# Using bitext mining to identify translated material:

## practical assessment and new applications

*Zhilu Tu, Minghao Wang,*

*Mark Shuttleworth, and Zhiwen Hua*
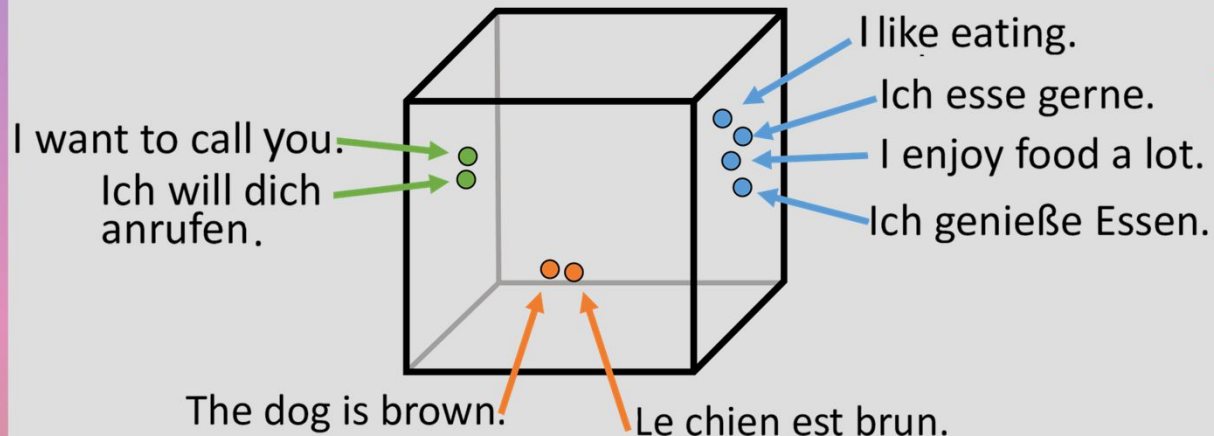
Hong Kong Baptist University

TIIS @HKBU

WikiAligner

# Introduction

## *Bilingual text mining*

SENTENCE LENGTH

LEXICON RECOURSES

I want to call you.
Ich will dich anrufen.

I like eating.
Ich esse gerne.
I enjoy food a lot.
Ich genieße Essen.

The dog is brown.
Le chien est brun.
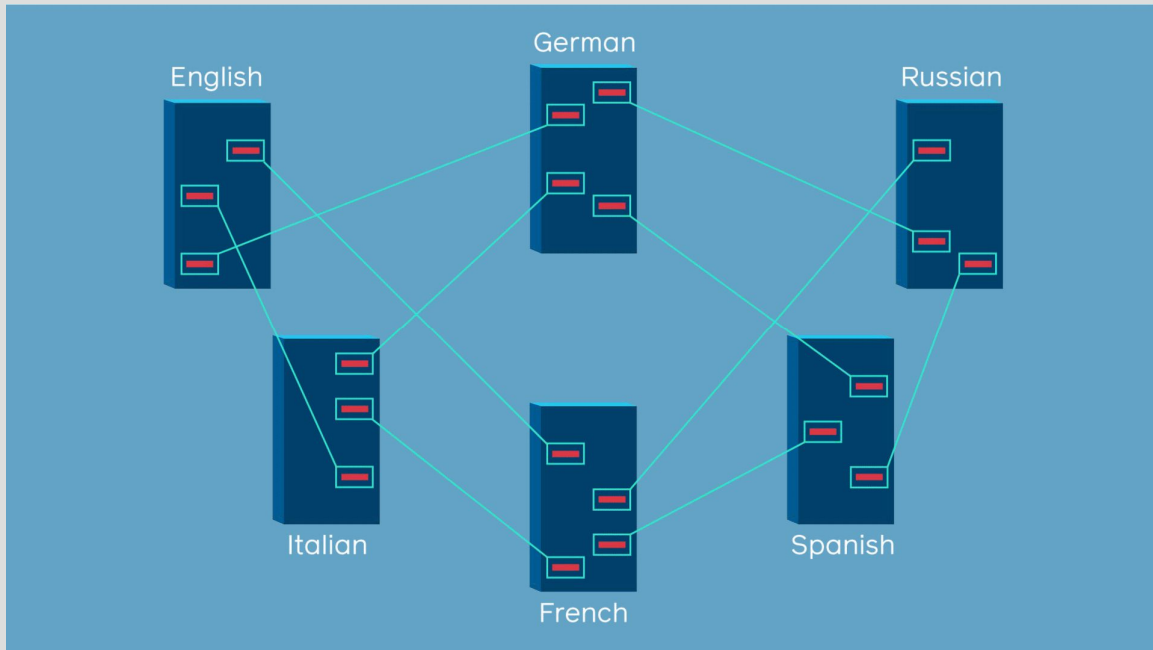
(Schwenk, 2019; Thompson & Zhang, 2019)

- Traditional
  - Sentence length-based algorithm
  - Bilingual cross-reference resources
  - Monotonic - can't locate scattered sentences

- Conventional
  - Sentence-embedding
  - Pre-trained model – LASER; LaBSE
  - Non-monotonic

# Objectives

## *Towards practicality*



CCMatrix (Joulin & Schwenk, 2020)

*Russian State against our country. Or the Russian government lost control of this potentially catastrophically damaging nerve agent and allowed it to get into the hands of others.*

правительства против нашей страны, либо российское правительство потеряло контроль над распространением потенциально катастрофически опасного нервного токсина, и он попал в чужие руки.

- A variety of big data applications have now been created using bitext mining technology.

- The paper investigates the use of bitext mining within a small-scale dataset setting as an example of how they might be exploited in other practical scenarios

# Objectives

## *Example*

• Translation largely remains Wikipedia's "dark matter"; this has led to discrepancies and changes in points of view (Shuttleworth, 2017).

*Russian State against our country. Or the Russian government lost control of this potentially catastrophically damaging nerve agent and allowed it to get into the hands of others.*

правительства против нашей страны, либо российское

правительство потеряло контроль над распространением

потенциально катастрофически опасного нервного токсина, и он

попал в чужие руки.

[... or the Russian government lost control of the distribution of a potentially catastrophically dangerous nerve toxin and it got into the hands of others]

# Objectives

## *Hypothesis*

*Russian State against our country. Or the Russian government lost control of this potentially catastrophically damaging nerve agent and allowed it to get into the hands of others.*

правительства против нашей страны, либо российское правительство потеряло контроль над распространением потенциально катастрофически опасного нервного токсина, и он попал в чужие руки.

[... or the Russian government lost control of the distribution of a potentially catastrophically dangerous nerve toxin and it got into the hands of others]

- Aligning materials manually affected the efficiency of the research and the possibility of enlarging the study scale

- Bitext mining can help reduce the researcher's workload while opening up the possibility of analyzing more data

# Methodology

*Our tool*

- WikiAligner: A Browser/Server Wikipedia bitext mining solution integrating back-end pipeline and front-end re-presentation of the data.

# Methodology

## *Demo - Minghao Wang*

# Methodology

## *WikiAligner workflow*



1. Front-end UI input
   - Keywords for the articles
   - SL title & TL title

2. Back-end data pipeline
   - Segmented sentences
   - LaBSE(Feng et al., 2022) embedding
   - Similarity search

3. Front-end display
   - Highlight
   - Tags & Threshold
   - Right-click translation

# Discussion

*Assessments*

WikiAligner

## Performance

- All the alignments that have been extensively analyzed by Shuttleworth (2018) are also located by using our tool.

## Accessibility

- As with other web-based services, potential users do not need to attend the coding of the bitext tool before using

## Manifestations

- Researchers will be able to scrutinize every translation pair and their distribution by using the highlight feature

# Discussion

*New applications*

# Conclusion

- Standing on the shoulders of giants from NLP, we present a method to utilize bitext mining in the context of a small-scale dataset.

- WikiAligner can be an entry point for applying bitext mining to more real-world scenarios with new features.

# References

- Feng, F., Yang, Y., Cer, D., Arivazhagan, N., & Wang, W. (2022). Language-agnostic BERT Sentence Embedding. *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 878–891. https://doi.org/10.18653/v1/2022.acl-long.62

- Joulin, A., & Schwenk, H. (2020, February). *CCMatrix: A billion-scale bitext dataset for training translation models*. https://ai.facebook.com/blog/ccmatrix-a-billion-scale-bitext-data-set-for-training-translation-models/

- Schwenk, H. (2019, January 22). *LASER natural language processing toolkit - Engineering at Meta*. https://engineering.fb.com/2019/01/22/ai-research/laser-multilingual-sentence-embeddings/

- Shuttleworth, M. (2017). Locating foci of translation on Wikipedia. *Translation Spaces*, *6*(2), 310–332. https://doi.org/10.1075/TS.6.2.07SHU

- Shuttleworth, M. (2018). Translation and the Production of Knowledge in "Wikipedia": Chronicling the Assassination of Boris Nemtsov. *Alif: Journal of Comparative Poetics*, *38*, 231–263. https://www.jstor.org/stable/26496376

- Thompson, B., & Zhang, S. (2019). *thompsonb/vecalign: Improved Sentence Alignment in Linear Time and Space*. https://github.com/thompsonb/vecalign

# Q&A

- wikialigner@outlook.com