

# Language Scaling: Applications, Challenges and Approach



Lecture Tutorial for KDD 2021

August 14-18, 2021, Singapore

<https://languagescalingkdd.github.io/>

# Outline

- Introduction [Dixin Jiang]
  - Motivating examples in Microsoft products
  - Problem description
  - Categorization of applications
  - Challenges
- Methodology [Dixin Jiang]
  - Model Transfer
  - Data Transfer
- Applications\*
  - Dependency Parsing [Xiubo Geng]
  - Machine Reading Comprehension [Ming Gong]
  - Grammar Error Correction [Linjun Shou]
- Summary & Future directions [Jian Pei]



Dixin Jiang

Software Technology Center at Asia (STCA) of Microsoft



Linjun Shou

Software Technology Center at Asia (STCA) of Microsoft



Xiubo Geng



Ming Gong



Jian Pei

Simon Fraser University

\*For more applications, please refer to our tutorial at  
The Web Conference 2021

# Why language matters (society)

[Why Languages Matter | SIL International](#)

## Eradicate extreme poverty and hunger

Higher literacy rates often result in higher per capita incomes.



Photo: Rodney Ballard

## Achieve universal primary education

Primary education programs that begin in the mother tongue help students gain literacy and numeracy skills more quickly.



Photo: Rodney Ballard

## Promote gender equality and empower women

Nearly two-thirds of the world's 875 million illiterate people are women.



Photo: Marc Ewell

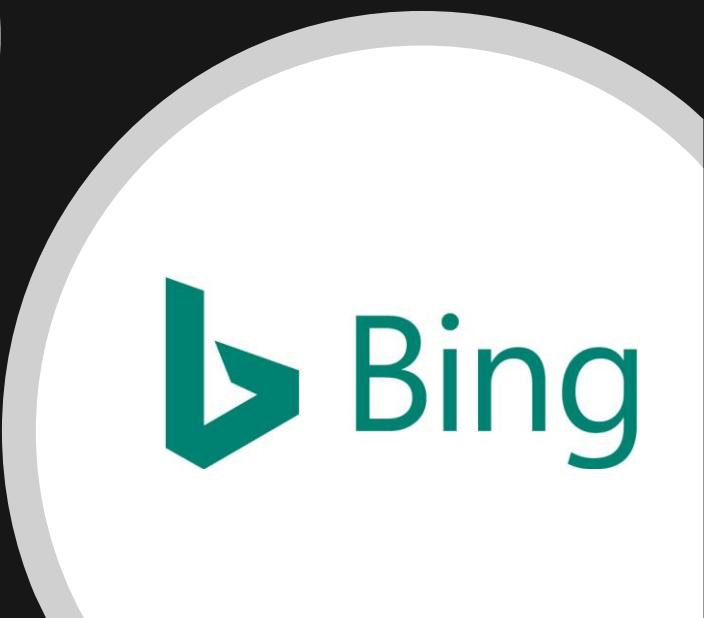
## Reduce child mortality

The mortality rate for children under five years of age is reduced when vital health information about disease prevention and treatment is available in local languages.



# Why Language Matters (Industry)

- Microsoft's mission statement: to empower *every person and every organization on the planet* to achieve more.
- NLP is widely adopted in Microsoft products and services

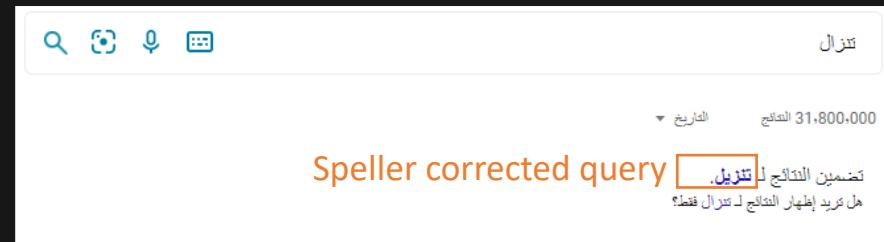


# Speller | Web Relevance | Suggested Replies | Natural Language Understanding

ar-EG

{تنزال}

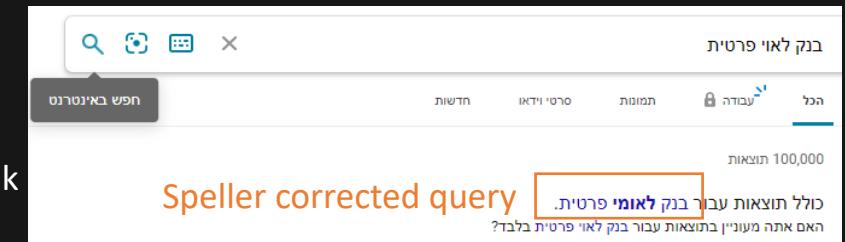
English Translation:  
download



he-IL

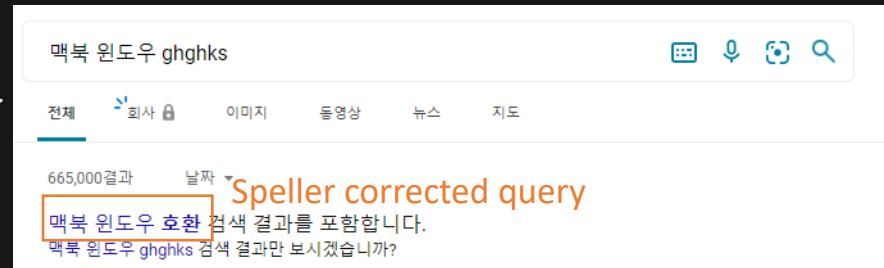
בנק לאומי

English Translation:  
Private National Bank



ko-KR

{맥북 윈도우 ghghks}  
English Translation:  
MacBook Windows  
Compatible



ja-JP

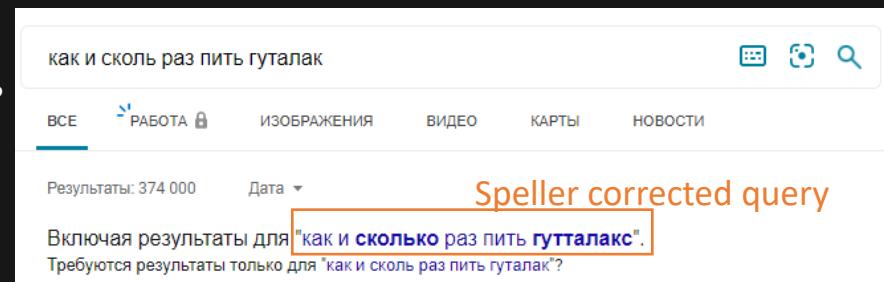
{田中みどらら}

English Translation:  
Ms. Tanaka



ru-RU

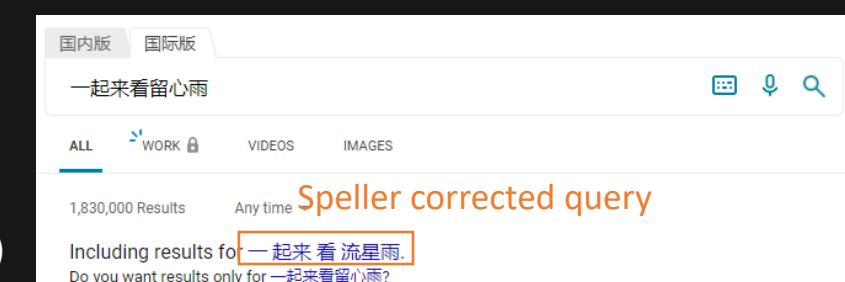
{как и сколь раз пить  
гуталак}  
English Translation:  
how many times to  
drink guttalax



zh-CN

{一起来看流星雨}

English Translation:  
Watch meteor  
together (TV series)



# Scalability vs. Specificity

Scalability: scaling out to many languages

Specificity: language specific features

## Japanese Speller

- Four scripts (Kanji, Hiragana, Katakana and Romaji)
- Two IME (Romaji on PC, and Kana on mobile)



Example: みやさき vs. 宮崎

	typo	correction
Kana	みや <span style="color:red">さ</span> き	みや <span style="color:green">ざ</span> き
Romaji	miyasaki	miyazaki

1 keystroke difference

- Need language-specific training data
- Trade off between scalability and specificity to meet the highest ROI (Return Over Invest)

Speller | **Web Relevance** | Suggested Replies | Natural Language Understanding

# ▼ English

the crown has had its scandals

ALL WORK IMAGES VIDEOS MAPS NEWS SHOPPING

15,100,000 Results Any time ▾

**'The Crown' Has Had Its Scandals, but There's Nothing Like ...**  
<https://www.nytimes.com/2020/11/12/arts/television/the-crown-princess-diana.html>  
Nov 12, 2020 · 'The Crown' Has Had Its Scandals, but There's Nothing Like Diana The face that launched a thousand tabloid stories – and books and documentaries and ...

**'The Crown' Has Had Its Scandals, however There's Nothing ...**  
[https://lightlynnews.com/2020/11/12/arts/the-crown-has... ▾](https://lightlynnews.com/2020/11/12/arts/the-crown-has...)  
Nov 12, 2020 · 'The Crown' Has Had Its Scandals, however There's Nothing Like Diana by Lightlynnews.com · On November 12, 2020 · In Arts / Television When we first glimpse her, minutes into...

**The Crown has had its scandals, but there's nothing like ...**  
[https://celebrityml.com/entertainment/the-crown-has-had... ▾](https://celebrityml.com/entertainment/the-crown-has-had...)  
Nov 21, 2020 · The Crown has had its scandals, but there's nothing like Diana 11/21/2020 The face that launched a thousand tabloid stories – and books and documentaries and Instagram feeds – takes cent...

**'The Crown' Has Had Its Scandals, however There's Nothing ...**  
[https://www.thinkipos.com/the-crown-has-had-its... ▾](https://www.thinkipos.com/the-crown-has-had-its...)  
'The Crown' Has Had Its Scandals, however There's Nothing Like Diana November 12, 2020 November 12, 2020 - by Kumar - Leave a Comment Once we first glimpse her, minutes into Season four of "The Crown,"...

**'The Crown' Has Had Its Scandals, however There's Nothing ...**  
[https://www.thinkipos.com/the-crown-has-had-its-scandals... ▾](https://www.thinkipos.com/the-crown-has-had-its-scandals...)  
'The Crown' Has Had Its Scandals, however There's Nothing Like Diana November 15, 2020 November 15, 2020 - by Kumar - Leave a Comment Once we first glimpse her, minutes into Season four of "The Crown,"...

# Spanish

espejos tocador con luces comprar	grid icon	microphone icon	camera icon			
ALL	WORK	IMAGES	VIDEOS			
MAPS	NEWS	SHOPPING				
1,860,000 Results	Any time ▾					
<a href="#">Amazon.es: tocador maquillaje con luz</a>						
<a href="https://www.amazon.es/tocador-maquillaje-luz/s?k=tocador+maquillaje+con+luz">https://www.amazon.es/tocador-maquillaje-luz/s?k=tocador+maquillaje+con+luz</a> ▾						
BEAUTME Espejo de tocador con Luces, tocador de Maquillaje Iluminado o Espejos de Belleza montados en la Pared con atenuador, Espejo cosmético Hollywood con 15 Bombillas LED (Plateado 99,99 € 99,99)						
<a href="#">Espejo Tocador Con Luces en Mercado Libre México</a>						
<a href="https://listado.mercadolibre.com.mx/espejo-tocador-con-luces">https://listado.mercadolibre.com.mx/espejo-tocador-con-luces</a> ▾						
Encuentra Espejo Tocador Con Luces en Mercado Libre México. Descubre la mejor forma de comprar online.						
<a href="#">Mejor espejo tocador con luces: ofertas y reseñas 2020</a>						
<a href="https://quecomprar.org/espejo-tocador-con-luzes">https://quecomprar.org/espejo-tocador-con-luzes</a> ▾						
Que Comprar ha desarrollado una rápida selección de los mejores modelos de espejo tocador con luces disponibles en línea a los mejores precios. Índice de contenidos  Top 5 mejores espejos tocadores con luces: bestseller en Amazon						

processo histórico car

TUDO TRABALHO

1.170.000 Resultados

## História do Canadá

<https://www.infoescol>

História do Canadá. Os pri

(algonquinos, esquimós, ir

Ásia. Teriam atravessado de

Opanadé - História

## Canada - History

<https://brasilescola.uol.com.br>

Canada Historia da Americ

Autor: Rainer Goncalves So

História do Canadá

<https://pt.wikipedia.org>

<https://pt.wikipedia.org>

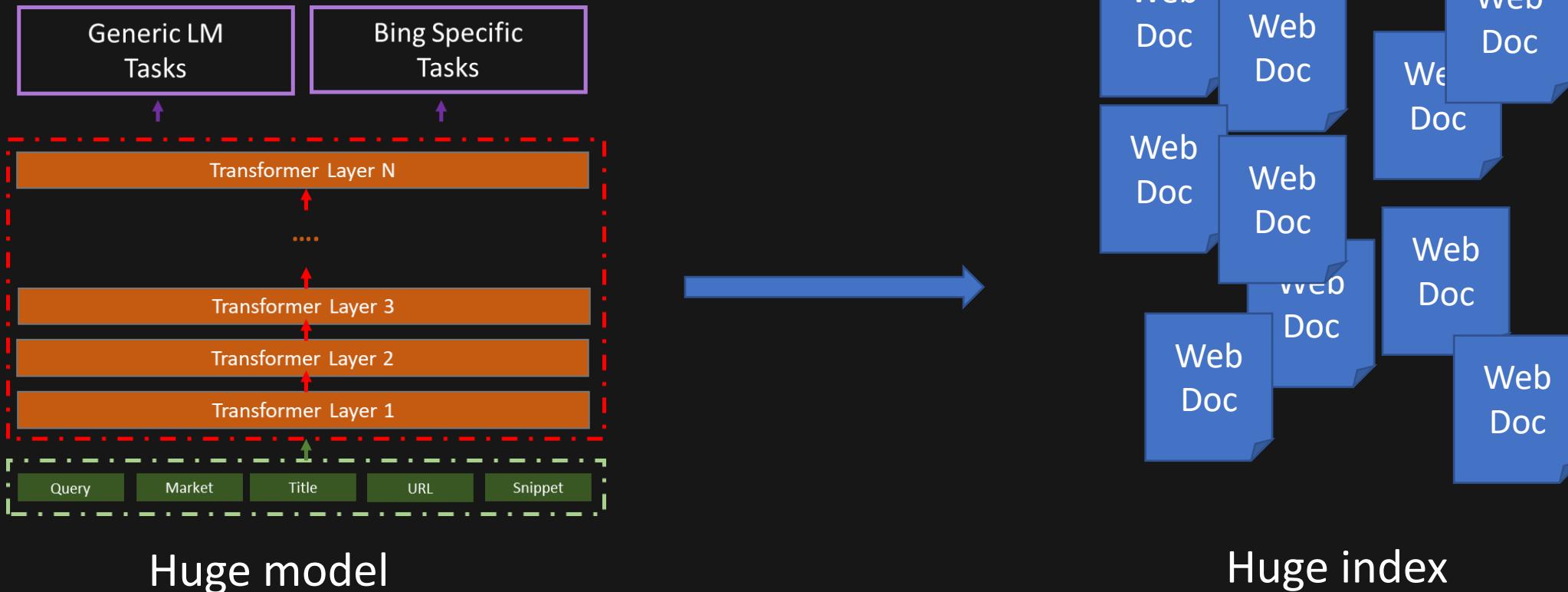
20/02/2005 - A História do

ameríndios entre os quais

entanto os primeiros habita

For more information about the study, please contact Dr. Michael J. Hwang at (319) 356-4550 or via email at [mhwang@uiowa.edu](mailto:mhwang@uiowa.edu).

# Efficiency



## English

Mary Smith <wrbutpt1013@gma  
il.com>     

Tue 2020-11-17 17:06  
To: Melody Wang

Thanks for the heads up. Did your friends wonder why you had invited the paper boy to your party?

 Are the suggestions above helpful?

[Reply](#) | [Forward](#)

## Spanish

Mary Smith <wrbutpt1013@gma  
il.com>     

Tue 2020-11-17 21:34  
To: Melody Wang

Necesita los originales. Dámelos hoy en la oficina.  
Gustavo

 Are the suggestions above helpful?

[Reply](#) | [Forward](#)

# Privacy



Data Center 1



Data Center 2



Data Center K

Universal model to support  
all languages

According to compliance restrictions, data cannot be moved out of the data center  
How to use the data in all data centers to train a universal model?

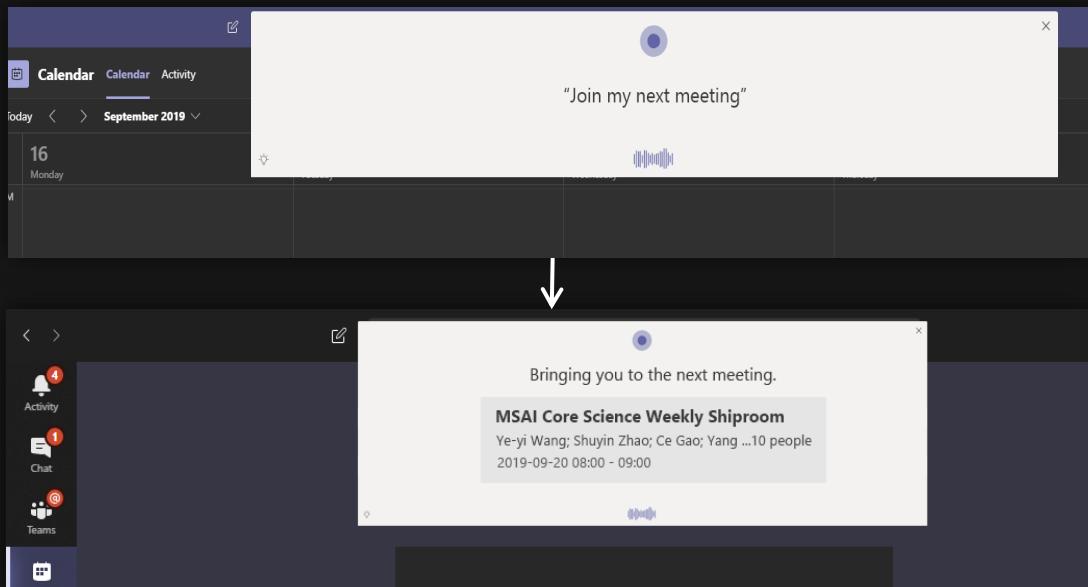
**Teams Voice Skill:**
*Join my next meeting*

Domain: Calendar

Intent: Connect\_to\_meeting

Slot: &lt;start\_time&gt;next&lt;/start\_time&gt;

## Desktop Experience



Start the presentation

Quando è il mio incontro con Tom

Message à Carlos que je serai en retard

打开未读消息

Agrega Alice a este chat

次の会議に参加する

Compartilhe o arquivo com Bob

Füge Alice zum Cortana - Kanal hinzu

## Mobile Experience

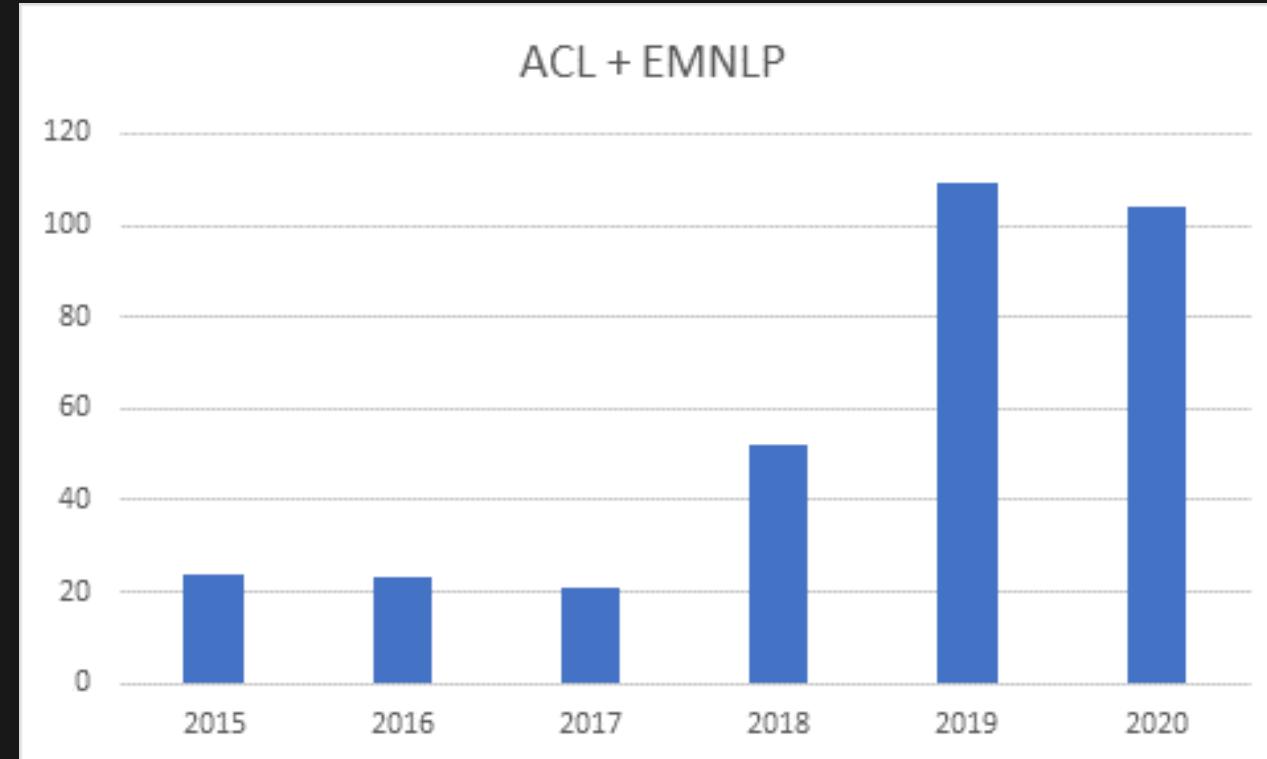


# M Domains × N Languages

- Need large amounts of training data
- Need to handle many models
  - Domain classification model, intent detection model, slot tagging model
  - Training, deployment, maintenance

# Why Language Matters (Research)

Searching for  
“multilingual”,  
“cross-lingual” and  
“bilingual” in the  
ACL anthology  
(ACL+EMNLP)



# Problem Description for Language Scaling



Given an NLP application (such as spelling correction, Web search, suggested reply or NLU)

1. Usually have relatively rich English training data, and well-trained English model
2. Scale out the model to **100+ languages effectively and efficiently**

Source: Both Q and D are in English

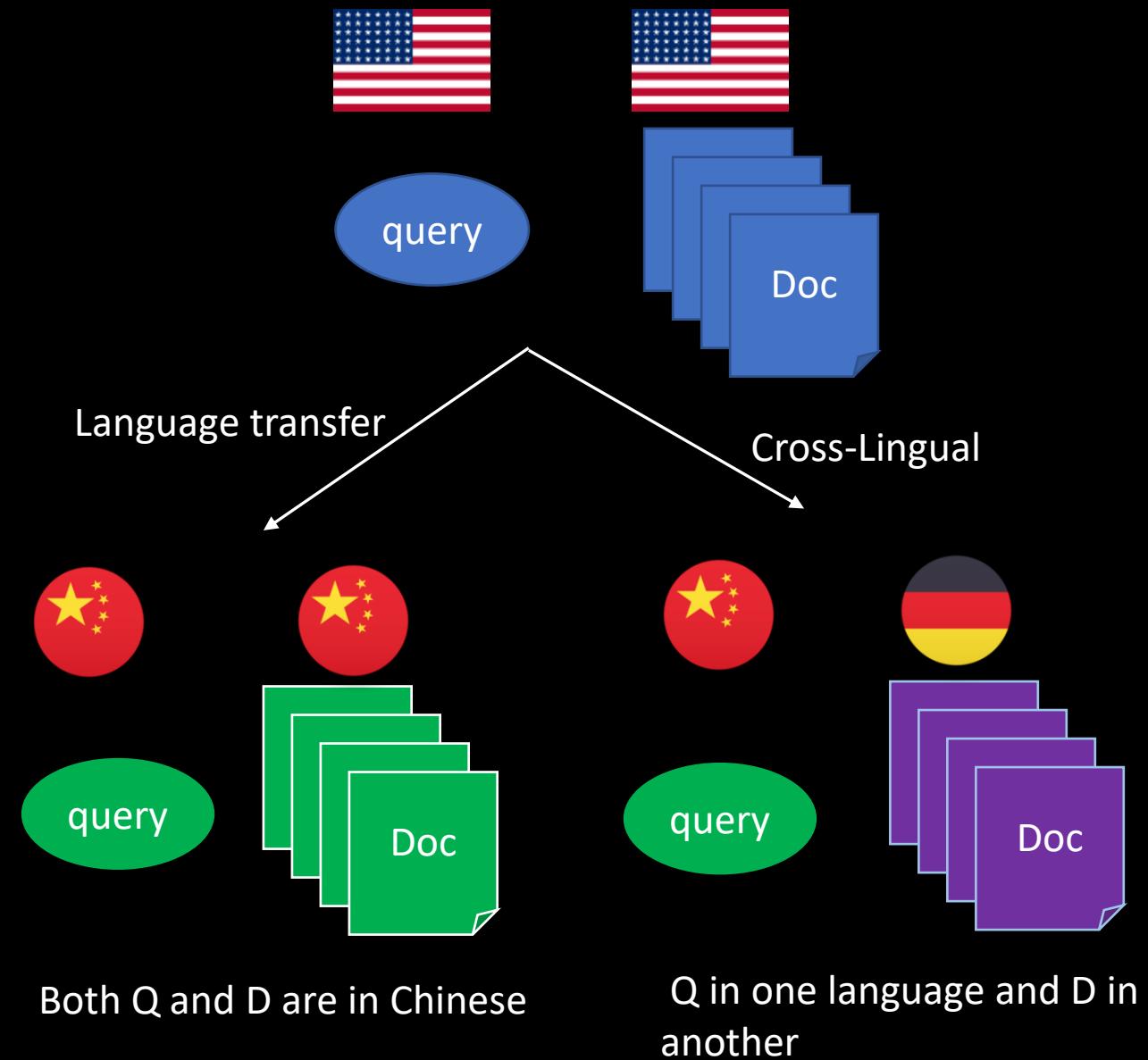
# A Related Problem

- **Language transfer**

- Given English training data, transfer to other languages

- **Cross-lingual**

- e.g., in Web search, query is in one language and document is in another language
- Special case for languages transfer
  - Similar multi-lingual representation and cross-lingual alignment



# NLP Tasks

Type	Category	Sub Category	Example
NLU	Text Classification	Single text 	Domain identification, Intent detection, Sentiment classification
		Text pair 	Information retrieval, Natural language inference
	Sequence Labeling	Single text 	Named entity recognition, Slot tagging
		Text pair 	Machine reading comprehension
	Structure Prediction	Single text 	Dependency parsing, constituency parsing, semantic role labeling
	Text Generation	Token level 	Spelling correction, Sentence auto completion
		Sentence level	Machine translation, Conversation, Question generation



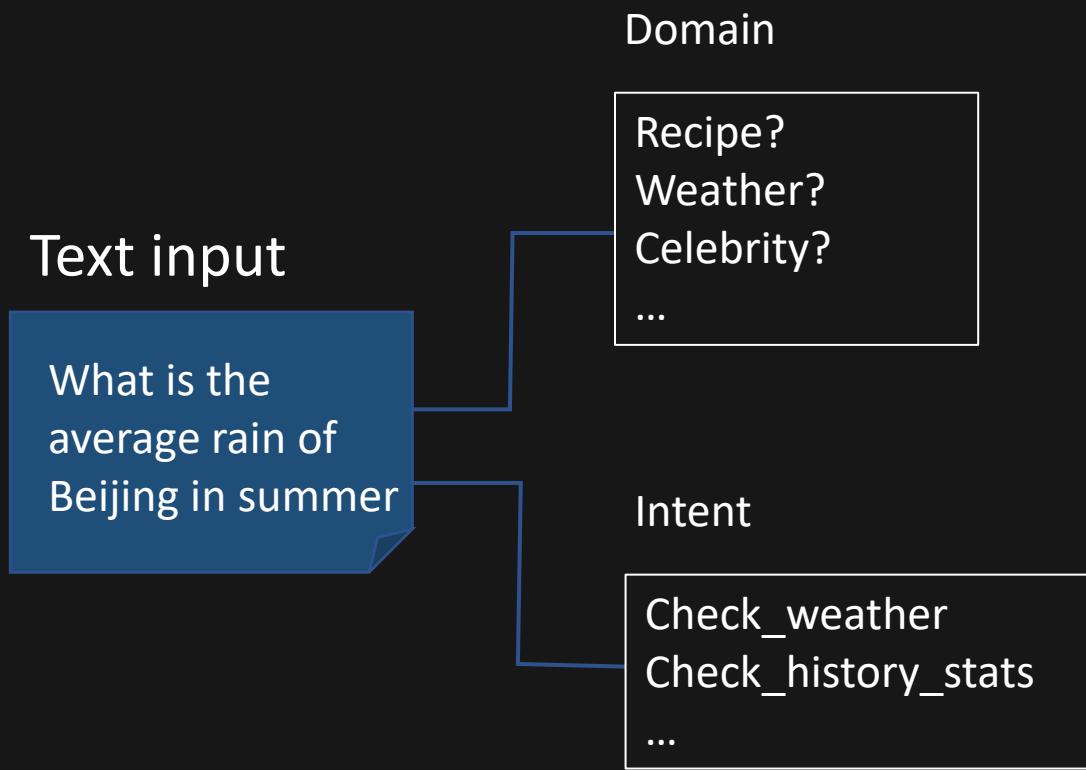
Will be presented in this tutorial



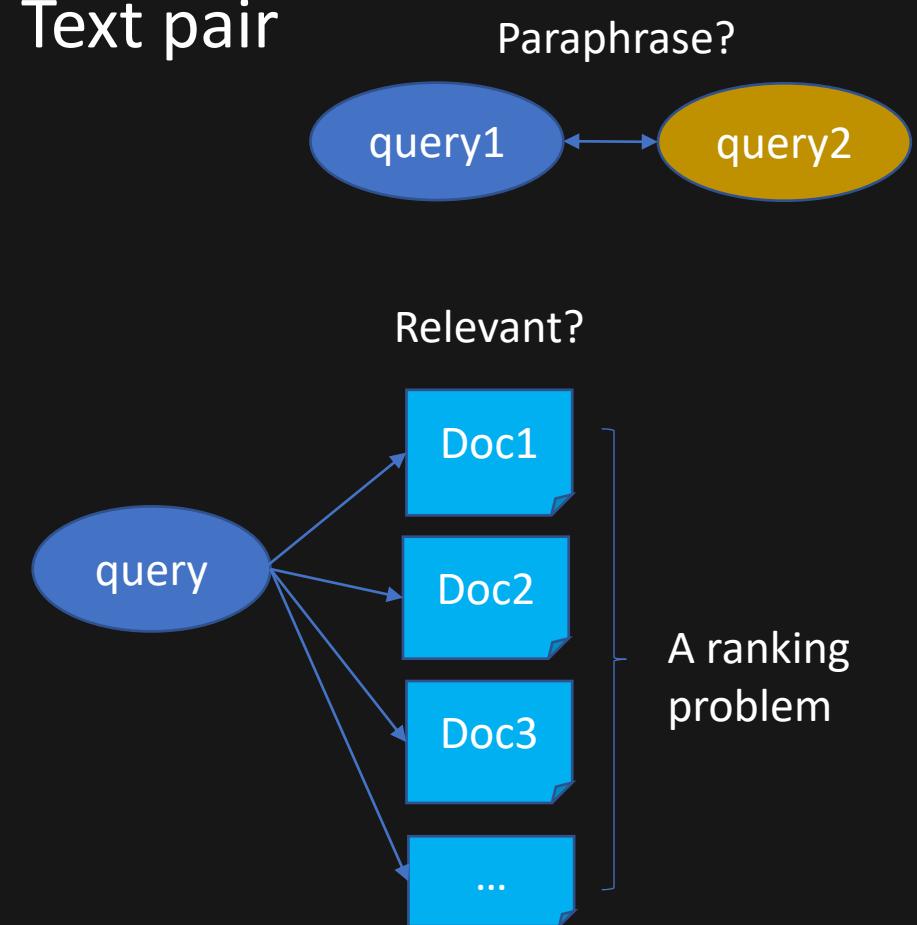
Has been presented in a related tutorial in The Web Conference 2021

# Text Classification

- Single text



- Text pair

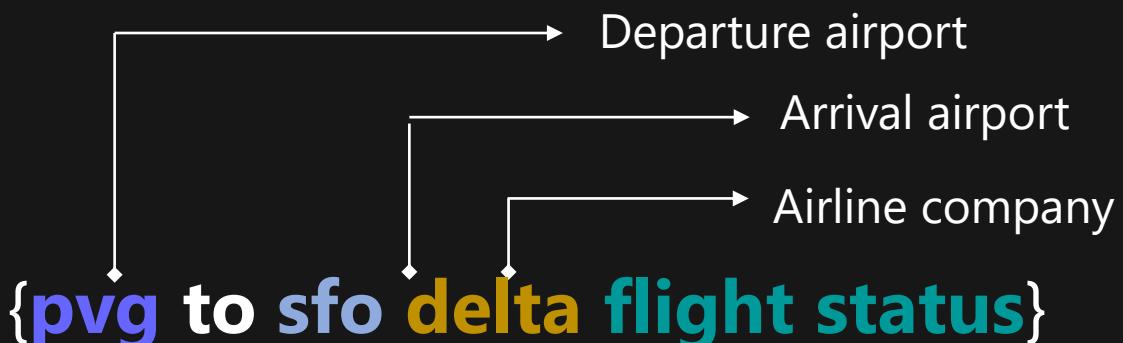


# Sequence Labeling

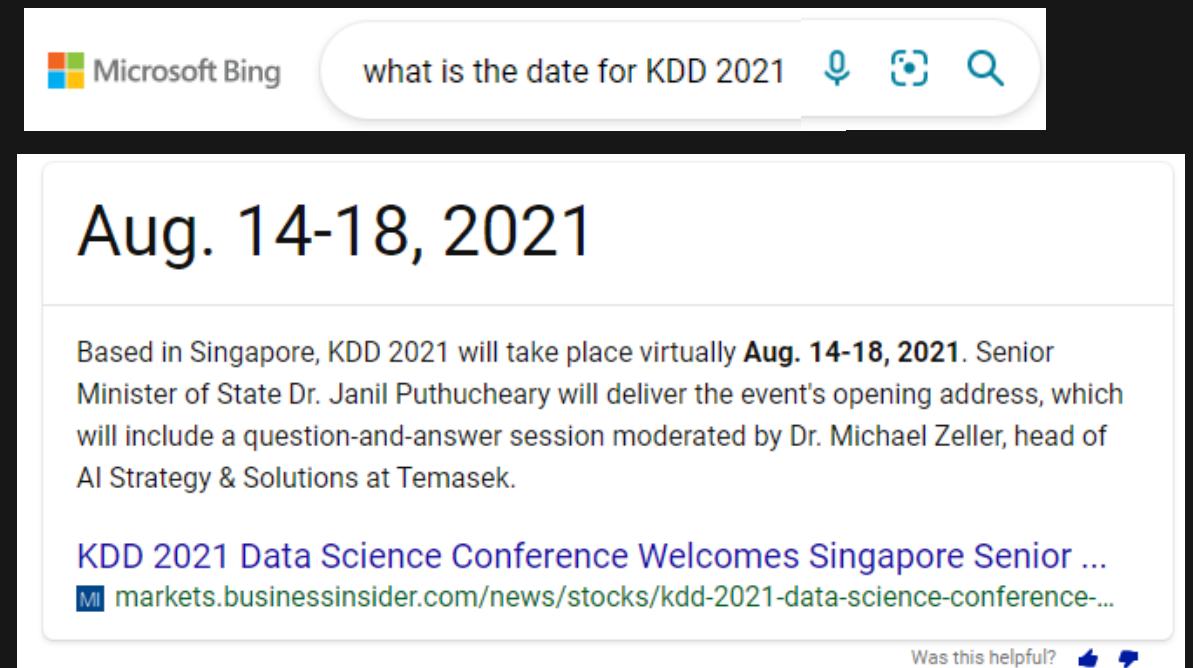
- Single text

Luke Rawlence PERSON joined Aiimi ORG as a data scientist in Milton Keynes PLACE, after finishing his computer science degree at the University of Lincoln. ORG

Example from <https://www.aiimi.com/>



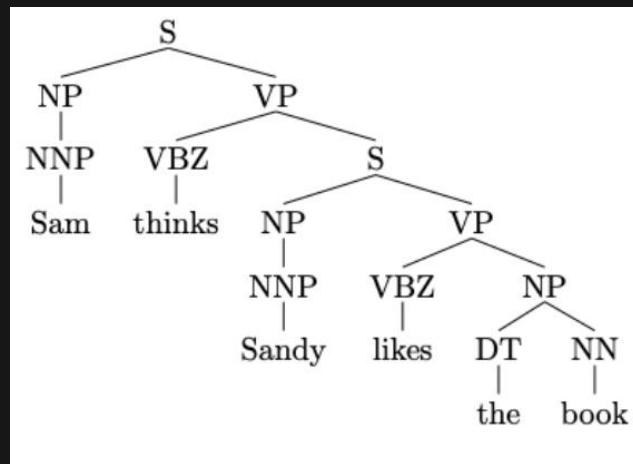
- Text pair



# Structure Prediction

- Given a sentence, parse the syntactic structure of the sentence
  - e.g., constituency parsing, dependency parsing

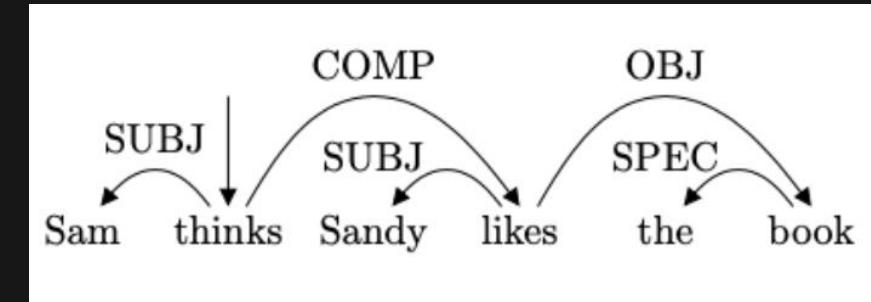
Constituency parsing: how words group into nested constituents



The structure of a sentence is represented by a tree, which is based on the phrase structure grammar.

Examples are from: <https://www.cs.princeton.edu/courses/archive/fall19/cos484/lectures/lec10.pdf>

Dependency parsing: how words depend on each other



Each word has a directed edge towards its head, and there is a label to describe the dependency type.

## Sentence level

# Text Generation

Token level

### Spelling Correction

Britnay spears vidios

Britney speaks videos

### Sentence Auto Completion

The web conference

The web conference is a yearly international

The Web Conference is a yearly international academic conference on the topic of the future development

### Machine Translation

The Web Conference is a yearly international academic conference on the topic of the future direction of the World Wide Web.

La Conferencia Web es una conferencia académica internacional anual sobre el tema de la futura dirección de la World Wide Web.

### Conversation

A: I've been hearing some strange noises around the house at night.

B: oh no! That's scary! What do you think it is?

A: I don't know, that's what's making me anxious.

Response: ???

\*EMPATHETICDIALOGUES

# Challenges for Language Scaling

## Lack of Training Data

- Examples of labeling items
  - Passage-QA: **millions** of QA pairs labeled for English
  - Web relevance: **millions** of query-document pairs labeled for English
- Unrealistic to label so much data for each language

## Too Many Models

- Suppose we have M applications and N languages, we need  $O(M*N)$  models
- In reality,  $M=O(100)$ ,  $N=O(100)$
- The cost for model building, serving, and maintenance is huge

## Model Size and Latency

- Recent deep learning and pre-trained models are getting larger and larger
  - TULRv2: 270M parameters
  - GPT-3: 175B parameters
- Unrealistic to serve such models for online service

- Language-specific features, data privacy issues, cold start problem are all related to this challenge
- Major focus on this tutorial

- Multi-task learning
- Adapters

- Hardware acceleration
- Computation library optimization
- Model compression

# Outline

- Introduction [Dixin Jiang]
  - Motivating examples in Microsoft products
  - Problem description
  - Categorization of applications
  - Challenges
- Methodology [Dixin Jiang]
  - Model Transfer
  - Data Transfer
- Applications\*
  - Dependency Parsing [Xiubo Geng]
  - Machine Reading Comprehension [Ming Gong]
  - Grammar Error Correction [Linjun Shou]
- Summary & Future directions [Jian Pei]



Dixin Jiang

Software Technology Center at Asia (STCA) of Microsoft



Linjun Shou

Software Technology Center at Asia (STCA) of Microsoft



**Xiubo Geng**



Ming Gong



Jian Pei

Simon Fraser University

\*For more applications, please refer to our tutorial at  
The Web Conference 2021

# Overview of Approaches

## Tasks

NLU

Classification

Single text

Text Pair

Sequence labeling

Single text

Text Pair

Structure prediction

## NLG

Token-level

Sentence-level

## Approaches

Model transfer

Word-based

Contextual-based

Data Transfer

Translation

Semi-supervised

Weakly-supervised

Data generation

To align the representations of different languages

## Data

Mono-lingual corpus

Unlabeled data

Tree bank

Translated data

Parallel data

Bilingual dict

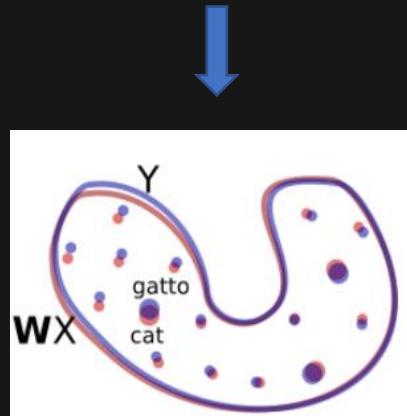
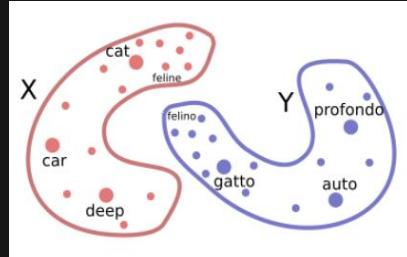
Which method to choose?

- In general, there is no single “winning” method
- Similar tasks, similar methods
- Method selection is often constrained by data availability

# Approaches: Model Transfer and Data Transfer

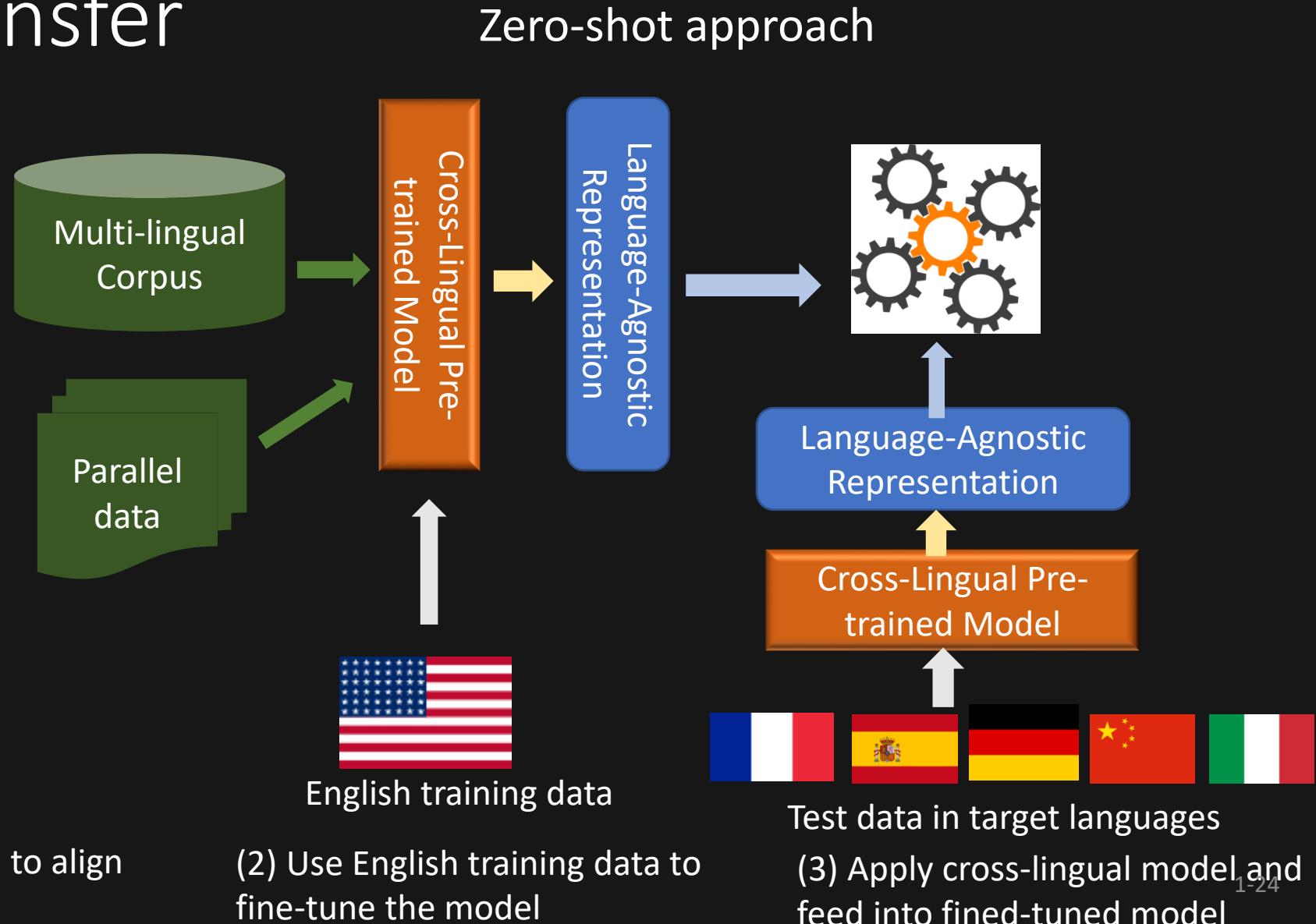
## (1) Model Transfer

Cross-lingual model



From MUSE

(1) Pre-train a cross-lingual model to align different languages.



(2) Use English training data to fine-tune the model

(3) Apply cross-lingual model and feed into fined-tuned model

# Cross-Lingual Models

- Goal: **represent** different languages in a shared vector space, such that the texts with similar meaning are **aligned** close to each other, no matter in which languages they are expressed
- Cross-lingual ***word*** embedding
- Cross-lingual ***contextual*** embedding

# Cross-Lingual Word Embedding

- Mapping-based methods

$$J = \underbrace{\mathcal{L}(X^s) + \mathcal{L}(X^t)}_1 + \underbrace{\Omega(\underline{X^s}, \underline{X^t}, W)}_2$$

$X^t \approx WX^s$



$\underline{X^s}, \underline{X^t}$ : word embedding matrix for source and target languages, respectively

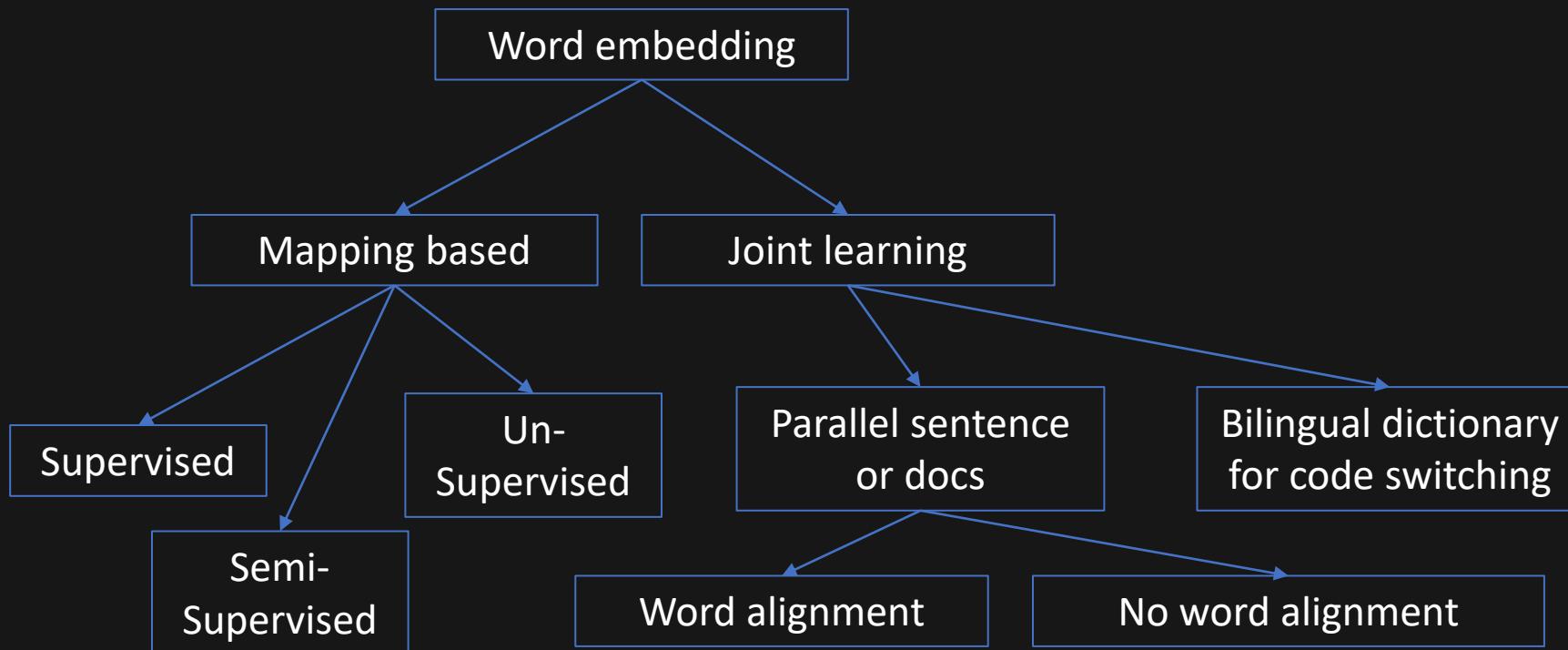
$W$ : mapping matrix from  $X^s$  to  $X^t$

- Joint-learning methods

$$J = \mathcal{L}(X^s) + \mathcal{L}(X^t) + \Omega(X^s, X^t)$$

$\mathcal{L}, \Omega, J$ : Loss functions for individual language embedding, cross-lingual mapping, and total loss, respectively

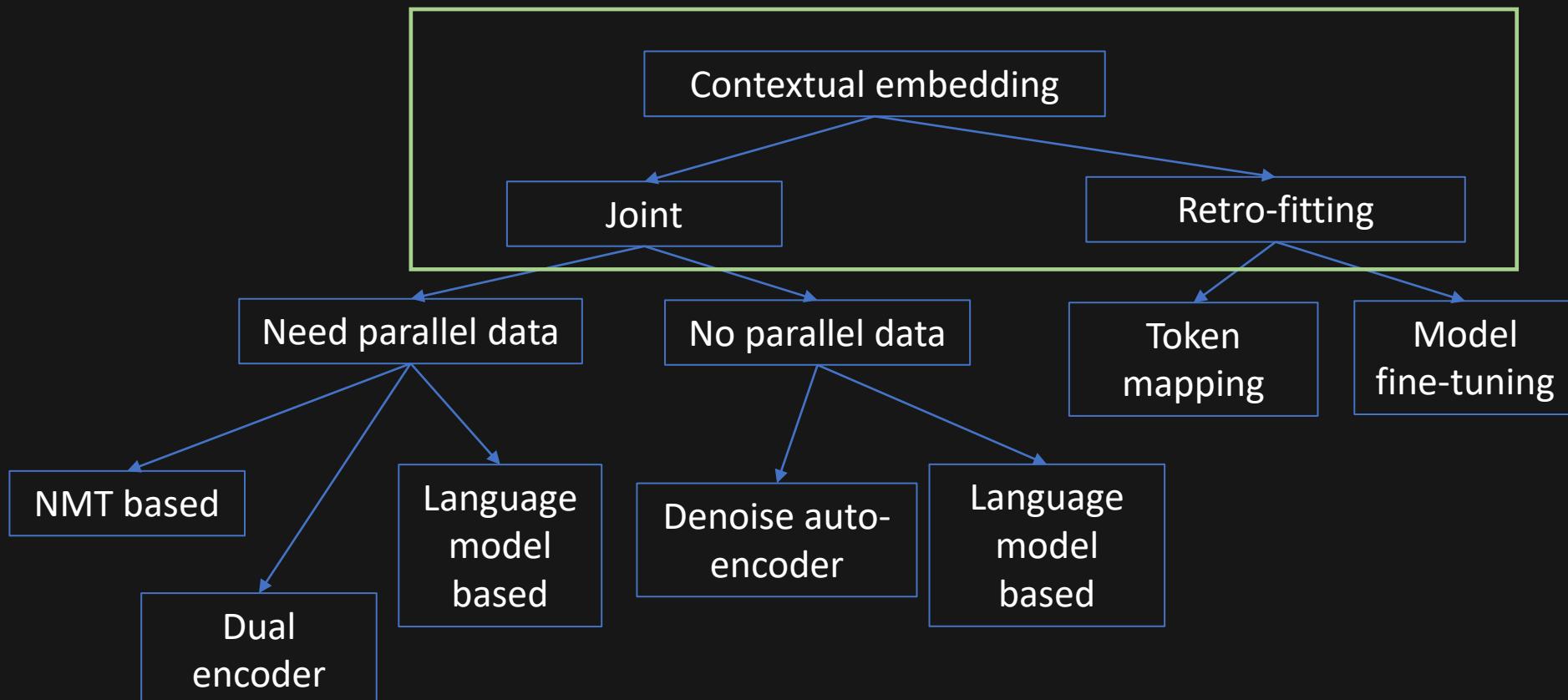
# Cross-Lingual Word Embedding



Sebastian Ruder, Ivan Vulić, Anders SØgaard

- A survey of cross-lingual word embedding models,. Journal of Artificial Intelligence Research, May 2019.
- 2019 ACL Tutorial [Unsupervised Cross-lingual Representation Learning \(ruder.io\)](https://ruder.io/unsupervised-cross-lingual-representation-learning/)

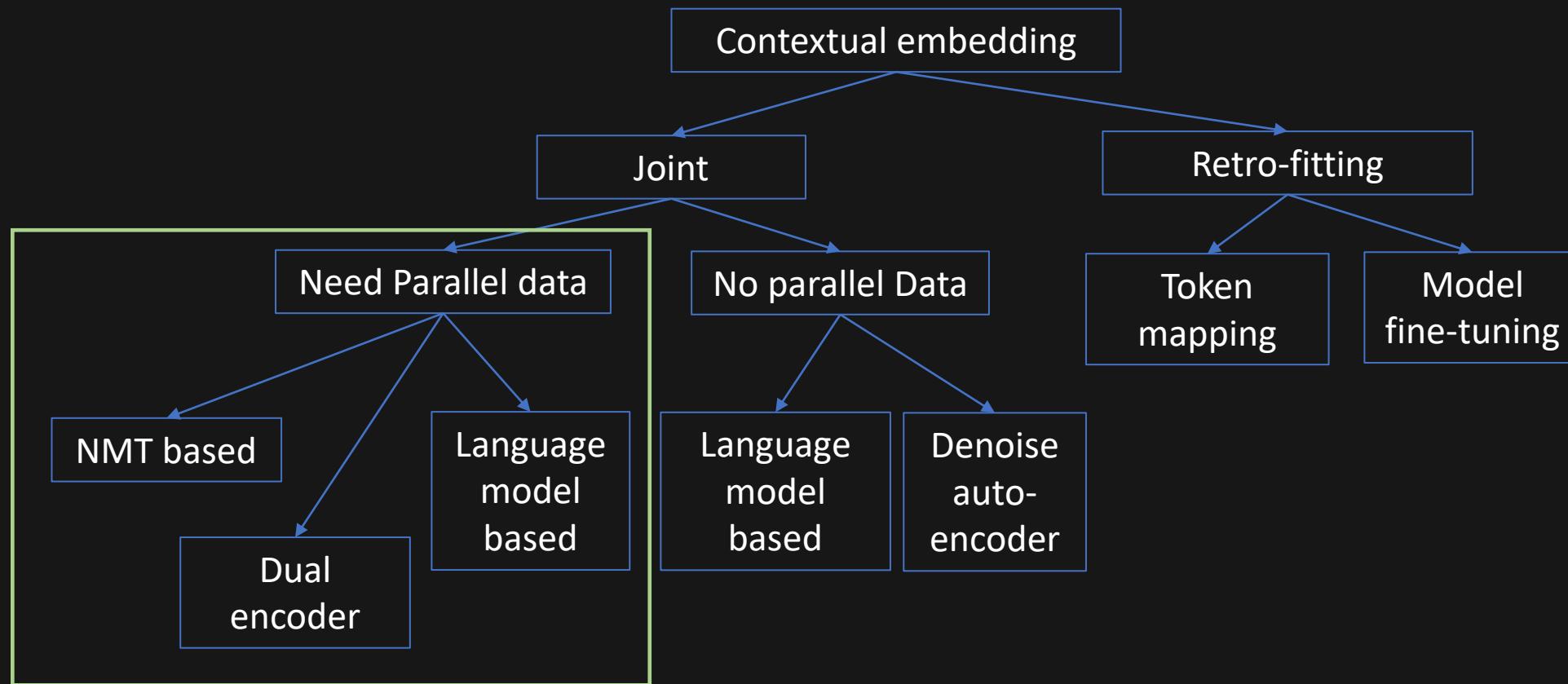
# Cross-Lingual Contextual Embedding



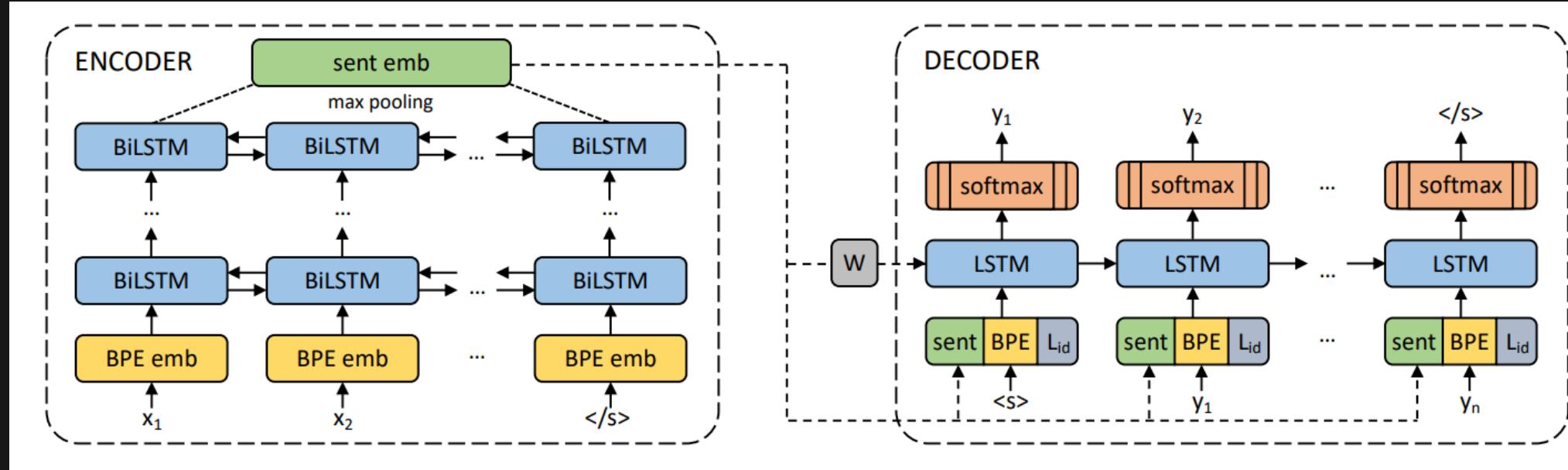
Joint or Retro-fitting?

- *Joint*: the alignment between languages happens simultaneously with the learning of contextual embedding
- *Retro-fitting*: the alignment happens after the contextual embedding for individual languages

# Joint Learning with Parallel Data



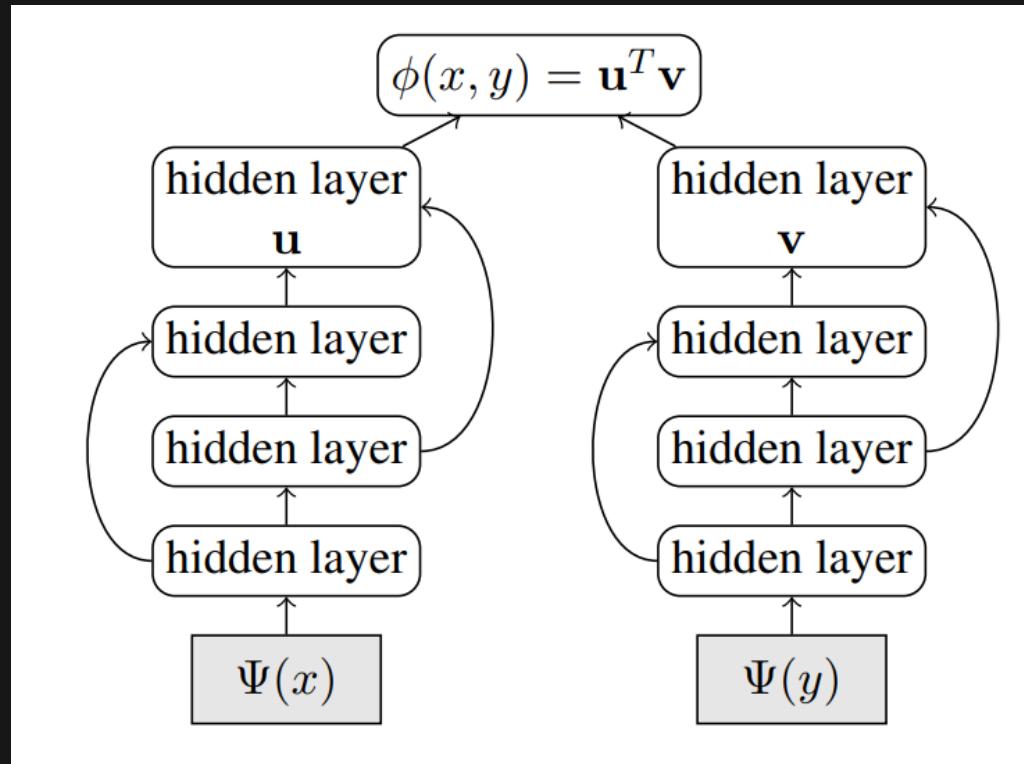
# Neural Machine Translation (NMT) based Method



Mikel Artetxe, Holger Schwenk: Massively Multilingual Sentence Embeddings for Zero-Shot Cross-Lingual Transfer and Beyond. Trans. Assoc. Comput. Linguistics 7: 597-610 (2019)

- Encoder-decoder architecture for machine translation
- M-to-N translation: no indicator to distinguish the input languages => implicitly enforces different languages to have a uniform encoding

# Dual Encoder



- $x$  is in one language, and  $y$  is in another language
- Using parallel data to align the encoding  $u$  and  $v$
- Dual encoder is explicit alignment, while NMT based approach is implicit alignment

Mandy Guo, Qinlan Shen, Yinfei Yang, et al.

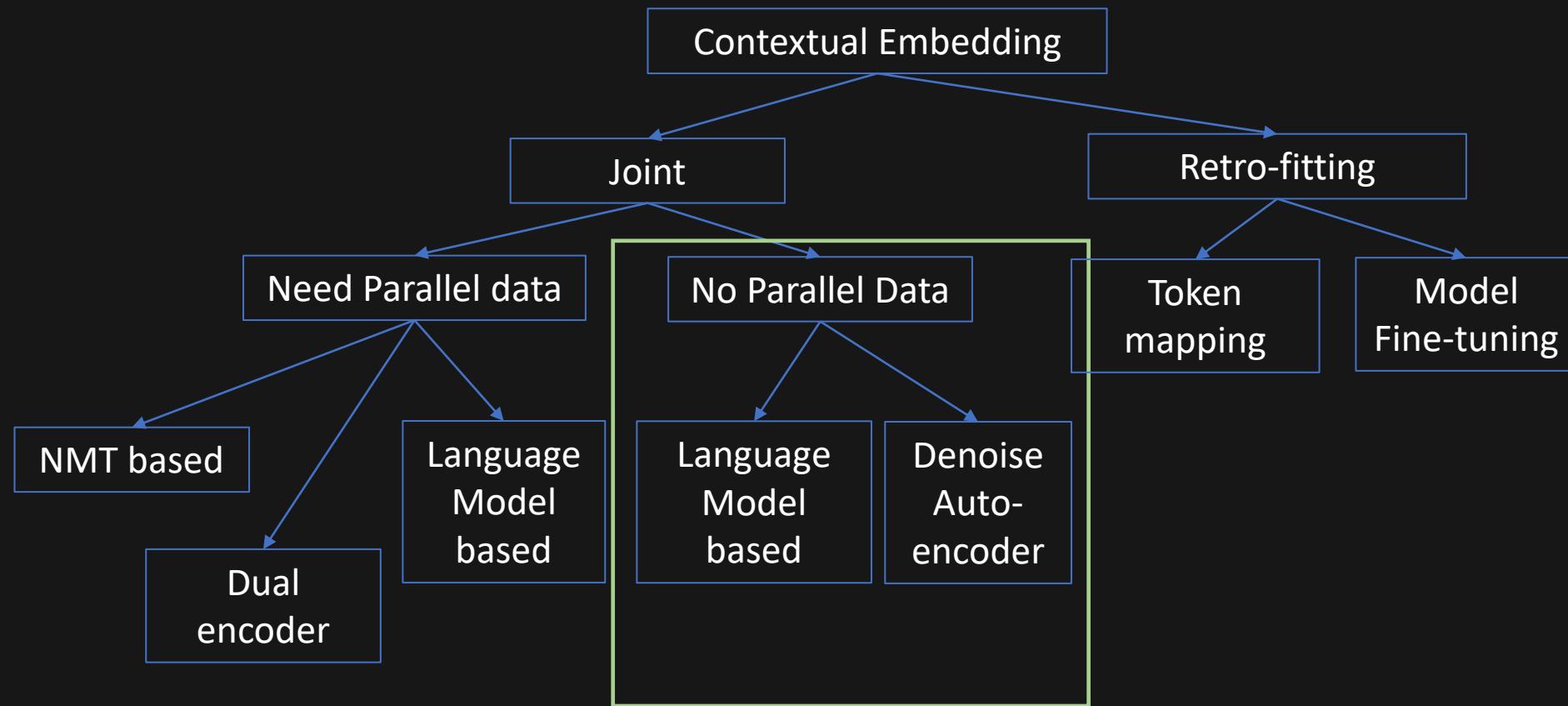
Effective Parallel Corpus Mining using Bilingual Sentence Embeddings. WMT 2018: 165-176

# Language Model based Methods

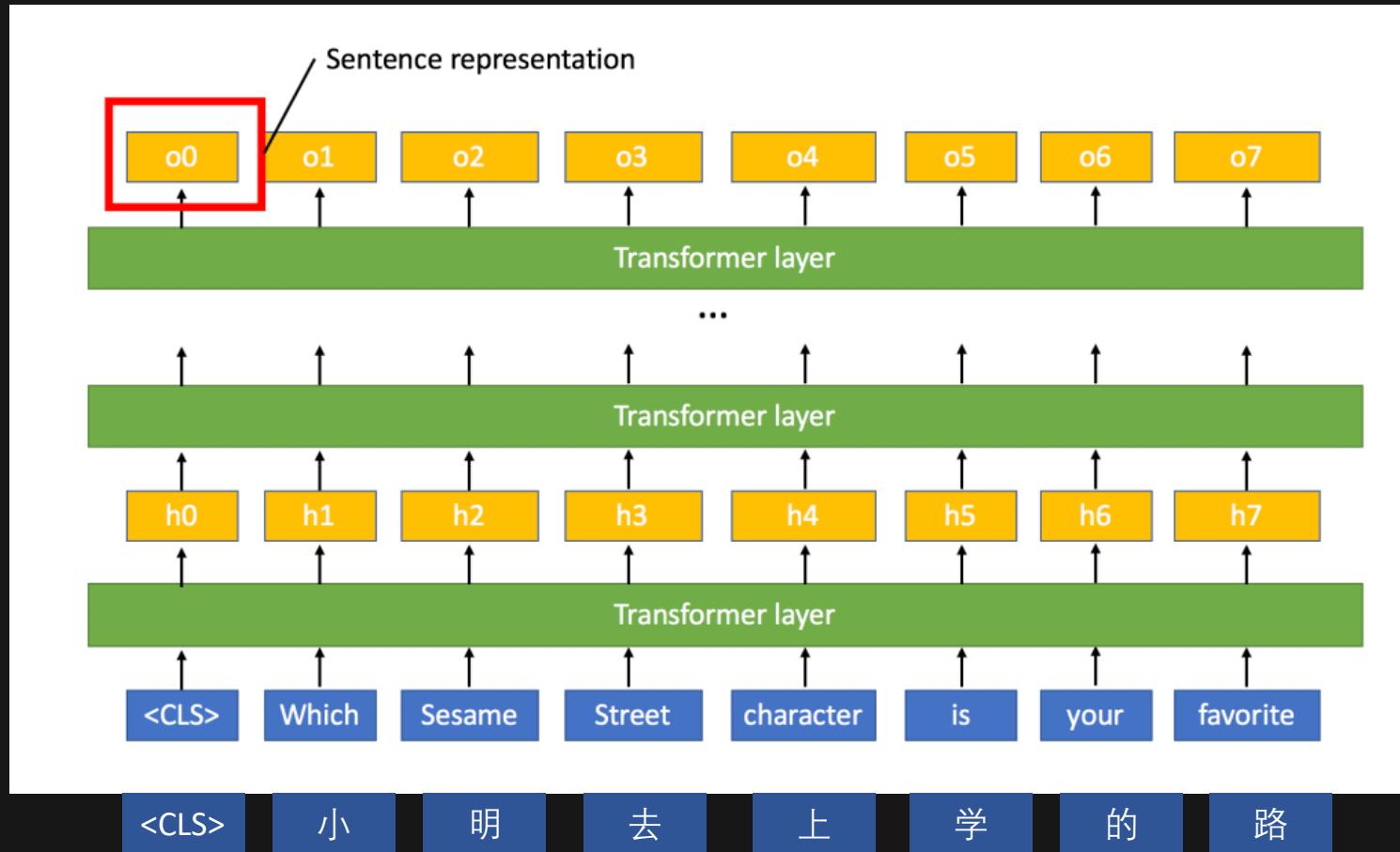


- Transformer based language model
- Cross-lingual representation at **both sentence level and token level**

# Joint Learning without Parallel Data



# Multi-Lingual BERT

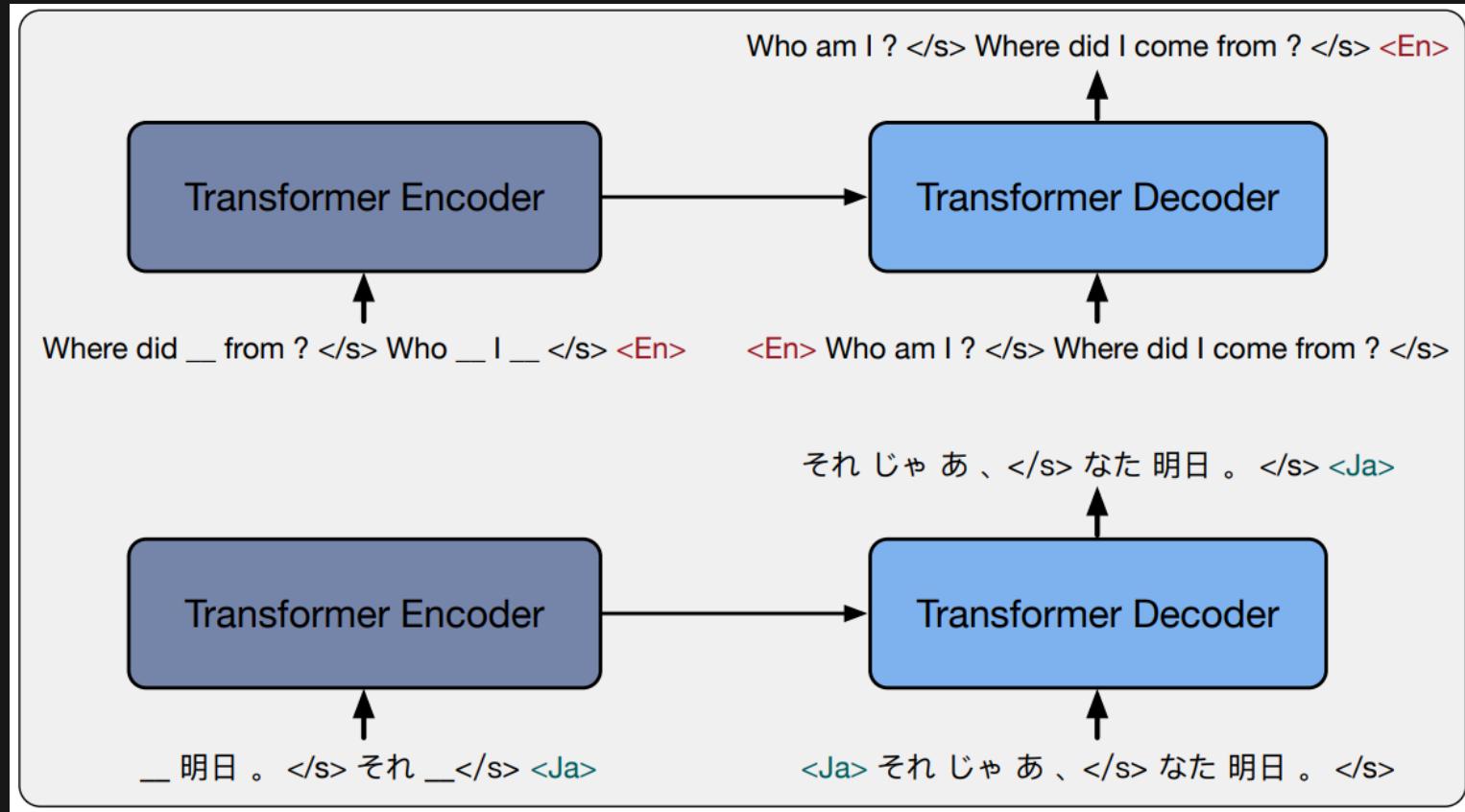


- Share the same architecture and training scheme with monolingual BERT
- 104 languages Wikipedia pages
- 110k shared WordPiece vocabulary

# Why mBERT has Cross-Lingual Effectiveness

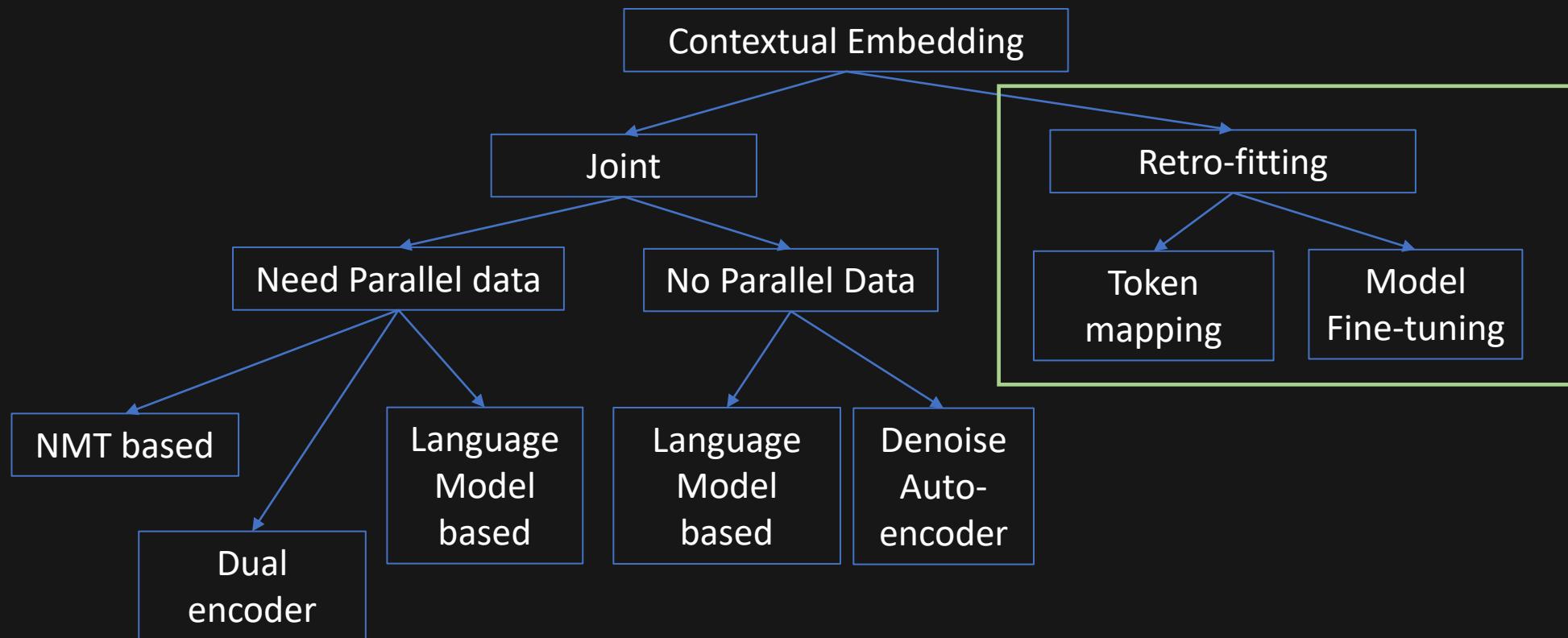
- Vocabulary overlap in different languages
    - Cognate words
    - Foreign words
    - However, it has been shown that mBERT still has cross-lingual effectiveness when there is zero vocabulary overlap
  - Common structure in different languages
    - Word order
    - Most frequent words
- 
- Pires T. et al. How multilingual is Multilingual BERT? ACL'19.
  - Wu, S. and Dredze M. Beto, Bentz, Becas: The Surprising Cross-Lingual Effectiveness of BERT, EMNLP'19.
  - Artetxe M. et al. On the Cross-lingual Transferability of Monolingual Representations, ACL'20.
  - Conneau, A. et al. Emerging Cross-lingual Structure in Pretrained Language Models, ACL'20.
  - Karthikeyan, K. et al. Cross-Lingual Ability of Multilingual BERT: An Empirical Study, ICLR'20.
  - Libovický, J. et al. On the Language Neutrality of Pre-trained Multilingual Representations, EMNLP Findings'20.

# Multi-Lingual Bidirectional and Auto-Regressive Transformers (mBART)

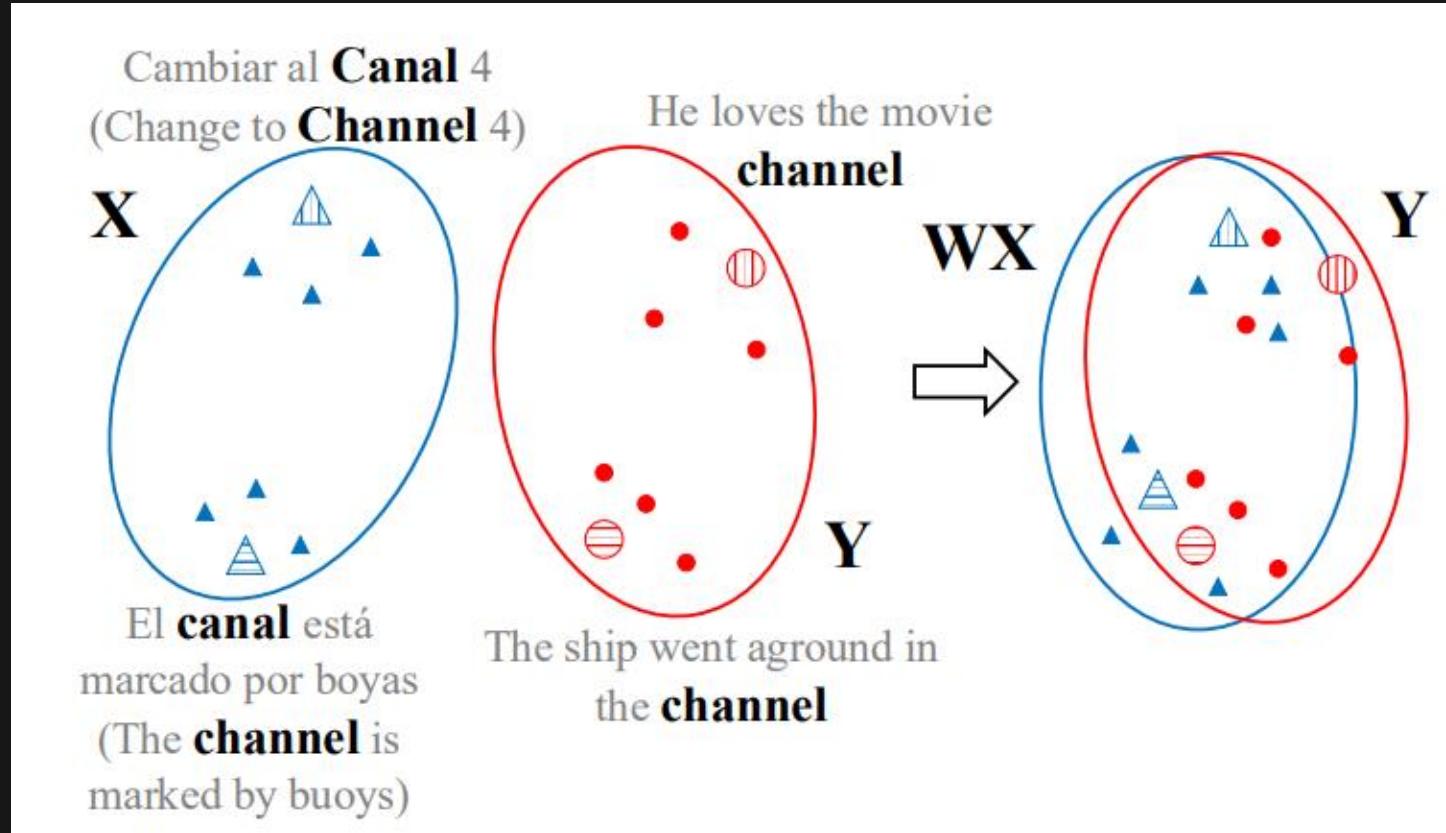


- Same architecture and training scheme with monolingual BART
- 25 languages crawled web pages
- 250K shared Sentence Piece vocabulary

# Joint learning without Parallel Data

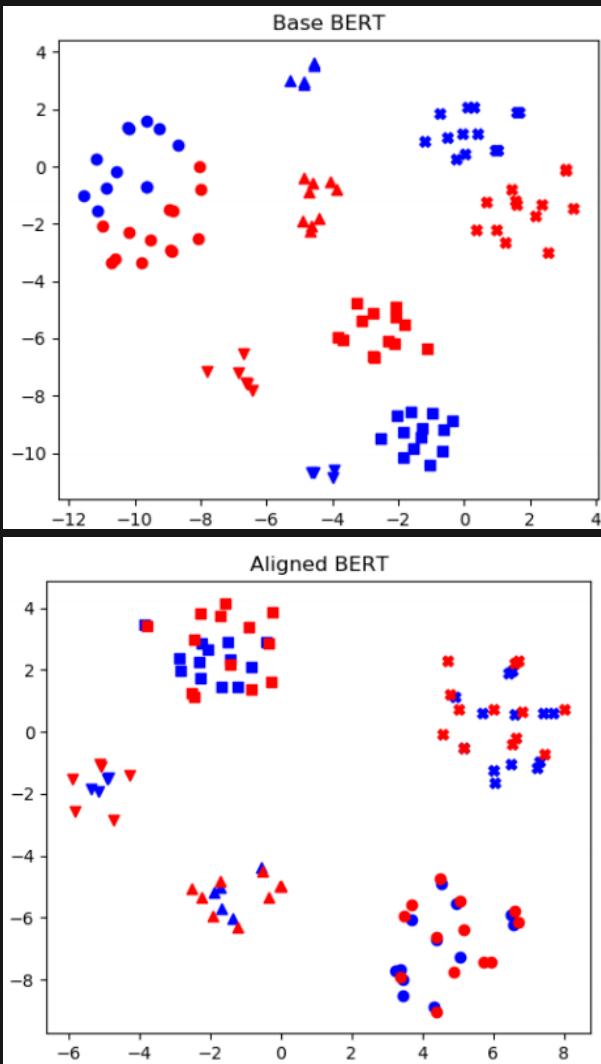


# Contextual Token Mapping



- Learn the mapping matrix  $W$   
$$\min_W \sum_{i=1}^n \|Wx_i - y_i\|^2$$
 where  $W^\top W = I$
- Instead of using static words in dictionary, use parallel sentences to align the word occurrences in sentences

# Model Fine-Tuning



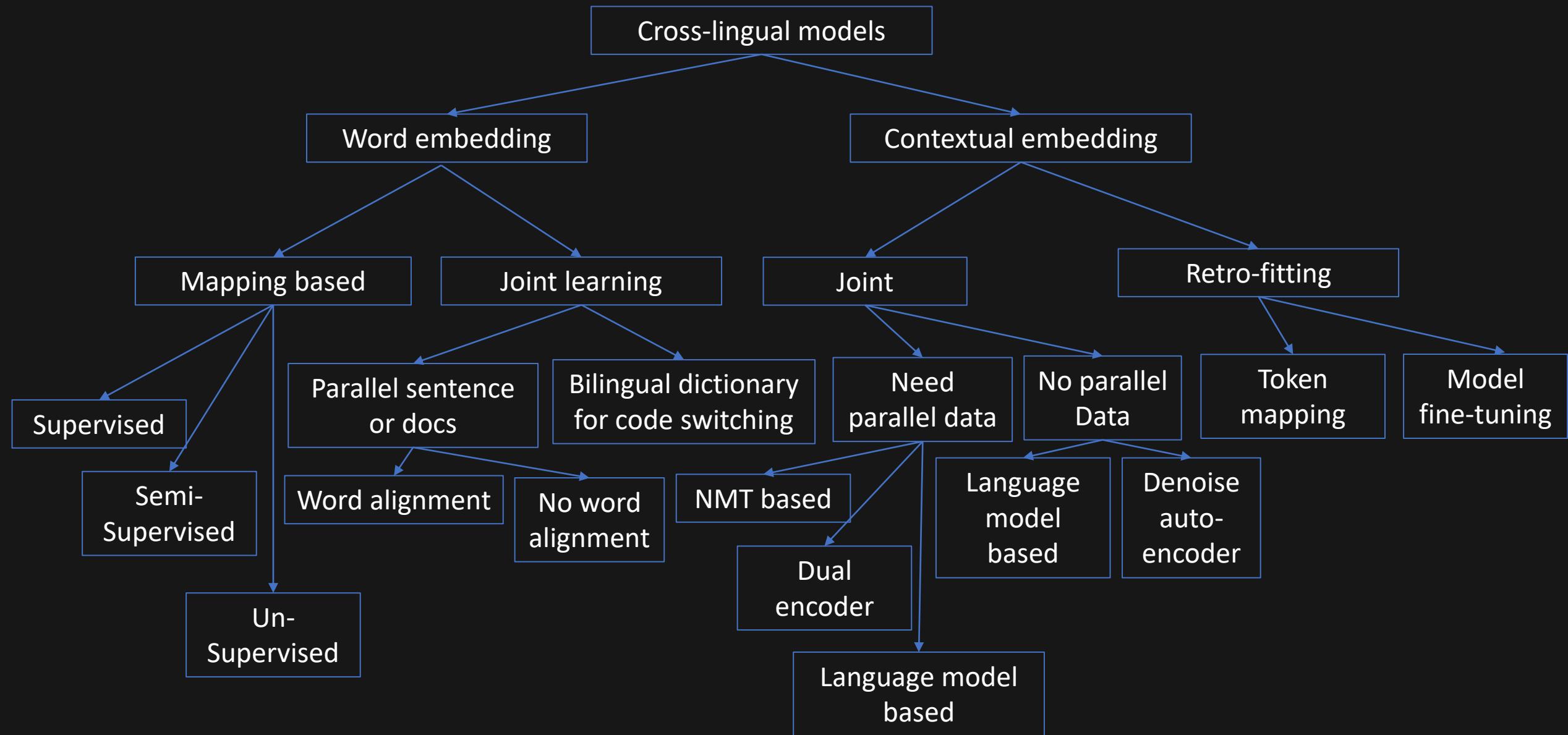
- Upper figure: mBERT roughly aligns the words
- Lower figure: effect of fine-tuning
- After mBERT training, fine-tuning by

$$\mathcal{L}(f, C) = - \sum_{(s,t) \in C} \sum_{(i,j) \in a(s,t)} \|f(i,s) - f(j,t)\|^2$$

( $s, t$ ) are parallel sentences, ( $i, j$ ) are aligned words,  $f()$  is the encoding function

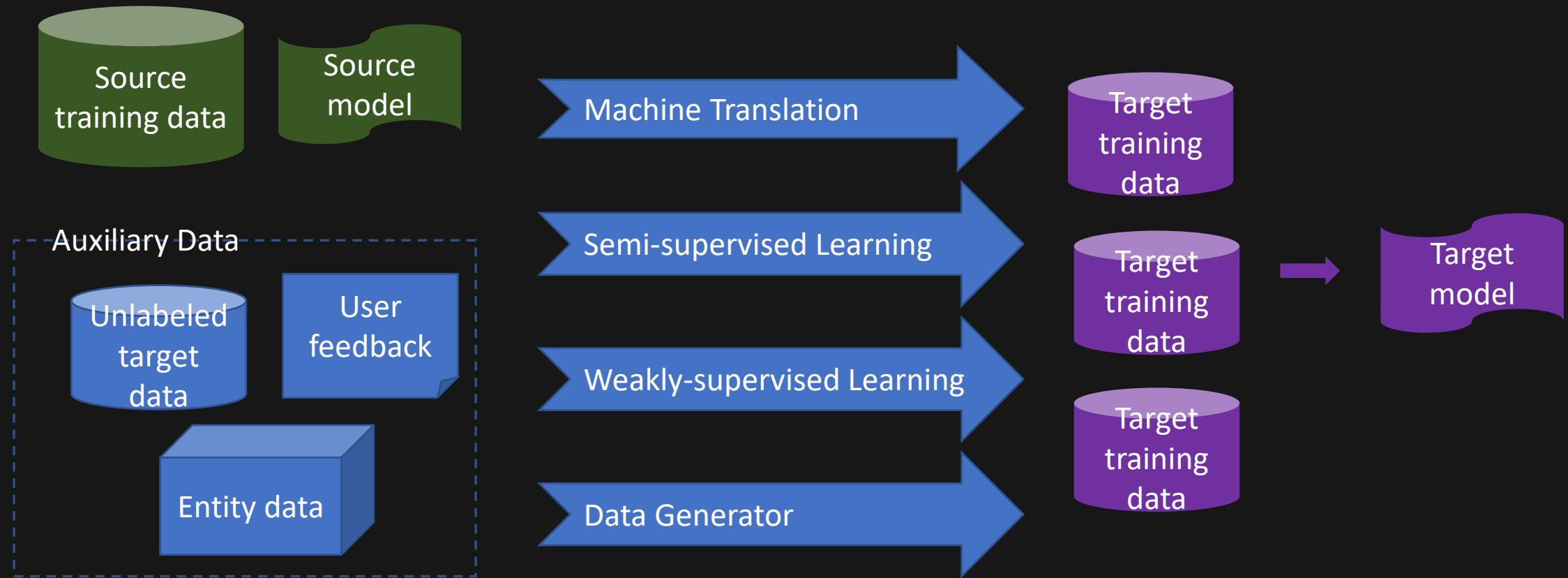
$$R(f, C) = \sum_{t \in C} \sum_{i=1}^{\text{len}(t)} \|f(i, t) - f_0(i, t)\|^2$$

To avoid trivial alignment (e.g.,  $f$  is constant), add regularization



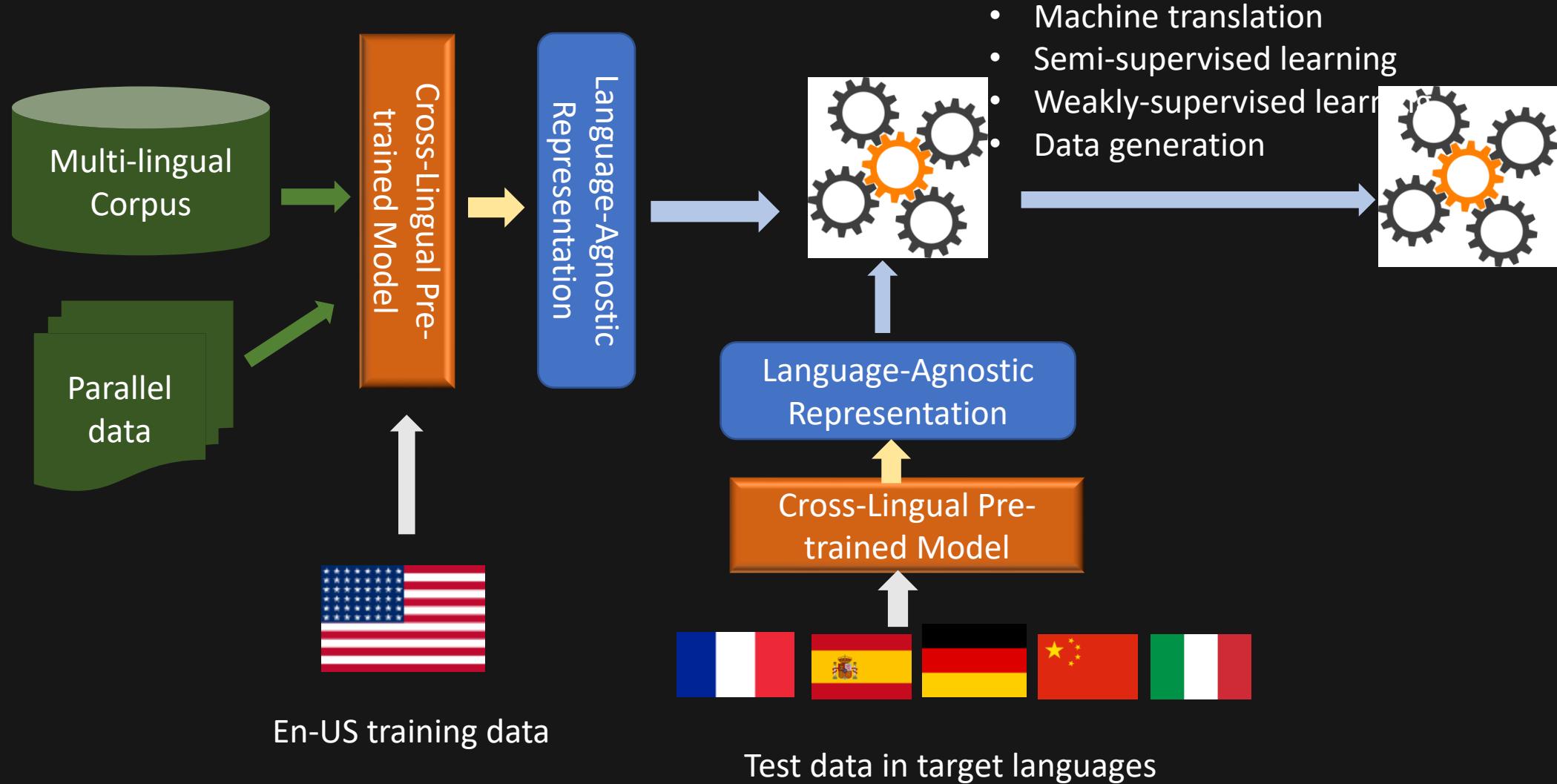
# Approaches: Model Transfer and Data Transfer

## (2) Data Transfer



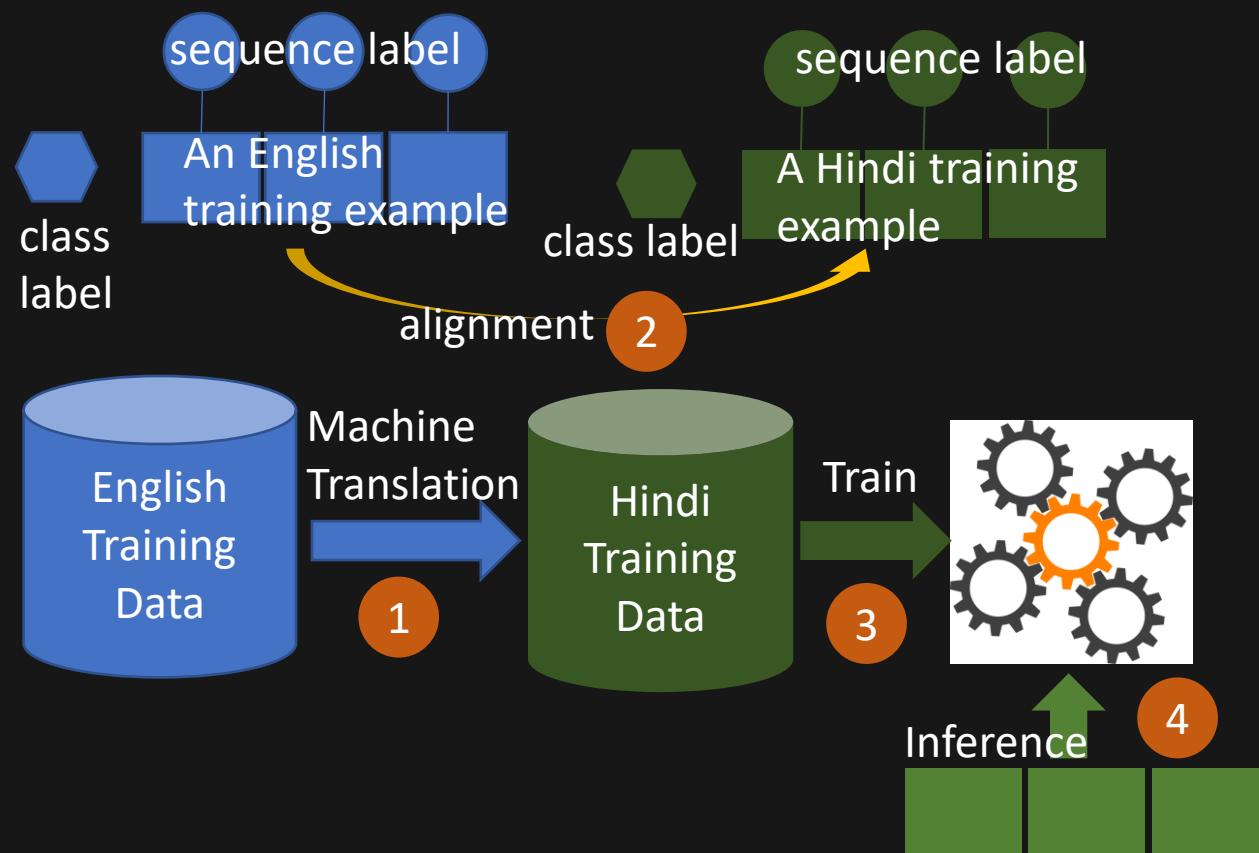
# Approaches: Model Transfer and Data Transfer

## (1) Model Transfer + (2) Data Transfer

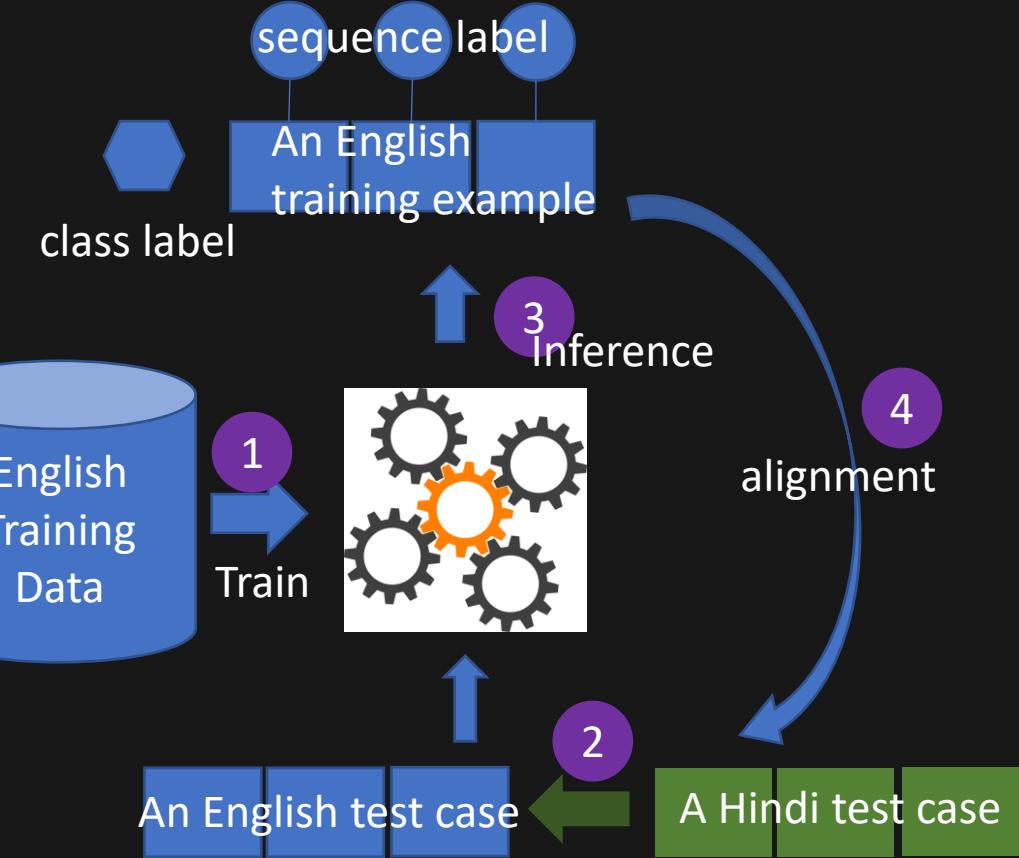


# Machine Translation

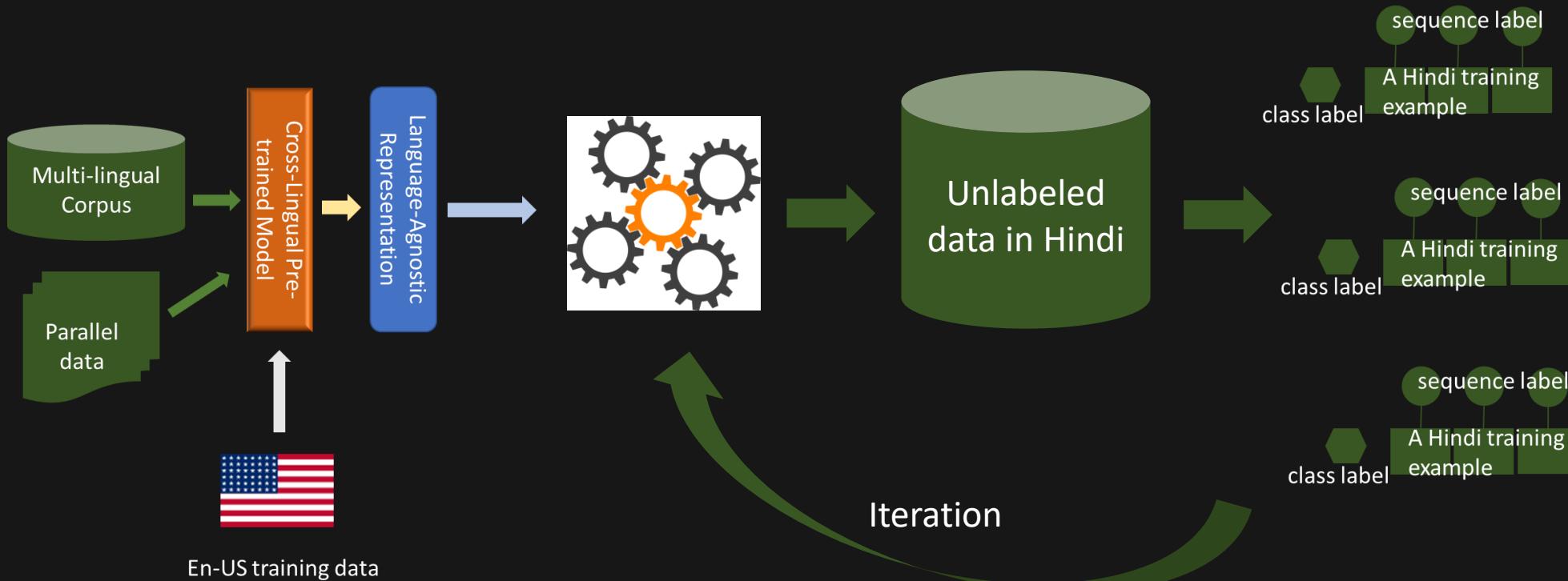
- Translate-on-train



- Translate-on-test



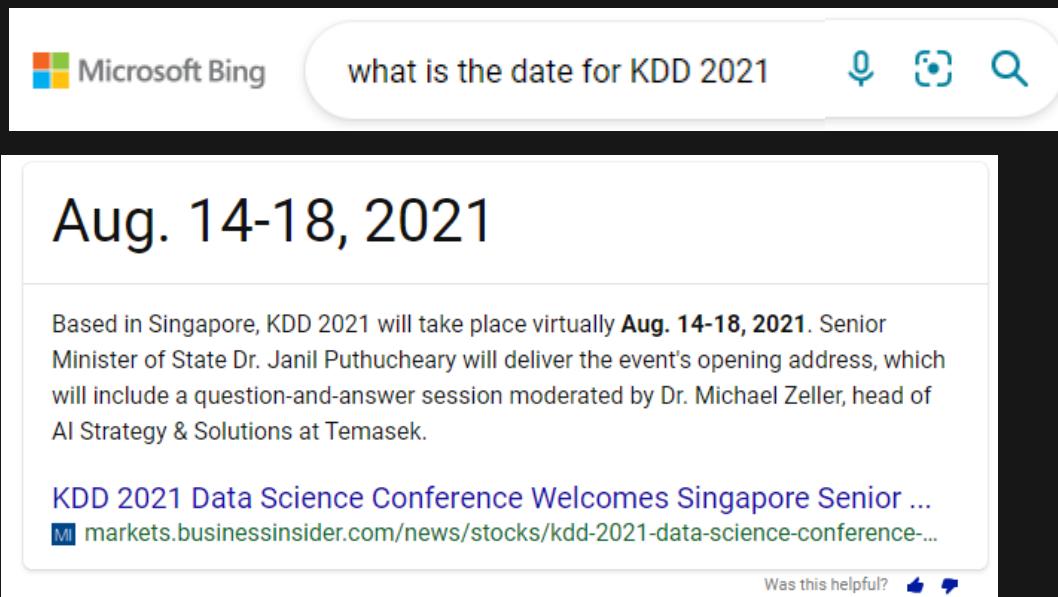
# Semi-Supervised Learning



1. Build an initial model by model transfer approach
2. Apply the initial model on unlabeled data to get labels
3. Use the labeled data to improve the initial model
4. Iterate the above steps 2-3 for several rounds (self learning, model ensemble, reinforcement learning)

# Weakly-Supervised Learning

- Collect auxiliary data
  - E.g., user feedback, knowledge data

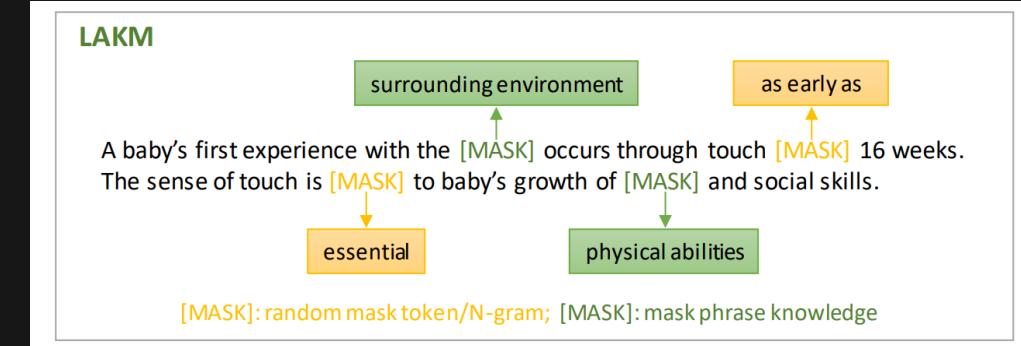


- Rich user behaviors in search logs
- Any single behavior is a weak signal, but more powerful when combined in a ML model
- Generally applicable to all languages

**[Question]:** who were the kings of the southern kingdom  
**[Passage]:** In the southern kingdom there was only one dynasty, that of king David, except usurper Athaliah from the northern kingdom, who by marriage, []  
**[Answer - ground truth]:** king David  
**[Answer - model predication]:** David, except usurper Athaliah

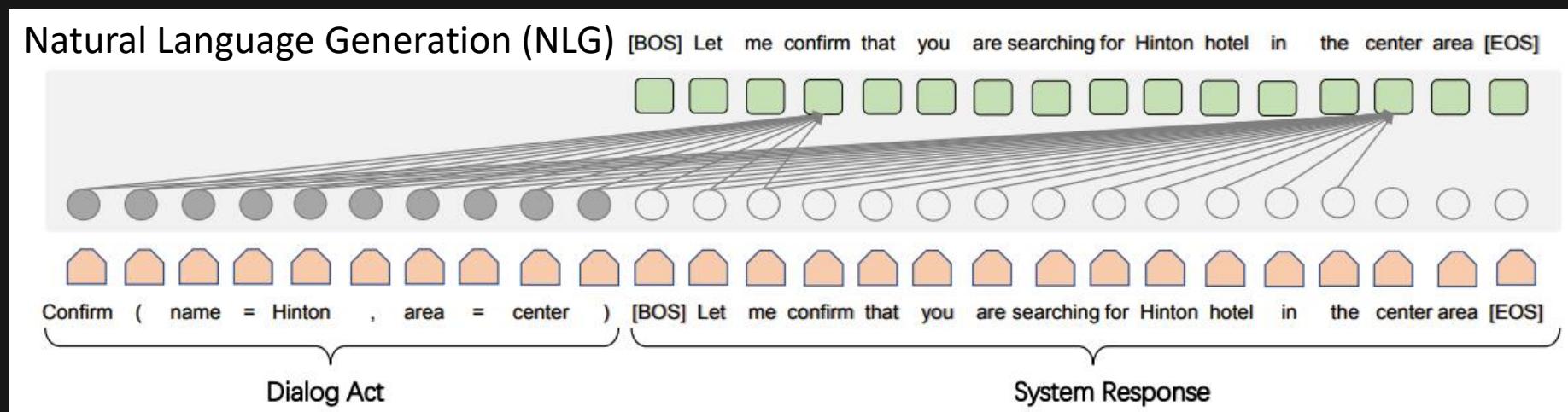
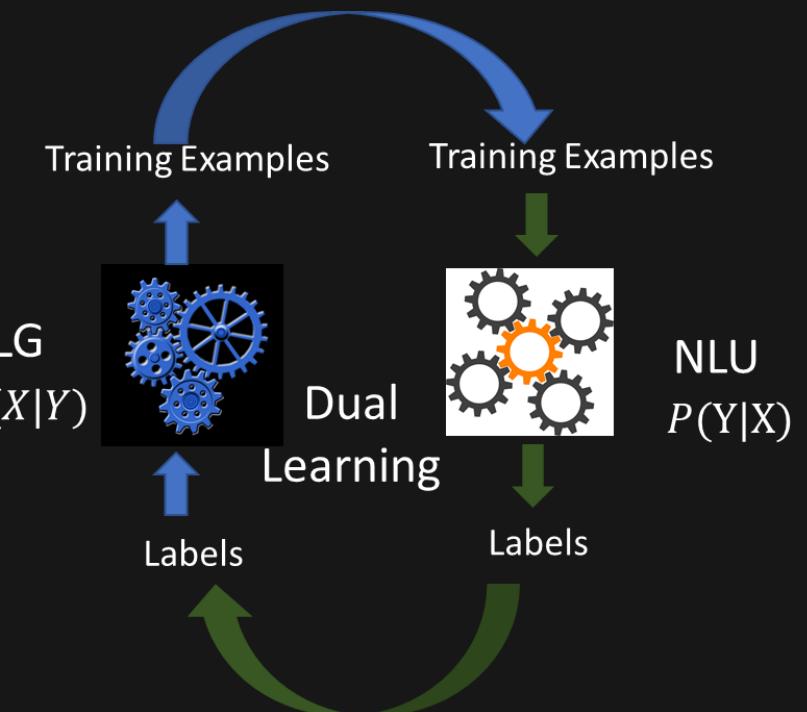
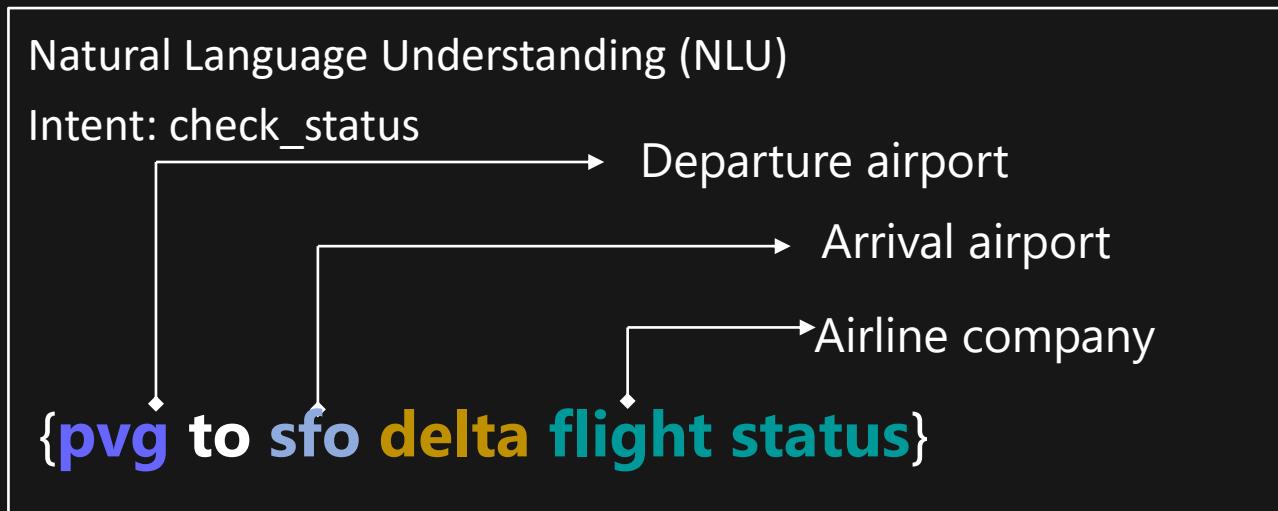
**[Question]:** What is the suggested initial does dosage of chlordiazepoxide  
**[Passage]:** If the drug is administered orally, the suggested initial dose is 50 to 100 mg, to be followed by repeated doses as needed until agitation is controlled up to 300 mg per day. []  
**[Answer - ground truth]:** 50 to 100 mg  
**[Answer - model predication]:** 100 mg

## Examples of boundary detection error



**Mask** entities and phrases in sentences and let the model to **recover** them as a pre-training task

# Data Generation



# Summary: Data Transfer Methods

The central idea of data transfer is to produce pseudo training data ( $X, Y$ ) in the target language, where  $X$  is the training instance, and  $Y$  is the label

Methods	Data instances ( $X$ )	How to derive labels ( $Y$ )
Machine translation	Translated	Alignment
Semi-supervised learning	Real data	Self-learning with iteration
Weakly-supervised learning	Real data	Implicitly/partially derived from auxiliary data
Data generation	Synthesized	Generated

# Outline

- Introduction [Dixin Jiang]
  - Motivating examples in Microsoft products
  - Problem description
  - Categorization of applications
  - Challenges
- Methodology [Dixin Jiang]
  - Model Transfer
  - Data Transfer
- Applications\*
  - Dependency Parsing [Xiubo Geng]
  - Machine Reading Comprehension [Ming Gong]
  - Grammar Error Correction [Linjun Shou]
- Summary & Future directions [Jian Pei]



Dixin Jiang

Software Technology Center at Asia (STCA) of Microsoft



Linjun Shou

Software Technology Center at Asia (STCA) of Microsoft



**Xiubo Geng**



Ming Gong



Jian Pei

Simon Fraser University

\*For more applications, please refer to our tutorial at  
The Web Conference 2021

# Cross Lingual Dependency Parsing

Xiubo Geng

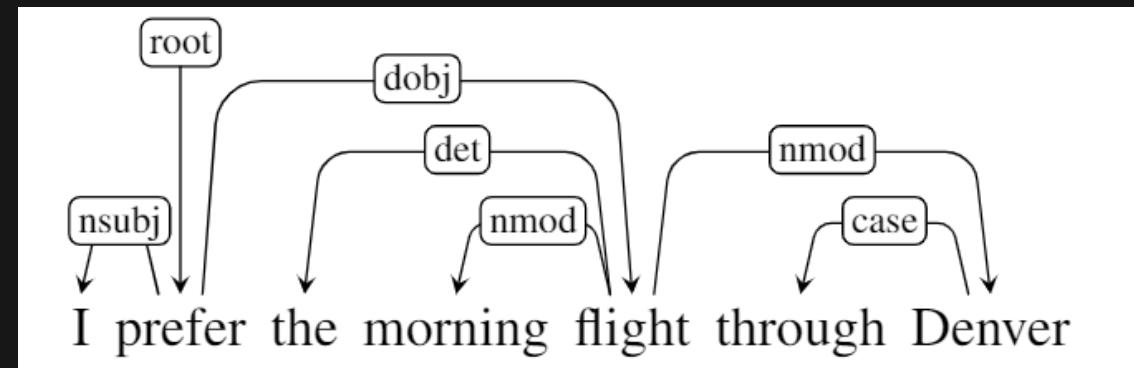
[xigeng@microsoft.com](mailto:xigeng@microsoft.com)

# Cross-lingual Dependency Parsing

- Background
- Model Transfer
  - Delexicalized model and improvements
  - Multilingual pre-trained models
- Data Transfer
  - Translation-Based Approach
    - Treebank Translation (source training data + translator)
    - Annotation Projection (parallel data + source model)
  - Unsupervised Approach
    - Source reordering
    - Direct unsupervised adaption

# Dependency Parsing

- Given a sentence, parse it into a dependency tree
- A dependency tree is a directed graph  $G = (V, A)$ , where
  - The nodes in  $V$  are words  $w_1, \dots, w_n$  of a sentence
  - The arcs in  $A$  have the form  $(w_i, l, w_j)$ , where  $l$  is label



# Dependency Parsing

Type	Category	Sub Category	Example
NLU	Text Classification	Single text	Domain identification, Intent detection, Sentiment classification
		Text pair	Information retrieval, Natural language inference
	Sequence Labeling	Single text	Named entity recognition, Slot tagging
		Text pair	Machine reading comprehension
NLG	Structure Prediction	Single text	Dependency parsing, constituency parsing, semantic role labeling
	Text Generation	Token level	Spelling correction, Sentence auto completion
		Sentence level	Machine translation, Conversation, Question generation

# Cross-lingual Dependency Parsing

- Given treebanks of rich-resource languages (e.g. English), derive dependency parsing models for low-resource languages
- Benchmarks – Universal Dependencies (UD)
  - A framework for consistency annotation of grammar (POS, syntactic dependencies, morphological features)
  - Nearly 200 treebanks in over 100 languages

# Dependency Parsing

- Evaluation Metrics
  - Labeled attachment score (LAS) =  
percentage of words that get the correct head and label
  - Unlabeled attachment score (UAS) =  
percentage of words that get the correct head

# Monolingual Dependency Parsing

- Graph-based Parser
  - Define a space of candidate dependency trees for a sentence
  - Learn to score dependency trees, factored into subgraphs
  - Parse by finding the highest-scoring dependency tree
- Arc-factored Model
  - Every arc is a separate subgraph

$$s(G = (V, A)) = \sum_{(w_i, l, w_j) \in A} s(w_i, l, w_j)$$

score of the labeled arc

- First parse an unlabeled tree, then label each arc
- So we need two scoring functions
  - $s^{(arc)}(w_i, w_j)$  to construct the tree
  - $s^{(label)}(l|w_i, w_j)$  to label each predicted arc

# Monolingual Model

- Biaffine attention scoring for arc prediction

$$h_i^{arc-dep} = MLP^{arc-dep}(r_i)$$

$$h_j^{arc-head} = MLP^{arc-head}(r_j)$$

$$s_{i \leftarrow j}^{arc} = \text{BiAffine}\left(h_i^{arc-dep}, h_j^{arc-head}\right)$$

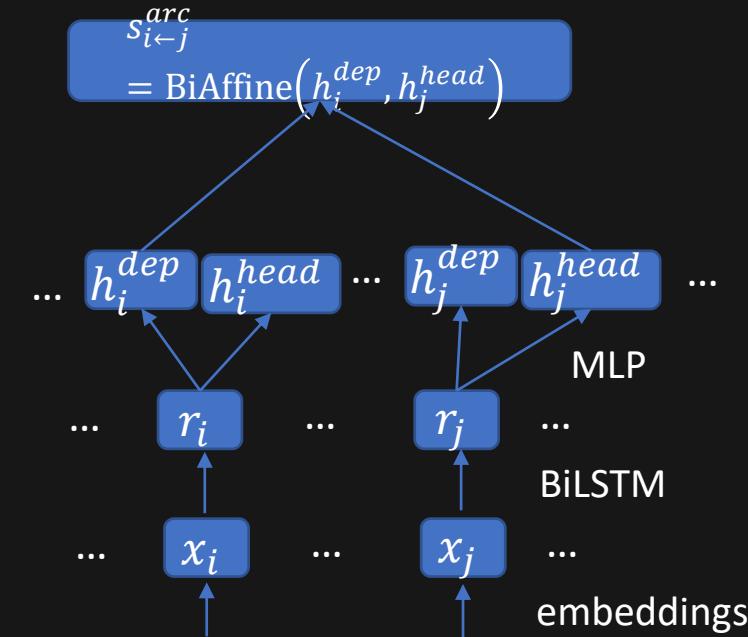
$$\text{BiAffine}(h_i, h_j) = h_j^T U^1 h_i + (h_j, h_i)^T U^2 + b$$

- Biaffine classifier to predict dependency labels

$$h_i^{label-dep} = MLP^{label-dep}(r_i)$$

$$h_j^{label-head} = MLP^{label-head}(r_j)$$

$$s_{i \leftarrow j}^{label} = \text{BiAffine}(h_i^{label-dep}, h_j^{label-head})$$

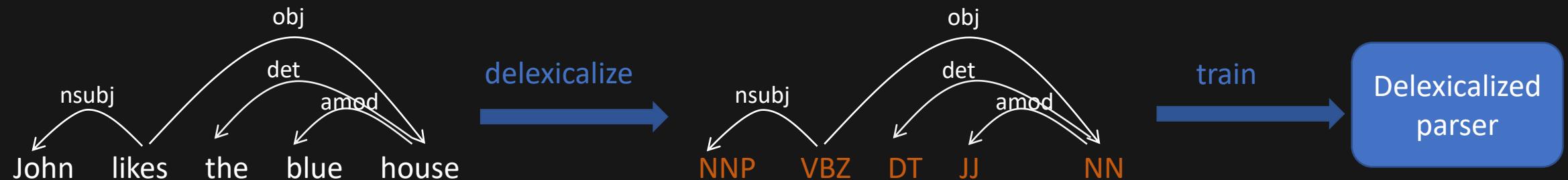


# Cross-lingual Dependency Parsing

- Background
- Model Transfer
  - Delexicalized model and improvements
  - Multilingual pre-trained models
- Data Transfer
  - Translation-Based Approach
    - Treebank Translation (source training data + translator)
    - Annotation Projection (parallel data + source model)
  - Unsupervised Approach
    - Source reordering
    - Direct unsupervised adaption

# Traditional Approaches

- Delexicalized model
  - Using shared POS features across languages
  - Success for closely related languages



# Traditional Approaches

- Fine-Grained Features

Feature	Description
81A	Order of Subject, Object and Verb
85A	Order of Adposition and Noun
86A	Order of Genitive and Noun
87A	Order of Adjective and Noun
88A	Order of Demonstrative and Noun
89A	Order of Numeral and Noun

Typological Features

Cluster	Lang.	Sample words
60	EN	was, wasn't, was'nt, wasn"t, hasnít, doesn't, ...
60	ES	estaba, estaráen, estubo, fúe, quedaba, ...
101	EN	very, mildly, wholly, terribly, gloriously, ...
101	ES	muchomás, fuerte, fuertes, duro, duros, poco, ...
153	EN	chicken, bird, ostriches, beef, pork, burger, steak, ...
153	ES	pollo, achote, manzana, tortugas, marsupiales, ...
195	EN	The, ...
195	ES	El, La, Los, Las, LoS, ...
236	EN	dry, wet, moist, lifeless, dullish, squarish, limpid, ...
236	ES	seco, secos, semiseco, semisecos, mojado, humedo, ...

Cross-lingual word clusters

# Modern Approach

- Pre-trained multilingual embeddings (mBERT) performs relatively well for zero-shot model transfer
- There is still a large gap between zero-shot transfer and supervised learning, especially for languages with large distance to English.

mBERT(S) : supervised learning

Baseline(Z): dictionary supervised cross-lingual word embeddings

mBERT(Z) : zero-shot transfer from English training data

mBERT(Z+POS) : gold POS is used as input

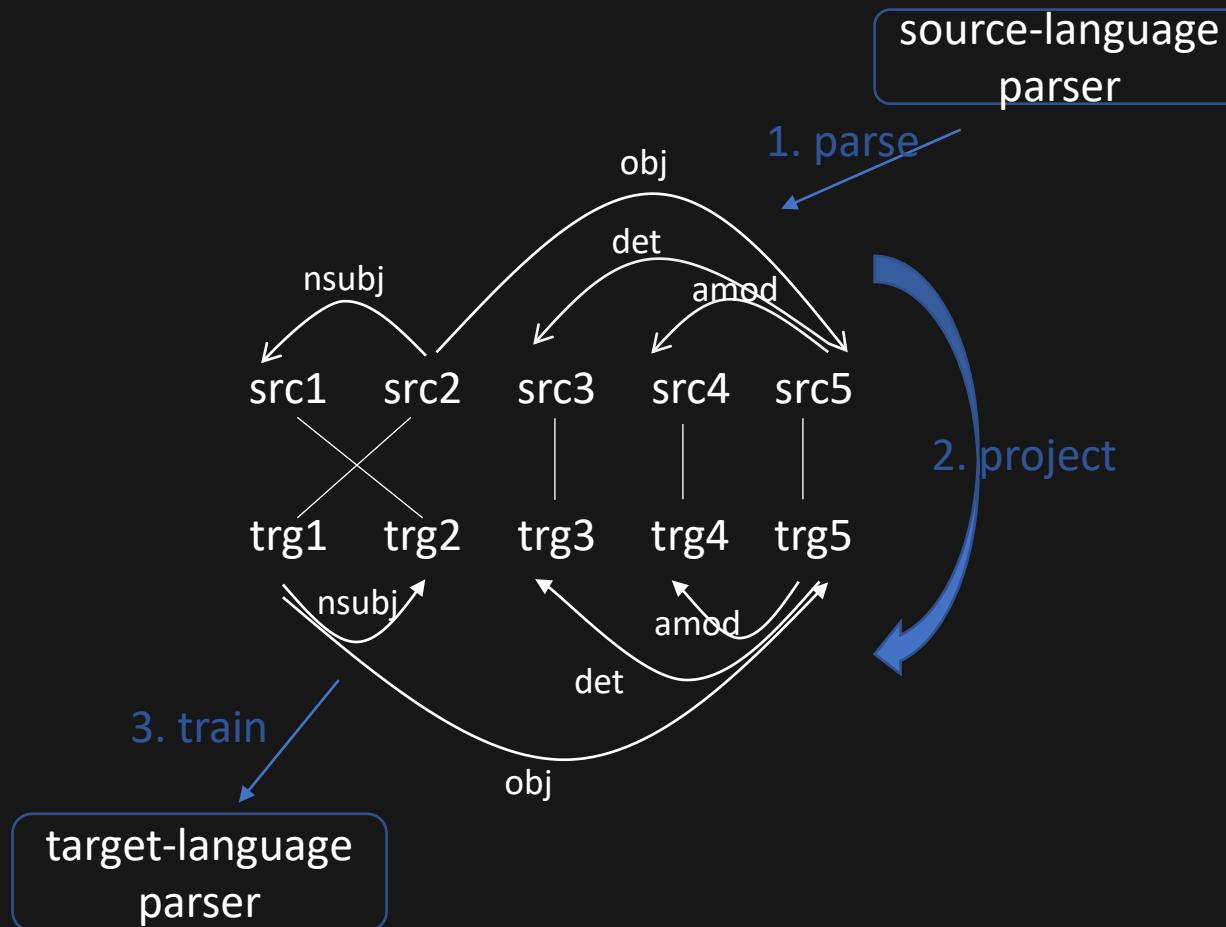
Languages are sorted according to word order distance to English

	Dist	mBERT(S)	Baseline(Z)	mBERT(Z)	mBERT(Z+POS)
en	0.00	91.5/81.3	90.4/ <b>88.4</b>	<u>91.5</u> /81.3	<b>91.8</b> / <u>82.2</u>
no	0.06	93.6/85.9	<u>80.8</u> / <b>72.8</b>	80.6/68.9	<b>82.7</b> / <u>72.1</u>
sv	0.07	91.2/83.1	81.0/ <u>73.2</u>	<u>82.5</u> /71.2	<b>84.3</b> / <u>73.7</u>
fr	0.09	91.7/85.4	77.9/ <u>72.8</u>	<u>82.7</u> /72.7	<b>83.8</b> / <u>76.2</u>
pt	0.09	93.2/87.2	76.6/ <b>67.8</b>	<u>77.1</u> /64.0	<b>78.3</b> / <u>66.9</u>
da	0.10	89.5/81.9	76.6/ <u>67.9</u>	<u>77.4</u> /64.7	<b>79.3</b> / <u>68.1</u>
es	0.12	92.3/86.5	74.5/ <u>66.4</u>	<u>78.1</u> /64.9	<b>79.0</b> / <u>68.9</u>
it	0.12	94.8/88.7	80.8/ <u>75.8</u>	<u>84.6</u> /74.4	<b>86.0</b> / <u>77.8</u>
ca	0.13	94.3/89.5	73.8/ <u>65.1</u>	<u>78.1</u> /64.6	<b>79.0</b> / <u>67.9</u>
hr	0.13	92.4/83.8	61.9/52.9	<b>80.7</b> / <u>65.8</u>	<u>80.4</u> / <u>68.2</u>
pl	0.13	94.7/79.9	74.6/ <u>62.2</u>	<u>82.8</u> /59.4	<b>85.7</b> / <u>65.4</u>
sl	0.13	88.0/77.8	68.2/ <u>56.5</u>	<u>72.6</u> /51.4	<b>75.9</b> / <u>59.2</u>
uk	0.13	90.6/83.4	60.1/52.3	<b>76.7</b> / <u>60.0</u>	<u>76.5</u> / <u>65.5</u>
bg	0.14	95.2/85.5	79.4/ <b>68.2</b>	<u>83.3</u> /62.3	<b>84.4</b> / <u>68.1</u>
cs	0.14	94.2/86.6	63.1/53.8	<u>76.6</u> / <u>58.7</u>	<b>77.4</b> / <u>63.6</u>
de	0.14	86.1/76.5	71.3/61.6	<u>80.4</u> /66.3	<b>83.5</b> / <u>71.2</u>
he	0.14	91.9/83.6	55.3/48.0	<b>67.5</b> / <u>48.4</u>	<u>67.0</u> / <u>54.3</u>
nl	0.14	94.0/85.0	68.6/60.3	<u>78.0</u> / <u>64.8</u>	<b>79.9</b> / <u>67.1</u>
ru	0.14	94.7/88.0	60.6/51.6	<b>73.6</b> / <u>58.5</u>	<u>73.2</u> / <u>61.5</u>
ro	0.15	92.2/83.2	65.1/54.1	<b>77.0</b> / <u>58.5</u>	<u>76.9</u> / <u>62.6</u>
id	0.17	86.3/75.4	49.2/43.5	<b>62.6</b> / <u>45.6</u>	<u>59.8</u> / <u>48.6</u>
sk	0.17	93.8/83.3	66.7/58.2	<u>82.7</u> / <u>63.9</u>	<b>82.9</b> / <u>67.8</u>
lv	0.18	87.3/75.3	<b>70.8</b> / <u>49.3</u>	66.0/41.4	<u>70.4</u> / <u>48.5</u>
et	0.20	88.8/79.7	<u>65.7</u> / <u>44.9</u>	<u>66.9</u> /44.3	<b>70.8</b> / <u>50.7</u>
fi	0.20	91.3/81.8	66.3/48.7	68.4/47.5	<b>71.4</b> / <u>52.5</u>
zh*	0.23	88.3/81.2	42.5/25.1	<b>53.8</b> / <u>26.8</u>	<u>53.4</u> / <u>29.0</u>
ar	0.26	87.6/80.6	38.1/28.0	<u>43.9</u> / <u>28.3</u>	<b>44.7</b> / <u>32.9</u>
la	0.28	85.2/73.1	<u>48.0</u> / <b>35.2</b>	47.9/26.1	<b>50.9</b> / <u>32.2</u>
ko	0.33	86.0/74.8	34.5/16.4	<b>52.7</b> / <u>27.5</u>	<u>52.3</u> / <u>29.4</u>
hi	0.40	94.8/86.7	35.5/26.5	<u>49.8</u> / <u>33.2</u>	<b>58.9</b> / <u>44.0</u>
ja*	0.49	94.2/87.4	28.2/ <u>20.9</u>	<u>36.6</u> /15.7	<b>41.3</b> / <u>30.9</u>
AVER	0.17	91.3/82.6	64.1/53.8	<u>71.4</u> / <u>54.2</u>	<b>73.0</b> / <u>58.9</u>

# Cross-lingual Dependency Parsing

- Background
- Model Transfer
  - Delexicalized model and improvements
  - Multilingual pre-trained models
- Data Transfer
  - Translation-Based Approach
    - Treebank Translation
    - Annotation Projection
  - Unsupervised Approach
    - Source reordering
    - Direct unsupervised adaption

# Annotation Projection

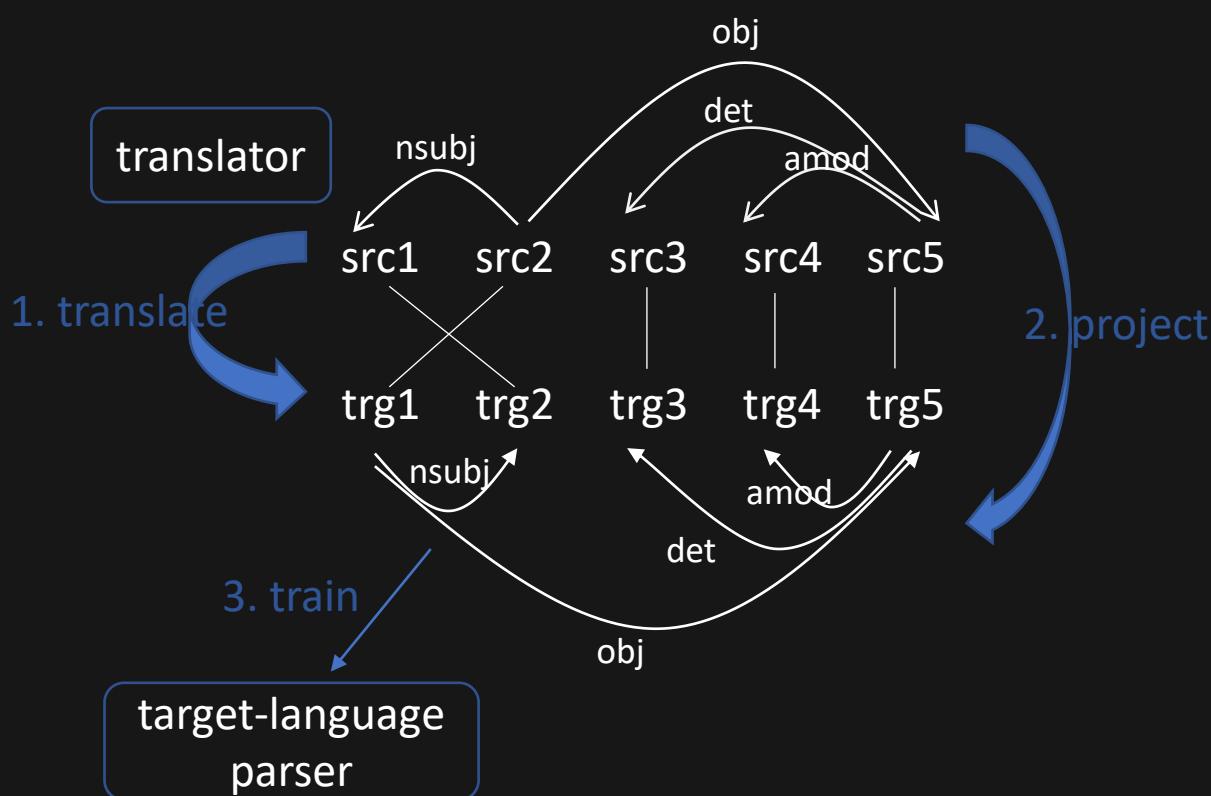


- Given a word-aligned parallel corpus and a source-language parser
  - Annotate the source sentences with the parser
  - Use word alignment to map the source annotation to target sentences
  - Train a target-language parser with the projected annotation

Yarowsky, D., Ngai, G., & Wicentowski, R. (2001). Inducing Multilingual Text Analysis Tools via Robust Projection Across Aligned Corpora. In Proceedings of HLT.

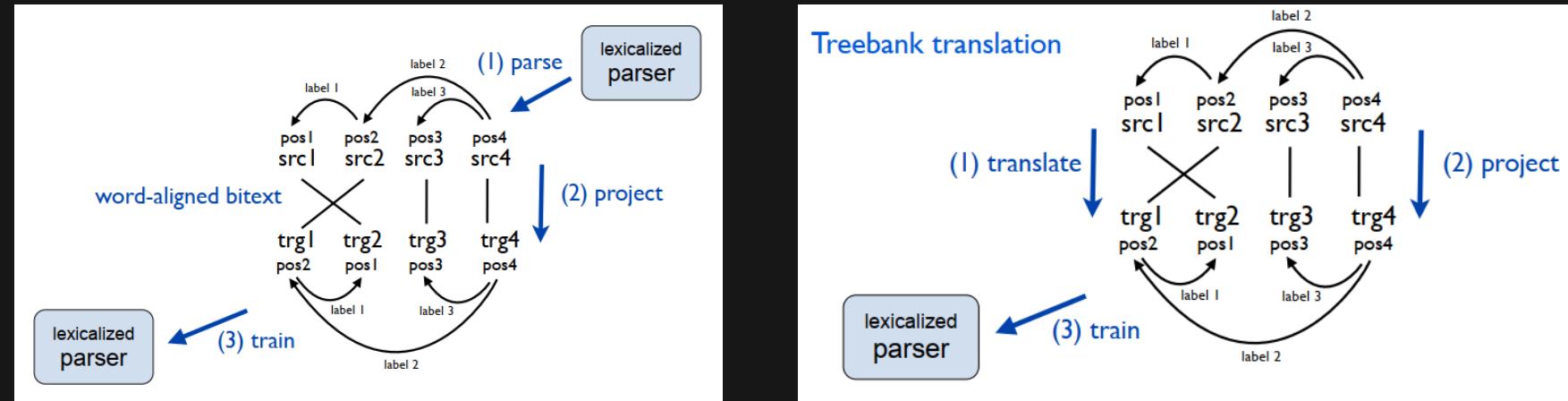
Hwa, R., Resnik, P., Weinberg, A., Cabezas, C., & Kolak, O. (2005). Bootstrapping Parsers via Syntactic Projection across Parallel Texts. Natural Language Engineering.

# Treebank Translation



- Given a source-language treebank corpus and a source-to-target translator
  - Translate source sentences into target sentences with word-to-word or phrase-to-phrase alignment
  - Map source annotation to target sentences
  - Train a target-language parser with the projected annotation

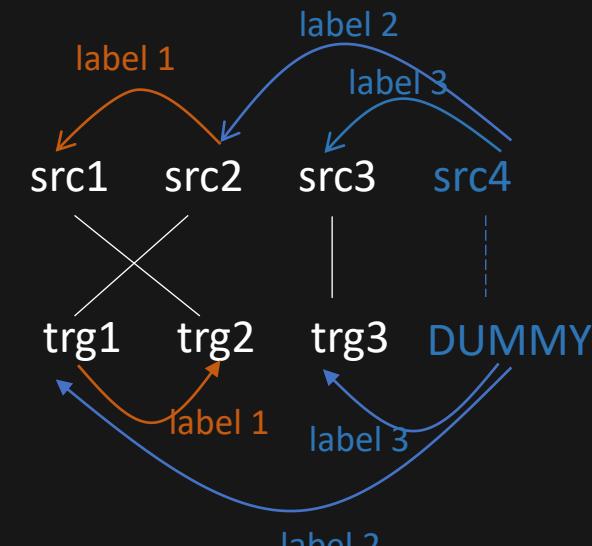
# Discussion



	Annotation Projection	Treebank Translation
Data Required	Source-Target sentence pairs	Annotated source treebank
Model Required	Source annotator	Translator
Source-language annotation	By an annotation model	By human annotators
Target-language sentence	By human	By a translation model
Target-language annotation	Projected from source annotation	Projected from source annotation

# Projecting Strategy

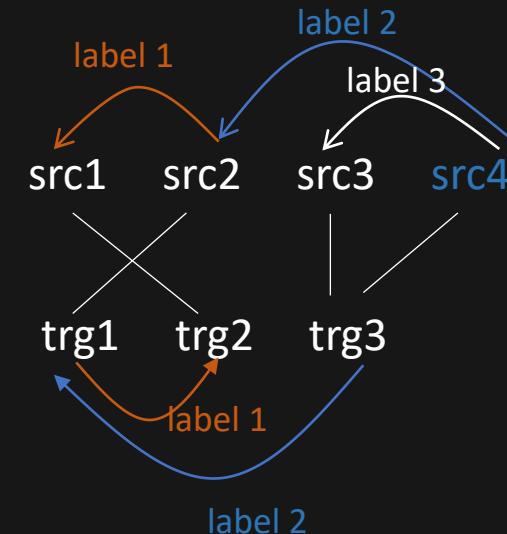
- One-to-one
  - Directly copy dependency relations
  - Unaligned source language words are covered by additional DUMMY nodes
  - Unaligned target language words are deleted



one-to-one

# Projecting Strategy

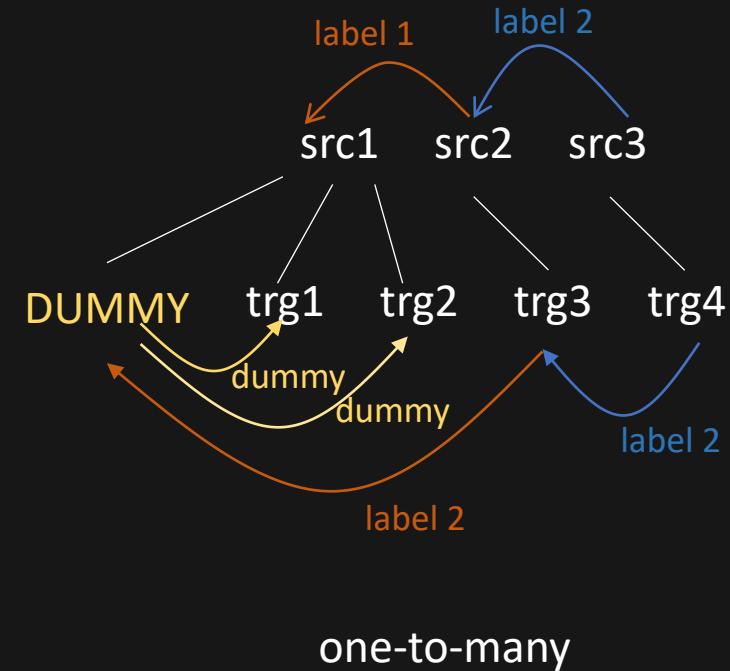
- One-to-one
- Many-to-one
  - Only keep the link to the head of the aligned source language words
  - Delete all other links



many to one

# Projecting Strategy

- One-to-one
- Many-to-one
- One-to-many
  - Introduce an additional DUMMY node
  - The DUMMY node acts as the immediate parent in the target language
  - The DUMMY node captures the dependency relation of the source side annotation



# Projecting Strategy

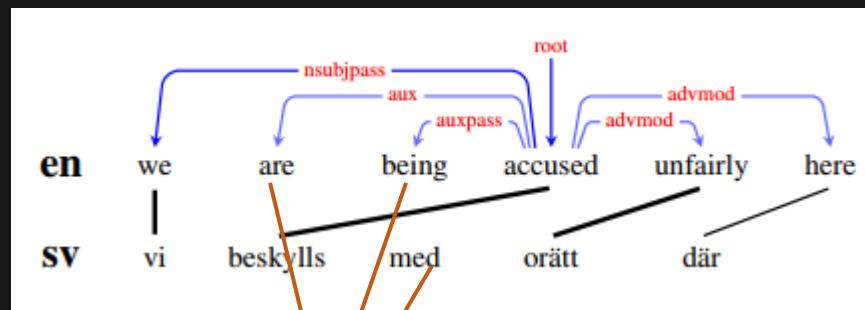
- One-to-one
- Many-to-one
- One-to-many
- Many-to-many
  - First apply one-to-many, then apply many-to-one

# Discussion

- Advantages
  - Synthesize target-language data for model training
  - Competitive and complementary to cross-lingual models

# Discussion

- Advantages
  - Synthesize target-language data for model training
  - Competitive and complementary to cross-lingual models
- Limitations
  - Non-isomorphic alignment between source and target language



The English words “are” and “being”, and the Swedish word “med”, do not have corresponding word-level translation

# Discussion

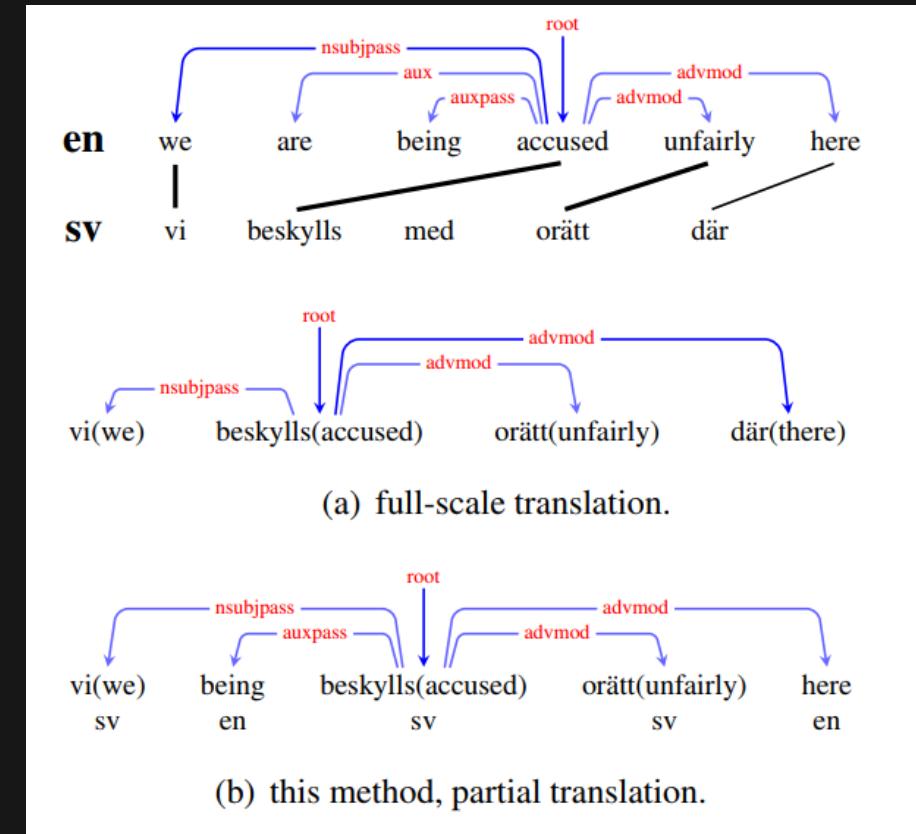
- Advantages
  - Synthesize target-language data for model training
  - Competitive and complementary to cross-lingual models
- Limitations
  - Non-isomorphic alignment between source and target language
  - Translation error or parsing error

# Discussion

- Advantages
  - Synthesize target-language data for model training
  - Competitive and complementary to cross-lingual models
- Limitations
  - Non-isomorphic alignment between source and target language
  - Translation error or parsing error
  - Conflicting annotation

# Code-Mixed Treebank

- Translate a source treebank into a code-mixed treebank
  - Source words with highly confident target alignments are translated into target words
  - Continuous spans of target words are reordered according to translation
- Make the best use of source syntax
  - With the minimum noise introduced
  - Capture the structure of target-language sentences as much as possible



# Code-Mixed Treebank

Lang.	Delex		PartProj		Src		Tgt		Src + Tgt		Mix		<b>Src + Mix</b>	
	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS	UAS	LAS
DE	64.10	53.77	69.90	61.28	66.87	57.46	70.84	62.30	72.41	63.74	71.41	63.46	<b>72.78</b>	<b>64.38</b>
ES	71.53	63.33	75.81	66.83	75.63	65.85	76.49	67.39	77.00	67.95	81.18	71.80	<b>81.44</b>	<b>71.66</b>
FR	75.13	67.26	75.54	67.63	78.13	70.63	76.91	69.39	78.75	71.17	83.20	76.32	<b>83.77</b>	<b>76.48</b>
IT	77.71	69.27	77.71	69.27	81.11	72.83	79.30	71.65	81.56	74.09	85.30	77.43	<b>86.13</b>	<b>78.38</b>
PT	74.03	67.70	79.44	71.30	77.37	69.36	78.32	70.67	79.73	71.84	83.54	75.34	<b>84.05</b>	<b>75.89</b>
AVG	72.50	64.27	75.68	67.26	75.82	67.23	76.37	68.28	77.89	69.76	80.93	72.87	<b>81.63</b>	<b>73.36</b>

- All lexicalized model outperform Delex, showing effectiveness of lexicalized features
- The proposed code-mixed approach (Mix) increases the accuracy

# Cross-lingual Dependency Parsing

- Background
- Model Transfer
  - Delexicalized model and improvements
  - Multilingual pre-trained models
- Data Transfer
  - Translation-Based Approach
    - Treebank Translation (source training data + translator)
    - Annotation Projection (parallel data + source model)
  - Unsupervised Approach
    - Source reordering
    - Direct unsupervised adaption

# Source Reordering

- Given a source-language treebank and unannotated target-language corpus
- Rearrange the order of source sentences to meet the word order in target language
- Under the guidance of a POS language model



Source-language tree

Reorder according to a  
POS language model



Reordered source-language tree

# Source Reordering

- Given a source-language treebank and unannotated target-language corpus
- Rearrange the order of source sentences to meet the word order in target language
- Under the guidance of a POS language model

$$x^* = \arg \max_{x' \in R(x)} p_{\tau}(x')$$

Set of all possible permutations  
of the words in  $x$

The POS language model learned  
from the target-language corpus

# Unsupervised Adaptation

- Leverage the unsupervised parser based on CRF autoencoder framework
  - Applicable for both supervised and unsupervised learning
- Approach
  - Train a source parser
  - Use it to initialize and regularize a target parser
  - Unsupervised training on unannotated target data

# Unsupervised Adaptation

- Supervised Learning

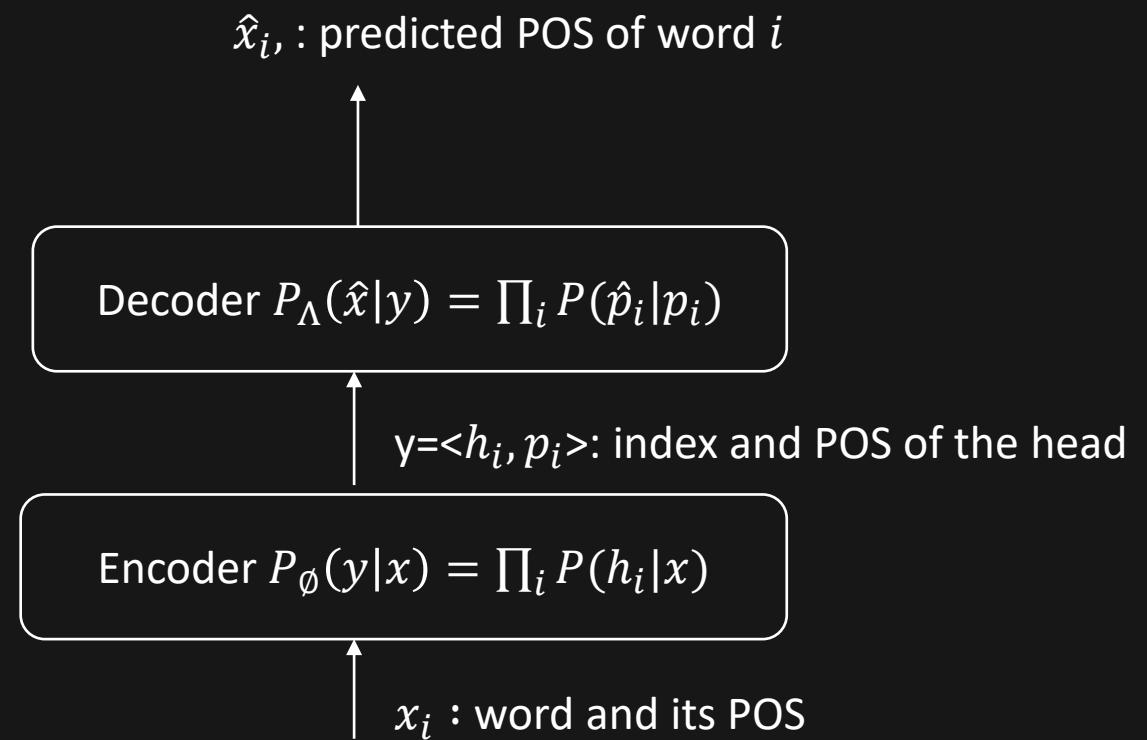
$$\mathcal{L} = -\sum \log P(y^*, \hat{x}|x)$$

- Unsupervised Learning

$$\mathcal{L} = -\sum \max_y \log P(y, \hat{x}|x)$$

- Cross-lingual adaptation

- Train a CRF autoencoder in a supervised way with multilingual BERT
- Use the source model to initialized a second CRF autoencoder
- Train it in an unsupervised way on unannotated target data with regularization by source model



# Data Transfer: Summary

	<b>Required resource</b>	<b>Input sentence(X)</b>	<b>How to derive labels (Y)</b>	<b>Training data languages</b>
Treebank translation	Source treebank + translator	translated	Alignment according to source treebank	Target language
Annotation projection	Parallel data + source annotator	from parallel data	Alignment according to source treebank annotated by source model	Target language
Source reordering	Source treebank +unannotated target sentences	Reordered from source sentence	Direct copy from source treebank	Source language
Direct unsupervised adaptation	Source treebank + unannotated target sentences	Real data	No label	Target language

# Key Takeaway

- Dependency parsing parses a sentence into a dependency tree
  - Predicts the head and arc label for each token
- Model Transfer
  - Traditional approach leverages delexicalized features
  - Model approach relies on multilingual embeddings
- Data Transfer
  - Translation-based approach aligns annotation from source language to target language
  - Unsupervised learning approach directly uses unannotated target-language corpus

# Outline

- Introduction [Dixin Jiang]
  - Motivating examples in Microsoft products
  - Problem description
  - Categorization of applications
  - Challenges
- Methodology [Dixin Jiang]
  - Model Transfer
  - Data Transfer
- Applications\*
  - Dependency Parsing [Xiubo Geng]
  - Machine Reading Comprehension [Ming Gong]
  - Grammar Error Correction [Linjun Shou]
- Summary & Future directions [Jian Pei]



Dixin Jiang

Software Technology Center at Asia (STCA) of Microsoft



Linjun Shou

Software Technology Center at Asia (STCA) of Microsoft



**Xiubo Geng**



Ming Gong



Jian Pei

Simon Fraser University

\*For more applications, please refer to our tutorial at  
The Web Conference 2021

# Cross Lingual Sequence Labelling

- Machine Reading Comprehension (MRC) as Example

Ming Gong

[migon@microsoft.com](mailto:migon@microsoft.com)

# Outline

- Sequence Labeling Application Overview
- CLMRC Definition & Benchmark Datasets
- CLMRC Baseline Approach & Challenges
- Advanced Approaches for CLMRC
- Key Take-aways

# Cross-lingual Sequence Labelling

Type	Category	Sub Category	Example
NLU	Text Classification	Single text	Domain identification, Intent detection, Sentiment classification
		Text pair	Information retrieval, Natural language inference
	Sequence Labeling	Single text	Named entity recognition, Slot tagging
		Text pair	Extractive machine reading comprehension (MRC)
	Structure Prediction	Single text	Dependency parsing, semantic role labeling
NLG	Text Generation	Token level	Spelling correction, Sentence auto completion, Grammar error correction
		Sentence level	Machine translation, Conversation, Question generation

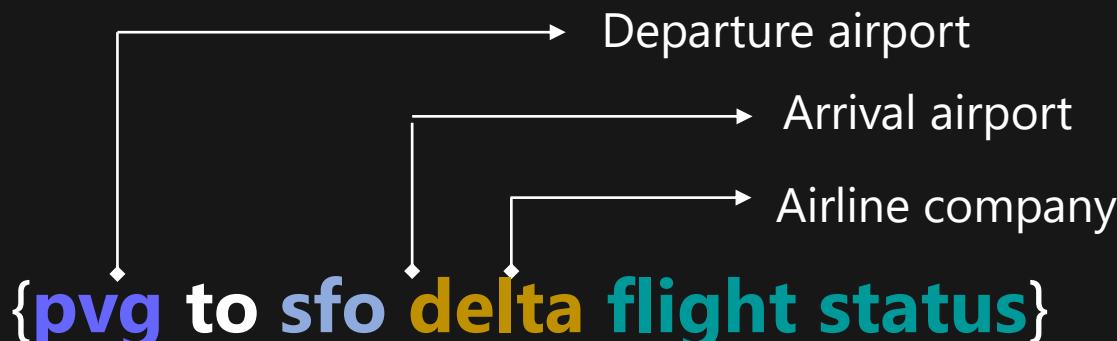
# Sequence Labeling Example Tasks

- Single text

- Given a text and a set of labels, predict the label for each token in the text

Luke Rawlence PERSON joined Aiimi ORG as a data scientist in Milton Keynes PLACE, after finishing his computer science degree at the University of Lincoln. ORG

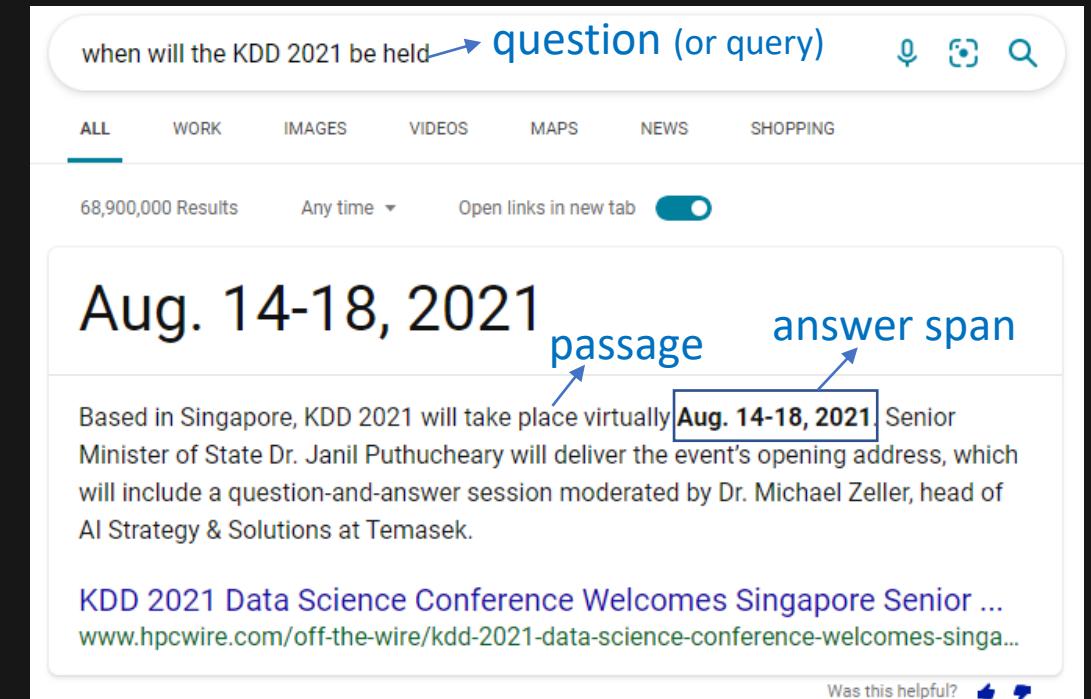
Application1: Named Entity Recognition (NER)



Application2: Slot Tagging

- Text pair

- Given two texts, label the tokens in one text w.r.t. the other text.



Application3: Machine reading comprehension (MRC)

# Cross-lingual MRC (CLMRC): Definition & Dataset

- Benchmark Datasets

- MLQA*: Patrak et al. MLQA: Evaluating Cross-lingual Extractive Question Answering. ACL 2020.
- TydiQA*: Clark, J. et al. TyDi QA: A Benchmark for Information-Seeking Question Answering in Typologically Diverse Languages. TACL 2020.
- XQuAD*: Artetxe, M. et al. On the Cross-lingual Transferability of Monolingual Representations. ACL 2020.

En	During what time period did the Angles migrate to Great Britain?
The name "England" is derived from the Old English name Englaland [...] The Angles were one of the Germanic tribes that settled in Great Britain during the <b>Early Middle Ages</b> . [...] The Welsh name for the English language is "Saesneg"	
De	Während welcher Zeitperiode migrierten die Angeln nach Großbritannien?
Der Name England leitet sich vom altenglischen Wort Englaland [...] Die Angeln waren ein germanischer Stamm, der das Land im <b>Frühmittelalter</b> besiedelte. [...] ein Verweis auf die weißen Klippen von Dover.	
Ar	في أي حقبة زمنية هاجر الأنجل إلى بريطانيا الظماء؟
والتي تعني "أرض الأنجل". والأنجل كانت واحدة، Englaland يشتهر اسم "الإنجلترا" من الكلمة الإنجليزية القديمة من القبائل الجرمانية التي استقرت في إنجلترا <b>خلال أوائل العصور الوسطى</b> . [...] وقد سماها العرب قديماً الإنكلترا	
Vi	Trong khoảng thời gian nào người Angles di cư đến Anh?
Tên gọi của Anh trong tiếng Việt bắt nguồn từ tiếng Trung. [...] Người Angle là một trong những bộ tộc German định cư tại Anh trong <b>Thời đầu Trung Cổ</b> . [...] dường như nó liên quan tới phong tục gọi người German tại Anh là Angli Saxones hay Anh - Sachsen.	
En	What are the names given to the campuses on the east side of the land the university sits on?
The campus is in the residential area of Westwood [...] The campus is informally divided into <b>North Campus and South Campus</b> , which are both on the eastern half of the university's land. [...] The campus includes [...] a mix of architectural styles.	
Es	¿Cuáles son los nombres dados a los campus ubicados en el lado este del recinto donde se encuentra la universidad?
El campus incluye [...] una mezcla de estilos arquitectónicos. Informalmente está dividido en <b>Campus Norte y Campus Sur</b> , ambos localizados en la parte este del terreno que posee la universidad. [...] El Campus Sur está enfocado en la ciencias físicas [...] y el Centro Médico Ronald Reagan de UCLA.	
Zh	位于大学占地东半部的校园名称是什么？
整个校园被不正式地分为 <b>南北两个校园</b> ，这两个校园都位于大学占地的东半部。北校园是原校园的中心，建筑以意大利文艺复兴时代建筑闻名，其中的威尔图书馆 (Powell Library) 成为好莱坞电影的最佳拍摄场景。 [...] 这个广场曾在许多电影中出现。	
Hi	विश्वविद्यालय जहाँ स्थित है, उसके पूर्वी दिशा में बने परिसरों को क्या नाम दिया गया है?
जब 1919 में यूटीएलए ने अपना नया परिसर खोला, तब इसमें चार इमारतें थी। [...] परिसर अनोपचारिक रूप से <b>उत्तरी परिसर और दक्षिणी परिसर</b> में विभाजित है, जो दोनों विश्वविद्यालय की जमीन के पूर्वी हिस्से में स्थित हैं। [...] दक्षिणी परिसर में भौतिक विज्ञान, जीव विज्ञान, इंजीनियरिंग, मनोविज्ञान, गणितीय विज्ञान, सभी रसायन से संबंधित क्षेत्र और यूएलसीए मेडिकल सेंटर स्थित हैं।	

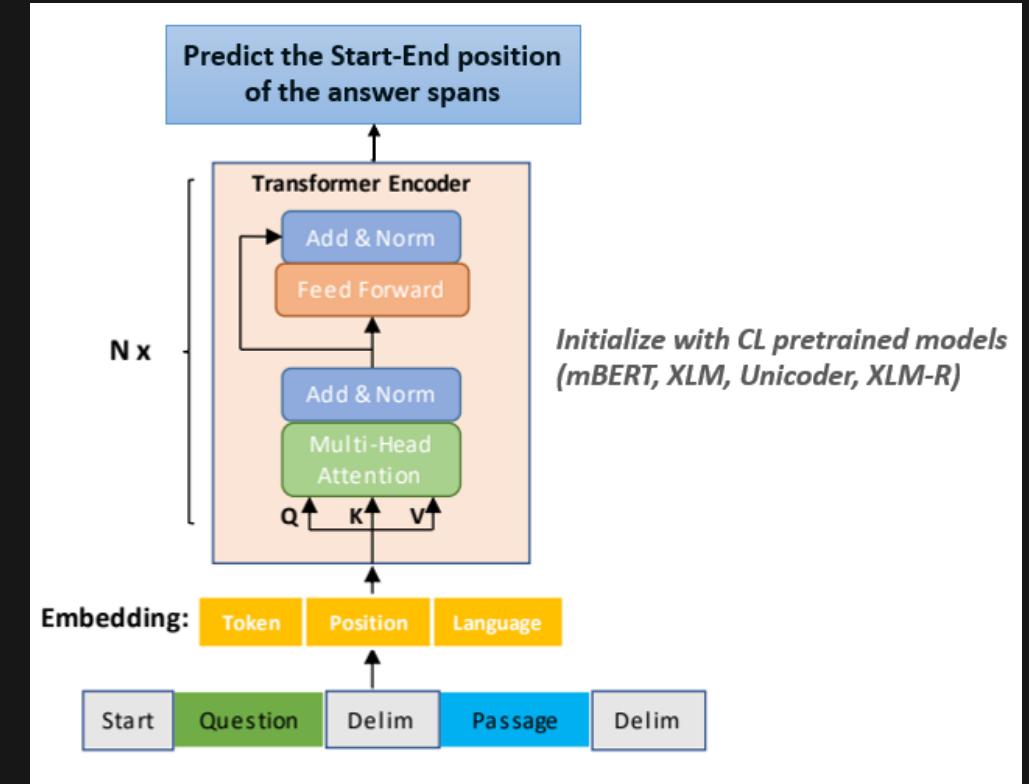
Examples from MLQA

## Typical Experiment Setting

- There are only English labeled data for training.
- For non-EN languages, there are only test sets in specific languages.

# Baseline Approach for CLMRC

- Zero-Shot Train
  - CL Pretrained Model + *English Labeled Data*
- Translation Train
  - CL Pretrained Model + *Translated Data in Target Languages*
  - CL Pretrained Model + *English Data + Translated Data in Target Languages*



# Challenges of Baseline Approach for CLMRC

**A**

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

**B**

Machine Translation Data Quality

1. MT noises impact more for pair-wise tasks
2. Answer span alignment is challenging after MT

**C**

Machine translation data is not available in some cases or quality is bad for tail languages

**D**

Answer span boundary detection leads to the major errors of CLMRC model prediction

# Challenges of Baseline Approach for CLMRC (A)

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

## Machine Translation Data Quality

- MT noises impact more for pair-wise tasks
- Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction

Classification vs Sequence Labeling using CL Pretrained models

Language	EM	MRC	NLI	
		Gap to en		
en	62.4	—	85.0	—
es	49.8	-12.6	78.9	-6.1
de	47.6	-14.8	77.8	-7.2
ar	36.3	-26.1	73.1	-11.9
hi	27.3	-35.1	69.6	-15.4
vi	41.8	-20.6	76.1	-8.9
zh	39.6	-22.8	76.5	-8.5

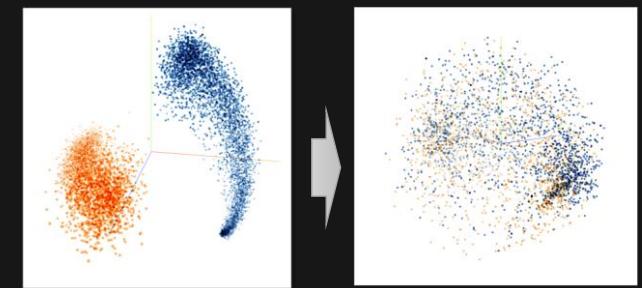


Table 1: The gap between target languages and English on Machine Reading Comprehension (MRC) (Lewis et al., 2019) is significantly larger than sentence level classification task like Natural Language Inference (NLI) (Conneau et al., 2018). In this experiment, we fine-tune XLM (Conneau and Lample, 2019) on English and directly test on other languages.

*For phrase boundary task, CL pretraining models **don't have sufficient transfer capability**  
→ **Require better semantic alignment***

Fei Yuan, Linjun Shou, Xuanyu Bai, Ming Gong, Yaobo Liang, Nan Duan, Yan Fu, Dixin Jiang. Enhancing Answer Boundary Detection for Multilingual Machine Reading Comprehension. ACL, 2020.

# Challenges of Baseline Approach for CLMRC (B-1)

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

Machine Translation Data Quality

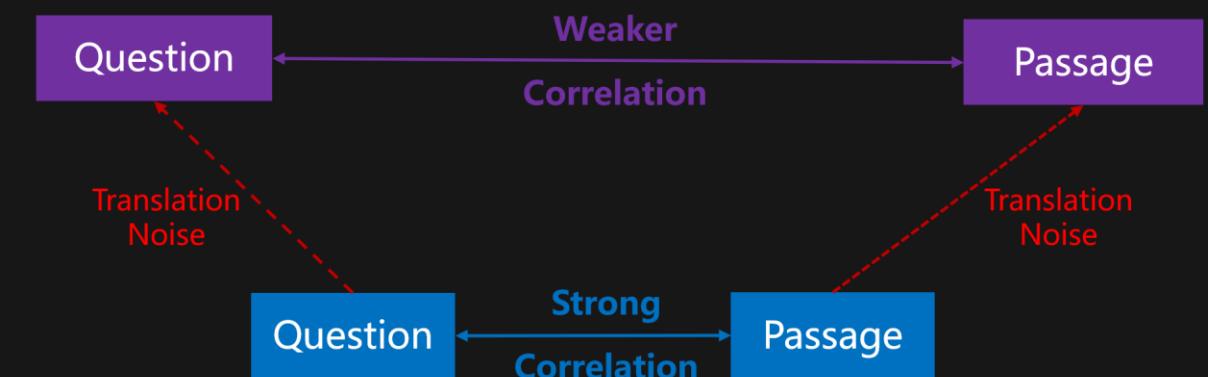
1. **MT noises impact more for pair-wise tasks**
2. Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction



# Challenges of Baseline Approach for CLMRC (B-2)

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

## Machine Translation Data Quality

1. MT noises impact more for pair-wise tasks
2. **Answer span alignment is challenging after MT**

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction

Question (EN): {Where is the Earth during the full moon?}

### Passage (EN)

The full moon is the lunar phase when the Moon appears fully illuminated from Earth's perspective. This occurs when Earth is located between the Sun and the Moon.

### Passage (Translate to DE)

Der Vollmond ist die Mondphase, in der der Mond vollständig aus der Perspektive der Erde erleuchtet erscheint. Dies geschieht, wenn sich die Erde zwischen Sonne und Mond befindet.

Answer span translate to DE:  
befindet sich zwischen der Sonne und dem Mond

Miss alignment with the correct span in DE passage

# Challenges of Baseline Approach for CLMRC (**C**)

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

Machine Translation Data Quality

- MT noises impact more for pair-wise tasks
- Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction

# Challenges of Baseline Approach for CLMRC (D)

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

## Machine Translation Data Quality

- MT noises impact more for pair-wise tasks
- Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction

Error analysis on MLQA for zero-shot CLMRC model using XLM (Lewis et al., 2019) showed that major errors come from answer spans *partially overlap with golden span*.

[Question]	who were the kings of the southern kingdom
[Passage]	In the southern kingdom there was only one dynasty, that of king David, except usurper Athaliah from the northern kingdom, who by marriage, []
[Answer - ground truth]	king David
[Answer - model predication]	David, except usurper Athaliah
[Question]	What is the suggested initial does dosage of chlordiazepoxide
[Passage]	If the drug is administered orally, the suggested initial dose is 50 to 100 mg, to be followed by repeated doses as needed until agitation is controlled up to 300 mg per day. []
[Answer - ground truth]	50 to 100 mg
[Answer - model predication]	100 mg

Table 2: Bad answer boundary detection cases of multilingual MRC model.

	#Test	#Error	Boundary error
es	4, 054	955	66.4%
de	5,390	1,648	75.2%

Fei Yuan, Linjun Shou, Xuanyu Bai, Ming Gong, Yaobo Liang, Nan Duan, Yan Fu, Dixin Jiang. Enhancing Answer Boundary Detection for Multilingual Machine Reading Comprehension. ACL, 2020.

# Advanced Approaches for CLMRC

## Challenges of the Baseline Approach

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks



B

Machine Translation Data Quality

1. MT noises impact more for pair-wise tasks
2. Answer span alignment is challenging after MT



C

Machine translation data is not available in some cases or quality is bad for tail languages



D

Answer span boundary detection leads to the major errors of CLMRC model prediction



## Advanced Approaches

### Model transfer

- New Mix-MRC task
- Dual encoders
- Adversarial training



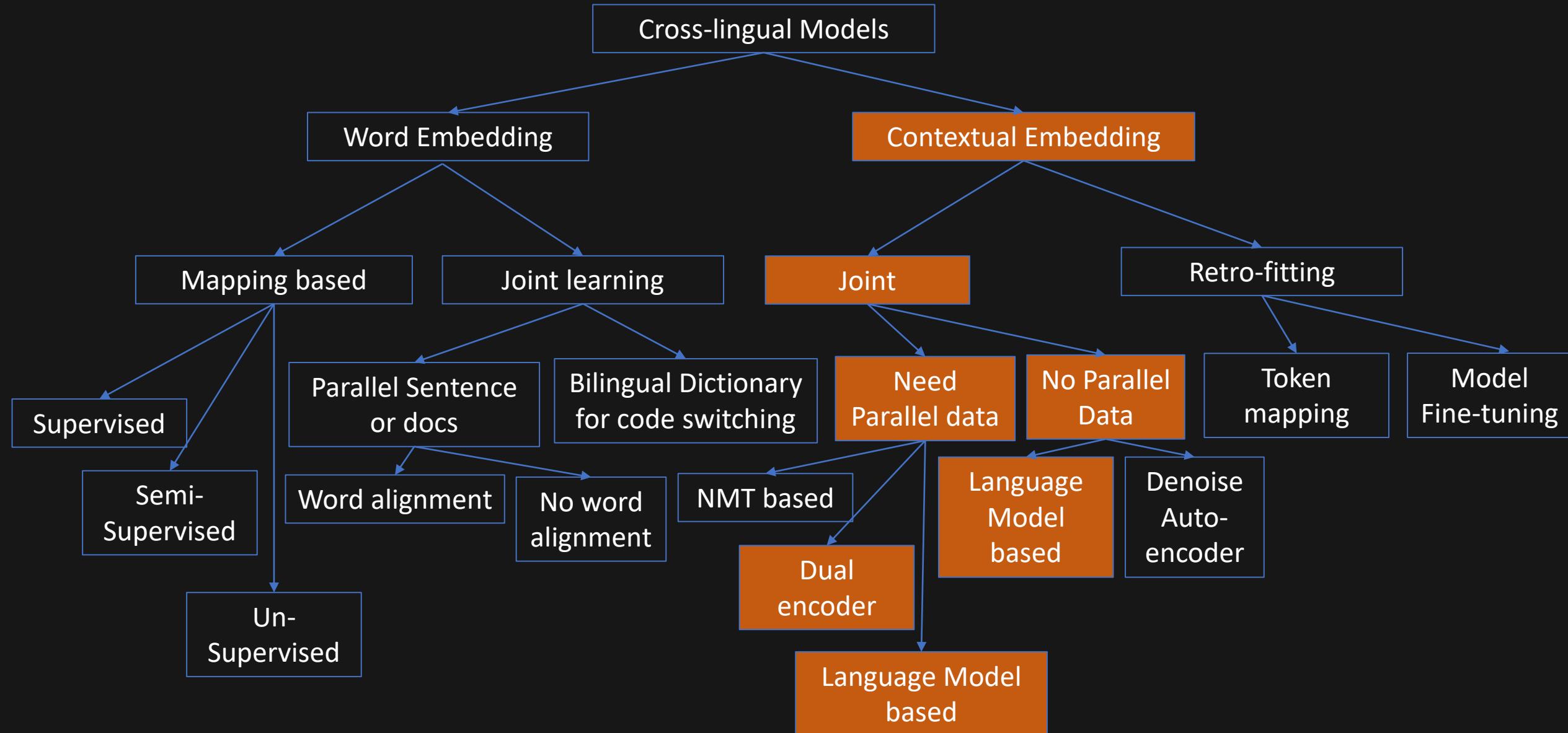
MT data

### Data transfer

- Multiple teacher models + Language branch MT data
- Instance level best teacher selection through reinforcement learning
- Back-translation & Verification to improve answer span alignment

- Unlabeled language specific data
- Generated language specific data

- New LAKM task + Language specific phrase knowledge data
- Calibration network for primary boundary correction



# Advanced Approaches for CLMRC

## Challenges of the Baseline Approach

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks



B

### Machine Translation Data Quality

- MT noises impact more for pair-wise tasks
- Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction

## Advanced Approaches

### Model transfer

- New Mix-MRC task
- Dual encoders
- Adversarial training



MT data

### Data transfer

# Mixed MRC Task with MT Data for Semantic Transfer

- **Motivation:** if human truly understand two languages, they could well perform MRC task in mixed languages. Could we enhance model semantic transfer capability using the MixMRC task?

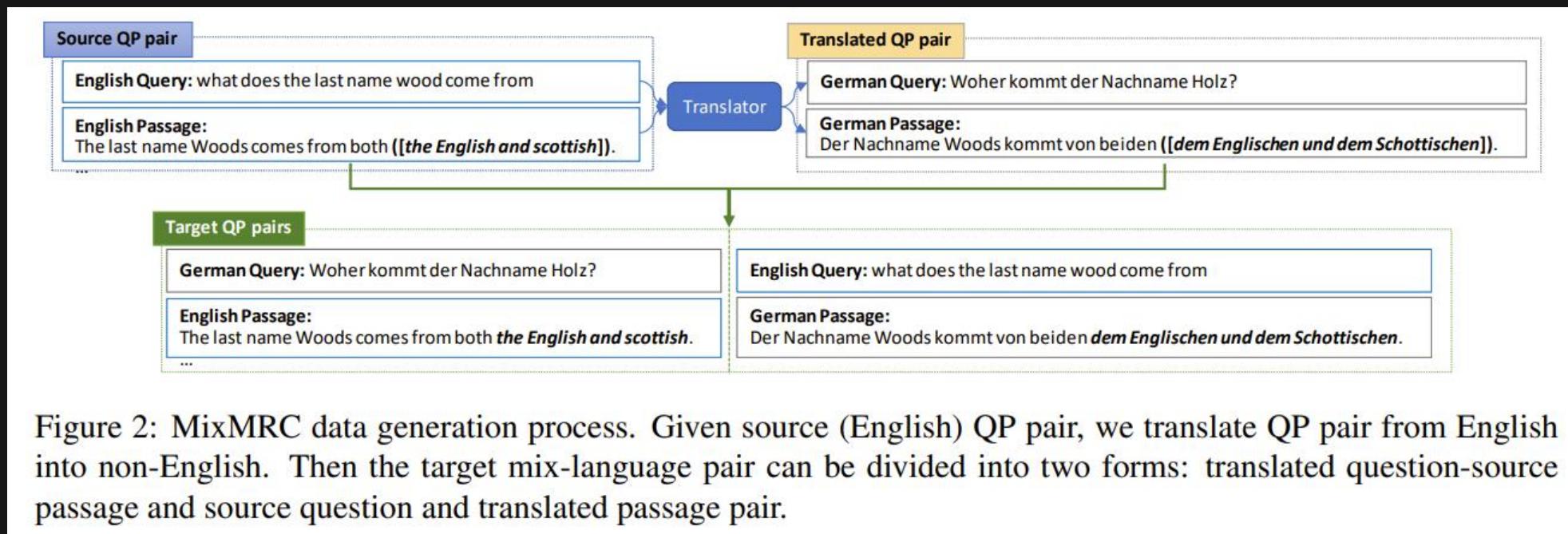


Figure 2: MixMRC data generation process. Given source (English) QP pair, we translate QP pair from English into non-English. Then the target mix-language pair can be divided into two forms: translated question-source passage and source question and translated passage pair.

# Mixed MRC Task with MT Data for Semantic Transfer

- Translation Train Results

<b>Model</b>	<b>Methods</b>	<b>MLQA (EM / F1)</b>			<b>MTQA (EM / F1)</b>		
		en	es	de	en	fr	de
M-BERT	Lewis et al. (2019)	65.2 / 77.7	37.4 / 53.9	47.5 / 62.0	-	-	-
	Baseline	65.4 / 79.0	50.4 / 68.5	46.2 / 60.6	67.0 / 86.9	52.9 / 78.2	59.8 / 81.4
	LAKM	<b>66.9</b> / 80.1	51.5 / 69.5	49.9 / 64.4	<b>68.8</b> / 87.6	56.8 / 78.8	62.4 / 81.9
	mixMRC	65.4 / 79.4	50.5 / 69.1	49.1 / 64.0	67.9 / 86.8	56.4 / 77.8	62.4 / 81.0
	mixMRC + LAKM	64.7 / 79.2	<b>52.1</b> / 70.4	<b>50.9</b> / 65.6	68.6 / 87.0	<b>57.5</b> / 78.5	<b>62.9</b> / 81.3
XLM	Lewis et al. (2019)	62.4 / 74.9	47.8 / 65.2	46.7 / 61.4	-	-	-
	Baseline	64.1 / 77.6	50.4 / 68.4	47.4 / 62.0	67.1 / 86.8	51.5 / 75.8	61.6 / 81.3
	LAKM	<b>64.6</b> / 79.0	52.2 / 70.2	50.6 / 65.4	<b>68.3</b> / 87.3	52.5 / 75.9	61.9 / 81.2
	mixMRC	63.8 / 78.0	52.1 / 69.9	49.8 / 64.8	66.5 / 85.9	52.9 / 75.0	62.1 / 80.5
	mixMRC + LAKM	64.4 / 79.1	<b>52.2</b> / 70.3	<b>51.2</b> / 66.0	68.2 / 86.8	<b>53.6</b> / 75.9	<b>62.5</b> / 80.9

Table 7: Experimental results on MLQA and MTQA dataset under translation condition (%).

# Dual Language Encoders for Semantic Transfer - DualBERT

- Simultaneously model the training data in both source and target language to better exploit the relations between two languages

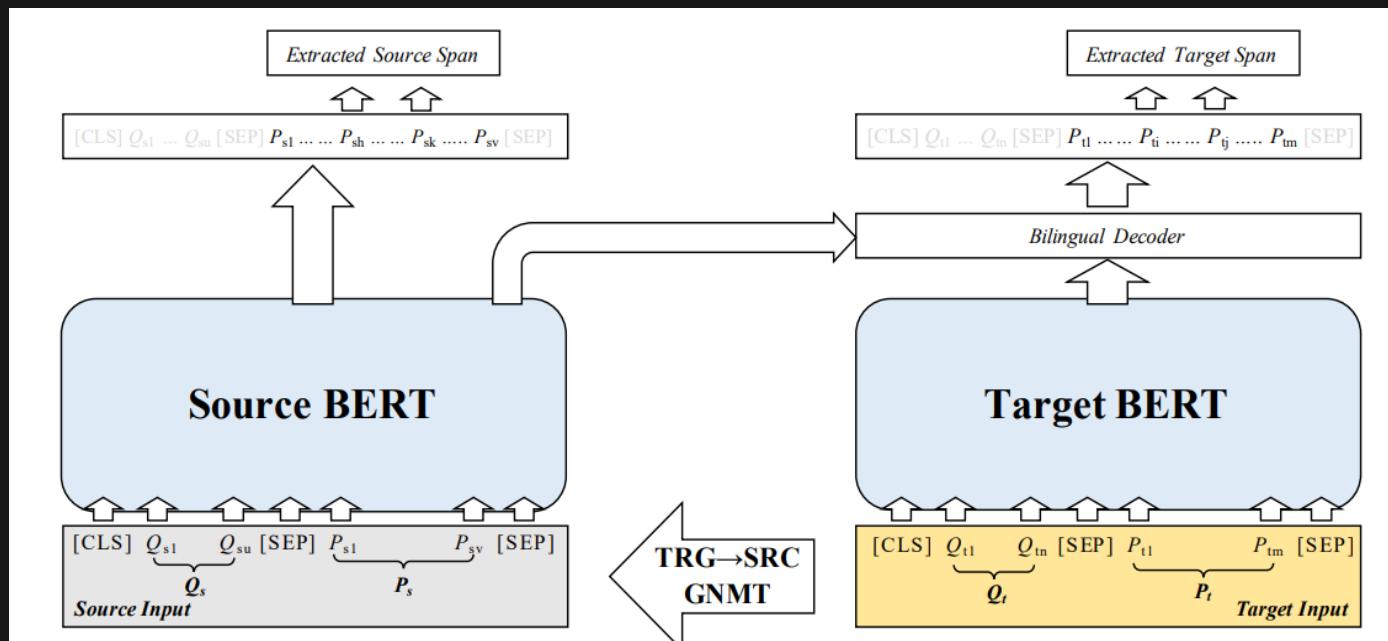
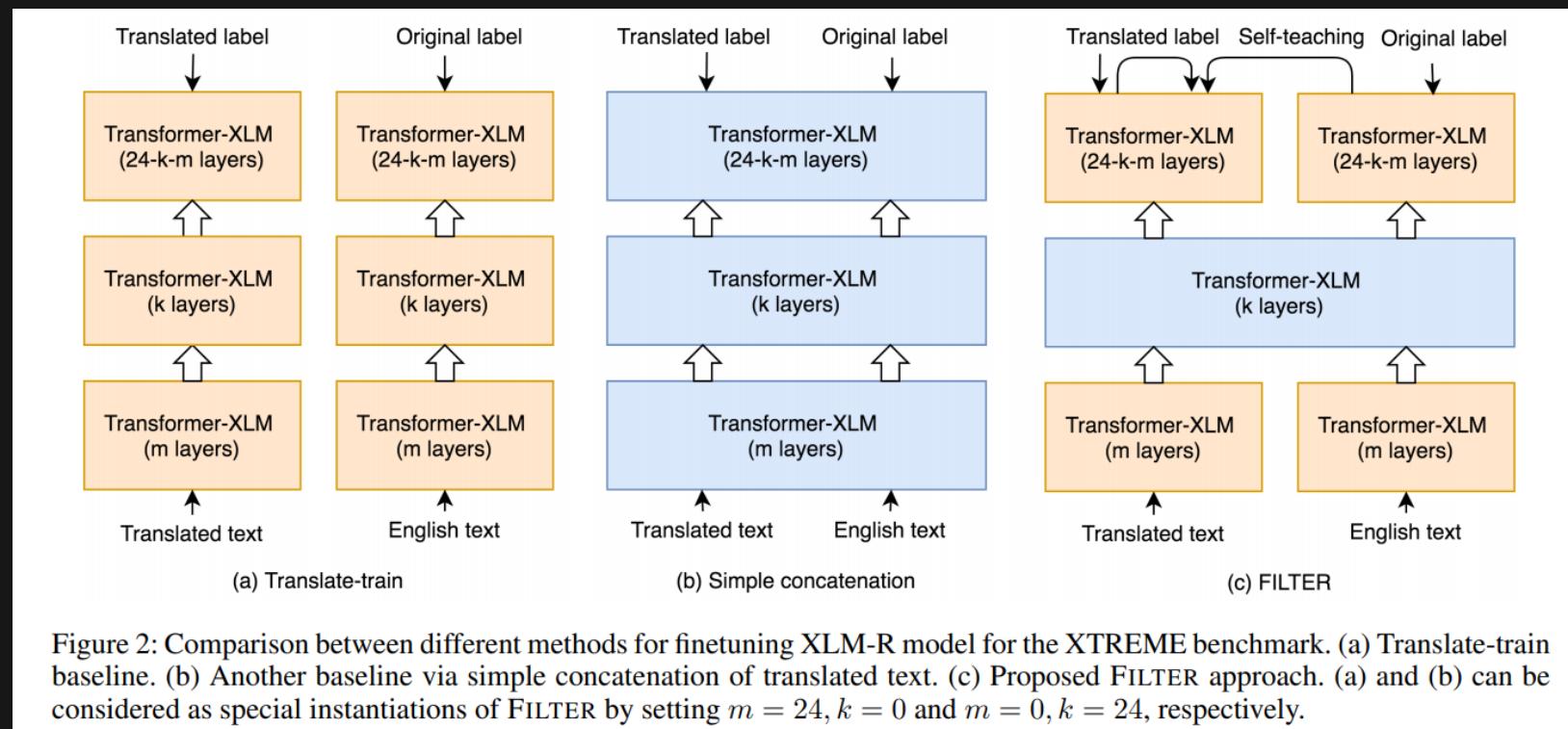


Figure 2: System overview of the Dual BERT model for cross-lingual machine reading comprehension task.

#	System	CMRC 2018				DRCD			
		Dev	Test	Dev	Test	Dev	Test	Dev	Test
		EM	F1	EM	F1	EM	F1	EM	F1
6	BERT <sub>B<sub>cn</sub></sub>	63.6	83.9	67.8	86.0	18.4	42.1	83.4	90.1
7	BERT <sub>B<sub>mul</sub></sub>	64.1	84.4	68.6	86.8	18.6	43.8	83.2	89.9
8	<b>Dual BERT</b>	<b>65.8</b>	<b>86.3</b>	<b>70.4</b>	<b>88.1</b>	<b>23.8</b>	<b>47.9</b>	<b>84.5</b>	<b>90.8</b>

# Dual Language Encoders for Semantic Transfer - FILTER

- Simultaneously model the training data in both source and target language to better exploit the relations between two languages



# Dual Language Encoders for Semantic Transfer - FILTER

Model	Pair sentence		Structured prediction		XQuAD	Question answering	
	XNLI	PAWS-X	POS	NER		MLQA	TyDiQA-GoldP
Metrics	Acc.	Acc.	F1	F1	F1 / EM	F1 / EM	F1 / EM
<i>Cross-lingual zero-shot transfer (models are trained on English data)</i>							
mBERT	65.4	81.9	70.3	62.2	64.5 / 49.4	61.4 / 44.2	59.7 / 43.9
XLM	69.1	80.9	70.1	61.2	59.8 / 44.3	48.5 / 32.6	43.6 / 29.1
XLM-R	79.2	86.4	72.6	65.4	76.6 / 60.8	71.6 / 53.2	65.1 / 45.0
InfoXLM	81.4	-	-	-	- / -	73.6 / 55.2	- / -
X-STILTs	80.0	87.9	74.4	64.0	78.7 / 63.3	72.4 / 53.7	<b>76.0 / 59.5</b>
<i>Translate-train (models are trained on English training data and its translated data on the target language)</i>							
mBERT	74.0	86.3	-	-	70.0 / 56.0	65.6 / 48.0	55.1 / 42.1
mBERT, multi-task	75.1	88.9	-	-	72.4 / 58.3	67.6 / 49.8	64.2 / 49.3
XLM-R, multi-task (Ours)	82.6	90.4	-	-	80.2 / 65.9	72.8 / 54.3	66.5 / 47.7
FILTER (Ours)	83.6	91.2	75.5	66.7	82.3 / 67.8	75.8 / 57.2	68.1 / 49.7
FILTER + Self-Teaching (Ours)	<b>83.9</b>	<b>91.4</b>	<b>76.2</b>	<b>67.7</b>	<b>82.4 / 68.0</b>	<b>76.2 / 57.7</b>	68.3 / 50.9

# Multi-lingual Adversarial Training for Semantic Transfer

- Adversarial Training to bring the embeddings of different languages closer to each other to achieve effective cross-lingual transfer

$$\mathcal{L}_{QA} + \mathcal{L}_{adv} = \mathcal{L}_{QA} - \sum_{l=1}^L KL[U(\mathbf{g}_l) || \log \mathbf{p}_l]$$
$$\mathcal{L}_D = - \sum_{l=1}^L \mathbb{1}(\mathbf{g}_l) \underline{\log \mathbf{p}_l} \text{ (Discriminator)}$$

Adversarial Loss

# Multi-lingual Adversarial Training for Semantic Transfer

- Bringing the embeddings of different languages closer to each other to achieve effective cross-lingual transfer

Model	Method	MLQA Languages (G-XLT)							G-XLT	XLT
		ar	de	en	es	hi	vi	zh		
MBERT <sub>QA</sub>	ZS	46.9	51.4	60.2	55.0	47.0	52.0	49.7	51.7 ( $\pm 0.4$ )	61.7 ( $\pm 0.3$ )
Trans	T(Q)	<b>53.8</b>	60.8	<b>73.5</b>	65.4	<b>53.2</b>	63.2	56.7	60.9 ( $\pm 0.2$ )	<b>64.9</b> ( $\pm 0.2$ )
	T(C)	44.8	51.7	62.0	57.6	42.7	55.8	50.8	52.2 ( $\pm 1.0$ )	58.5 ( $\pm 0.9$ )
	T(Q+C)	48.9	58.4	70.4	63.6	46.8	61.4	54.4	57.7 ( $\pm 0.1$ )	64.3 ( $\pm 0.0$ )
	T(ALL)	52.6	<b>61.1</b>	<b>73.8</b>	<b>66.3</b>	50.6	<b>64.8</b>	<b>58.2</b>	<b>61.1</b> ( $\pm 0.1$ )	64.2 ( $\pm 0.2$ )
AT	(en-zh)	50.5	56.7	68.0	60.8	50.7	57.4	51.5	56.5 ( $\pm 0.1$ )	62.8 ( $\pm 0.1$ )
	(en-all)	<b>54.1</b>	<b>61.1</b>	<b>73.6</b>	<b>65.5</b>	<b>54.2</b>	<b>63.4</b>	<b>56.8</b>	<b>61.2</b> ( $\pm 0.1$ )	<b>65.2</b> ( $\pm 0.1$ )

# Advanced Approaches for CLMRC

## Challenges of the Baseline Approach

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

Machine Translation Data Quality

1. MT noises impact more for pair-wise tasks
2. Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction



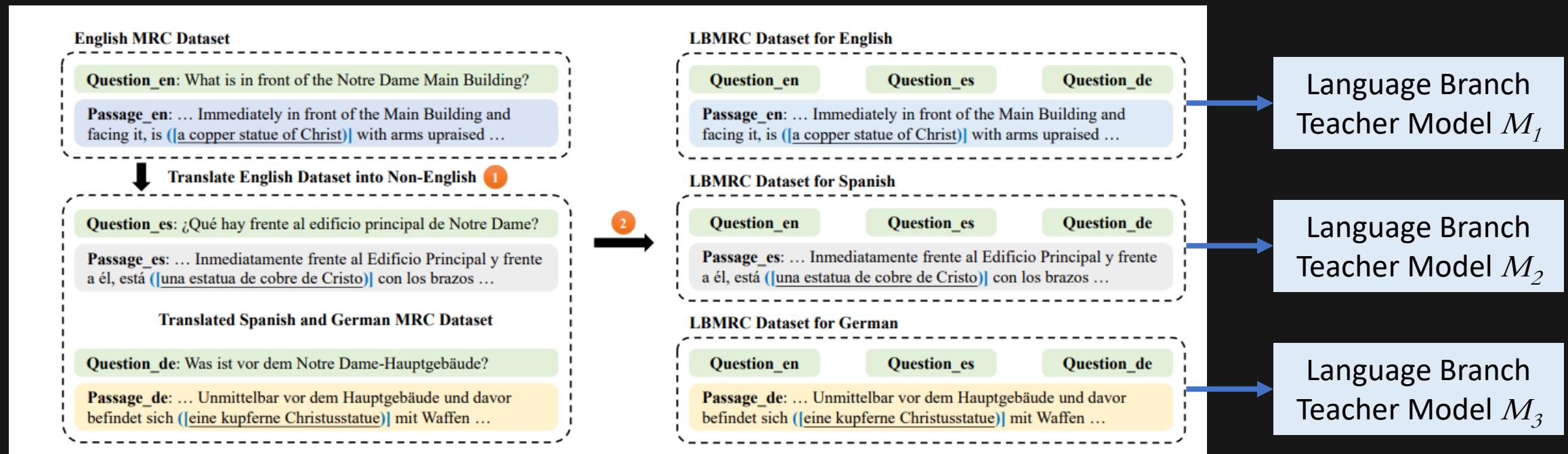
### Model transfer

- Multiple teacher models + Language branch MT data
- Instance level best teacher selection through reinforcement learning
- Back-translation & Verification to improve answer span alignment

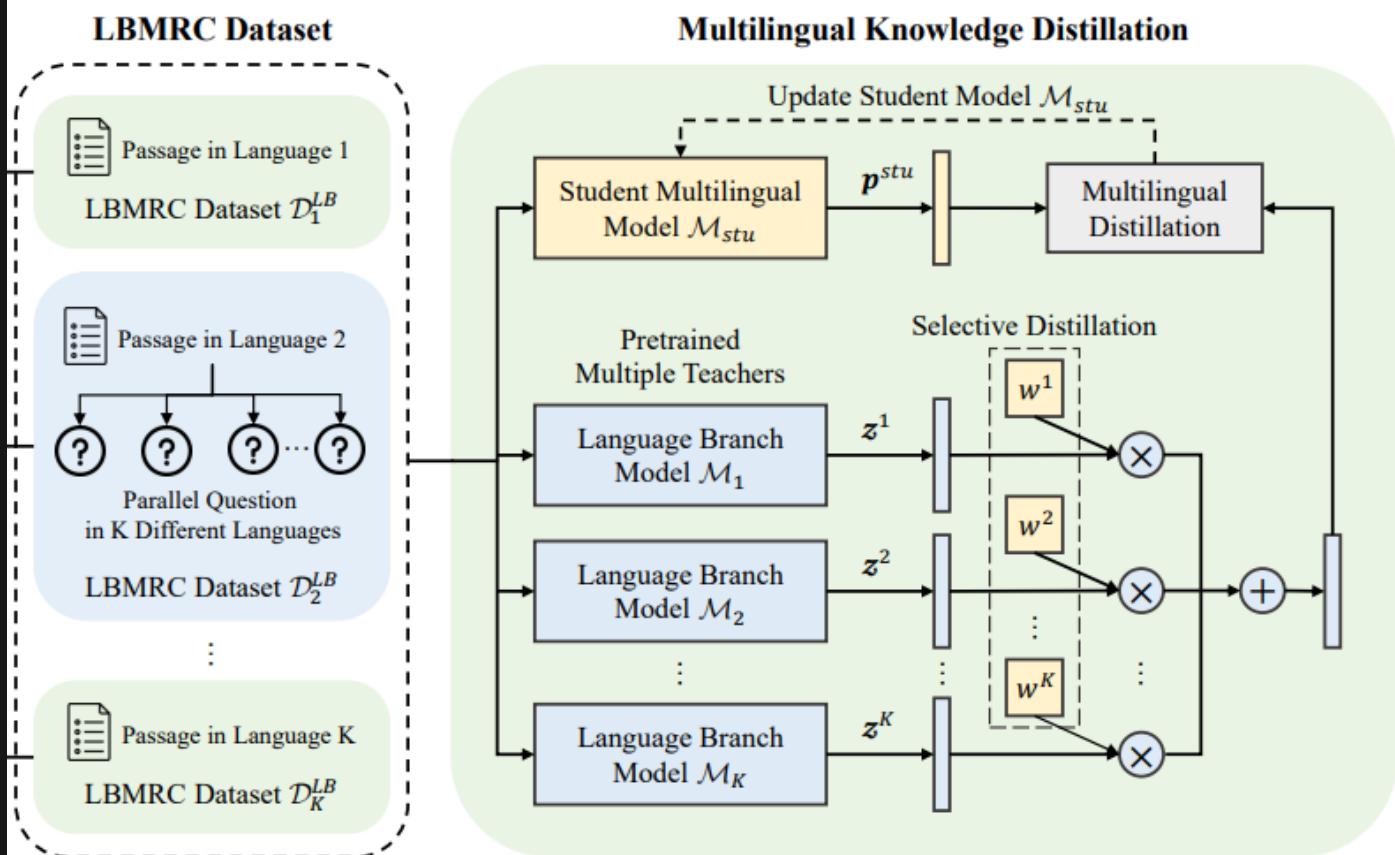
### Data transfer

# LBMRC: Language Branch Knowledge Distillation for Reducing MT Noise

- **Motivation:** How to reduce the negative impact of the noisy MT data?  
→ Leverage multiple teacher models to overcome the MT noise



# LBMRC: Language Branch Knowledge Distillation for Reducing MT Noise



Methods	MLQA (EM / F1)					
	en	es	de	ar	hi	vi
Lewis <sup>¶</sup>	62.4 / 74.9	47.8 / 65.2	46.7 / 61.4	34.4 / 54.0	33.4 / 50.7	39.4 / 59.3
Baseline	63.4 / 77.3	49.7 / 68.2	48.5 / 63.7	37.5 / 56.9	37.0 / 54.3	42.4 / 63.5
LAKM	64.6 / 79.0	52.2 / 70.2	50.6 / 65.4	-	-	-
mixMRC	63.8 / 78.0	52.1 / 69.9	49.8 / 64.8	38.5 / 58.4	40.1 / 57.1	45.2 / 66.2
mixMRC + LAKM	64.4 / 79.1	52.2 / 70.3	51.2 / 66.0	-	-	-
Ours-hyper	<b>64.8 / 79.3</b>	53.9 / 71.8	52.1 / 66.8	<b>40.4 / 60.0</b>	42.8 / 59.8	46.1 / 67.2
Ours-imp	64.7 / 79.2	<b>54.3 / 72.0</b>	<b>52.4 / 66.9</b>	40.1 / 59.9	<b>42.9 / 59.9</b>	<b>46.5 / 67.5</b>

Table 1: EM and F1 score of 6 languages on the MLQA dataset. The left 3 languages (en, es, de) are under translation condition while the right part (ar, hi, vi) results are under the zero-shot transfer method. The results with <sup>¶</sup> are adopted from Lewis et al. (2019).

Methods	XQuAD (EM / F1)					
	ar	hi	vi	el	ru	tr
Baseline	43.2 / 62.6	46.0 / 63.1	48.7 / 70.4	49.4 / 68.5	55.2 / 72.3	44.1 / 63.1
mixMRC	42.4 / 63.6	50.0 / 66.2	52.7 / 72.6	51.1 / 72.1	58.7 / 75.9	47.8 / 65.8
Ours-hyper	<b>44.5 / 65.0</b>	52.0 / 67.4	55.5 / 74.6	52.2 / 73.1	59.3 / 76.6	<b>50.8 / 68.3</b>
Ours-imp	44.0 / 64.6	<b>52.5 / 67.9</b>	<b>55.6 / 74.9</b>	<b>52.4 / 73.3</b>	<b>59.6 / 76.6</b>	50.2 / 67.7

Table 3: EM and F1 score of 6 languages on the XQuAD dataset under the zero-shot transfer setting.

# Reinforced Multi-Teacher Selection for Knowledge Distillation (RL-KD)

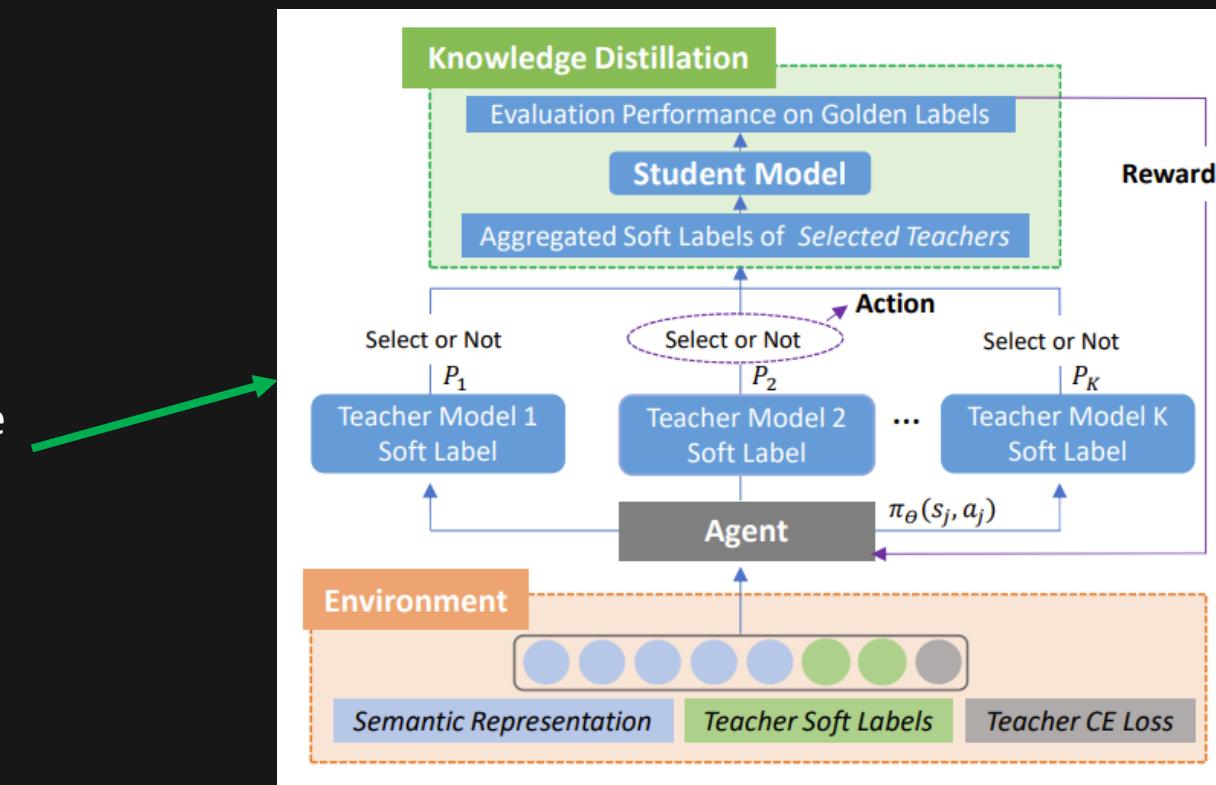
- **Motivation:** we have multiple teachers, how to select best teacher model for each case to generate soft label for better KD?

## Baseline Methods:

- Vanilla KD from teacher ensemble with static weights

## New Approach:

- RL method to select from multiple teacher models at instance level



# Reinforced Multi-Teacher Selection for Knowledge Distillation (RL-KD)

- **Motivation:** we have multiple teachers, how to select best teacher model for each case to generate soft label for better KD?

## Baseline Methods:

- Vanilla KD from teacher ensemble with static weights

## New Approach:

- RL method to select from multiple teacher models at instance level

Teacher	Student	Strategy	QQP	MRPC	MNLI-(mm/m)	RTE	QNLI	SST-2	AVG.
-	BERT <sub>3</sub>	FT	88.3	72.1	74.6 / 74.8	60.7	83.3	85.9	77.1
BERT <sub>12</sub>	BERT <sub>3</sub>	V-KD	88.4	74.0	75.3 / 75.6	56.7	83.7	87.5	77.3
Robert <sub>12</sub>	BERT <sub>3</sub>	V-KD	85.3	73.0	71.9 / 71.7	55.2	81.9	86.0	75.0
XLNet <sub>12</sub>	BERT <sub>3</sub>	V-KD	85.4	74.3	72.0 / 71.6	55.2	81.4	87.4	75.3
ALBERT <sub>12</sub>	BERT <sub>3</sub>	V-KD	88.4	74.0	76.4 / 75.6	56.7	83.7	87.5	77.5
Rand-Single-Ensemble	BERT <sub>3</sub>	V-KD	87.3	70.6	75.0 / 74.5	51.6	83.7	86.0	75.5
W-Ensemble	BERT <sub>3</sub>	V-KD	85.3	74.8	72.2 / 72.0	56.7	82.4	87.5	75.8
LR-Dev-Ensemble	BERT <sub>3</sub>	V-KD	88.6	71.8	75.4 / 75.4	54.9	84.2	86.8	76.7
Best-Single-Ensemble	BERT <sub>3</sub>	V-KD	88.6	77.2	75.1 / 74.7	56.0	84.1	87.0	77.5
Our Method ( <i>reward</i> <sub>1</sub> )	BERT <sub>3</sub>	RL-KD	<b>89.1</b>	76.0	76.9 / <b>76.8</b>	61.4	<b>85.4</b>	<b>89.1</b>	79.2
Our Method ( <i>reward</i> <sub>2</sub> )	BERT <sub>3</sub>	RL-KD	<b>89.1</b>	76.2	<b>77.4</b> / 76.3	63.5	84.8	88.5	79.4
Our Method ( <i>reward</i> <sub>3</sub> )	BERT <sub>3</sub>	RL-KD	89.0	76.7	76.7 / 75.7	<b>64.6</b>	85.3	<b>89.1</b>	<b>79.6</b>

(1)*reward*<sub>1</sub>: use the minus of ground-truth loss (CE) of student model as the reward function for teacher model selection.

(2)*reward*<sub>2</sub>: besides (1), also introduce the minus of knowledge distillation loss (DL) into the reward function.

(3)*reward*<sub>3</sub>: besides (2), also take the Accuracy metric of student model on dev set into account for better generalization

# Improve Span Alignment for MT: Back-translation + Verification

- **Motivation:** the answer spans after translation may not exist in translated passage or have boundary errors.

- Translate source span to target span and further improve boundary by:
  1. *Simple Match*
  2. *Answer Aligner*
  3. *Answer Verifier*

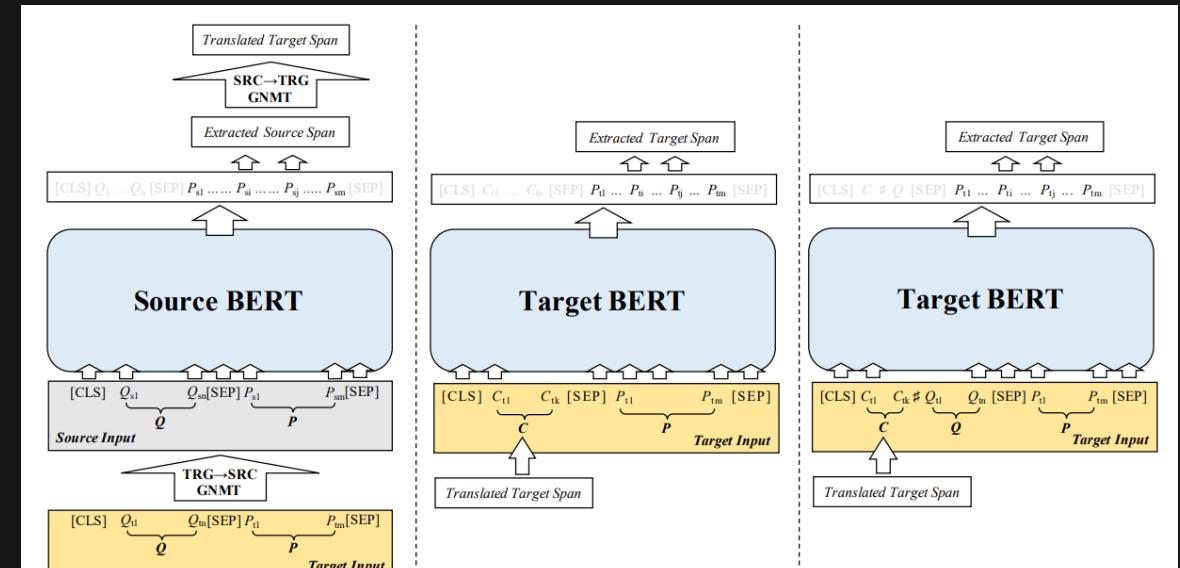


Figure 1: Back-translation approaches for cross-lingual machine reading comprehension (Left: GNMT, Middle: Answer Aligner, Right: Answer Verifier)

# Improve Span Alignment for MT: Back-translation + Verification

Experiment Results using different answer span alignment approaches

#	System	CMRC 2018						DRCD			
		Dev		Test		Challenge		Dev		Test	
		EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
Without any alignment approaches	1 GNMT+BERT <sub>SQ-B<sub>en</sub></sub>	15.9	40.3	20.8	45.4	4.2	20.2	28.1	50.0	26.6	48.9
	2 GNMT+BERT <sub>SQ-L<sub>en</sub></sub>	16.8	42.1	21.7	47.3	5.2	22.0	28.9	52.0	28.7	52.1
	3 GNMT+BERT <sub>SQ-L<sub>en</sub></sub> +SimpleMatch	26.7	56.9	31.3	61.6	9.1	35.5	36.9	60.6	37.0	61.2
	4 GNMT+BERT <sub>SQ-L<sub>en</sub></sub> +Aligner	46.1	66.4	49.8	69.3	16.5	40.9	60.1	70.5	59.5	70.7
	5 GNMT+BERT <sub>SQ-L<sub>en</sub></sub> +Verifier	64.7	84.7	68.9	86.8	20.0	45.6	83.5	90.1	82.6	89.6

# Advanced Approaches for CLMRC

## Challenges of the Baseline Approach

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

### Machine Translation Data Quality

- MT noises impact more for pair-wise tasks
- Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages



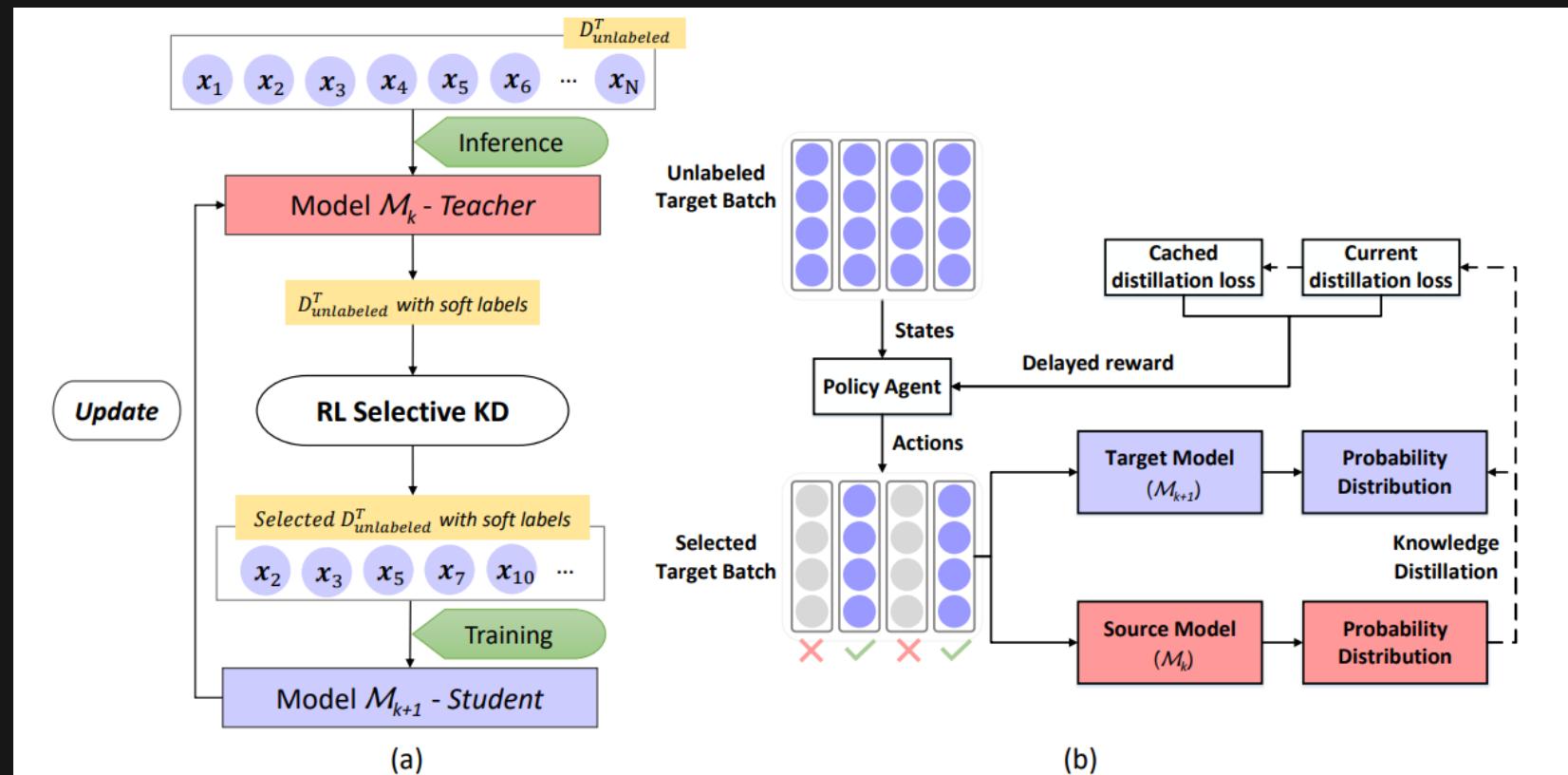
### Data transfer

- Unlabeled language specific data
- Generated language specific data

Answer span boundary detection leads to the major errors of CLMRC model prediction

# Unlabeled Language Specific Data + Iterative Knowledge Distillation

- **Motivation:** real language specific data (e.g. <question, passage> pairs) without label is relatively easy to get in many application scenarios. Could we generate good labels (e.g. answer spans) for them?



# Unlabeled Language Specific Data + Iterative Knowledge Distillation

- Experiment Results

(a) Results on CoNLL.

Method	es	nl	de	Average
Täckström [41]	59.30	58.40	40.40	52.70
Tsai et al. [6]	60.55	61.56	48.12	56.74
Ni et al. [14]	65.10	65.40	58.50	63.00
Mayhew et al. [25]	65.95	66.50	59.11	63.85
Xie et al. [16]	72.37	71.25	57.76	67.13
Wu and Dredze [43] <sup>†</sup>	74.50	79.50	71.10	75.03
Moon [36] <sup>†</sup>	75.67	80.38	71.42	75.82
Wu et al. [21]	76.75	80.44	73.16	76.78
Wu et al. [22] (mBERT)	76.94	80.89	73.22	77.02
<b>RIKD (mBERT)</b>	<b>77.84<sup>¶</sup></b>	<b>82.46<sup>¶</sup></b>	<b>75.48<sup>¶</sup></b>	<b>78.59<sup>¶</sup></b>
Wu et al. [22] (XLM-R)	78.77	80.99	74.67	78.14
<b>RIKD (XLM-R)</b>	<b>79.46<sup>¶</sup></b>	<b>81.40<sup>¶</sup></b>	<b>78.40<sup>¶</sup></b>	<b>79.75<sup>¶</sup></b>
Täckström [41] <sup>‡</sup>	61.90	59.90	36.40	52.73
Moon [36] <sup>‡</sup>	76.53	83.35	72.44	77.44
Wu et al. [22] <sup>‡</sup>	78.00	81.33	75.33	78.22

(b) Results on WikiAnn.

Method	ar	hi	zh	Average
Wu and Dredze [44]	42.30	67.60	52.90	54.27
Wu et al. [22] (mBERT)	43.12	69.54	48.12	53.59
<b>RIKD (mBERT)</b>	<b>45.96<sup>¶</sup></b>	<b>70.28<sup>¶</sup></b>	<b>50.40<sup>¶</sup></b>	<b>55.55<sup>¶</sup></b>
Pfeiffer et al. [15]	41.80	-	20.50	-
Wu and Dredze [44]	45.50	66.60	43.90	52.00
Wu et al. [22] (XLM-R)	50.91	72.48	31.14	51.51
<b>RIKD (XLM-R)</b>	<b>54.46<sup>¶</sup></b>	<b>74.42<sup>¶</sup></b>	<b>37.48<sup>¶</sup></b>	<b>55.45<sup>¶</sup></b>

# Generated Language Specific Data for Model Enhancement

- **Motivation:** Language specific plain text data (i.e. passages) is easier to get.  
Could we generate both questions and labels (i.e. answer spans) for them?

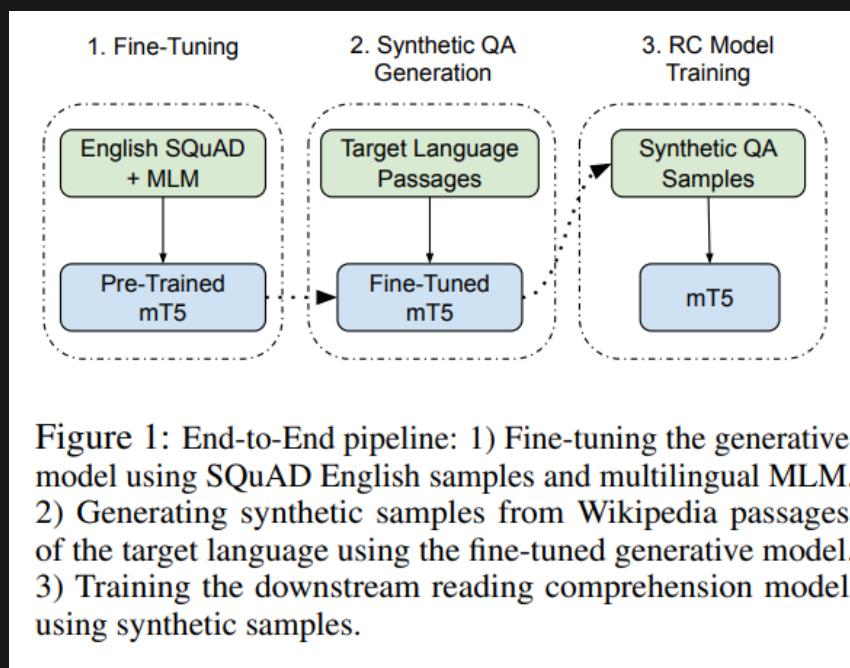


Figure 1: End-to-End pipeline: 1) Fine-tuning the generative model using SQuAD English samples and multilingual MLM. 2) Generating synthetic samples from Wikipedia passages of the target language using the fine-tuned generative model. 3) Training the downstream reading comprehension model using synthetic samples.

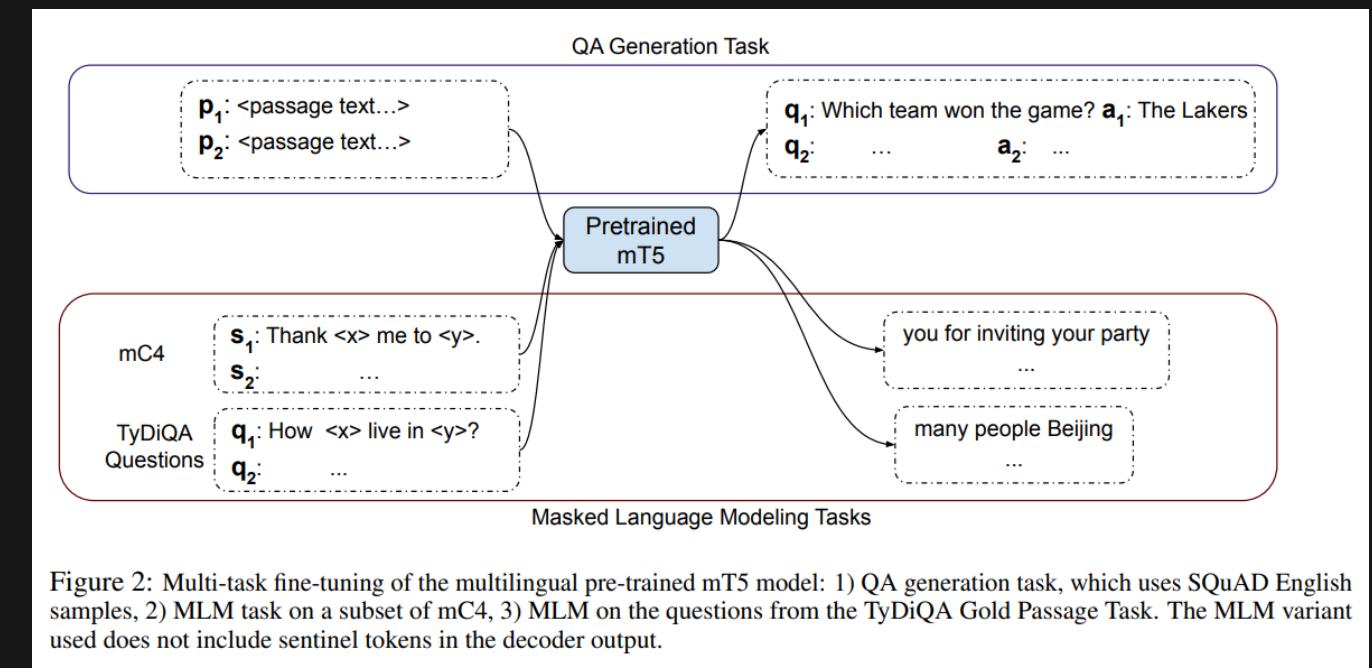


Figure 2: Multi-task fine-tuning of the multilingual pre-trained mT5 model: 1) QA generation task, which uses SQuAD English samples, 2) MLM task on a subset of mC4, 3) MLM on the questions from the TyDiQA Gold Passage Task. The MLM variant used does not include sentinel tokens in the decoder output.

- Siamak Shakeri, Noah Constant, Mihir Sanjay Kale, Linting Xue. Towards Zero-Shot Multilingual Synthetic Question and Answer Generation for Cross-Lingual Reading Comprehension. Arxiv 2021.
- Arij Riabi, Thomas Scialom, Rachel Keraron, Benoît Sagot, Djamel Seddah, Jacopo Staiano. Synthetic Data Augmentation for Zero-Shot Cross-Lingual Question Answering, arxiv 2020

# Generated Language Specific Data for Model Enhancement

- Experiment Results

Dataset	en	ar	de	el	es	hi	ru	th	tr	vi	zh	avg
SQuAD en (paper)	88.4	75.2	80.0	77.5	81.8	73.4	74.7	73.4	76.5	79.4	75.9	77.8
SQuAD en (ours)	88.6	75.0	80.4	76.5	81.6	73.9	74.1	73.8	76.2	80.1	76.4	77.4
SQuAD en + ru	88.2	76.6	81.2	79.1	82.6	76.1	77.6	72.1	75.1	78.4	77.4	<b>78.6</b>
SQuAD en + hi	<b>88.9</b>	76.7	81.5	79.4	<b>82.9</b>	73.4	78.7	<b>74.1</b>	75.0	79.1	78.0	<b>78.6</b>
SQuAD en + de	88.0	72.7	79.7	73.0	82.0	73.6	76.4	71.6	74.8	78.7	76.2	76.6
SQuAD en + ar	88.0	73.3	81.2	78.8	82.4	75.1	78.5	71.4	75.6	77.3	<b>78.2</b>	77.8
SQuAD en + es	88.2	<b>77.2</b>	<b>81.8</b>	<b>79.9</b>	81.3	<b>76.4</b>	<b>79.2</b>	72.3	75.8	79.5	77.7	<b>78.7</b>
Supervised	87.3	79.4	82.7	81.8	83.8	78.0	81.9	74.7	80.2	80.4	83.2	81.2

Table 5: Performance of fine-tuned mT5 Large models on XQuAD. *Supervised* refers to training on SQuAD en + translate-train dataset of the target language.

# Advanced Approaches for CLMRC

## Challenges of the Baseline Approach

A

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks

B

### Machine Translation Data Quality

- MT noises impact more for pair-wise tasks
- Answer span alignment is challenging after MT

C

Machine translation data is not available in some cases or quality is bad for tail languages

D

Answer span boundary detection leads to the major errors of CLMRC model prediction

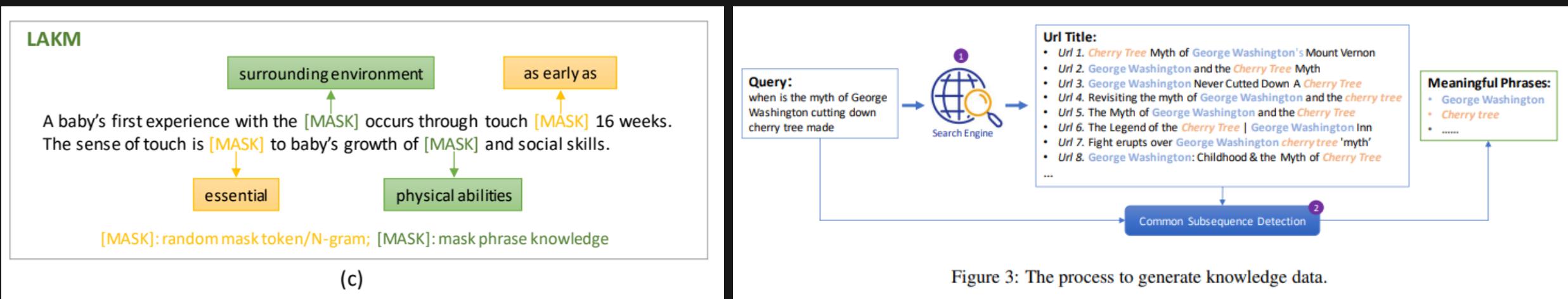
Model transfer

Data transfer

- 
- New LAKM task + Language specific phrase knowledge data
  - Calibration network for primary boundary correction

# Enhance boundary detection: LAKM task + Language specific phrase data

- **Language Agnostic Knowledge Masking** task (LAKM) targets introducing language specific boundary prior knowledge into the model.



- As an auxiliary task of the main MRC task in the fine-tuning stage (i.e. multi-task training)
- Phrase Candidates Generation (for recall)
- Phrase Filtering (for precision)

Fei Yuan, Linjun Shou, Xuanyu Bai, Ming Gong, Yaobo Liang, Nan Duan, Yan Fu, Dixin Jiang. Enhancing Answer Boundary Detection for Multilingual Machine Reading Comprehension. ACL, 2020.

# Enhance boundary detection: LAKM task + Language specific phrase data

- Experiment Results

		MLQA (EM / F1)		
		en	es	de
Baseline		65.2 / 77.7	46.6 / 64.3	44.3 / 57.9
LAKM		<b>66.8 / 80.0</b>	<b>48.0 / 65.9</b>	<b>45.5 / 60.5</b>

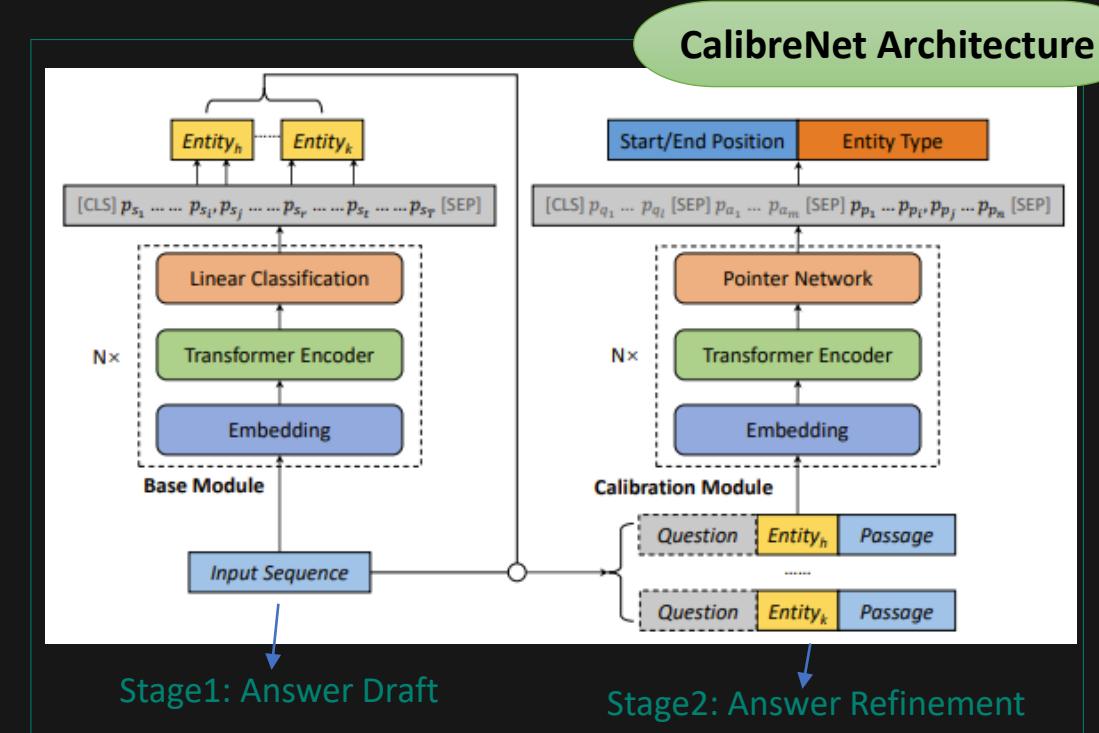
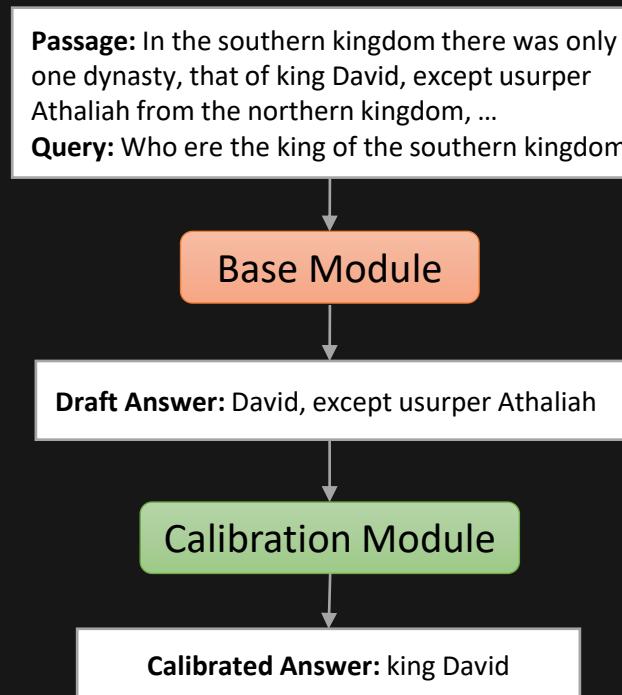
  

		MTQA (EM / F1)		
		en	fr	de
Baseline		65.8 / 86.6	41.3 / 70.9	50.7 / 76.2
LAKM		<b>67.8 / 87.2</b>	<b>44.6 / 72.1</b>	<b>54.5 / 77.8</b>

Table 9: Zero Shot experimental results on MLQA and MTQA datasets (%). We only use English MRC training data and don't use translation data.

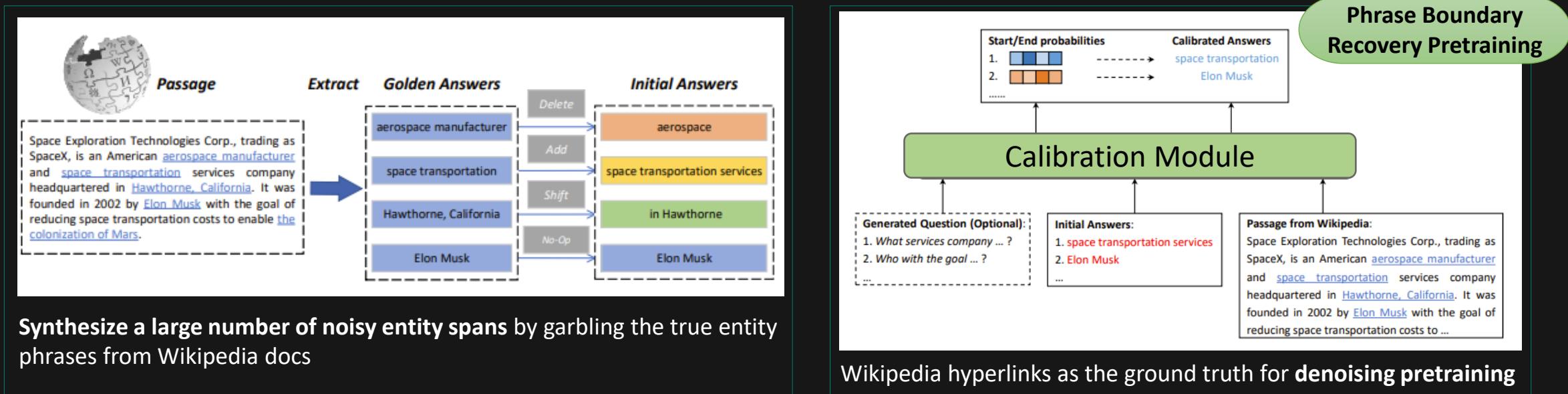
# Post-processing for Boundary Refinement: CalibreNet

- Propose a Two-stage Approach inspired by how human write articles
  - Answer draft (*base module*) → Answer refinement (*calibration module*)



# Post-processing for Boundary Refinement: CalibreNet

- Design a novel **Phrase Boundary Recovery task (PBR)** to pretrain the calibration module on *large-scale data synthesized from Wikipedia docs* in multiple languages



Results on MLQA  
(F1/EM)

Model	Methods	en	es	de	ar	hi	vi	Avg.
XLM R <sub>base</sub>	MLQA <sub>org</sub>	77.80 / 65.30	67.20 / 49.70	60.80 / 47.10	53.00 / 34.70	57.90 / 41.70	63.10 / 43.10	63.30 / 46.90
	MLQA <sub>re-imp</sub>	79.21 / 65.96	67.73 / 49.94	60.92 / 46.31	55.76 / 36.89	59.58 / 42.92	66.30 / 45.64	64.92 / 47.94
	CalibreNet (MLQA <sub>re-imp</sub> )	79.68 / 66.51	68.04 / 50.77	61.66 / 47.55	56.14 / 37.83	59.97 / 43.84	66.92 / 46.59	65.40 / 48.84

# Key Take-aways

- MRC is more challenging in language scaling and **the baseline approaches (CL pretrained model + MT data)** don't work well due to

**A**

CL pretrain models don't have sufficient transfer capability for phrase boundary tasks



## Model transfer

- New Mix-MRC task
- Dual encoders
- Adversarial training



## Data transfer

MT data

**B**

### Machine Translation Data Quality

- MT noises impact more for pair-wise tasks
- Answer span alignment is challenging after MT



- Multiple teacher models + Language branch MT data
- Instance level best teacher selection through reinforcement learning
- Back-translation & Verification to improve answer span alignment

**C**

Machine translation data is not available in some cases or quality is bad for tail languages



- Unlabeled language specific data
- Generated language specific data

**D**

Answer span boundary detection leads to the major errors of CLMRC model prediction



## Advanced Approaches

- New LAKM task + Language specific phrase knowledge data
- Calibration network for primary boundary correction

# Outline

- Introduction [Dixin Jiang]
  - Motivating examples in Microsoft products
  - Problem description
  - Categorization of applications
  - Challenges
- Methodology [Dixin Jiang]
  - Model Transfer
  - Data Transfer
- Applications\*
  - Dependency Parsing [Xiubo Geng]
  - Machine Reading Comprehension [Ming Gong]
  - Grammar Error Correction [Linjun Shou]
- Summary & Future directions [Jian Pei]



Dixin Jiang

Software Technology Center at Asia (STCA) of Microsoft



Linjun Shou

Software Technology Center at Asia (STCA) of Microsoft



Xiubo Geng



Ming Gong



Jian Pei

Simon Fraser University

\*For more applications, please refer to our tutorial at  
The Web Conference 2021

# Grammar Error Correction & Question Generation

Linjun Shou  
[lijsho@microsoft.com](mailto:lijsho@microsoft.com)

# Cross-lingual GEC & QG

Type	Category	Sub Category	Example
NLU	Text Classification	Single text	Domain identification, Intent detection, Sentiment classification
		Text pair	Information retrieval, Natural language inference
	Sequence Labeling	Single text	Named entity recognition, Slot tagging
		Text pair	Machine reading comprehension
	Structure Prediction	Single text	Dependency parsing, semantic role labeling
NLG	Text Generation	Token level	Spelling correction, Sentence auto completion, <b>Grammar error correction</b>
		Sentence level	Machine translation, Conversation, <b>Question generation</b>

# Grammar Error Correction

Grammatical Error Correction (GEC) is the task of correcting different kinds of errors in text such as spelling, punctuation, grammatical, and word choice errors.

Input (Erroneous)	Output (Corrected)
She see Tom is catched by policeman in park at last night.	She saw Tom caught by a policeman in the park last night.

# Question Generation

Question Generation (QG) is the task of generating a question based on passage context and an input answer.

<b>Input Passage</b>	... In case of European Union law which should have been transposed into the laws of member stages, such Directives, the European Commission can take proceedings against the member stage under the <i>Treaty on the Functioning of the European Union</i> , ...
<b>Input Answer</b>	Treaty on the Functioning of the European Union
<b>Reference</b>	Under what treaty can the European Commission take actions against member states?

# Approaches for Cross-lingual Transfer

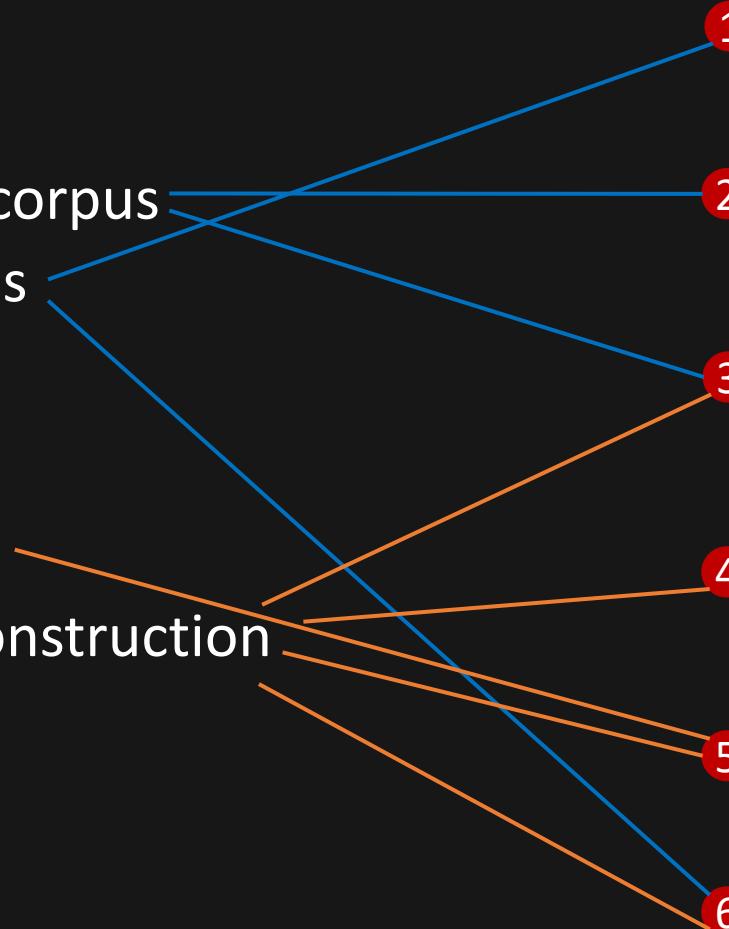
## Representative Works

- **Model Transfer**

- w/ monolingual corpus
- w/ parallel corpus

- **Data Transfer**

- Back translation
- Synthetic data construction



### Cross-Lingual Natural Language Generation via Pre-Training

Zewen Chi<sup>†\*</sup>, Li Dong<sup>‡</sup>, Furu Wei<sup>‡</sup>, Wenhui Wang<sup>‡</sup>, Xian-Ling Mao<sup>‡</sup>, Heyan Huang<sup>‡</sup>  
<sup>†</sup>Beijing Institute of Technology  
<sup>‡</sup>Microsoft Research  
{czw,maoxl,hhy63}@bit.edu.cn  
{lidong1,fwei,Wenhui.Wang}@microsoft.com

### Multilingual Denoising Pre-training for Neural Machine Translation

Yinhan Liu<sup>†\*</sup>, Jiatao Gu<sup>†\*</sup>, Naman Goyal<sup>†\*</sup>, Xian Li<sup>†</sup>, Sergey Edunov<sup>†</sup>,  
Marjan Ghazvininejad<sup>†</sup>, Mike Lewis<sup>†</sup>, and Luke Zettlemoyer<sup>‡</sup>

### mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer

Linting Xue<sup>\*</sup> Noah Constant<sup>\*</sup> Adam Roberts<sup>\*</sup>  
Mihir Kale Rami Al-Rfou Aditya Siddhant Aditya Barua Colin Raffel  
Google Research

### A Simple Recipe for Multilingual Grammatical Error Correction

Sascha Rothe Google rothe@google.com	Jonathan Mallinson Google jonmall@google.com	Eric Malmi Google emalmi@google.com
Sebastian Krause Google bastik@google.com	Aliaksei Severyn Google severyn@google.com	

### An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

Shun Kiyono<sup>1,2</sup> Jun Suzuki<sup>2,1</sup> Masato Mita<sup>1,2</sup> Tomoya Mizumoto<sup>1,2\*</sup> Kentaro Inui<sup>2,1</sup>  
<sup>1</sup> RIKEN Center for Advanced Intelligence Project <sup>2</sup> Tohoku University  
{shun.kiyono, masato.mita, tomyo.mizumoto}@riken.jp;  
{jun.suzuki, inui}@ecei.tohoku.ac.jp

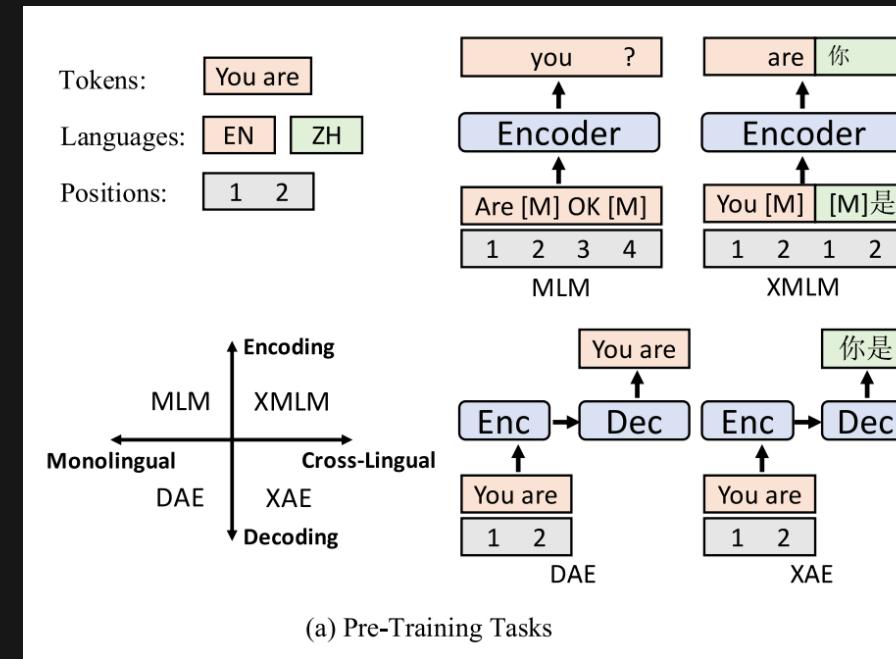
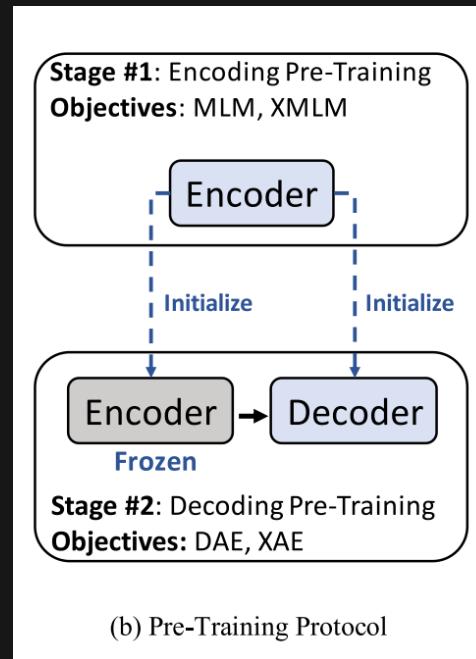
### Cross-lingual Transfer Learning for Grammatical Error Correction

Ikumi Yamashita Satoru Katsumata<sup>\*</sup> Masahiro Kaneko  
Aizhan Imankulova Mamoru Komachi  
Tokyo Metropolitan University, Japan  
yamashita-ikumi@ed.tmu.ac.jp  
satoru.katsumata@retrieva.jp  
{kaneko-masahiro, imankulova-aizhan}@ed.tmu.ac.jp  
komachi@tmu.ac.jp

# Model Transfer - XNLG

- **Model Transfer**

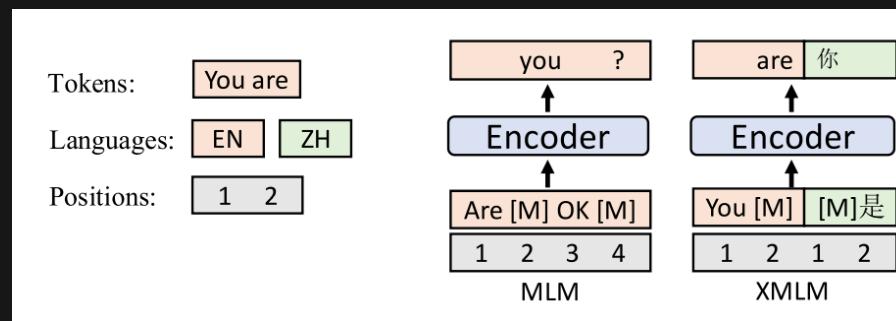
- w/ monolingual corpus
- w/ parallel corpus



# Model Transfer - XNLG

- **Encoder Pretraining**

- Monolingual data
- Parallel data



$$\mathcal{L}_{\text{MLM}}^{(x)} = - \sum_{i \in M_x} \log p(x_i | x_{\setminus M_x})$$

$$\mathcal{L}_{\text{XMLM}}^{(x,y)} = - \sum_{i \in M_x} \log p(x_i | x_{\setminus M_x}, y_{\setminus M_y}) \quad (2)$$

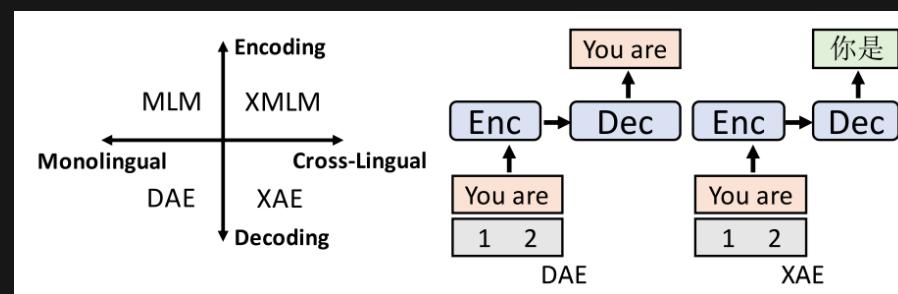
$$- \sum_{i \in M_y} \log p(y_i | x_{\setminus M_x}, y_{\setminus M_y}) \quad (3)$$

where  $M_x, M_y$  represent the masked positions of  $x$  and  $y$ , respectively.

# Model Transfer - XNLG

- **Decoder Pretraining (w/ Encoder fixed)**

- Monolingual data
- Parallel data



$$\mathcal{L}_2 = \sum_{(x,y) \in \mathcal{D}_p} \mathcal{L}_{\text{XAE}}^{(x,y)} + \sum_{x \in \mathcal{D}_m} \mathcal{L}_{\text{DAE}}^{(x)}$$

$$\mathcal{L}_{\text{XAE}}^{(x,y)} = -\log p(y|x) - \log p(x|y)$$

- Word order is locally shuffled
- Randomly drop tokens with probability = 0.1
- substitute tokens with the special padding token [P] with a probability of 0.1.

# Model Transfer - XNLG

- Results on Zero-shot QG

Models	BL-4	MTR	RG-L
XLM	0.25	0.62	2.56
PIPELINE (XLM)	4.42	9.59	21.22
w/ Google Translator	9.95	14.92	29.37
XNLG	<b>16.37</b>	<b>18.74</b>	<b>34.93</b>

Table 2: Evaluation results of zero-shot Chinese-Chinese question generation. Same shorthands apply as in Table 1.

Models	Rel	Flu	Corr
XLM	0	0	0
PIPELINE (XLM)	0.50	0.80	0.03
w/ Google Translator	1.31	<b>1.43*</b>	0.69
XNLG	<b>1.68*</b>	1.29	<b>0.89*</b>

Table 3: Human evaluation results of zero-shot Chinese-Chinese question generation. Rel is short for relatedness, Flu for fluency, and Corr for correctness. “\*” indicates the improvements are significant at  $p < 0.05$ .

# Model Transfer – mBART

- Model Transfer
  - w/ monolingual corpus

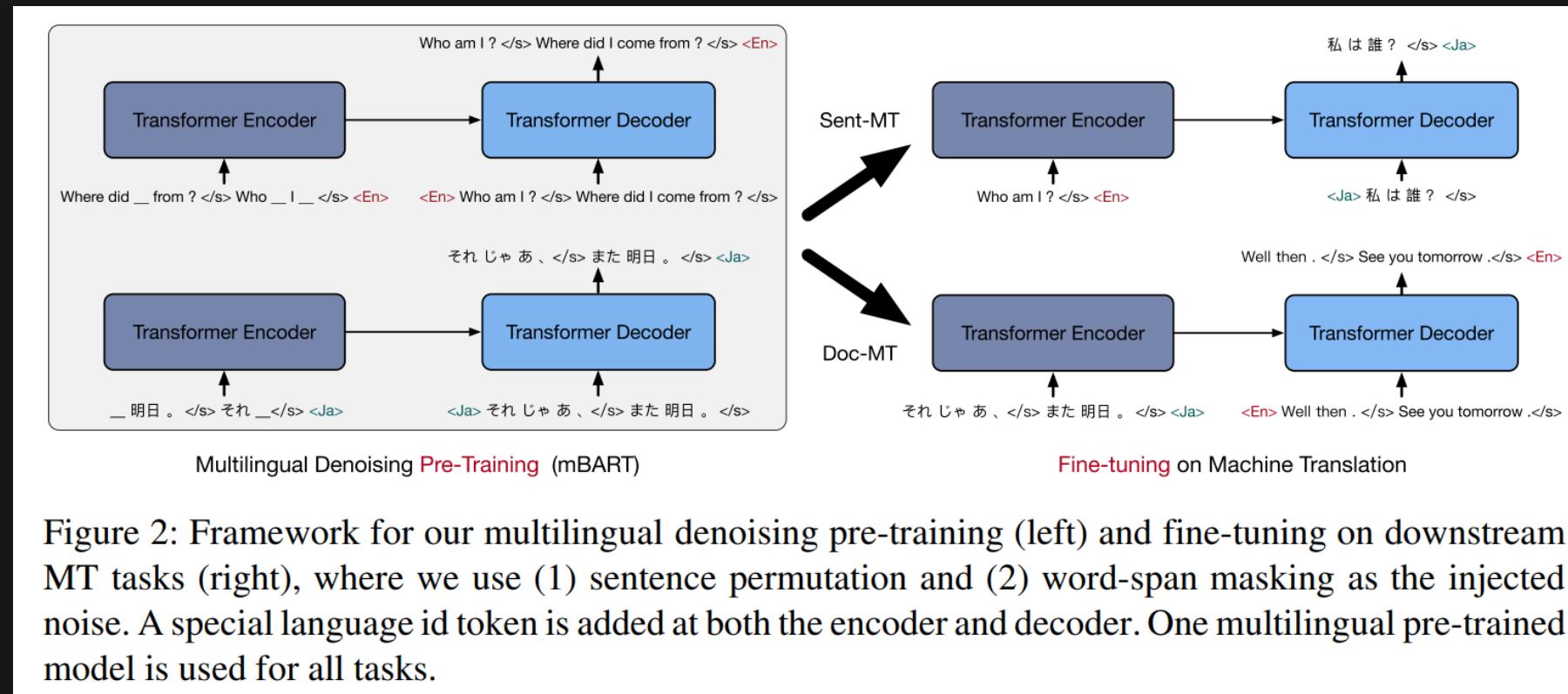


Figure 2: Framework for our multilingual denoising pre-training (left) and fine-tuning on downstream MT tasks (right), where we use (1) sentence permutation and (2) word-span masking as the injected noise. A special language id token is added at both the encoder and decoder. One multilingual pre-trained model is used for all tasks.

# Model Transfer – mBART

- Model Transfer
  - w/ monolingual corpus

$$\mathcal{L}_\theta = \sum_{\mathcal{D}_i \in \mathcal{D}} \sum_{X \in \mathcal{D}_i} \log P(X|g(X); \theta)$$

Noise function  $g(X)$ :

- mask 35% of the words in each instance by random sampling a span length according to a Poisson distribution ( $\lambda = 3.5$ )
- Permute the order of sentences within each instance.

Code	Language	Tokens/M	Size/GB
En	English	55608	300.8
Ru	Russian	23408	278.0
Vi	Vietnamese	24757	137.3
Ja	Japanese	530 (*)	69.3
De	German	10297	66.6
Ro	Romanian	10354	61.4
Fr	French	9780	56.8
Fi	Finnish	6730	54.3
Ko	Korean	5644	54.2
Es	Spanish	9374	53.3
Zh	Chinese (Sim)	259 (*)	46.9
It	Italian	4983	30.2
Nl	Dutch	5025	29.3
Ar	Arabic	2869	28.0
Tr	Turkish	2736	20.9
Hi	Hindi	1715	20.2
Cs	Czech	2498	16.3
Lt	Lithuanian	1835	13.7
Lv	Latvian	1198	8.8
Kk	Kazakh	476	6.4
Et	Estonian	843	6.1
Ne	Nepali	237	3.8
Si	Sinhala	243	3.6
Gu	Gujarati	140	1.9
My	Burmese	56	1.6

Table 1: Languages and Statistics of the CC25 Corpus. A list of 25 languages ranked with monolingual corpus size. Throughout this paper, we replace the language names with their ISO codes for simplicity. (\*) Chinese and Japanese corpus are not segmented, so the tokens counts here are sentences counts

# Model Transfer – mT5

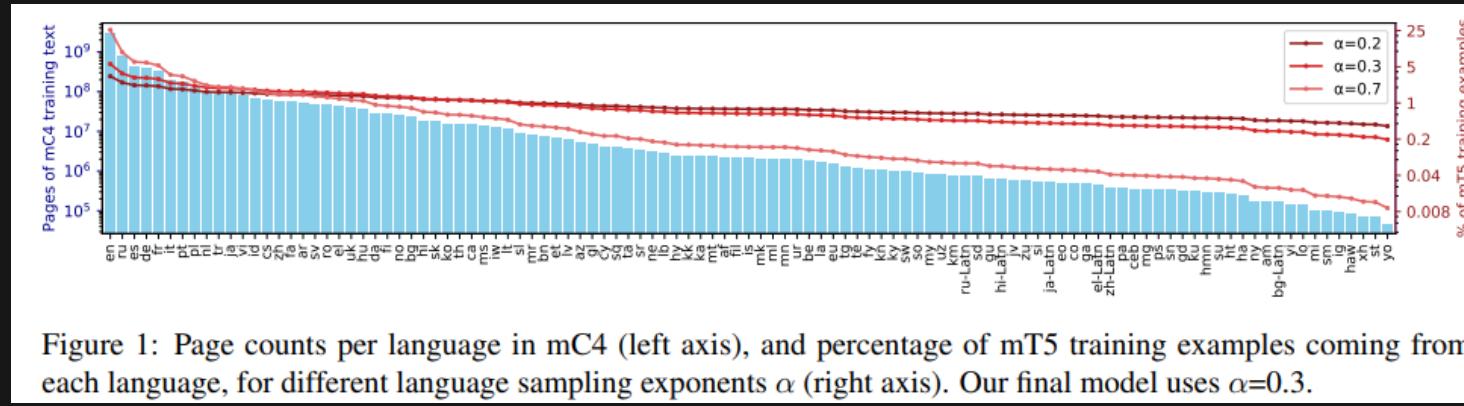
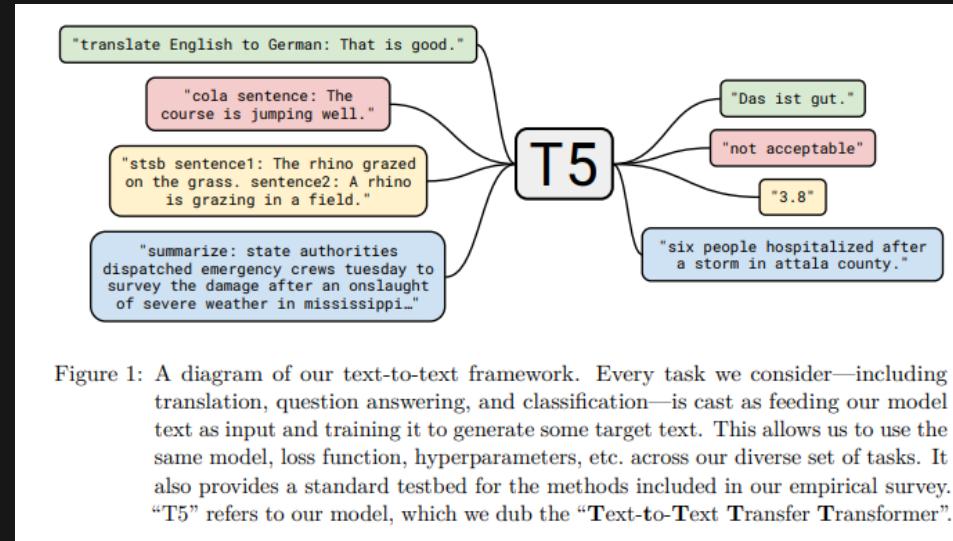
- Model Transfer
  - w/ monolingual corpus

Model	Architecture	Parameters	# languages	Data source
mBERT (Devlin, 2018)	Encoder-only	180M	104	Wikipedia
XLM (Conneau and Lample, 2019)	Encoder-only	570M	100	Wikipedia
XLM-R (Conneau et al., 2020)	Encoder-only	270M – 550M	100	Common Crawl (CCNet)
mBART (Lewis et al., 2020b)	Encoder-decoder	680M	25	Common Crawl (CC25)
MARGE (Lewis et al., 2020a)	Encoder-decoder	960M	26	Wikipedia or CC-News
mT5 (ours)	Encoder-decoder	300M – 13B	101	Common Crawl (mC4)

Table 1: Comparison of mT5 to existing massively multilingual pre-trained language models. Multiple versions of XLM and mBERT exist; we refer here to the ones that cover the most languages. Note that XLM-R counts five Romanized variants as separate languages, while we ignore six Romanized variants in the mT5 language count.

# Model Transfer – mT5

- Multi-lingual version of T5
- mC4 in 101 languages



# Model Transfer – mT5

## Pretraining Tasks:

Objective	Inputs	Targets
Prefix language modeling BERT-style <a href="#">Devlin et al. (2018)</a>	Thank you for inviting Thank you <M> <M> me to your party apple week .	me to your party last week . <i>(original text)</i>
Deshuffling MASS-style <a href="#">Song et al. (2019)</a>	party me for your to . last fun you inviting week Thank Thank you <M> <M> me to your party <M> week .	<i>(original text)</i> <i>(original text)</i>
I.i.d. noise, replace spans I.i.d. noise, drop tokens	Thank you <X> me to your party <Y> week . Thank you me to your party week .	<X> for inviting <Y> last <Z> for inviting last
Random spans	Thank you <X> to <Y> week .	<X> for inviting me <Y> your part

Objective	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
Prefix language modeling	80.69	18.94	77.99	65.27	<b>26.86</b>	39.73	<b>27.49</b>
BERT-style ( <a href="#">Devlin et al., 2018</a> )	<b>82.96</b>	<b>19.17</b>	<b>80.65</b>	<b>69.85</b>	<b>26.78</b>	<b>40.03</b>	<b>27.41</b>
Deshuffling	73.17	18.59	67.61	58.47	26.11	39.30	25.62

Table 4: Performance of the three disparate pre-training objectives described in Section 3.3.1.

# Model Transfer – mT5

- Results

Model	Sentence pair		WikiAnn NER	Question answering		
	XNLI	PAWS-X		XQuAD	MLQA	TyDiQA-GoldP
Metrics	Acc.	Acc.	F1	F1 / EM	F1 / EM	F1 / EM
<i>Cross-lingual zero-shot transfer (models fine-tuned on English data only)</i>						
mBERT	65.4	81.9	62.2	64.5 / 49.4	61.4 / 44.2	59.7 / 43.9
XLM	69.1	80.9	61.2	59.8 / 44.3	48.5 / 32.6	43.6 / 29.1
InfoXLM	81.4	-	-	- / -	73.6 / 55.2	- / -
X-STILTs	80.4	87.7	64.7	77.2 / 61.3	72.3 / 53.5	76.0 / 59.5
XLM-R	79.2	86.4	65.4	76.6 / 60.8	71.6 / 53.2	65.1 / 45.0
VECO	79.9	88.7	65.7	77.3 / 61.8	71.7 / 53.2	67.6 / 49.1
RemBERT	80.8	87.5	<b>70.1</b>	79.6 / 64.0	73.1 / 55.0	77.0 / 63.0
mT5-Small	67.5	82.4	50.5	58.1 / 42.5	54.6 / 37.1	35.2 / 23.2
mT5-Base	75.4	86.4	55.7	67.0 / 49.0	64.6 / 45.0	57.2 / 41.2
mT5-Large	81.1	88.9	58.5	77.8 / 61.5	71.2 / 51.7	69.9 / 52.2
mT5-XL	82.9	89.6	65.5	79.5 / 63.6	73.5 / 54.5	75.9 / 59.4
mT5-XXL	<b>85.0</b>	<b>90.0</b>	69.2	<b>82.5 / 66.8</b>	<b>76.0 / 57.4</b>	<b>80.8 / 65.9</b>
<i>Translate-train (models fine-tuned on English data plus translations in all target languages)</i>						
XLM-R	82.6	90.4	-	80.2 / 65.9	72.8 / 54.3	66.5 / 47.7
FILTER + Self-Teaching	83.9	91.4	-	82.4 / 68.0	76.2 / 57.7	68.3 / 50.9
VECO	83.0	91.1	-	79.9 / 66.3	73.1 / 54.9	75.0 / 58.9
mT5-Small	64.7	79.9	-	64.3 / 49.5	56.6 / 38.8	48.2 / 34.0
mT5-Base	75.9	89.3	-	75.3 / 59.7	67.6 / 48.5	64.0 / 47.7
mT5-Large	81.8	91.2	-	81.2 / 65.9	73.9 / 55.2	71.1 / 54.9
mT5-XL	84.8	91.0	-	82.7 / 68.1	75.1 / 56.6	79.9 / 65.3
mT5-XXL	<b>87.8</b>	<b>91.5</b>	-	<b>85.2 / 71.3</b>	<b>76.9 / 58.3</b>	<b>82.8 / 68.8</b>
<i>In-language multitask (models fine-tuned on gold data in all target languages)</i>						
mBERT	-	-	89.1	-	-	77.6 / 68.0
mT5-Small	-	-	83.4	-	-	73.0 / 62.0
mT5-Base	-	-	85.4	-	-	80.8 / 70.0
mT5-Large	-	-	88.4	-	-	85.5 / 75.3
mT5-XL	-	-	90.9	-	-	87.5 / 78.1
mT5-XXL	-	-	<b>91.2</b>	-	-	<b>88.5 / 79.1</b>

Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou, R., Siddhant, A., Barua, A., & Raffel, C. (2021). mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer. NAACL.

# Model Transfer – mT5

## Zero shot generation Challenges – Illegal predictions

- Normalization
- Grammatical adjustment
- Accidental translation

Target	Prediction	Explanation
ຈຳນາມແລ້ວພາກ	ຈຳນາມແລ້ວພາກ	Decomposed Thai ຈຳ + ນາມ + ແລ້ວ + ພາກ
लोथर डे माइज़ियर	लोथर डे माइज़ियर	Decomposed Hindi लोथर + डे + माइज़ियर
27 - 30 %	27 - 30 %	Replaced full-width percent sign
12 . <sup>a</sup>	12 . a	Removed superscript
لِبَكْرِيَةُ الْأَمَوَانِيَّةُ	الْبَكْرِيَةُ الْأَمَوَانِيَّةُ	Arabic “for anaerobic bacteria” ⇒ “anaerobic bacteria”
строками битов	строки битов	Russian “bit strings (instrumental)” ⇒ “bit strings (nominative)”
seis años	six years	Translated from Spanish
Zweiten Weltkrieg	the Second World War	Translated from German
新英格兰爱国者队	New 英格兰爱国者队	Partially translated Chinese “New England Patriots”
хлоропласт	chloroplast	Partially translated Russian “chloroplast”

Table 5: Illegal mT5-XXL predictions on XQuAD zero-shot, illustrating normalization (top), grammatical adjustment (middle) and translation (bottom).

*During fine-tuning stage, combine with pretraining task for multi-task learning*

# Data Transfer –Domain Adaptation

## mT5 Pretraining

**Input:** A Simple [x] Multilingual Grammatical Error [y]  
**Target:** [x] Recipe for [y] Correction

- Span-prediction task for Seq2Seq pretraining
- mC4 data -> 50 billion documents in 101 languages

## GEC Pretraining

**Input:** Simple recipe for Multilingual Grammatical Correction Error  
**Target:** A Simple Recipe for Multilingual Grammatical Error Correction

- Adapt to GEC task
- Denoising pretraining
  - *Drop spans of tokens*
  - *Swap tokens*
  - *Drop spans of characters*
  - *Swap characters*
  - *Insert*
  - *Lower/Upper case of words*

## GEC Fine-tuning

# Data Transfer - Domain Adaptation

Models	<i>CoNLL-14</i>	<i>BEA test</i>	<i>Czech</i>	<i>German</i>	<i>Russian</i>
Omelianchuk et al.*	66.5	<b>73.6</b>	-	-	-
Lichtarge et al.*	<b>66.8</b>	73.0	-	-	-
Náplava and Straka	63.40	69.00	80.17	73.71	50.20
Katsumata and Komachi*	63.00	66.10	73.52	68.86	44.36
gT5 base	54.10	60.2	71.88	69.21	26.24
gT5 xxl	65.65	69.83	<b>83.15</b>	<b>75.96</b>	<b>51.62</b>

Table 2:  $F_{0.5}$  Scores. Models denoted with \* are ensemble models. We used the  $M^2$  scorer for CoNLL-14, Russian, Czech and German, and the ERRANT scorer (Bryant et al., 2019b) for BEA test.

# Data Transfer - An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

Q1: How to generate better pseudo data ( $X$ )

Q2: Seed Corpus Selection ( $Y$ )

- Clean data
- Grammatical complexity

Q3: Joint Training or Pretraining

$$\mathcal{L}(\mathcal{D}, \Theta) = -\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{X}, \mathbf{Y}) \in \mathcal{D}} \log(p(\mathbf{Y}|\mathbf{X}, \Theta)), \quad (1)$$

where  $p(\mathbf{Y}|\mathbf{X}, \Theta)$  denotes the conditional probability of  $\mathbf{Y}$  given  $\mathbf{X}$ .

# Data Transfer - An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

## Q1: How to Generate Better Pseudo Data:

1. Direct Noise
  - Inject noise directly
2. Back translation (Noisy)
  - Train a Seq2Seq model to generate noisy sentences, given grammatically correct sentence
  - Add noise during decoding stage
3. Back translation (sample)
  - Seq2Seq
  - Random Sample during decoding (beam search)

# Data Transfer - An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

## Q1: How to Generate Better Pseudo Data:

1. Direct Noise
2. Back translation (Noisy)
3. Back translation (sample)

Method	Prec.	Rec.	$F_{0.5}$
Baseline	46.6	23.1	38.8
BACKTRANS (SAMPLE)	44.6	27.4	39.6
BACKTRANS (NOISY)	42.5	<b>31.3</b>	39.7
DIRECTNOISE	<b>48.9</b>	25.7	<b>41.4</b>

Table 2: Performance of models on BEA-valid: a value in **bold** indicates the best result within the column. The seed corpus  $\mathcal{T}$  is SimpleWiki.

# Data Transfer - An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

## Q2: Seed Corpus Selection:

Corpus	#sent (pairs)	Notes
SimpleWiki	1,369,460	Simple grammar
Wikipedia	145,883,941	Complex grammar
Gigaword	131,864,979	Clean text

# Data Transfer - An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

## Q2: Seed Corpus Selection:

Method	Seed Corpus $\mathcal{T}$	Prec.	Rec.	$F_{0.5}$
Baseline	N/A	46.6	23.1	38.8
BACKTRANS (NOISY)	Wikipedia	43.8	30.8	40.4
BACKTRANS (NOISY)	SimpleWiki	42.5	31.3	39.7
BACKTRANS (NOISY)	Gigaword	43.1	33.1	40.6
DIRECTNOISE	Wikipedia	48.3	25.5	41.0
DIRECTNOISE	SimpleWiki	48.9	25.7	41.4
DIRECTNOISE	Gigaword	48.3	26.9	41.7

Table 3: Performance on BEA-valid when changing the seed corpus  $\mathcal{T}$  used for generating pseudo data ( $|\mathcal{D}_p| = 1.4M$ ).

# Data Transfer - An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

## Q3: Joint Training vs Pretraining:

Optimization Method		$ D_p $	Prec.	Rec.	$F_{0.5}$
N/A	Baseline	0	46.6	23.1	38.8
PRETRAIN	BACKTRANS (NOISY)	1.4M	49.6	24.3	41.1
PRETRAIN	DIRECTNOISE	1.4M	48.4	21.2	38.5
JOINT	BACKTRANS (NOISY)	1.4M	43.8	30.8	40.4
JOINT	DIRECTNOISE	1.4M	48.3	25.5	41.0
PRETRAIN	BACKTRANS (NOISY)	14M	50.6	30.1	44.5
PRETRAIN	DIRECTNOISE	14M	49.8	25.8	42.0
JOINT	BACKTRANS (NOISY)	14M	43.0	32.3	40.3
JOINT	DIRECTNOISE	14M	48.7	23.5	40.1

Table 4: Performance of the model with different optimization settings on BEA-valid. The seed corpus  $\mathcal{T}$  is Wikipedia.



PRETRAIN is crucial for utilizing extensive pseudo data.

# Data Transfer - An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

## Q3: Joint Training vs Pretraining:

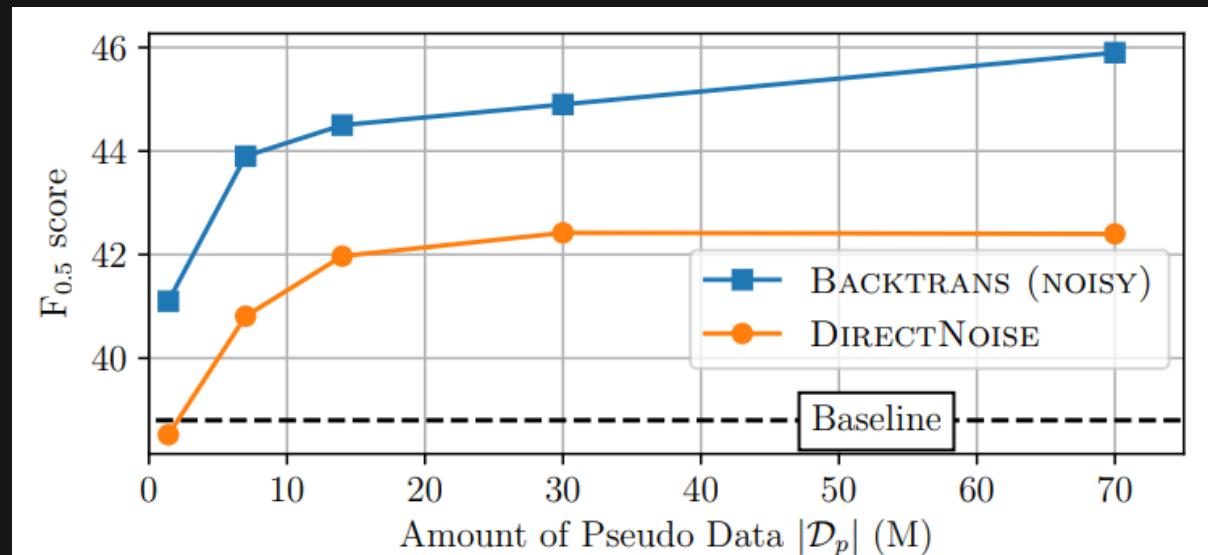


Figure 1: Performance on BEA-valid for different amounts of pseudo data ( $|\mathcal{D}_p|$ ). The seed corpus  $\mathcal{T}$  is Wikipedia.

# Data Transfer - Cross-lingual Transfer Learning for Grammatical Error Correction

- Effect of Language Similarity
- Effect of parallel corpus

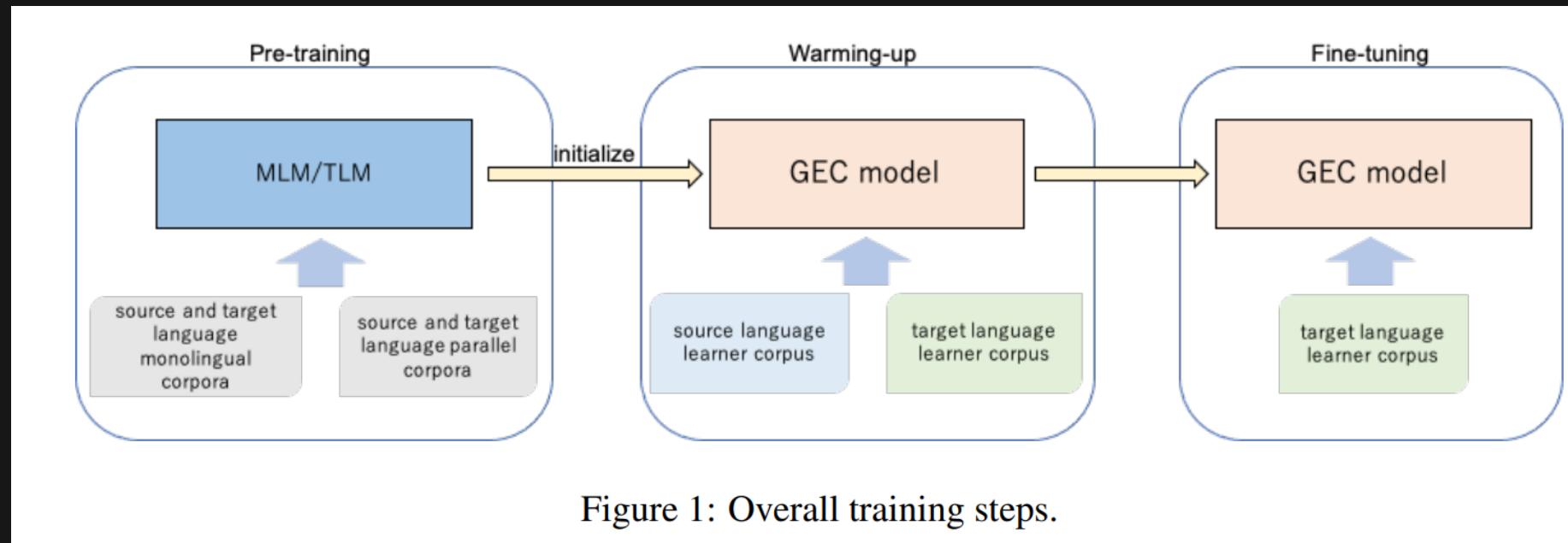


Figure 1: Overall training steps.

# Data Transfer - Cross-lingual Transfer Learning for Grammatical Error Correction

target	source similarity		
	high	moderate	low
Russian	Czech	English	Japanese
Czech	Russian	English	Japanese
English	German	Russian	Japanese

Table 2: Languages used in the experiment.

- Better transfer for similar languages
- Parallel data helps better transfer

	Model	P	R	F <sub>0.5</sub>
PLAIN	Ru-only	19.29	14.08	17.96
	Cs→Ru	19.05	12.78	17.35
	En→Ru	23.76	13.65	20.70
	Ja→Ru	20.70	13.57	18.73
MLM	Ru-only	19.95	<b>23.15</b>	20.52
	Cs→Ru	26.36	19.02	24.47
	En→Ru	26.02	19.74	24.47
	Ja→Ru	27.23	16.13	23.93
TLM	Cs→Ru	<b>28.51</b>	22.47	<b>27.06</b>
	En→Ru	27.60	22.18	26.31
	Ja→Ru	26.11	19.99	24.61

Table 4: Russian GEC results.

	Model	P	R	F <sub>0.5</sub>
PLAIN	Cs-only	52.05	39.29	48.88
	Ru→Cs	59.93	38.73	54.01
	En→Cs	61.35	39.01	55.05
	Ja→Cs	56.22	38.53	51.49
MLM	Cs-only	57.46	47.40	55.12
	Ru→Cs	63.58	47.15	59.43
	En→Cs	63.54	48.63	59.87
	Ja→Cs	62.15	47.35	58.50
TLM	Ru→Cs	<b>65.09</b>	<b>50.82</b>	<b>61.63</b>
	En→Cs	63.20	48.47	59.58
	Ja→Cs	63.84	45.70	59.15

Table 5: Czech GEC results.

# Key Takeaways

## Model Transfer

- Multi-lingual Denoising pretraining is effective for cross-lingual Seq2seq pretraining
- Parallel corpus helps language transfer
- Bigger model / Larger data lead to better performance

## Data transfer

- GEC pretraining could boost performance on top of general pretrained models

# Outline

- Introduction [Dixin Jiang]
  - Motivating examples in Microsoft products
  - Problem description
  - Categorization of applications
  - Challenges
- Methodology [Dixin Jiang]
  - Model Transfer
  - Data Transfer
- Applications\*
  - Dependency Parsing [Xiubo Geng]
  - Machine Reading Comprehension [Ming Gong]
  - Grammar Error Correction [Linjun Shou]
- • Summary & Future directions [Jian Pei]



Dixin Jiang

Software Technology Center at Asia (STCA) of Microsoft



Linjun Shou

Software Technology Center at Asia (STCA) of Microsoft



Xiubo Geng



Ming Gong



Jian Pei

Simon Fraser University

\*For more applications, please refer to our tutorial at  
The Web Conference 2021

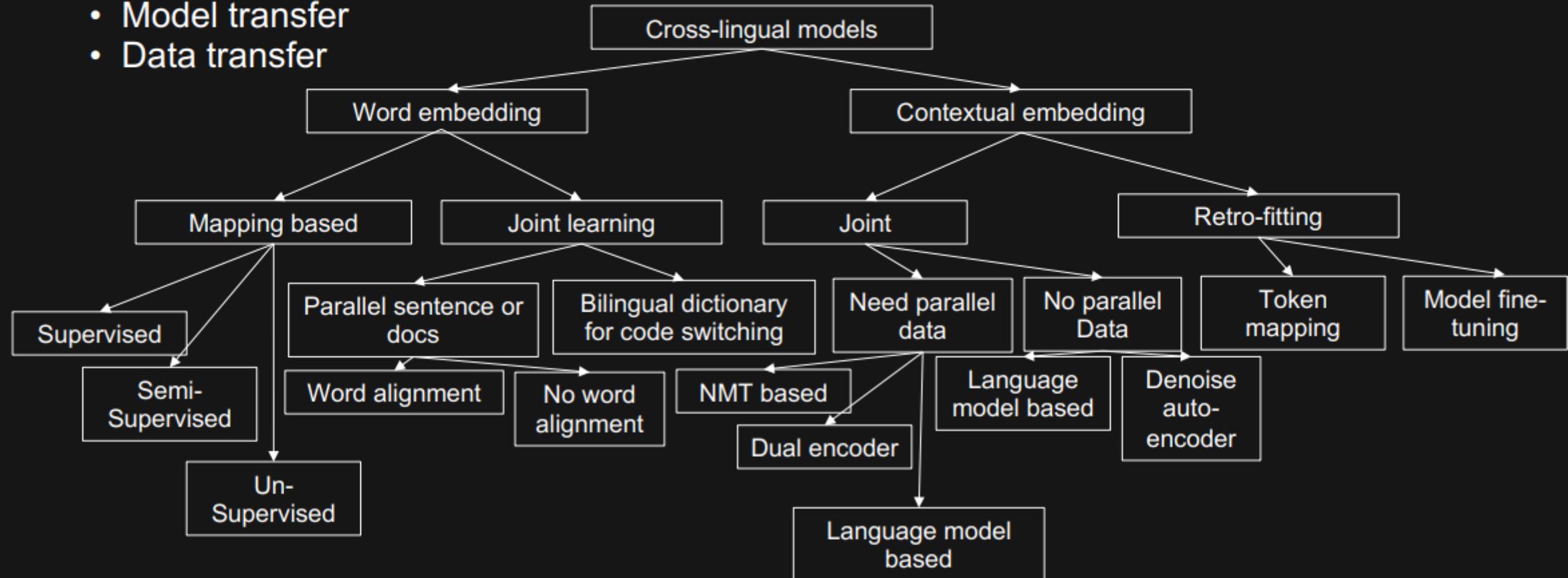
# Summary

- Language scaling is critical in business and many applications
  - Key technical challenge: lack of cross-lingual data

Type	Category	Sub Category	Example
NLU	Text Classification	Single text	Domain identification, Intent detection, Sentiment classification
		Text pair	Information retrieval, Natural language inference
	Sequence Labeling	Single text	Named entity recognition, Slot tagging
		Text pair	Machine reading comprehension
	Structure Prediction	Single text	Dependency parsing, constituency parsing, semantic role labeling
NLG	Text Generation	Token level	Spelling correction, Sentence auto completion
		Sentence level	Machine translation, Conversation, Question generation

# Summary

- Methodologies
  - Model transfer
  - Data transfer



Application examples:

- Dependency parsing
- Machine reading comprehension
- Grammar error correction

# Future Directions (1)

- Linguistic knowledge
  - Universal POS tagging and Dependency Parsing ([Universal Dependencies](#))
  - Typological databases such as the WALS (Dryer and Haspelmath, 2013),  
PHOIBLE (Moran et al., 2014), Ethnologue (Lewis et al., 2015), and Glottolog  
(Hammarstrom et al., 2015)
- Early works use topological databases to measure the distances  
between languages or select the intermediate transfer languages
- Look forward to systematic methods that integrate linguistic  
knowledge into transfer approaches

WALS: Matthew S. Dryer and Martin Haspelmath. 2013. The World Atlas of Language Structures Online. Max Planck Institute for Evolutionary Anthropology

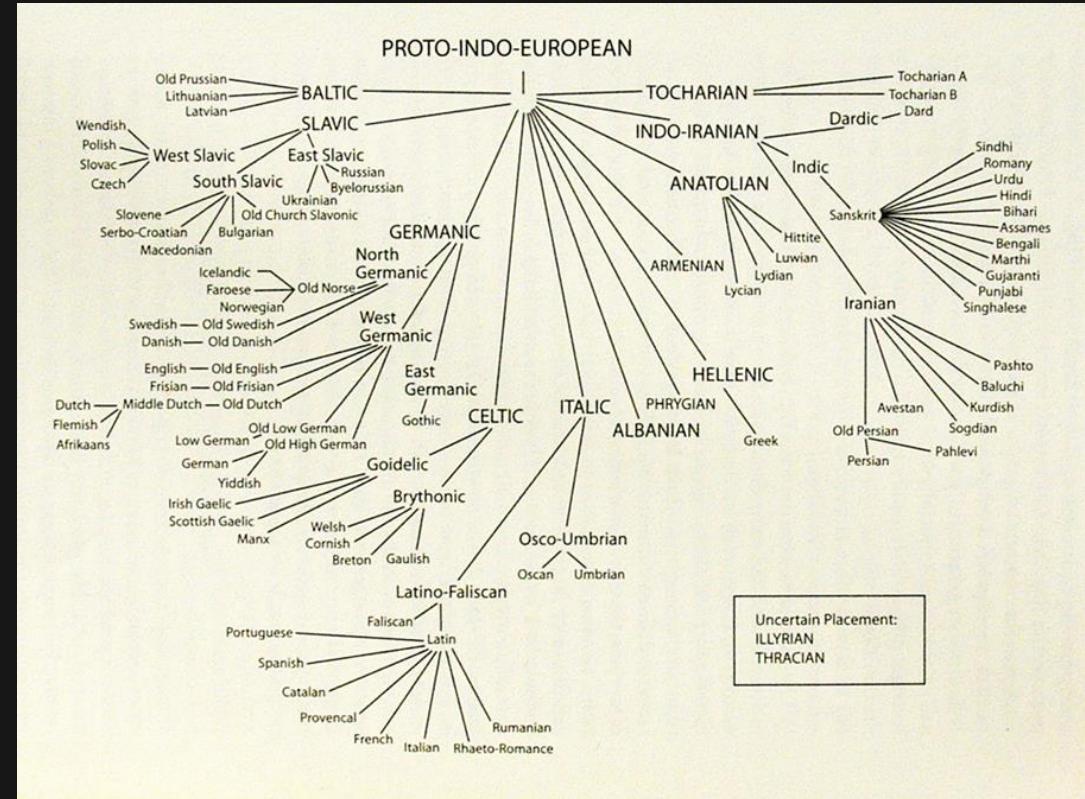
PHOIBLE: Steven Moran, Daniel McCloy, and Richard Wright. 2014. PHOIBLE Online. Max Planck Institute for Evolutionary Anthropology, Leipzig

Ethnologue: M. Paul Lewis, Gary F. Simons, and Charles D. Fennig. 2015. Ethnologue: Languages of the World, Eighteenth edition. SIL International, Dallas, Texas

Glottolog: Harald Hammarstrom, Robert Forkel, Martin Haspelmath, and Sebastian Bank. 2015. Glottolog 2.6. Max Planck Institute for the Science of Human History, Jena

# Future Directions (2)

- Best practice to train multi-lingual models
  - how to mitigate “catastrophic forgetting”?
  - do we align models with language families?
  - how to measure the “capacity” of models and the “size” of languages to determine the best combination of languages



Picture from [The English Cowpath: The Proto-Indo-European Homeland Puzzle](#)

thank you

danke 謝謝  
спасибо faafetai lava  
спасибо Dankie  
спасибо nandi  
спасибо kijitos  
спасибо dhanyavad  
спасибо bayatlaa  
спасибо mauruui  
спасибо namm  
спасибо dankie  
спасибо vinaka  
спасибо спасиби  
спасибо blagodaram  
спасибо kia ota  
спасибо barka  
спасибо welalin  
спасибо tack  
спасибо dank je  
спасибо misaotra  
спасибо matondo  
спасибо paldies  
спасибо grazzi  
спасибо mahalo  
спасибо tapadh leat  
спасибо xvala  
спасибо asante manana  
спасибо obrigada  
спасибо murakozé  
спасибо tenki  
спасибо chokkane  
спасибо mamnun  
спасибо djiere dieuf  
спасибо tau  
спасибо mochchakkeram  
спасибо dya'kuo  
спасибо go raibh maith agat  
спасибо arigatō takk  
спасибо dakujem trugarez  
спасибо merce  
спасибо merce  
спасибо xiexie  
спасибо ευχαριστώ  
спасибо 감사합니다  
спасибо sukriya  
спасибо kop khun krap  
спасибо terima kasih  
спасибо 감사합니다  
спасибо danke  
спасибо teşekkür ederim  
спасибо gracias  
спасибо merci  
спасибо ngiyabonga  
спасибо tesekkür ederim  
спасибо gracias  
спасибо thank you