：** "我来自我国最西部地理位置比较偏僻、经济不发达、教育水平和条件很低的以及医疗和医学落后的新疆，但是我一直以来坚持不懈的努力而向往美好的未来。"



Zh →Ug translation sample

# MT Performance Comparison

Ug→Zh translation sample



Zh →Ug translation sample

# Main Challenges

- **Omission** – Missing content or phrases

- **Over-translation** – Unnecessary or redundant content

- **Semantic Errors** – Incorrect meaning or misinterpretation

- **Syntactic Errors** – Grammatical or structural issues

- **Word-level Errors** – Inaccurate lexical choices

# Outline

- LLM Everywhere

- Machine Translation with LLMs

- Challenges of LLM-based MT on LRLs

- **Strategies for Text-only LLM-based MT**

- Findings on Multi-modal LLM-based MT

- Conclusion

# Strategies for Text-only LLM-based MT

Low-Resource Bottleneck: LLM performance constrained by scarce, low-quality instruction datasets

Proposed Solution: Refined Instruction Tuning — an automated pipeline for high-quality Chinese→X instruction corpora

# Constructing LRLs Instruction

- Experiment

  - Dataset

    - Constructed high-quality instruction datasets for 8 low-resource Chinese → X directions (with approximately 5k instances per direction), covering Uyghur (ug), Tibetan (bo), Persian (fa), Hebrew (he), Urdu (ur), Bengali (bn), Vietnamese (vi), and Indonesian (id).

| Region | Language | Size | AverageLength |
|---|---|---|---|
| China | zh→ug | 5k | 98.93 |
| | zh→ur | 5k | 106.36 |
| Middle East | zh→fa | 5k | 102.74 |
| | zh→he | 5k | 99.09 |
| South Asia | zh→ur | 5k | 104.03 |
| | zh→bn | 5k | 97.08 |
| Southeast Asia | zh→vi | 5k | 97.04 |
| | zh→id | 5k | 97.27 |

- Experiment

  - Main experiment

    - Main experimental results on the FLOWERS+, IWSLT, and CCMatrix test sets, with scores being the average of four evaluation metrics: Chrf++, COMET, XCOMET-XL, and BLEURT.
    - FLORES+ (8 languages): Achieved a SOTA macro-average of 61.78 (+5.57 over Direct baseline) .

| Method | FLORES+ (zh→xx) | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | ug | bo | ur | bn | fa | he | vi | id |
| Madlad | 42.99 | 40.64 | 54.10 | 56.85 | 60.31 | 61.74 | 70.58 | 75.08 |
| Direct | 35.52 | 40.73 | 52.94 | 55.53 | 64.25 | 54.29 | 74.30 | 72.11 |
| COD | 38.03 | **49.35** | 53.48 | 60.75 | 60.70 | 60.76 | 72.62 | 77.10 |
| MAPS | 37.15 | 48.22 | 55.28 | **63.21** | 66.15 | **63.06** | 74.95 | 78.44 |
| CompTrans | 34.37 | 46.58 | 52.10 | 53.10 | 62.92 | 58.51 | 73.59 | 77.17 |
| Ours | **49.66** | 45.84 | **58.11** | 57.45 | **66.24** | 62.29 | **75.57** | **79.07** |

- Motivation



Figure 1: The difference between machine and human translation in an English→Chinese example

- Inherent Limitations: LLMs struggle in low-resource translation via **Direct Prompting** due to a lack of explicit linguistic reasoning.

- Information Deficiency: Existing CoT methods suffer from coarse granularity, failing to systematically integrate the deep, **multi-source knowledge** required for precise translation.

- Dataset & Evaluation Metrics & Models

  - Dataset

    - The experiments utilize the official FLORES-200 benchmark alongside custom-curated test sets from IWSLT and CCMatrix, which were manually cleaned and filtered to ensure high-quality evaluation across diverse low-resource contexts.

| Dataset | Language Pairs | Size (per pair) |
|---------|----------------|-----------------|
| FLORES-200 | En $\leftrightarrow$ {Fa, Ur, Lo, Uz} | 1,012 |
| IWSLT | Zh $\rightarrow$ {Fa, He, Id, Vi} | 1,000 |
| CCMatrix | Zh $\rightarrow$ {Fa, He, Id, Bn, Vi} | 1,000 |

Table 1: Statistics of Evaluation Datasets

- Experiment

  - Main experiment

    - Achieved superior performance in most directions, notably reaching 88.10 on en-fa and 88.24 on en-uz. The model consistently outperforms 0-shot, MAPS and CompTra baselines across diverse language pairs.

    - Demonstrated substantial improvements in challenging low-resource directions, such as en-lo and en-uz. These results highlight strong cross-lingual transfer and generalization capabilities.

| Method | Middle East | | South Asia | | Southeast Asia | | Middle Asia | |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| | en-fa | fa-en | en-ur | ur-en | en-lo | lo-en | en-uz | uz-en |
| Madlad | 85.31 | 87.78 | 79.15 | 87.05 | **83.59** | **87.80** | N/A | N/A |
| 0-shot | 86.43 | 87.82 | 80.97 | 86.98 | 78.15 | 85.83 | 84.55 | 86.10 |
| MAPS | 87.59 | **88.38** | 82.24 | 87.61 | 79.86 | 80.54 | 87.25 | **87.11** |
| COD | 85.63 | 87.69 | 80.93 | 86.96 | 78.01 | 85.67 | 85.20 | 85.88 |
| CompTra | 84.77 | 87.26 | 78.79 | 86.30 | 72.75 | 84.61 | 84.00 | 85.59 |
| **Ours** | **88.10** | 88.27 | **82.91** | **87.62** | 82.14 | 86.57 | **88.24** | 86.96 |

Table 2: Main results across multiple low-resource language pairs on FLORES-200 using Qwen3-30B-A3B-Instruct-2507, evaluated with COMET.

- Motivation



- Although LLMs demonstrate promising translation and self-refinement capabilities , their performance remains severely constrained in low-resource scenarios.

- Leveraging external knowledge via RAG assists LLM translation , but relying solely on contextual parallel examples and single-round generation is insufficient to resolve diverse errors.

- # Experiment

  - ## Main experiment

    - Main experimental results on the FLORES-200, NTREX-128, and TICO-19 with XCOMET-XL and BLEURT-20 metrics.

    - The method consistently achieves the strongest or highly competitive performance on the multi-domain dataset FLORES-200.

| Method | Armenian | | Azerbaijani | | Hebrew | | Lao | |
|---|---|---|---|---|---|---|---|---|
| | XCOMET | BLEURT | XCOMET | BLEURT | XCOMET | BLEURT | XCOMET | BLEURT |
| 0-shot | 59.85 | 67.48 | 58.09 | 60.55 | 62.30 | 59.99 | 46.02 | 57.33 |
| Vanilla RAG | 64.25 | 70.19 | 61.49 | 62.06 | 65.71 | 62.17 | 52.86 | 65.23 |
| COD | 64.14 | 68.43 | 60.67 | 61.60 | 64.47 | 61.36 | 48.56 | 57.96 |
| MAPS | 67.88 | 72.08 | 64.41 | 64.53 | 69.03 | 64.45 | 52.92 | **63.77** |
| TEaR | 63.32 | 69.27 | 60.86 | 61.36 | 66.00 | 62.28 | 48.35 | 59.88 |
| CompTra | 55.37 | 58.29 | 59.86 | 61.29 | 65.26 | 62.01 | 48.95 | 52.27 |
| OURS | **68.32** | **71.80** | **65.44** | **64.68** | **71.45** | **65.83** | **54.01** | 64.44 |

| Method | Khmer | | Tamil | | Urdu | | Bengali | |
|---|---|---|---|---|---|---|---|---|
| | XCOMET | BLEURT | XCOMET | BLEURT | XCOMET | BLEURT | XCOMET | BLEURT |
| 0-shot | 49.51 | 52.69 | 49.78 | 70.25 | 62.37 | 55.22 | 63.01 | 69.48 |
| Vanilla RAG | 55.40 | 58.19 | 51.94 | 72.32 | 64.22 | 55.88 | 66.08 | 71.43 |
| COD | 50.39 | 52.01 | 48.26 | 67.96 | 58.44 | 54.01 | 62.63 | 69.46 |
| MAPS | 55.59 | 57.39 | **53.43** | 72.87 | 67.03 | **56.80** | 68.19 | 72.30 |
| TEaR | 51.50 | 53.59 | 51.25 | 71.53 | 63.93 | 55.27 | 65.72 | 70.63 |
| CompTra | 48.47 | 43.73 | 43.05 | 55.45 | 59.87 | 54.61 | 52.95 | 58.87 |
| OURS | **56.99** | **58.14** | 53.74 | **73.76** | **67.75** | 56.95 | **68.36** | **72.39** |

- Motivation



- Low-resource languages lack sufficient high-quality bilingual data, making traditional RAG hard to apply.

- LLMs are vulnerable to lexical ambiguity in cross-domain translation, which may cause semantic hallucinations.

- Existing RAG systems rely on coarse sentence- or document-level retrieval and cannot precisely disambiguate high-risk ambiguous words.

# Word Disambiguation

- ## Dataset

  - We evaluate Domain-aware RAG on three benchmarks: X-Bench , WMT , and IWSLT. As shown in Table, the X-Bench is a multi-domain benchmark that includes six Languages, such as Bengali (bn), Hungarian (hu), Urdu (ur), Persian (fa), Malay (ms), Indonesian (id), and seven domains. For each Language, we sample from public OPUS corpora.

| Lan. | OpenSubtitles | TED 2020 | QED | Tanzil | wikimedia | WikiMatrix | Europarl | WMT-News | TEP | Size |
|------|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| hu | ✓ | ✓ | ✓ | | | ✓ | ✓ | ✓ | | 1.2K |
| ms | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | 1.1K |
| ur | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | 0.8K |
| bn | ✓ | ✓ | ✓ | ✓ | | ✓ | | | | 0.9K |
| fa | | ✓ | ✓ | ✓ | | ✓ | | | ✓ | 0.9K |
| id | ✓ | ✓ | ✓ | ✓ | | ✓ | | | | 1.2K |

- ## Ablation study

  - ### The Influence of Decision Agent.

    - We sample 500 instances per direction from X-Bench, WMT, and IWSLT and use ChatGPT-4o for scoring.Qwen3-32B is the most consistent top performer and is adopted as the default decision agent.

# Outline

- LLM Everywhere

- Machine Translation with LLMs

- Challenges of LLM-based MT on LRLs

- Strategies for Text-only LLM-based MT

- Findings on Multi-modal LLM-based MT

- Conclusion

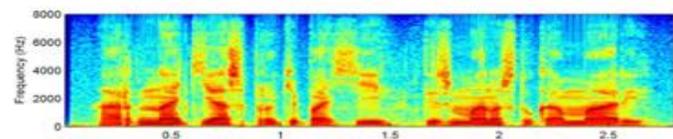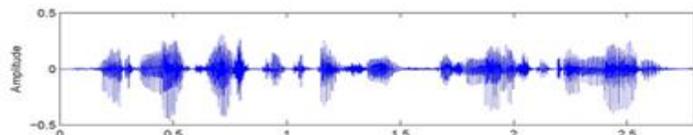# Findings on Multi-modal LLM-based MT

**Speech Processing vs NLP**

**From Continuous Signals to Discrete Tokens**

Prof. Liu CCMT2023

Image Captioning Model

There is a teacher teaching his students in the classroom.

有一个老师在教室里教他的学生。

Prof. Wang CCMT2023

# Constructing Multi-modal Datasets

- Motivation



**Traditional Pivot-based**

Pivot Translation — Error Accumulation → Noisy & Ungrounded

**SilkRoad-VL**

Visual Grounding → VLM & Ensemble → High-Quality & Visually Grounded

- Traditional pivot translation accumulates errors and lacks visual grounding, resulting in noisy and less accurate translations.

- We introduce visual information to directly support translation, improving translation quality and image-text consistency.

- ## Experiment

  - ### Dataset statistics

**Statistics of SilkRoad-VL Dataset**

| Language | Total | Short Count | Short Length | Long Count | Long Length |
|---|---|---|---|---|---|
| Uyghur(ug) | 8.7k | 3.7k | 15.7 | 5.0k | 34.7 |
| Kazakh(kk) | 26.1k | 13.0k | 14.8 | 13.2k | 31.6 |
| Kyrgyz(ky) | 23.2k | 12.8k | 15.4 | 10.3k | 33 |
| Tajik(tg) | 2.4k | 1.8k | 17 | 0.6k | 40 |
| Uzbek(uz) | 14.9k | 8.2k | 15.2 | 6.7k | 33 |
| Urdu(ur) | 9.6k | 5.1k | 23.3 | 4.5k | 56.6 |
| Total / Avg. | 84.9k | 44.6k | 16.9 | 40.3k | 38.2 |

- Experiment

  - Main experiment

Comparison with strong baselines on the SilkRoad testset using the average of **COMET-Kiwi, BERTScore, and CLIP.**

| Model | kk | ky | tg | ur | ug | uz |
|-------|----|----|----|----|----|----|
| Qwen3-VL | 52.69 | 49.93 | 46.61 | 55.11 | 47.14 | 54.53 |
| Qwen2.5-VL | 51.12 | 48.46 | 45.44 | 50.65 | 52.57 | 49.32 |
| LLaVA-v1.6 | 41.12 | 42.55 | 41.75 | 41.4 | 41.28 | 44.91 |
| LLaVA-1.5 | 49.93 | 48.18 | 44.82 | 43.49 | 42.38 | 47.06 |
| InternVL3 | 46.21 | 44.59 | 44.94 | 46.6 | 44.62 | 45.85 |
| Ours | **62.47** | **61.24** | **58.25** | **60.12** | **60.87** | **62.33** |

Comparison with strong baselines on the Multi30k testset using the average of **COMET-Kiwi and CLIP.**

| Model | kk | ky | tg | ur | ug | uz |
|-------|----|----|----|----|----|----|
| Qwen3-VL | 41.12 | 34.17 | 26.23 | 30.89 | 42.79 | 41.64 |
| Qwen2.5-VL | 34.07 | 30.15 | 24.78 | 36.43 | 34.77 | 33.55 |
| LLaVA-v1.6 | 22.92 | 23.88 | 24.48 | 23.38 | 23.43 | 28.3 |
| LLaVA-1.5 | 32.87 | 29.25 | 26.57 | 25.67 | 25.98 | 29.13 |
| InternVL3 | 29.16 | 27.23 | 24.36 | 30.81 | 28.49 | 28.38 |
| Ours | **50.36** | **49.35** | **40.2** | **48.41** | **48.98** | **49.06** |

- # Motivation

  - Directly injecting visual information into multimodal translation is often unstable and may introduce noise.

  - In low-resource settings, text-only retrieval is weak at disambiguating short or templated sentences.

  - Therefore, we proposes using Visual Neighbor Memory as auxiliary evidence for disambiguation rather than letting vision dominate generation.

# Vision Neighbor Memory

- ## Experiment

  - ## Main experiment

    - We conduct our main experiments on VG (zh→xx) with four low-resource target languages, where our method improves the average score from 51.02 to 57.81 (+6.79 over Direct). On Multi30K, our method further achieves the best overall average of 70.90, surpassing Direct by +1.34.

| Method | VG (en→xx) | | | | Multi30K(en→xx) | | |
|---|---|---|---|---|---|---|---|
| | bn | hi | ha | ma | de | fr | cs |
| Direct | 57.79 | 61.73 | 32.91 | 51.65 | 71.37 | 74.79 | 62.53 |
| COD | 58.85 | 62.27 | 34.91 | 55.20 | **72.39** | 75.05 | 64.51 |
| MAPS | 57.36 | 60.68 | 33.19 | 51.72 | 71.51 | 74.36 | 62.66 |
| CompTrans | 57.37 | 60.70 | 33.12 | 51.69 | 71.13 | 74.30 | 62.54 |
| Ours | **63.36** | **63.81** | **44.32** | **59.75** | 72.18 | **75.93** | **64.60** |

# Outline

- LLM Everywhere

- Machine Translation with LLMs

- Challenges of LLM-based MT on LRLs

- Strategies for Text-only LLM-based MT

- Findings on Multi-modal LLM-based MT

- Conclusion

# Conclusion

**1.Unified paradigm:** LLMs enable multilingual MT with reduced deployment cost.

**2.Text-only strategies:** Techniques such as CoT, RAG, and word sense disambiguation help compensate for limited parallel data.

**3.Multimodal MT:** Visual memory and cross-modal fusion provide additional support for low-resource translation.

**4.Data bottleneck:** High-quality instruction data remains critical; automated pipelines offer a scalable solution.

**5.Future directions:** Expand instruction corpora for more language pairs, integrate multimodal signals more effectively, and develop better evaluation metrics.

**6.Underlying cause:** The suboptimal performance of both **unimodal** and **multimodal** LLMs in **low-resource settings** ultimately stems from their **pre-trained foundation models**.
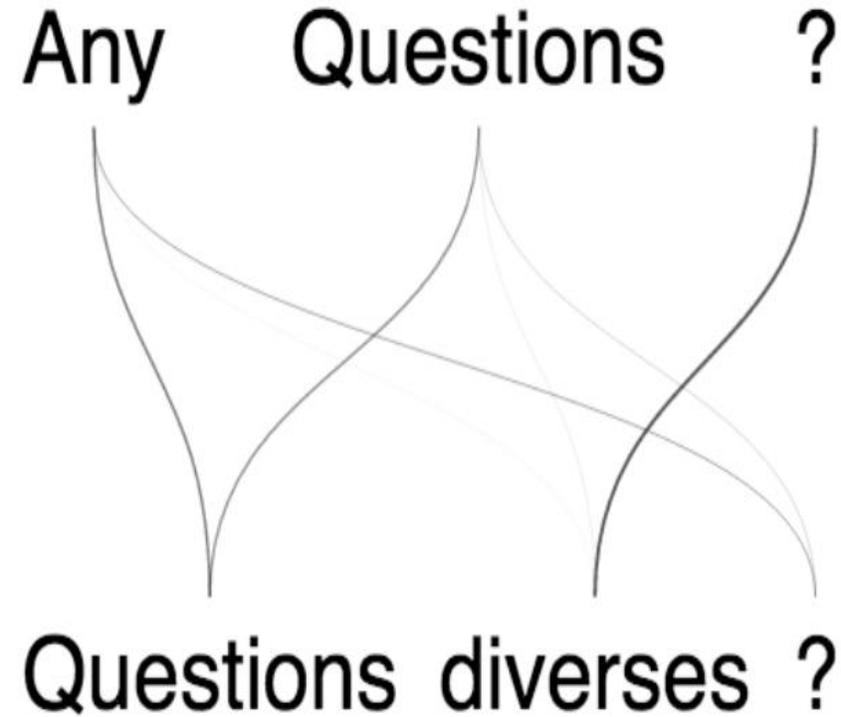
# Contributors

# Thank You!

This inspiration comes from Dzmitry Bahdanau @ ICLR2014

# Contact me~

Email: miradel_51@hotmail.com；miradeljan51@xju.edu.cn

- I'm hiring **self-motivated** master's and PhD candidates to join my team as a research intern (RI).

- I am open to collaborations both **online** and **on-site**.

- Computational resources：
  - H100  80G    50
  - 5090  32G     10
  - 4090  24G     10
  - 3090  16G      8

木拉丁
Haidian, Beijing

Scan QR code to add me