



OLSCOIL NA GAILLIMHE  
UNIVERSITY OF GALWAY

## Managing and Processing Multiple Languages

Jamal Nasir

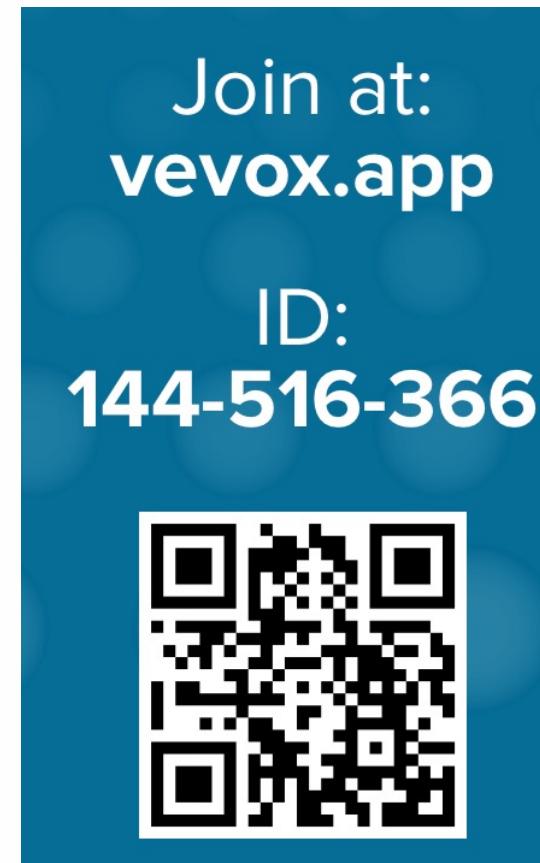
COST Action Training School in  
Computational Opinion Analysis

OP1N10N0P  
P1N10N0P1N1  
10N0P1N10N0  
N0P1N10N0P1  
N0P1N10N0

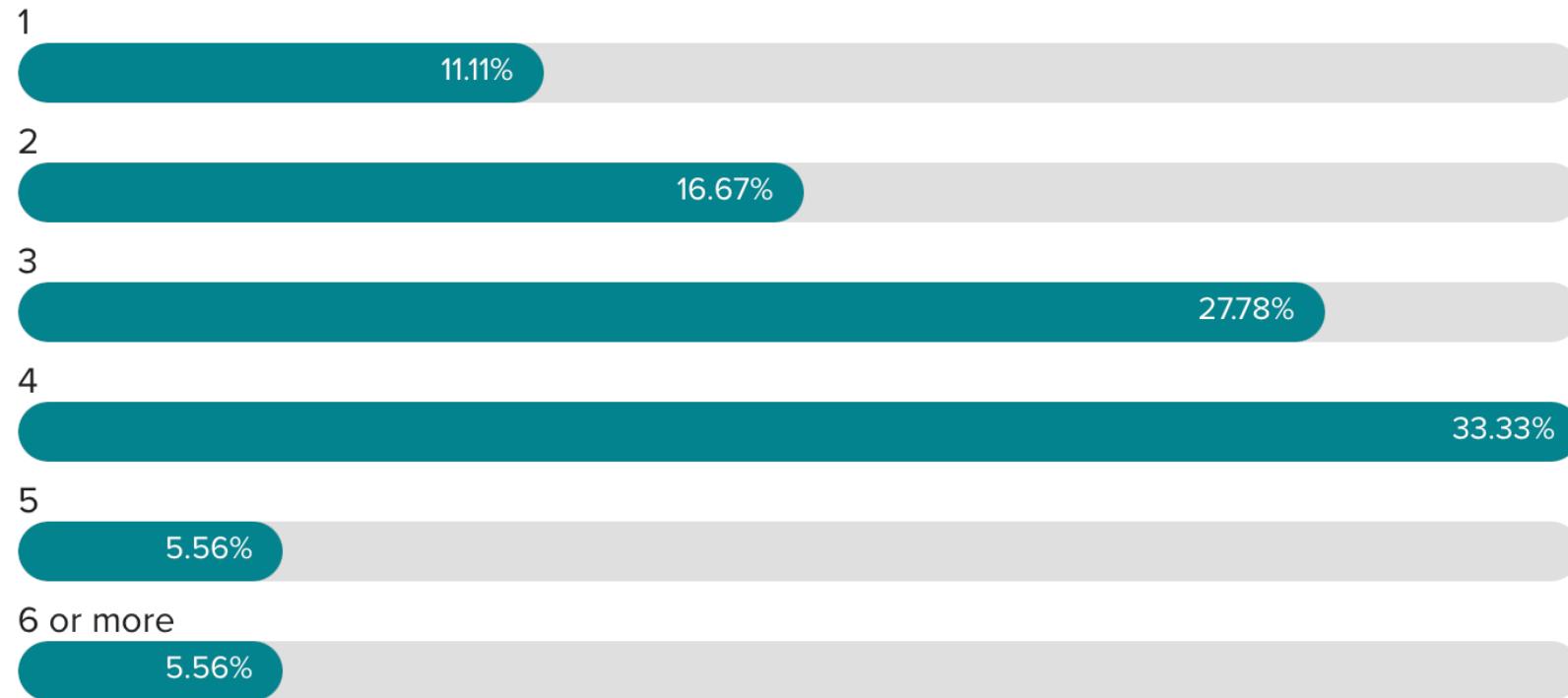


University  
*of*Galway.ie

How many languages can you speak or understand?



## How many languages can you speak or understand?





OLSCOIL NA GAILIMHE  
UNIVERSITY OF GALWAY

# Managing and Processing Multiple Languages: Data or Technology

University  
*of*Galway.ie

# How many Languages are there in the world?

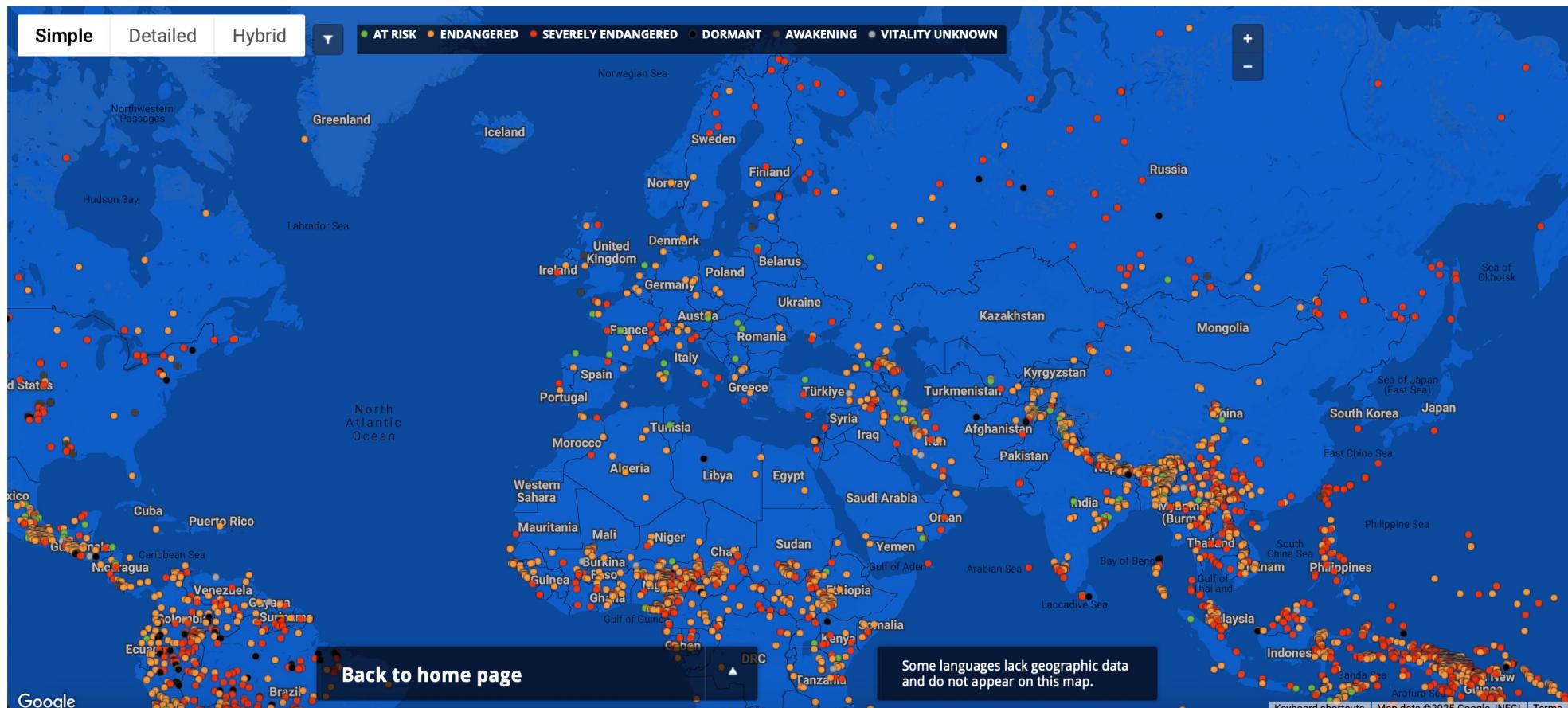
7,159 languages are in use today. □

That number is constantly in flux, because we're learning more about the world's languages every day. And beyond that, the languages *themselves* are in flux. They're living and dynamic, used by communities whose lives are shaped by our rapidly changing world. This is a fragile time: Roughly 44% of all languages are now [endangered](#), often with fewer than 1,000 users remaining. Meanwhile, the world's [20 largest languages](#) are the native tongue of more than 3.7 billion people total. That's just 0.3% of the world's languages accounting for nearly half of the world's population!



<https://www.ethnologue.com/insights/how-many-languages/>

# Endangered Languages



# Why Multilingual NLP?



Speaking **more languages** means communicating with **more people**...  
...and reaching **more users and customers**...

# Why Multilingual NLP?

- ❖ Decreasing the digital language divide



# Why Multilingual NLP?

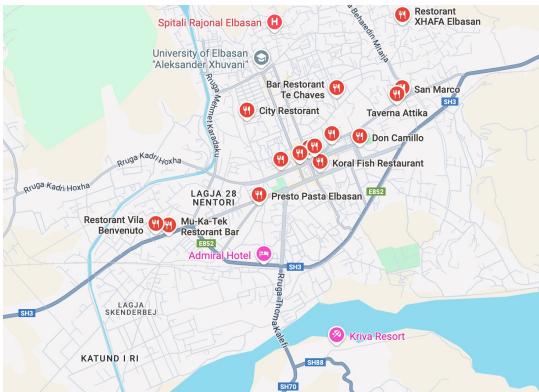
- ❖ Dealing with inequality of information (access)
- ❖ Mitigating cross-cultural biases
- ❖ Deploying language technology for underrepresented languages, dialects, minorities; societal impact
- ❖ Understanding cross-linguistic differences

“95% of all languages in use today will never gain traction online” (Andras Kornai)

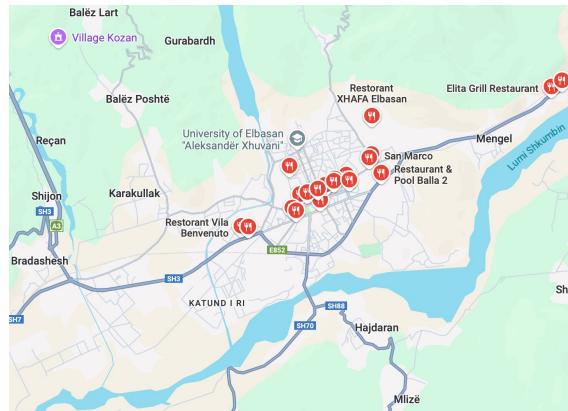
# Why Multilingual NLP?

- ❖ Inequality of information and representation can also affect how we understand places, events, processes...

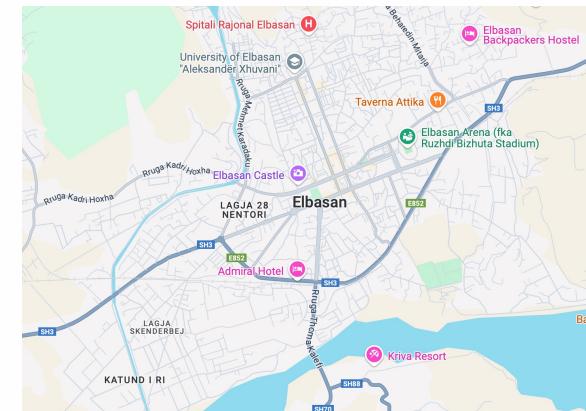
**We're in Elbasan searching for...**



**... restorantet (Al)**



**...Restaurant (En)**



**...éttermek (Hu)**

# How to Cope Multilingual?

## Better Data:

- ❖ – every piece of relevant data can help - be resourceful!
- ❖ – make data if necessary - be connected!

## Better Models or Algorithms:

- ❖ – sophisticated modelling/training methods - know NLP/ML!
- ❖ – linguistically informed methods - know linguistics!

## Better Deployment:

- ❖ – different situations require different solutions - be aware!

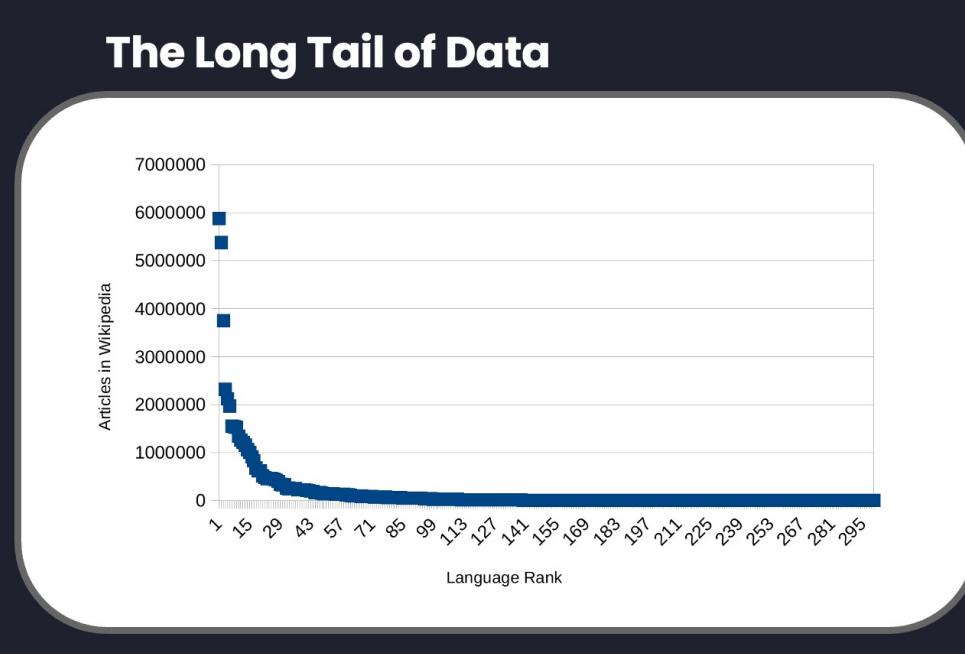


OLSCOIL NA GAILIMHE  
UNIVERSITY OF GALWAY

# Data

University  
*of*Galway.ie

# Issues: The Long Tail of Data



Even getting “raw” unannotated data is problematic for many languages...

Natural Language Processing for Multilingual Task-Oriented Dialogue  
Tutorial Abstract

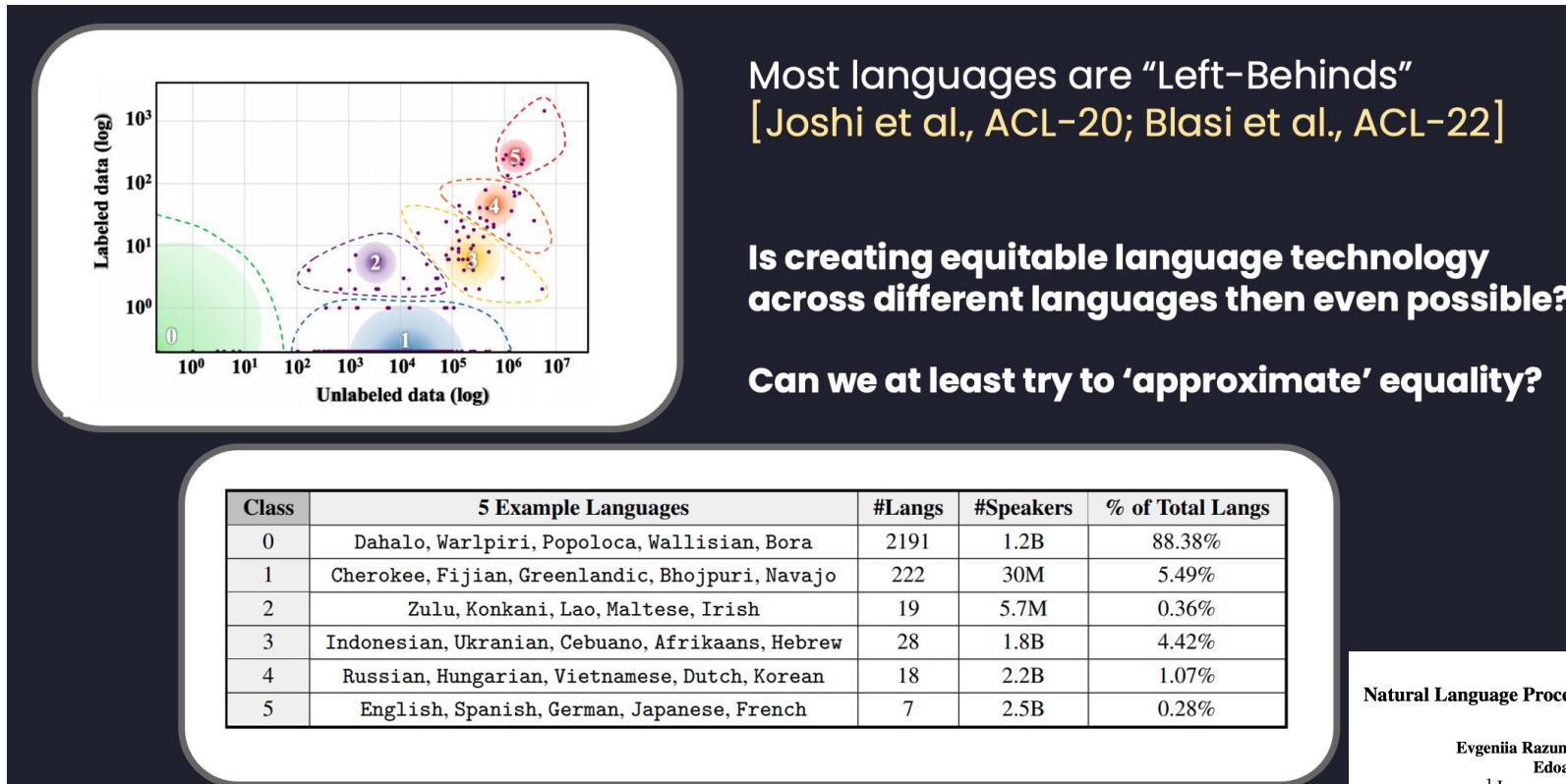
Evgenia Razumovskaya<sup>1</sup>, Goran Glavaš<sup>2</sup>, Olga Majewska<sup>1</sup>  
Edoardo Maria Ponti<sup>4,5</sup>, Ivan Vulić<sup>1,6</sup>

<sup>1</sup> Language Technology Lab, University of Cambridge

<sup>2</sup> Center for Artificial Intelligence and Data Science, University of Würzburg

<sup>4</sup>Mila – Quebec Artificial Intelligence Institute <sup>5</sup>McGill University <sup>6</sup>PolyAI Limited  
ter563,cm304,iv2501@cam.ac.uk goraninformatik.uni-mannheim.de  
edoardo-maria.ponti@mila.quebec edoardo-maria.ponti@mila.quebec

# Are all languages created equal?



Natural Language Processing for Multilingual Task-Oriented Dialogue  
Tutorial Abstract

Evgeniia Razumovskai<sup>1</sup>, Goran Glava<sup>2</sup>, Olga Majewska<sup>1</sup>  
Edoardo Maria Ponti<sup>4,5</sup>, Ivan Vulic<sup>1,6</sup>

<sup>1</sup> Language Technology Lab, University of Cambridge

<sup>2</sup> Center for Artificial Intelligence and Data Science, University of Würzburg  
<sup>4</sup>Mila – Quebec Artificial Intelligence Institute   <sup>5</sup>McGill University   <sup>6</sup>PolyAI Limited  
{er563,om304,iv250}@cam.ac.uk goran@informatik.uni-mannheim.de  
edoardo-maria.ponti@mila.quebec

# What we can do in Multilingual Setup?

- ❖ Many NLP tasks share common knowledge about language (e.g. linguistic representations, structural similarities)
- ❖ Languages share common structure (on the lexical, syntactic, and semantic level)
- ❖ Annotated data is rare, make use of as much supervision as available

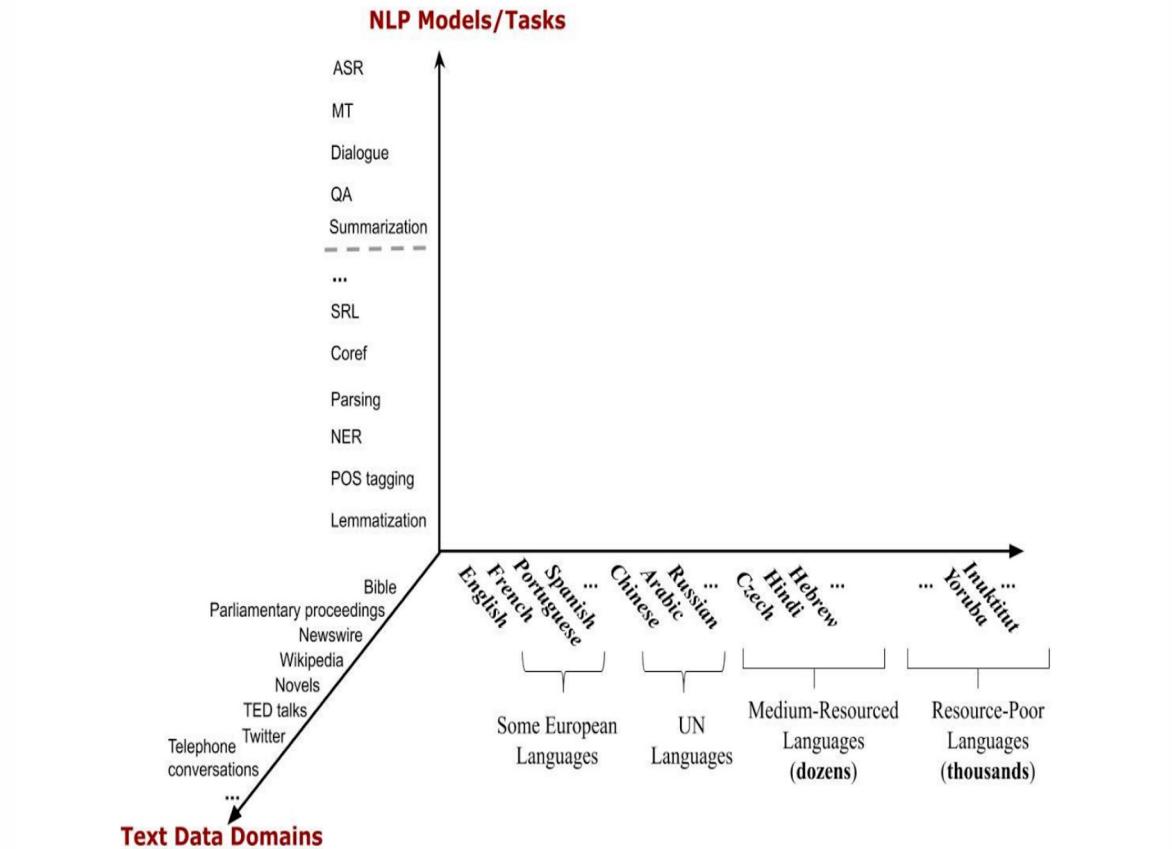


Image courtesy of *Yulia Tsvetkov*

# Data Creation/Curation

- ❖ What types of data? (monolingual? multilingual? annotated?)
- ❖ Where can we get it? (annotated data sources? curated text collections? scraping?)
- ❖ Can we create data? (efficient, high-quality creation strategies)
- ❖ How do we deal with the ethical issues? (working with communities, language ownership)

# Parallel Corpora

## CLASSIC SOUPS

Sm. Lg.

清 嫩 雞 湯	57.	House Chicken Soup (Chicken, Celery, Potato, Onion, Carrot) .....	1.50	2.75
雞 飯 湯	58.	Chicken Rice Soup .....	1.85	3.25
雞 麵 湯	59.	Chicken Noodle Soup .....	1.85	3.25
廣 東 雲 吞	60.	Cantonese Wonton Soup.....	1.50	2.75
蕃 茄 雲 湯	61.	Tomato Clear Egg Drop Soup .....	1.65	2.95
雲 吞 湯	62.	Regular Wonton Soup .....	1.10	2.10
酸 辣 湯	63.	Hot & Sour Soup .....	1.10	2.10
蛋 花 湯	64.	Egg Drop Soup.....	1.10	2.10
雲 蛋 湯	65.	Egg Drop Wonton Mix.....	1.10	2.10
豆 腐 菜 湯	66.	Tofu Vegetable Soup .....	NA	3.50
雞 玉 米 湯	67.	Chicken Corn Cream Soup .....	NA	3.50
蟹 肉 玉 米 湯	68.	Crab Meat Corn Cream Soup.....	NA	3.50
海 鮮 湯	69.	Seafood Soup.....	NA	3.50

# Do We have enough Parallel data?

<b>Parallel Corpus</b>	<b>Sentences</b>	<b>Parallel Corpus</b>	<b>Sentences</b>
Romanian-English	399,375	Greek-English	1,235,976
Bulgarian-English	406,934	Swedish-English	1,862,234
Slovene-English	623,490	Italian-English	1,909,115
Hungarian-English	624,934	German-English	1,920,209
Polish-English	632,565	Finnish-English	1,924,942
Lithuanian-English	635,146	Portuguese-English	1,960,407
Latvian-English	637,599	Spanish-English	1,965,734
Slovak-English	640,715	Danish-English	1,968,800
Czech-English	646,605	Dutch-English	1,997,775
Estonian-English	651,746	French-English	2,007,723

Europarl parallel data: <http://www.statmt.org/europarl/>

# The Open Parallel corpus

[opus.nlpl.eu](https://opus.nlpl.eu)

The screenshot shows the OPUS website interface. At the top, there is a navigation bar with links for Contribute, Publications, Corpora (which is highlighted in purple), and Dashboard. Below the navigation bar, there is a news section with 19 items. The main content area features the OPUS logo and a search form for finding corpora. The search form includes dropdowns for Source language (Afar) and Target language (Abkhazian), and a Search button. Below the search form, there is a section titled "An overview of the OPUS collection" which displays various statistics: 1,212 corpora, 58,851,021,412 total sentence pairs, and 747 languages available. A note states that the table displays 100 corpora, which make up 94.85% of the entire OPUS collection. To the right of this text is a table showing the top corpora in the OPUS collection, including their names, sentence counts, and percentages of the total.

Corpus	Sentences	% of OPUS
OpenSubtitles	20B	34.34
NLLB	13B	22.10
CCMatrix	11B	18.46
MultiCCAligned	2.2B	3.80864
ParaCrawl	1.5B	2.54692
DGT	1.1B	1.85690
MultiHPLT	897M	1.52502
XLEnt	883M	1.50013
MultiParaCrawl	789M	1.34012

# What if don't have parallel data?

<https://arxiv.org/pdf/1804.07755.pdf>

## Phrase-Based & Neural Unsupervised Machine Translation

<b>Guillaume Lample<sup>†</sup></b> Facebook AI Research Sorbonne Universités glample@fb.com	<b>Myle Ott</b> Facebook AI Research myleott@fb.com	<b>Alexis Conneau</b> Facebook AI Research Université Le Mans aconneau@fb.com
<b>Ludovic Denoyer<sup>†</sup></b> Sorbonne Universités ludovic.denoyer@lip6.fr	<b>Marc'Aurelio Ranzato</b> Facebook AI Research ranzato@fb.com	

<https://arxiv.org/pdf/1711.00043.pdf>

## UNSUPERVISED MACHINE TRANSLATION USING MONOLINGUAL CORPORA ONLY

**Guillaume Lample †‡ , Alexis Conneau † , Ludovic Denoyer ‡ , Marc'Aurelio Ranzato †**  
† Facebook AI Research,  
‡ Sorbonne Universités, UPMC Univ Paris 06, LIP6 UMR 7606, CNRS  
{gl, aconneau, ranzato}@fb.com, ludovic.denoyer@lip6.fr

<https://arxiv.org/pdf/1901.07291.pdf>

## Cross-lingual Language Model Pretraining

<b>Guillaume Lample*</b> Facebook AI Research Sorbonne Universités glample@fb.com	<b>Alexis Conneau*</b> Facebook AI Research Université Le Mans aconneau@fb.com
--	---

<https://arxiv.org/pdf/1710.11041.pdf>

## UNSUPERVISED NEURAL MACHINE TRANSLATION

**Mikel Artetxe, Gorka Labaka & Eneko Agirre**  
IXA NLP Group  
University of the Basque Country (UPV/EHU)  
{mikel.artetxe, gorka.labaka, e.agirre}@ehu.eus

**Kyunghyun Cho**  
New York University  
CIFAR Azrieli Global Scholar  
kyunghyun.cho@nyu.edu

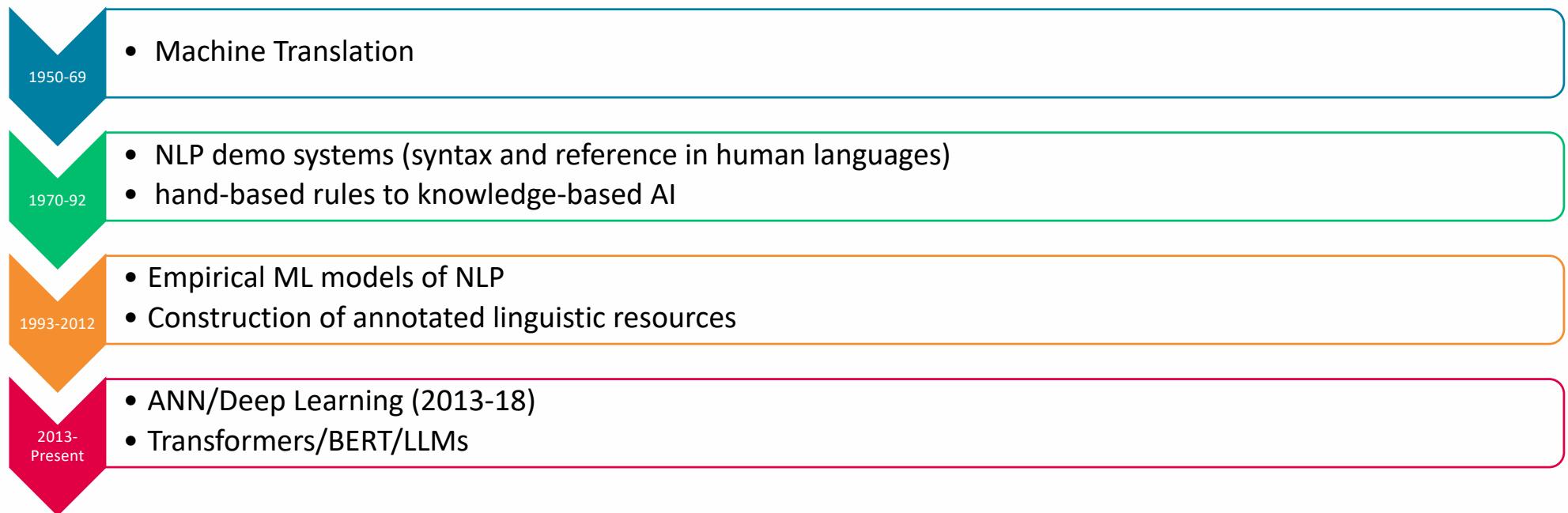


OLSCOIL NA GAILLIMHE  
UNIVERSITY OF GALWAY

# Technology

University  
*of*Galway.ie

# NLP History



# NLP Applications

- ❖ Text classification
- ❖ Token classification
- ❖ Question answering
- ❖ Causal language modelling
- ❖ Masked language modelling
- ❖ Translation
- ❖ Summarization
- ❖ Multiple choice

# Language Models

$$p(w_1, w_2, w_3, \dots, w_N) = p(w_1) p(w_2|w_1) p(w_3|w_1, w_2) \times \dots \times p(w_N|w_1, w_2, \dots, w_{N-1})$$

Conditional probability

Sentence: "the cat sat on the mat"

$$\begin{aligned} P(\text{the cat sat on the mat}) &= P(\text{the}) * P(\text{cat}|\text{the}) * P(\text{sat}|\text{the cat}) \\ &\quad * P(\text{on}|\text{the cat sat}) * P(\text{the}|\text{the cat sat on}) \\ &\quad * P(\text{mat}|\text{the cat sat on the}) \end{aligned}$$

Implicit order

# Language Models

$$\begin{aligned} P(w_1, w_2, \dots, w_n) &= p(w_1)p(w_2|w_1)p(w_3|w_1, w_2)\dots p(w_n|w_1, w_2, \dots, w_{n-1}) \\ &= \prod_{i=1}^n p(w_i|w_1, \dots, w_{i-1}) \end{aligned}$$

$S = \text{Where are we going}$

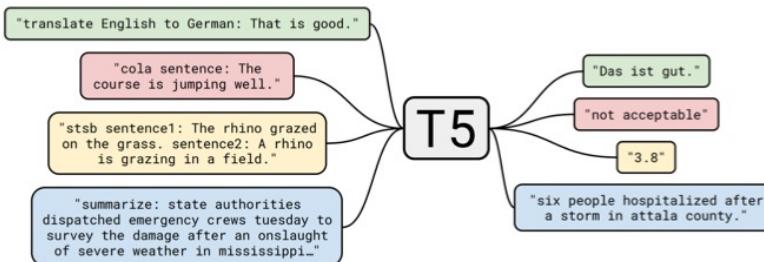
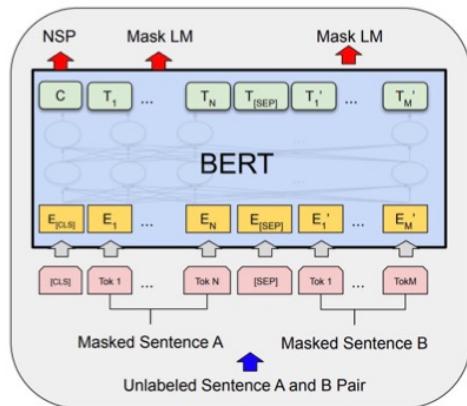
The diagram shows the sentence "S = Where are we going". A bracket above the words "are we" is labeled "Previous words (Context)". An arrow points from this bracket to the words "are we". Another arrow points from the word "going" to the word "going", which is highlighted in pink.

Previous words (Context)      Word being predicted

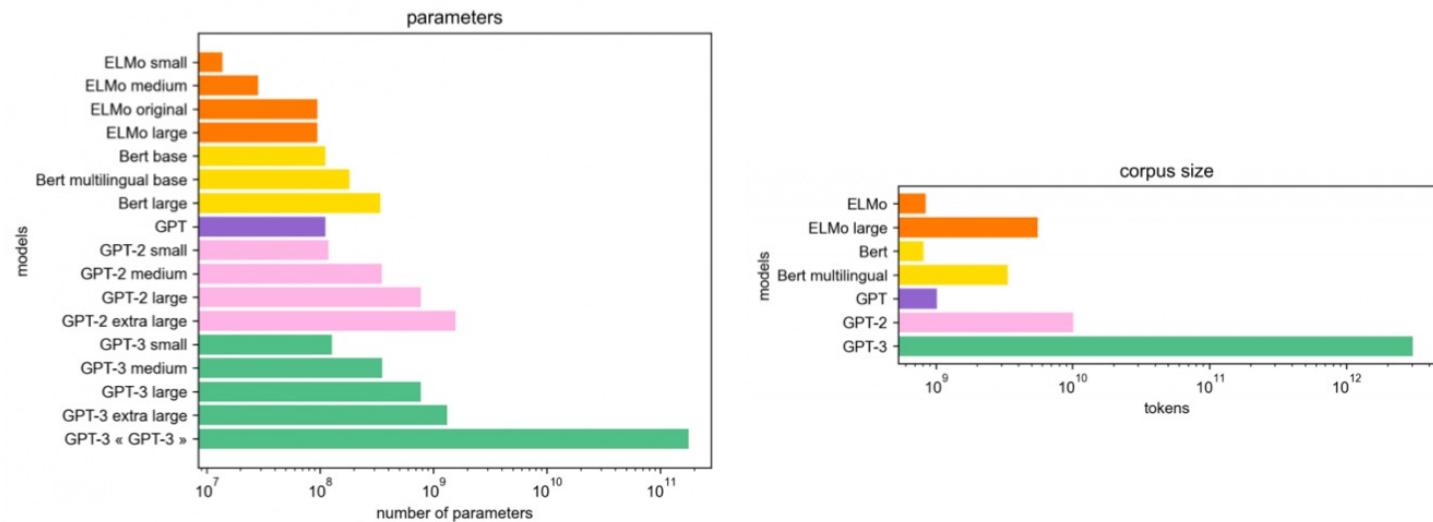
$$P(S) = P(\text{Where}) \times P(\text{are} \mid \text{Where}) \times P(\text{we} \mid \text{Where are}) \times P(\text{going} \mid \text{Where are we})$$

# Language Model

- Decoder-only models (GPT-x models)
- Encoder-only models (BERT, RoBERTa, ELECTRA)
- Encoder-decoder models (T5, BART)



# How large are ‘large’ LMs?



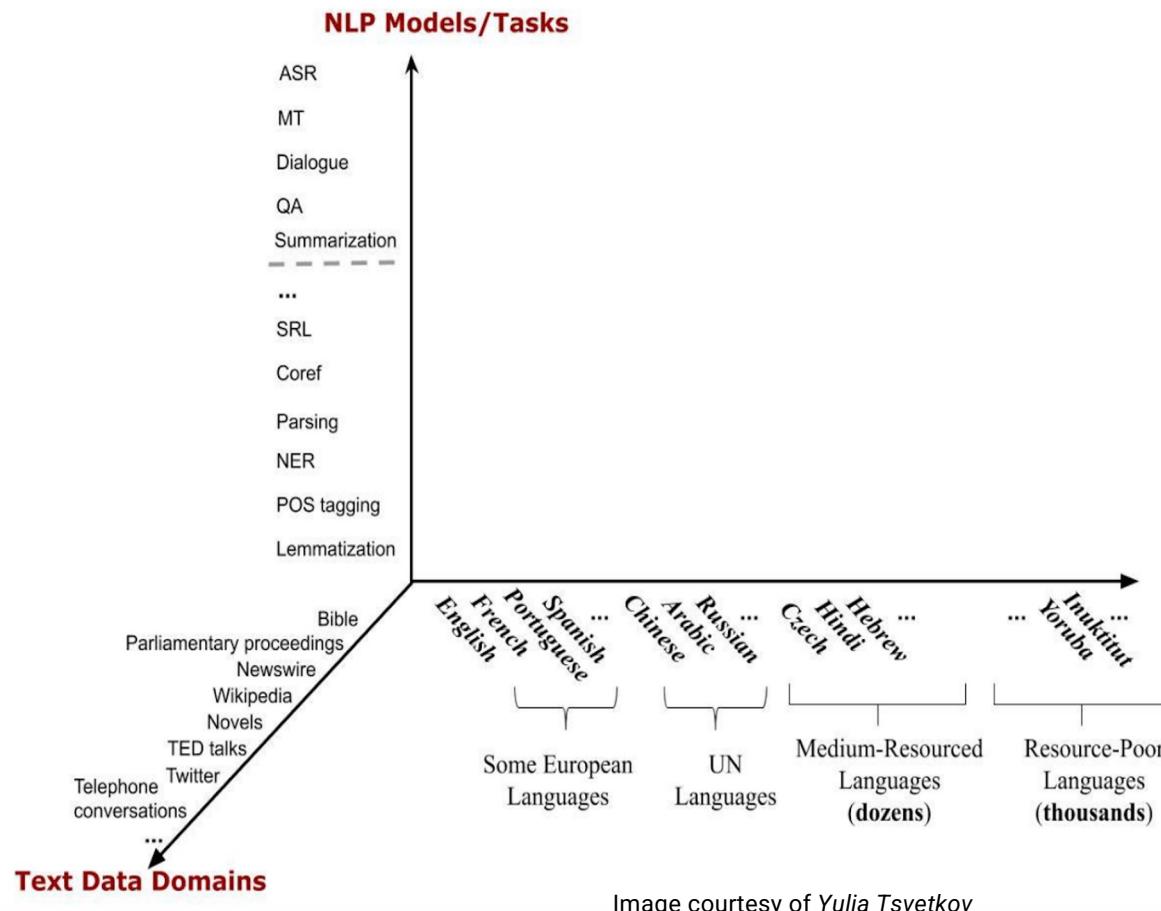
More recent models: PaLM (540B), OPT (175B), BLOOM (176B)...

Image source: <https://hellofuture.orange.com/en/the-gpt-3-language-model-revolution-or-evolution/>

# How large are ‘large’ LMs?

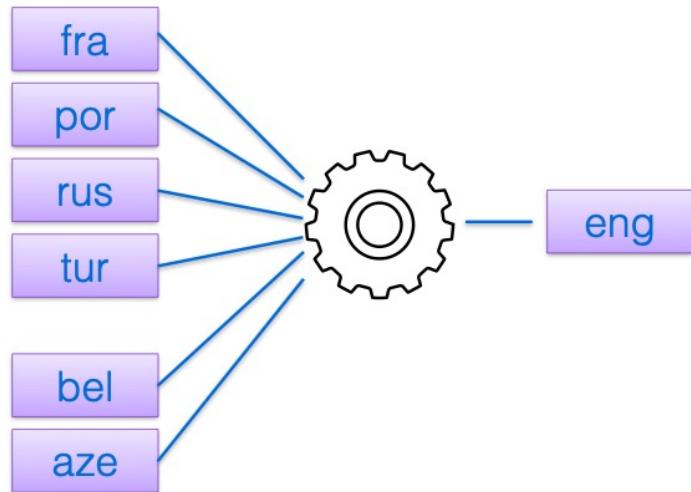
- Today, we mostly talk about two camps of models:
  - Medium-sized models: BERT/RoBERTa models (100M or 300M), T5 models (220M, 770M, 3B)
  - “Very” large LMs: models of 100+ billion parameters

# What we can do in Multilingual Setup?



# Multilingual training

- Train a large multi-lingual NLP system



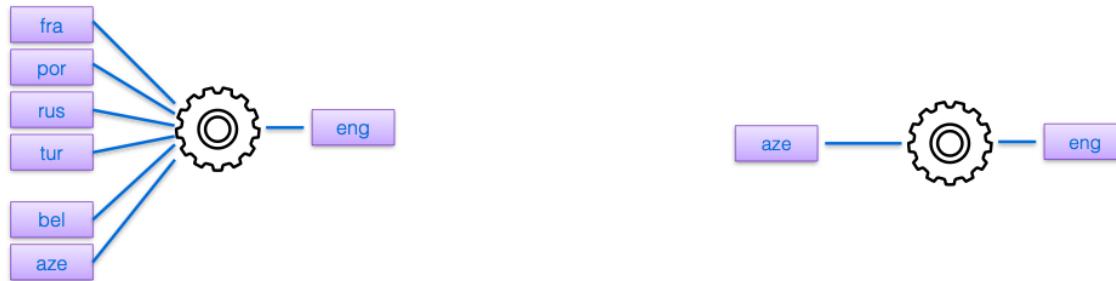
- Challenges: how to train effectively, how to ensure representation of low-resource languages

# Transfer Learning

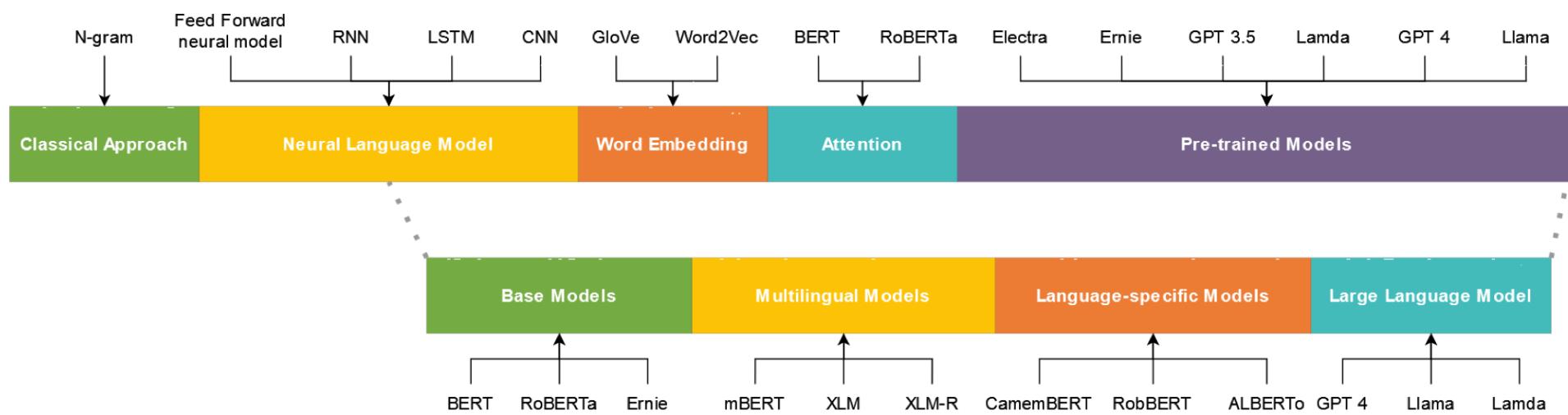
- Training on one (pair) language, transfer to another



- Train on many languages, transfer to another



# Evolution of Language Models in NLP

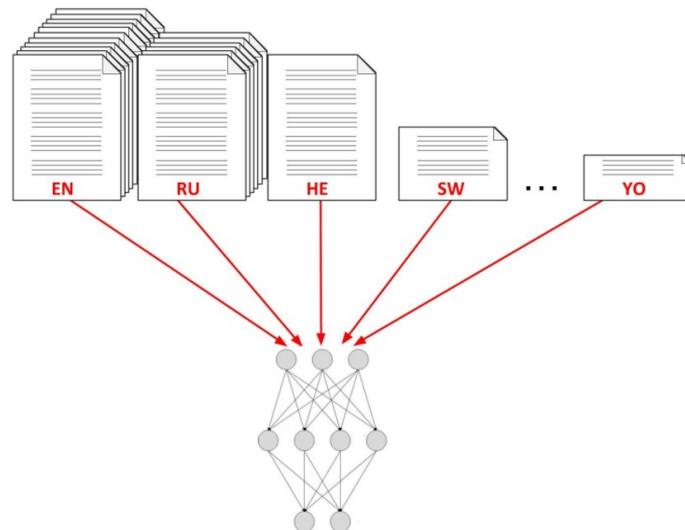


[https://hal.science/hal-04549672/file/Language\\_model\\_for\\_Multi\\_lingual\\_Tasks\\_a\\_survey\\_IJACSA\\_Journal\\_Version\\_Copy\\_.pdf](https://hal.science/hal-04549672/file/Language_model_for_Multi_lingual_Tasks_a_survey_IJACSA_Journal_Version_Copy_.pdf)

# Evolution of Language Models in NLP

Model	Type	Language	Year	Input Corpus Details
BERT [34]	Base model	English	2018	16GB of uncompressed text, BookCorpus (800M words), English Wikipedia (2500M words)
RoBERTa [35]	Base model	English	2019	160GB text: BookCorpus (800M words - 16GB) CC-News (63M English news articles - 76GB), OpenWebText (Web content extracted from URLs shared on Reddit - 38GB), Stories (subset of CommonCrawl data - 31GB)
ELECTRA [36]	Base model	English	2020	For experiments (Same Data as BERT): 3.3 billion tokens from Wikipedia and BooksCorpus. For Language model: extend the BERT dataset to 33B tokens by including data from ClueWeb; CommonCrawl; Gigaword
ERNIE [37]	Base model	English	2020	Processed Wikipedia Eng (4; 500M subwords and 140M entities)
ALBERT [38]	Base model	English	2020	16GB of uncompressed text consists of BookCorpus (800M words) English Wikipedia (2500M words)
UDify [39]	Base model	multilingual	2019	Full universal dependencies v2.3 corpus available on LINDAT, Arabic NYUAD, English ESL, Arabic NYUAD, French FTB, Hindi English HEINCS, Japanese BC-CWJ
XLNet [40]	Base model	English	2019	RACE Dataset, SQuAD, GLUE Dataset, ClueWeb09-B Dataset
mBERT	Multilingual Models	Cross-lingual	2018	Wikipedia, MultiUN, IIT Bombay corpus, OPUS, EUbookshop, OpenSubtitles, GlobalVoices, Kytea and PyThaiNLP5
XLM [41]	Multilingual Models	Cross-lingual	2019	Wikipedia, MultiUN, IIT Bombay corpus, OPUS, EUbookshop, OpenSubtitles, GlobalVoices, Kytea and PyThaiNLP5
CamemBERT [42]	Language-Specific model	French	2019	138GB of uncompressed text and 32.7B SentencePiece tokens consist of: French text extracted from CommonCrawlUnshuffled version of the French OSCAR corpus
RobBERT [43]	Language-Specific model	German	2020	39GB of uncompressed text consists of Dutch Section of OSCAR corpus (6.6B words - 39GB of texts)
BERTje [44]	Language-Specific model	Dutch	2019	Books: a collection of novels (4.4GB), TwNC a Dutch News Corpus (2.4GB), SoNaR-500 reference corpus (2.2GB), 4 Dutch news websites (1.6GB), Wikipedia dump (1.5GB), Total: 12 GB; 2.4B token
ALBERTo [45]	Language-Specific model	Italian	2019	TWITA:from twitter's official streaming API; 200M tweets and 191GB raw data
PhoBERT [46]	Language-Specific model	Vietnamese	2020	20GB texts: Vietnamese Wikipedia corpus (1GB)-(19GB) is a subset of a Vietnamese news corpus
BERT for Finnish [47]	Language-Specific model	Finnish	2019	Yle corpus, an archive of news and STT corpus of newswire articles
ParsBERT [48]	Language-Specific model	Persian	2021	In overall, more than 3M documents from Persian Wikipedia, BigBang Page, Chetor, Eligashm, Digikala, Ted Talks, books, Miras-Text
GPT-3.5	Large Language model	English	2022	vast amount of text data sourced from various publicly available sources on the internet including websites, books, articles, forums, and other forms of text content across different domains
Lamda [49]	Large Language model	English	2022	comprises 2.97 billion documents, 1.12 billion dialogues, and 13.39 billion dialogue utterances, totaling 1.56 trillion words.
Llama [50]	Large Language model	Multilingual	2023	English CommonCrawl, C4, Github, Wikipedia, Gutenberg and Books3, ArXiv, Stack Exchange
GPT 4 [51]	Large Language model	English	2023	vast amount of text data sourced from various publicly available sources on the internet including websites, books, articles, forums, and other forms of text content across different domains

# Joint Multilingual in a Nutshell



**Joint multilingual learning** – train a single model on a mix of datasets in all languages, to enable **data and parameter sharing** where possible

# Zero-shot vs Few Shot Learning

- Zero-shot learning – train a model in one domains and assume it generalizes more or less out-of-the-box in a low-resource domain
- Few shot learning – train a model in one domain and use only few examples from a low-resource domain to adapt it

# Zero-Shot transfer to (Low-Resource) Languages

**Step 1:**

**Train** a multilingual model.

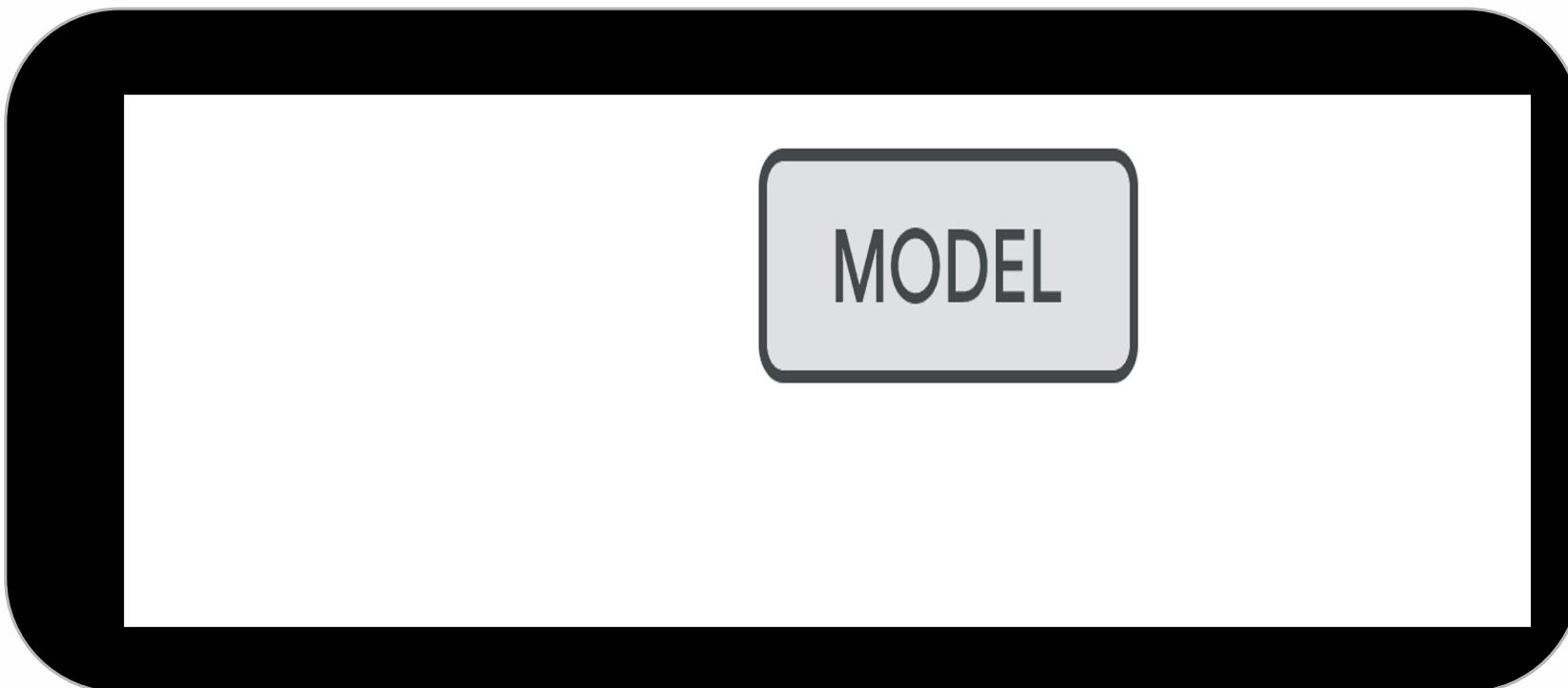
**Step 2:**

**Fine-tune** model on a **task** in a high resource **source language**.

**Step 3:**

Transfer and **evaluate** the model on a low resource **target language**.

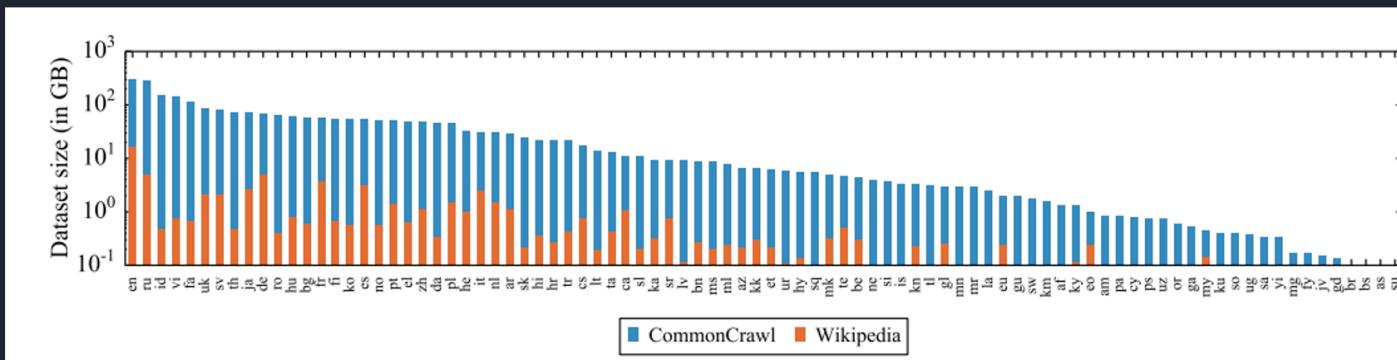
# Zero-Shot transfer to (Low-Resource) Languages



Animation: courtesy of Google Research

# Models

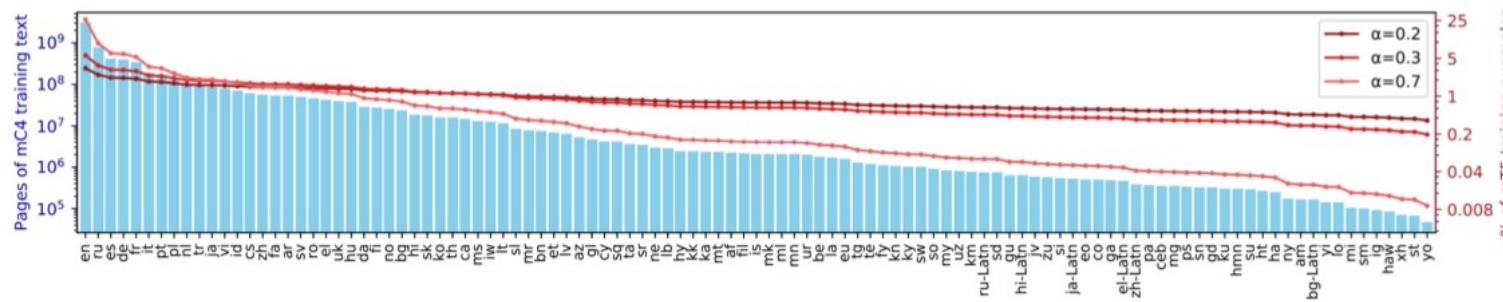
- Pretrained multilingual encoders are used as a baseline
  - mBERT ([Devlin et al., 2019](#)), XLM-R ([Conneau et al., 2020](#))
  - **Zero-shot cross-lingual transfer**
- The encoders are pre-trained on 100+ languages on the texts from Wikipedia and CC-100



Distribution of languages in CC-100 (Source: [Conneau et al., 2020](#))

# mC4 Dataset

- 107 languages, lower-resource languages upsampled based on their frequency in the dataset



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)

*mT5 paper, Sue et al., ACL 2021*



OLSCOIL NA GAILIMHE  
UNIVERSITY OF GALWAY

# Multilingual Translation

University  
*of*Galway.ie

# Why is it difficult to translate?

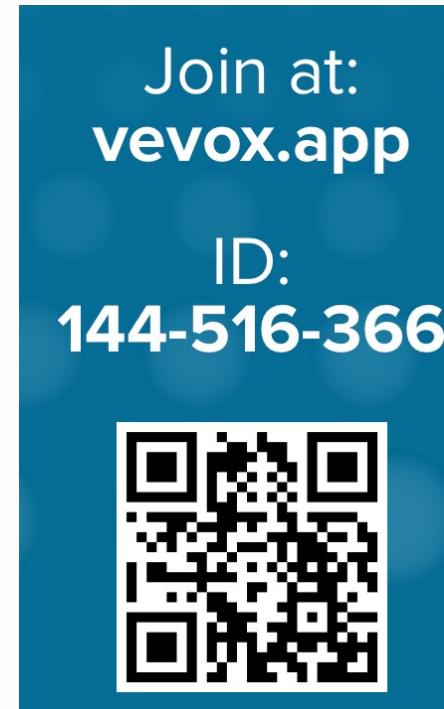
- ❖ Structural divergences
- ❖ Morphology
- ❖ Syntax

# Class Task

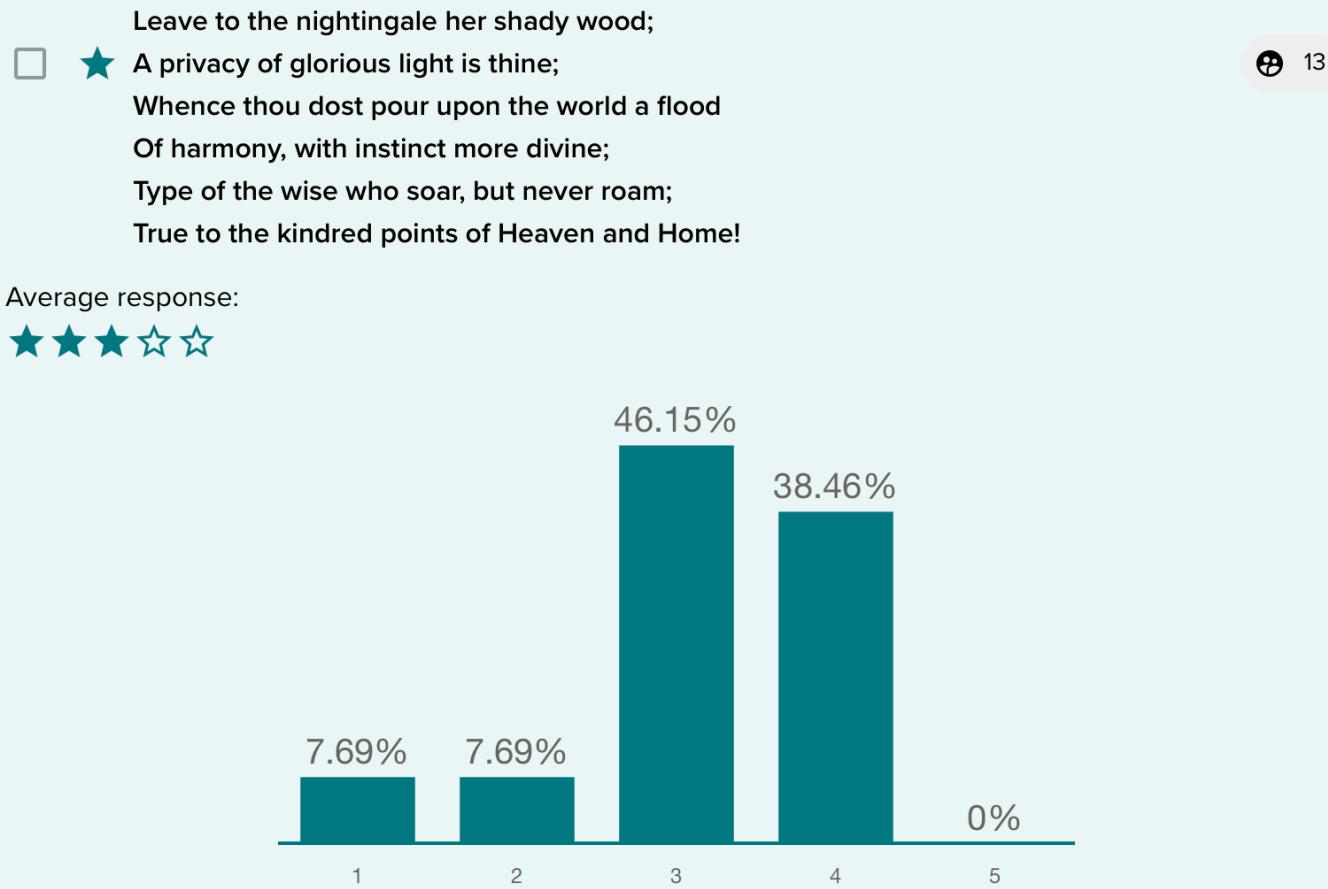
- ❖ Go to Google Translate
- ❖ Translate William Wordsworth's poetry into your mother tongue language



- ❖ Rate translation here



# Results



# Why is it difficult to translate?

This screenshot shows a translation interface. On the left, the source text "This cat is cute. Her name is Latte" is entered in English. On the right, the target language is set to Korean, and the translated text "이 고양이 너무 귀여워요. 이름은 라떼예요." is displayed, along with its phonetic transcription "i goyang-i neomu gwiyeowoyo. ileum-eun latteyeyo.". Both panels include a microphone icon for audio pronunciation, a "Send feedback" button, and other standard translation controls.

Detect language English German Spanish ↗

Korean Albanian German ↗

This cat is cute. Her name is Latte ×

이 고양이 너무 귀여워요. 이름은 라떼예요. ☆

i goyang-i neomu gwiyeowoyo. ileum-eun latteyeyo.

Send feedback

This screenshot shows a translation interface. On the left, the source text "이 고양이 너무 귀여워요. 이름은 라떼예요." is entered in Korean. On the right, the target language is set to English, and the translated text "This cat is so cute. His name is Latte." is displayed, along with its phonetic transcription "i goyang-i neomu gwiyeowoyo. ileum-eun latteyeyo.". Both panels include a microphone icon for audio pronunciation, a "Send feedback" button, and other standard translation controls.

Detect language Korean English German ↗

English Korean Albanian ↗

이 고양이 너무 귀여워요. 이름은 라떼예요. ×

This cat is so cute. His name is Latte. ☆

i goyang-i neomu gwiyeowoyo. ileum-eun latteyeyo.

Send feedback



OLSCOIL NA GAILIMHE  
UNIVERSITY OF GALWAY

# Multilingual Sentiment Analysis

University  
*of*Galway.ie

# Multilingual Sentiment Analysis



## Negative

I'm dissatisfied with your customer service.  
No one was able to help me with the  
problems I had with using your product.



## Neutral

The product has multiple features  
that are suitable for users with different  
levels of experience.

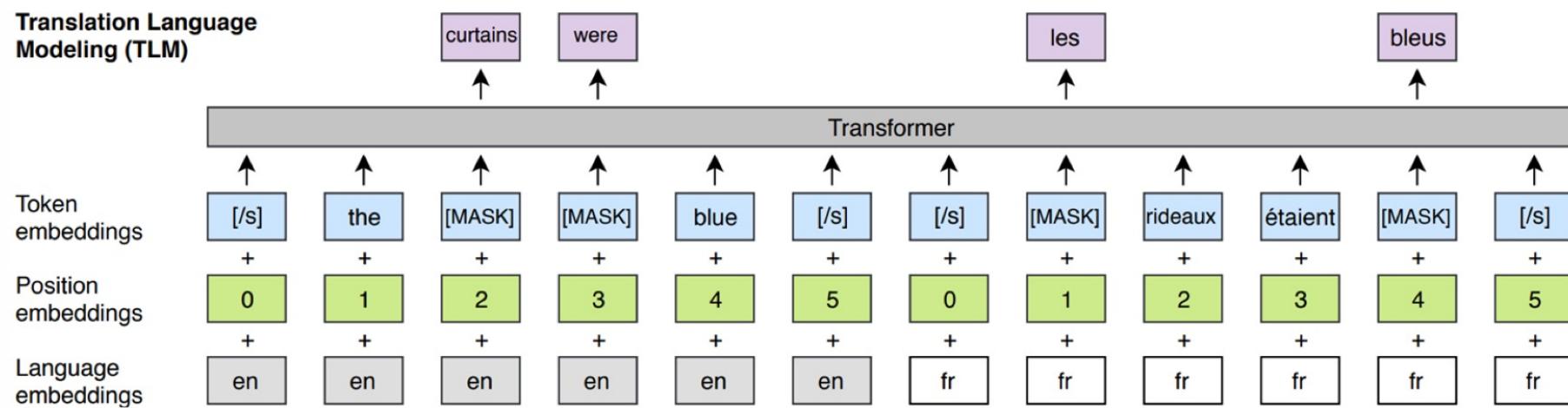


## Positive

I really enjoy how easy this product  
is to use and how it successfully helps  
my team complete their day-to-day tasks.

Source: Socialbakers

# XLM



# XLM-T

## **XLM-T: Multilingual Language Models in Twitter for Sentiment Analysis and Beyond**

**Francesco Barbieri<sup>♣</sup>, Luis Espinosa Anke<sup>◊</sup>, Jose Camacho-Collados<sup>◊</sup>**

<sup>♣</sup> Snap Inc., <sup>◊</sup> Cardiff NLP, School of Computer Science and Informatics, Cardiff University

<sup>♣</sup> Santa Monica, California, USA <sup>◊</sup> Cardiff, Wales, United Kingdom

fbarbieri@snap.com, {espinosa-ankel,camachocolladosj}@cardiff.ac.uk

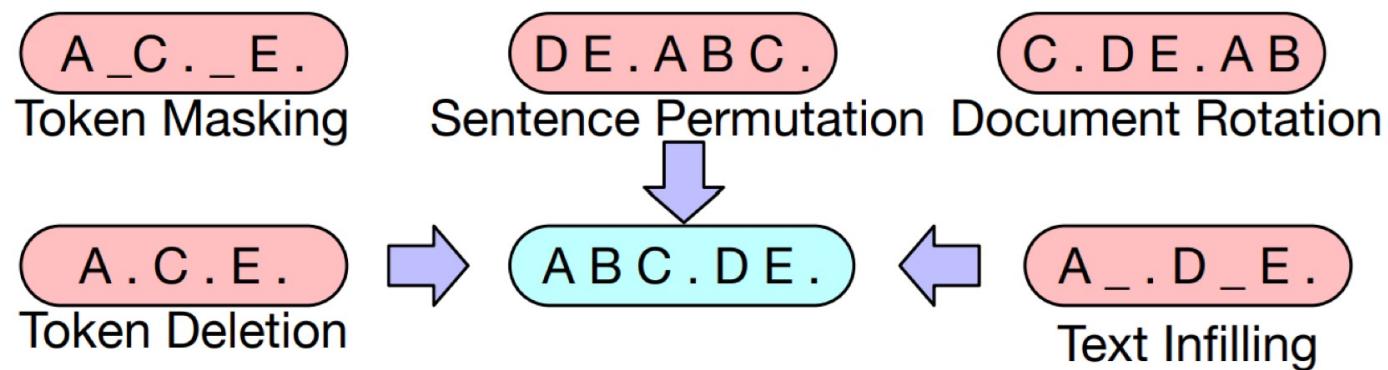
# XLM-RoBERTa

- ❖ This is a multilingual XLM-roBERTa-base model trained on ~198M tweets
- ❖ XLM-RoBERTa is a multilingual model trained on 100 different languages.
- ❖ Unlike some XLM multilingual models, it does not require lang tensors to understand which language is used, and should be able to determine the correct language from the input ids.

# XLM-RoBERTa

- ❖ Finetuned for sentiment analysis
- ❖ The sentiment fine-tuning was done on 8 languages
  - Ar → Arabic
  - En → English
  - Fr → French
  - De → German (Deutsch)
  - Hi → Hindi
  - It → Italian
  - Sp → Spanish (Español) \*(Note: The standard ISO 639-1 code for Spanish is "Es," but "Sp" is sometimes used informally.)\*
  - Pt → Portuguese (Português)

# BART



# BART

## mBART: Multi-Lingual BART

- Training on CC25 corpus
- Corpus of 25 languages
- A subset of *Common Crawl*
- A crawl of the internet

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

# Q&A

[jamal.nasir@universityofgalway.ie](mailto:jamal.nasir@universityofgalway.ie)

