

NATURAL LANGUAGE PROCESSING 2022/23

# Text Classification

MASSIVE: A 1M-Example Multilingual Natural Language Understanding  
Dataset with 51 Typologically-Diverse Languages

Margarida Vieira  
up201907907

Nuno Costa  
up201906272

Tiago Rodrigues  
up201907021

# Multiclass Classification Problem

This task's goal is to be able to **classify the intent of an utterance** - Intelligent Voice Assistant single-shot interactions.

As such, we'll be solving a **Multiclass Classification Problem**.

51 languages

55 slot types

60 intents

18 domains

**utterance:** what is the temperature in new york



**intent:** weather\_query

**MASSIVE** | a multilingual dataset

A **1M-example dataset** composed of realistic, **human-created virtual assistant utterance text**, mainly in the interrogative and imperative form.

It has been translated by professional translators and it spans 51 different languages.

# Results

Model	Multinomial NB			XGBoost Classifier			Logistic Regression		
	High	Low	Avg	High	Low	Avg	High	Low	Avg
TF-IDF	0.799 <b>ru-RU</b>	0.403 <b>my-MM</b>	0.739	0.789 <b>da-DK</b>	0.433 <b>my-MM</b>	0.752	0.817 <b>id-ID</b>	0.396 <b>my-MM</b>	0.760
Bag of Words	0.720 <b>pl-PL</b>	0.366 <b>my-MM</b>	0.655	0.782 <b>az-AZ</b>	0.429 <b>my-MM</b>	0.743	0.773 <b>en-US</b>	0.370 <b>my-MM</b>	0.711

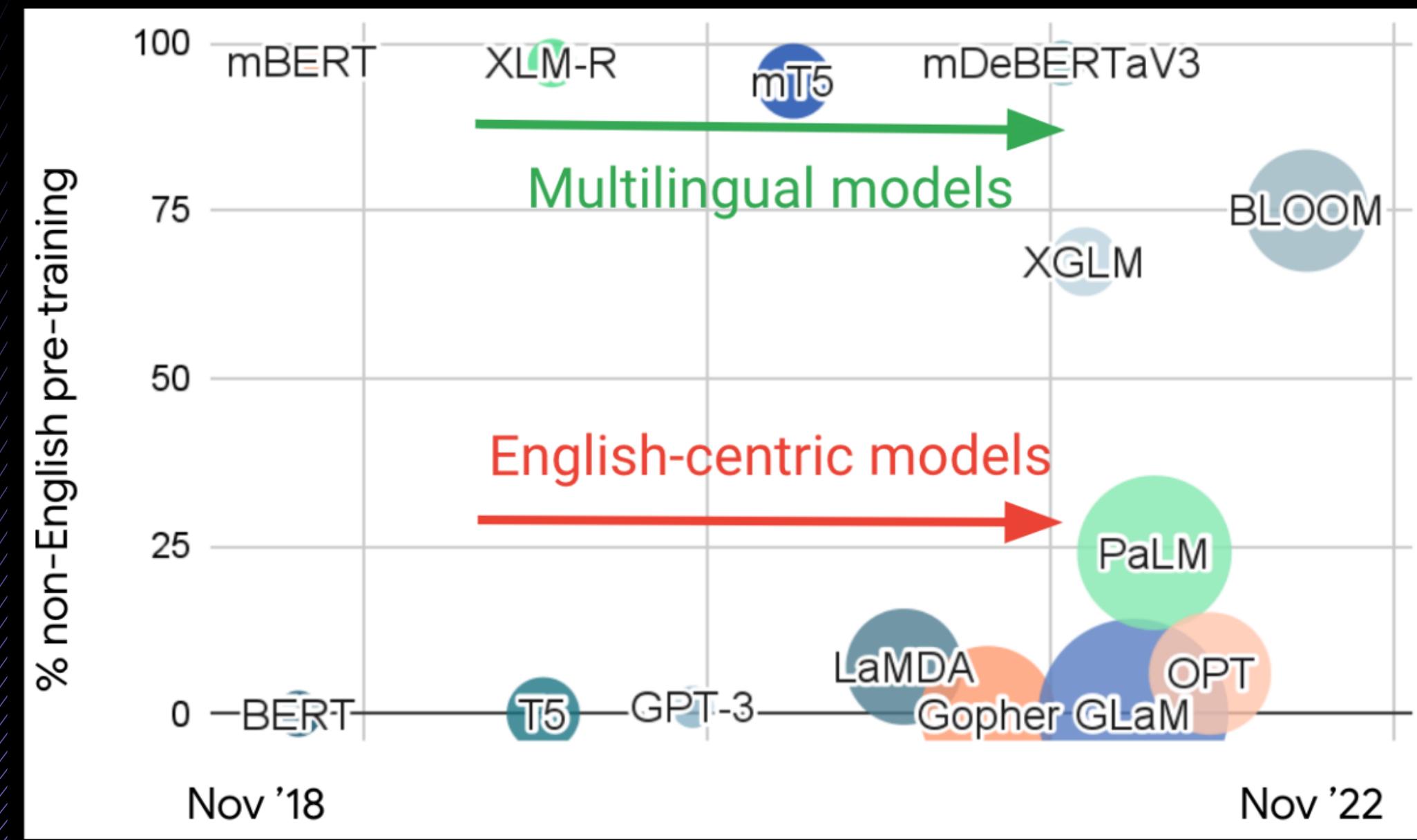
	High	Low	Avg
Our Baseline	0.745 <b>pl-PL</b>	0.119 <b>ja-JP</b>	0.638
XLM-R Base	0.883 <b>en-US</b>	0.772 <b>km-KH</b>	0.851

Used **Randomized Search** for *hyperparameter tuning*

Results are shown for the **best hyperparameters** with *lowercasing, tokenization, and stemming*

# Multilingual Models

## State-of-the-art



# Multilingual Models

## State-of-the-art

### XLM-RoBERTa (encoder)

multilingual version of RoBERTa,  
pre-trained on 2.5TB of filtered  
CommonCrawl data containing 100  
languages

### mDeBERTa (encoder)

pre-trained with the CC100 multilingual  
dataset, which comprises monolingual  
data for 100+ languages

### mBart (encoder-decoder)

sequence-to-sequence multilingual  
encoder-decoder model  
primarily intended for translation tasks  
unfeasible to use due to size

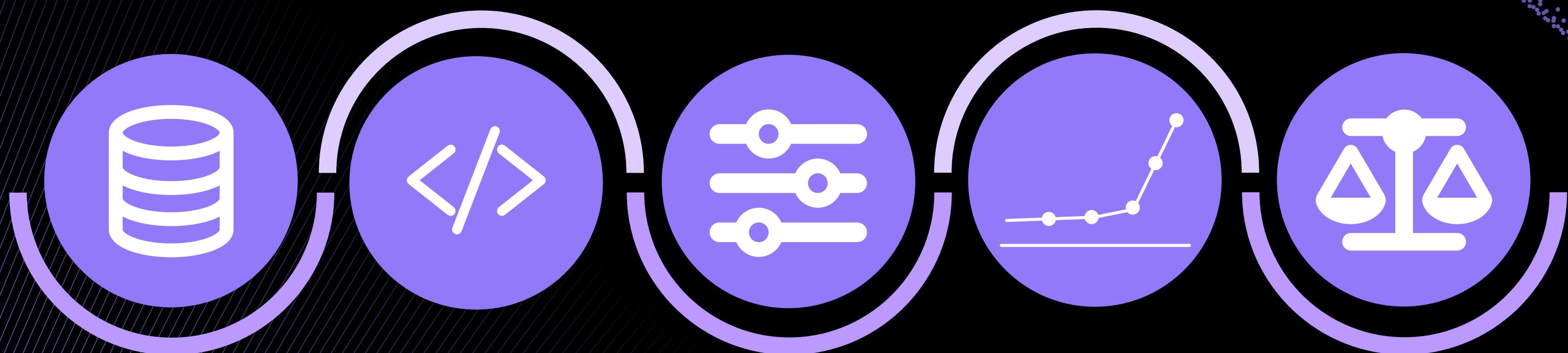
### BL00M (decoder)

largest multilingual open-source model to date  
autoregressive LLM with 176B parameters  
pre-trained on and able to generate text in 46 natural  
languages and 13 programming languages  
unfeasible to use due to size

### mT5 (encoder-decoder)

text-to-text transformer  
pre-trained on a new Common Crawl-based mC4 dataset,  
covering 101 languages  
unusual for classification tasks, since it is trained to  
generate the literal text of the class label instead of a class  
index

# Our Approach



## Small data preprocessing

punctuation removal  
lowercasing  
label encoding

## Tokenization

using the model specific tokenizer, taking advantage of AutoTokenizer from the **transformers** library

## Hyperparameter Definition

usage of **AdamW Optimizer** with:  
**learning rate** = 2e-5  
 $\epsilon$  = 1e-8  
and running with a **batch size** = 32

## Model Training

ran for **2 epochs** for most models using a **CrossEntropy loss** function, taking advantage of **DataParallelism** (T4x2 GPU)

## Model Evaluation

measure the model's **accuracy values** for the whole test set and also for each particular language to compare with the Paper's results

# Results

## XLM-RoBERTa `xlm-roberta-base`

	High	Low	Avg
<b>xlm-roberta-base</b> Our Model	0.875 <b>en-US</b>	0.564 <b>km-KH</b>	0.801
<b>xlm-roberta-base</b> Paper's Model	0.883 <b>en-US</b>	0.772 <b>km-KH</b>	0.851

## mDeBERTa `microsoft/mdeberta-v3-base`

	High	Low	Avg
<b>mdeberta-v3-base</b>	0.880 <b>en-US</b>	0.508 <b>km-KH</b>	0.803
<b>xlm-roberta-base</b> Paper's Model	0.883 <b>en-US</b>	0.772 <b>km-KH</b>	0.851

	Train Loss	Val. Loss	Val. Accuracy	Train Time	Val. Time
<b>Epoch 1</b>	0.861	0.674	0.832	3:22:03	0:11:30
<b>Epoch 2</b>	0.374	0.648	0.841	3:21:31	0:11:19

	Train Loss	Val. Loss	Val. Accuracy	Train Time	Val. Time
<b>Epoch 1</b>	0.837	0.721	0.830	3:48:51	0:13:21
<b>Epoch 2</b>	0.353	0.714	0.844	3:49:01	0:13:28

# Results

## DistilBERT

`distilbert-base-multilingual-cased`

	High	Low	Avg
<b>distilbert-base-multilingual-cased</b>	0.835 <b>en-US</b>	0.030 <b>am-ET</b>	0.756
<b>xlm-roberta-base</b> Paper's Model	0.883 <b>en-US</b>	0.772 <b>km-KH</b>	0.851

## InfoXLM

`microsoft/infoxlm-base`

	High	Low	Avg
<b>microsoft/infoxlm-base</b>	0.871 <b>en-US</b>	0.514 <b>km-KH</b>	0.796
<b>xlm-roberta-base</b> Paper's Model	0.883 <b>en-US</b>	0.772 <b>km-KH</b>	0.851

	Train Loss	Val. Loss	Val. Accuracy	Train Time	Val. Time
<b>Epoch 1</b>	1.221	0.850	0.775	1:37:05	0:05:06
<b>Epoch 2</b>	0.631	0.798	0.796	1:36:58	0:05:09
<b>Epoch 3</b>	0.475	0.796	0.803	1:38:02	0:05:11

	Train Loss	Val. Loss	Val. Accuracy	Train Time	Val. Time
<b>Epoch 1</b>	0.854	0.786	0.803	3:14:51	0:12:52
<b>Epoch 2</b>	0.457	0.736	0.828	3:15:03	0:12:47

# Conclusions

Hard to match the baseline results for the Paper given their significant processing power advantage to train XLM-Roberta especially when it comes to hyperparameter tuning and the number of epochs the model is trained for

---

Processing the text might affect the model's ability to accurately predict results. Besides this, by encoding the labels we may loose some helpful context that the model could use for guidance

---

The results clearly showcase the strength of Transformers when compared with more traditional and simpler models on more complex tasks such as multilingual and multilabel classification