



# The Road Towards Massively Multilingual Large Language Models

François Yvon - ISIR (Sorbonne Université, CNRS, INSERM)

BUCC @ LREC 2024, Torino

# Monolingual LLMs

---

Learning parameters on large **monolingual corpora** with **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

**Longtemps je me suis couché [mask]** - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

(eg. GPT\*, OPT, GPTJ, PALM\*, LLAMA\*, Mistral\*)

# Monolingual LLMs

---

Learning parameters on large **monolingual corpora** with **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

Longtemps je me suis couché [mask] - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

(eg. GPT\*, OPT, GPTJ, PALM\*, LLAMA\*, Mistral\*)

2. Predict missing word given bidirectional contexts : **pure encoder**

Longtemps je me suis couché [mask] bonne heure- unmask='de', train  $P_{\theta}(w_t | w_{-t})$

(eg. BERT, Roberta, CamemBERT, FlauBERT, etc)

# Monolingual LLMs

---

Learning parameters on large **monolingual corpora** with **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

Longtemps je me suis couché [mask] - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

(eg. GPT\*, OPT, GPTJ, PALM\*, LLAMA\*, Mistral\*)

2. Predict missing word given bidirectional contexts : **pure encoder**

Longtemps je me suis couché [mask] bonne heure- unmask='de', train  $P_{\theta}(w_t | w_{-t})$

(eg. BERT, Roberta, CamemBERT, FlauBERT, etc)

3. Denoising sequence to sequence : **encoder-decoder**

Longtemts je couché suis de bnone heur || Longtemps je me suis couché de bonne heure

train  $P_{\theta}(w | \tilde{w}) = \prod_t P(w_t | w_{<t}, \tilde{w})$  (eg. BART, T5, etc)

# Monolingual LLMs

---

Learning parameters on large **monolingual corpora** with **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

Longtemps je me suis couché [mask] - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

(eg. GPT\*, OPT, GPTJ, PALM\*, LLAMA\*, Mistral\*)

2. Predict missing word given bidirectional contexts : **pure encoder**

Longtemps je me suis couché [mask] bonne heure- unmask='de', train  $P_{\theta}(w_t | w_{-t})$

(eg. BERT, Roberta, CamemBERT, FlauBERT, etc)

3. Denoising sequence to sequence : **encoder-decoder**

Longtemts je couché suis de bnone heur || Longtemps je me suis couché de bonne heure

train  $P_{\theta}(w | \tilde{w}) = \prod_t P(w_t | w_{<t}, \tilde{w})$  (eg. BART, T5, etc)

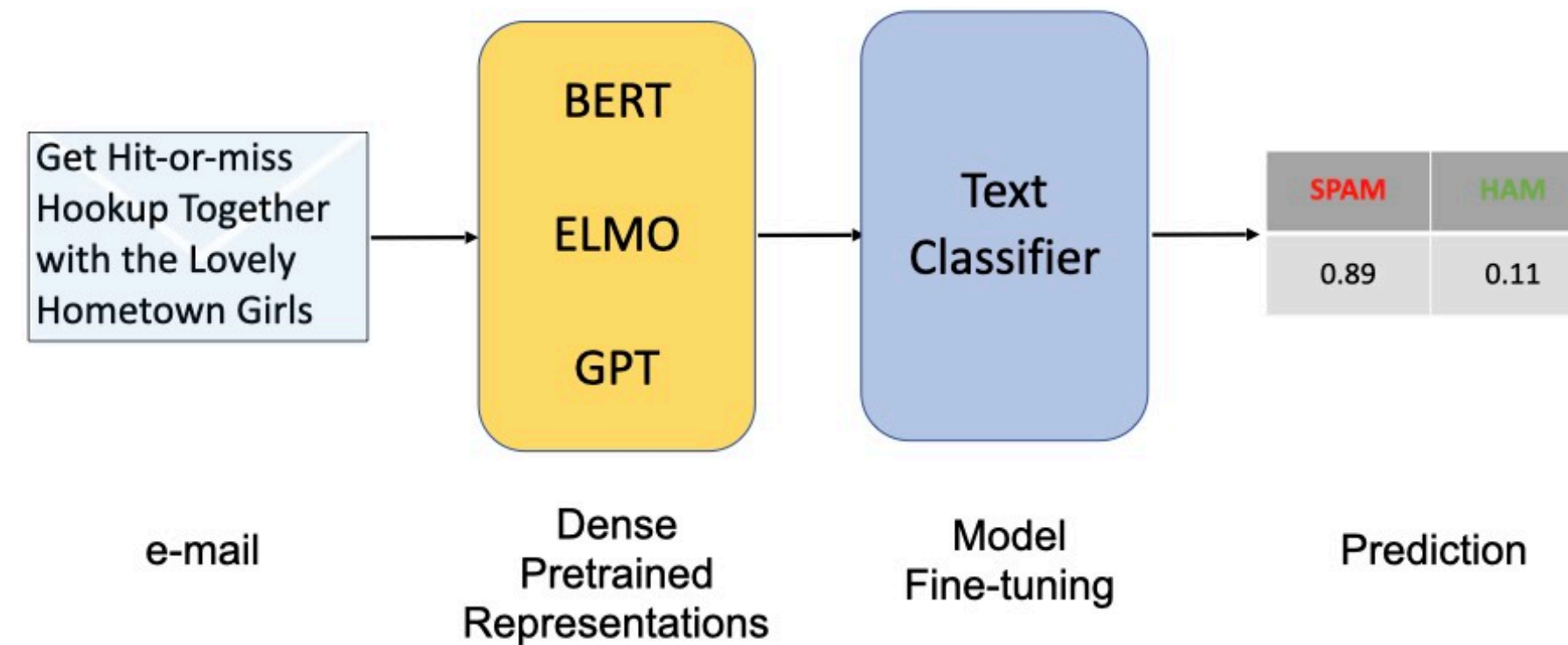
## † Alt. auxiliary tasks

- NSP
- span prediction (eg. SpanBERT)
- corruption detection (eg. Electra)

# Using Monolingual LLMs

---

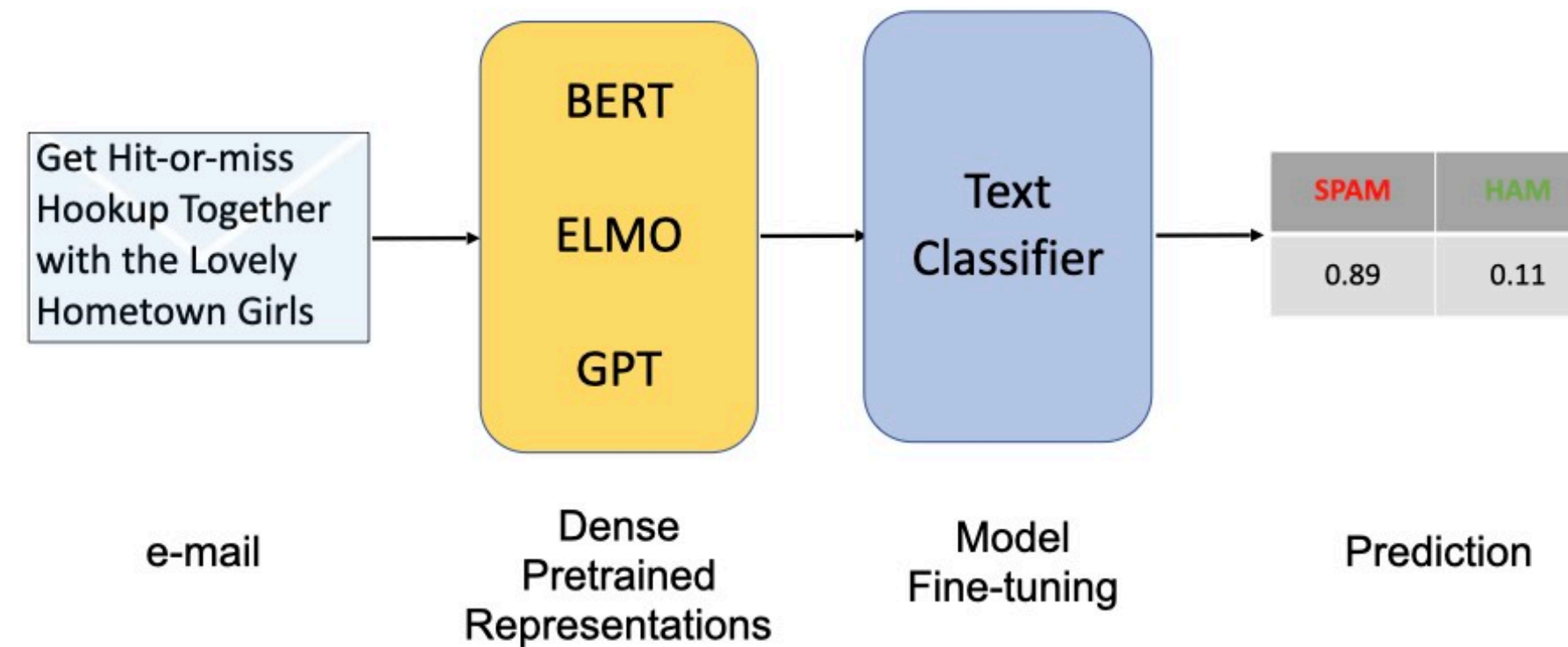
1. **Fine-tuned task-adapted model:**  $h_{\phi, \theta} = h_{\phi}(f_{\theta}(w); c)$



1+2. Multi-task fine-tuning with prompts a.k.a **instruction tuning**

# Using Monolingual LLMs

1. **Fine-tuned task-adapted model:**  $h_{\phi, \theta} = h_{\phi}(f_{\theta}(w); c)$



2. Multi-purpose text generation via **prompting**

**GEN** Of course. In Chorukor, Monday is ilopagar, Tuesday ilopager, ...  
Wednesday ilopagur, Thursday ilopagir ...

**Q&A** Answer this : What are the birth date and place of Ludvík Vaculík? ...  
23 July 1926, in Brumov, Moravia

**SA** "This Czech writer has written some the most wonderful French novels."  
is a positive comment? ... [Yes | No]

1+2. Multi-task fine-tuning with prompts a.k.a **instruction tuning**

# “Holistic” evaluation

---

## 81 models

AI21 Labs / J1-Jumbo v1 (178B)  
AI21 Labs / J1-Large v1 (7.5B)  
AI21 Labs / J1-Grande v1 (17B)  
AI21 Labs / J1-Grande v2 beta (17B)  
AI21 Labs / Jurassic-2 Jumbo (178B)  
AI21 Labs / Jurassic-2 Grande (17B)  
AI21 Labs / Jurassic-2 Large (7.5B)  
Aleph Alpha / Luminous Base (13B)  
Aleph Alpha / Luminous Extended (30B)  
Aleph Alpha / Luminous Supreme (70B)  
neurips / Local service  
Anthropic / Anthropic-LM v4-s3 (52B)  
Anthropic / Anthropic Claude 2.0  
Anthropic / Anthropic Claude v1.3  
Anthropic / Anthropic Claude Instant V1  
UC Berkeley / Koala (13B)  
BigScience / BLOOM (176B)  
BigScience / BLOOMZ (176B)  
BigScience / T0pp (11B)  
BigCode / SantaCoder (1.1B)  
BigCode / StarCoder (15.5B)  
Cerebras / Cerebras GPT (6.7B)  
Cerebras / Cerebras GPT (13B)  
Cohere / Cohere xlarge v20220609 (52.4B)  
Cohere / Cohere large v20220720 (13.1B)  
Cohere / Cohere medium v20220720 (6.1B)  
Cohere / Cohere small v20220720 (410M)  
Cohere / Cohere xlarge v20221108 (52.4B)  
Cohere / Cohere medium v20221108 (6.1B)  
Cohere / Cohere Command beta (6.1B)  
Cohere / Cohere Command beta (52.4B)

## 73 scenarios

Question answering

- MMLU
- BoolQ
- NarrativeQA
- NaturalQuestions (closed-book)
- NaturalQuestions (open-book)
- QuAC
- HellaSwag
- OpenbookQA
- TruthfulQA

Information retrieval

- MS MARCO (regular)
- MS MARCO (TREC)

Summarization

- CNN/DailyMail
- XSUM

Sentiment analysis

- IMDB

Toxicity detection

- CivilComments

Text classification

- RAFT

Aspirational scenarios

- Data-to-text generation
- Fact verification
- Copywriting
- Story generation

## 65 metrics

Accuracy

- none
- Quasi-exact match
- F1
- Exact match
- RR@10
- NDCG@10
- ROUGE-2
- Bits/byte
- Exact match (up to specified indicator)
- Absolute difference
- F1 (set match)
- Equivalent
- Equivalent (chain of thought)
- pass@1

Calibration

- Max prob
- 1-bin expected calibration error
- 10-bin expected calibration error
- Selective coverage-accuracy area
- Accuracy at 10% coverage
- 1-bin expected calibration error (after Platt scaling)
- 10-bin Expected Calibration Error (after Platt scaling)
- Platt Scaling Coefficient
- Platt Scaling Intercept

Robustness

- Quasi-exact match (perturbation: typos)
- F1 (perturbation: typos)
- Exact match (perturbation: typos)
- RR@10 (perturbation: typos)

## Multifacet Evaluation

- tasks, bias, fairness, openness, etc.

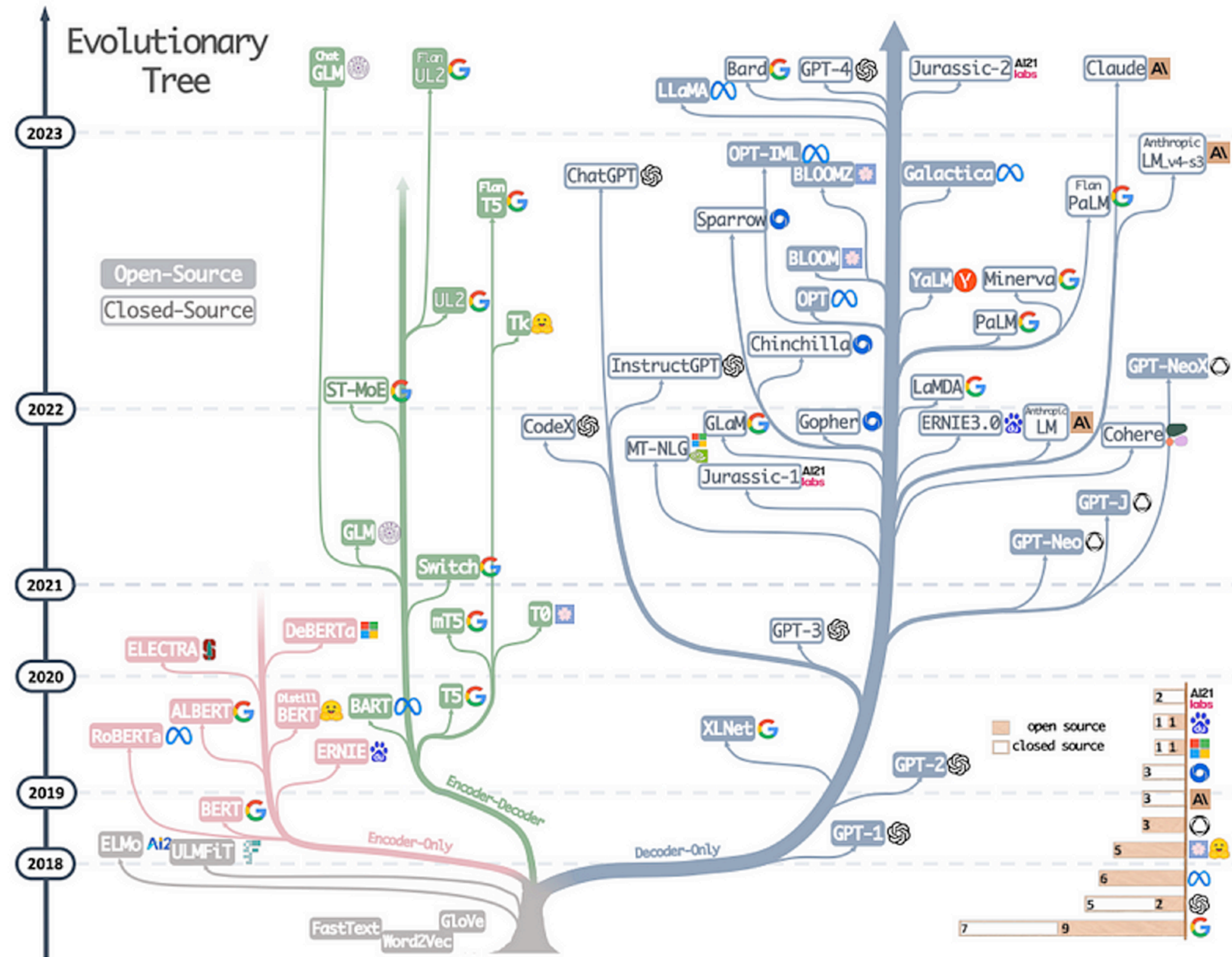
[Holistic Evaluation of Large Language Models \(Liang et al, 2022\)](#)





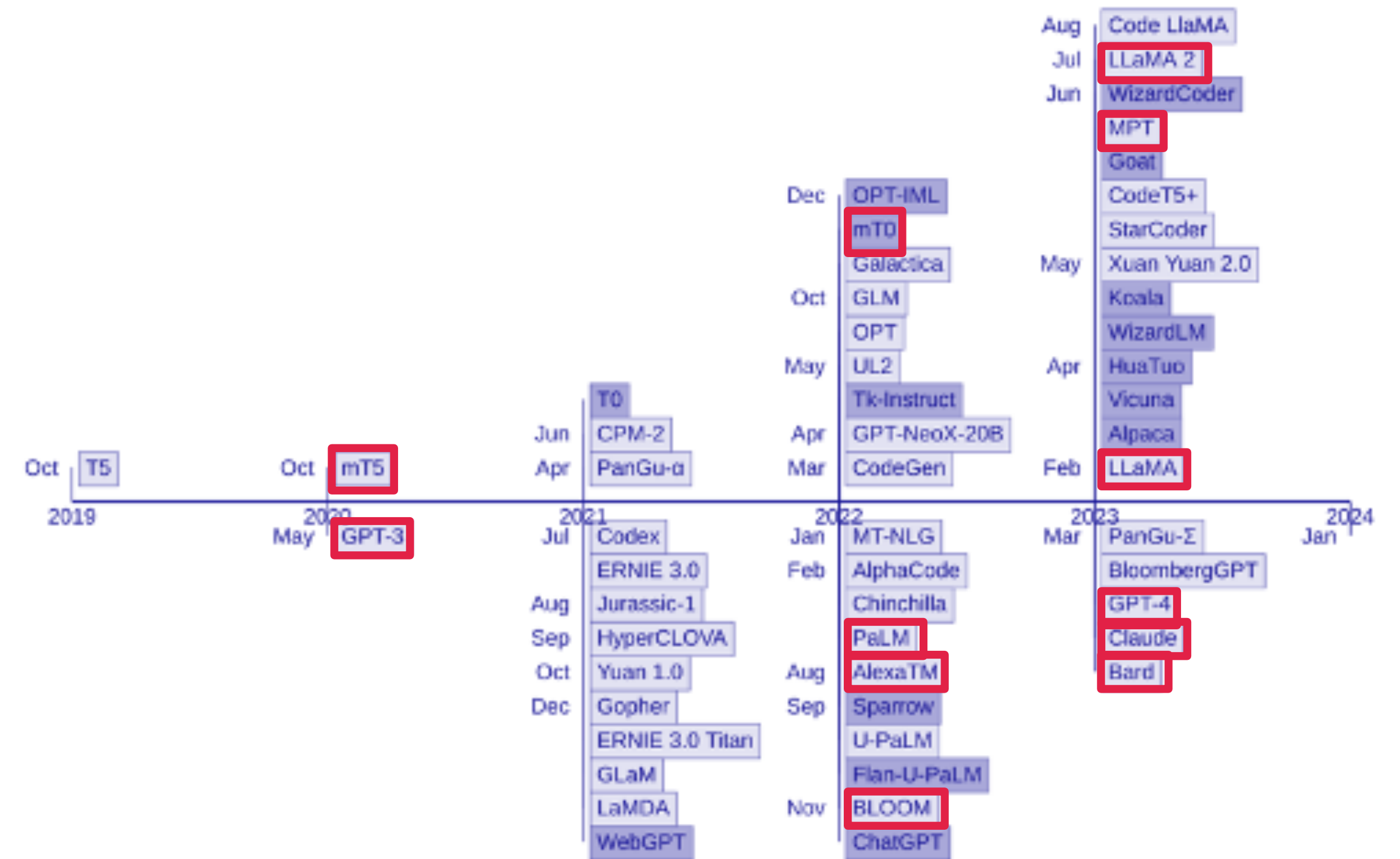


# Types of LLMs



source: <https://abiaryan.com/posts/intro-llms/>

Naveed, H., Khan, A. U., Qiu, S., Saqib, M., Anwar, S., Usman, M., ... & Mian, A. (2023). A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*.



## Also

- derivatives (fine-tuned, aligned)
- augmented (RAG, Tools, KBs)
- speech, image, video
- multimodal (text+image, +video)
- code, bio-chem, material, actions, ...

**This talk: multilingual (texts)**

# Multilingual LLMs: mLLMs

---

Learning parameters with **large multilingual corpora**, **auxiliary tasks** and **natural annotations**

# Multilingual LLMs: mLLMs

---

Learning parameters with **large multilingual corpora**, **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

Longtemps je me suis couché [mask] - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

Tutti ormai lo chiamavano don [mask] - unmask='Ciccio', train  $P_{\theta}(w_t | w_{<t})$

(eg. mGPT, XGLM, BLOOM, PALM-2, Falcon, LLM-jp, LLaMa2, ...)

# Multilingual LLMs: mLLMs

---

Learning parameters with **large multilingual corpora**, **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

Longtemps je me suis couché [mask] - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

Tutti ormai lo chiamavano don [mask] - unmask='Ciccio', train  $P_{\theta}(w_t | w_{<t})$

(eg. mGPT, XGLM, BLOOM, PALM-2, Falcon, LLM-jp, LLaMa2, ...)

2. Predict missing word given bidirectional contexts : **pure encoder**

(eg. mBERT, XLM-Roberta, etc)

# Multilingual LLMs: mLLMs

---

Learning parameters with **large multilingual corpora**, **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

Longtemps je me suis couché [mask] - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

Tutti ormai lo chiamavano don [mask] - unmask='Ciccio', train  $P_{\theta}(w_t | w_{<t})$

(eg. mGPT, XGLM, BLOOM, PALM-2, Falcon, LLM-jp, LLaMa2, ...)

2. Predict missing word given bidirectional contexts : **pure encoder**

(eg. mBERT, XLM-Roberta, etc)

3. Denoising sequence to sequence : **encoder-decoder**

(eg. mBART, mT5, ..., M2M)

# Multilingual LLMs: mLLMs

---

Learning parameters with **large multilingual corpora**, **auxiliary tasks** and **natural annotations**

1. Predict next word given prefix: **pure decoder**

Longtemps je me suis couché [mask] - unmask='de', train  $P_{\theta}(w_t | w_{<t})$

Tutti ormai lo chiamavano don [mask] - unmask='Ciccio', train  $P_{\theta}(w_t | w_{<t})$

(eg. mGPT, XGLM, BLOOM, PALM-2, Falcon, LLM-jp, LLaMa2, ...)

2. Predict missing word given bidirectional contexts : **pure encoder**

(eg. mBERT, XLM-Roberta, etc)

3. Denoising sequence to sequence : **encoder-decoder**

(eg. mBART, mT5, ..., M2M)

+ **Complementary objectives to align languages**

- parallel corpora (TLM loss, MT loss)
- bilingual dictionaries
- synthetic mixed-language data
- script normalization (romanization, transliteration)



# Using Multilingual LLMs

---

1. Fine-tune task-adapted models on L1 and process L2 with zero-shot model transfer

Only requires annotations in L1

# Using Multilingual LLMs

---

## 1. Fine-tune task-adapted models on L1 and process L2 with zero-shot model transfer

Only requires annotations in L1

## 2. Multi-purpose, multilingual text generation via prompting

Translate into English

“ By the end of the year, we will have seven new pharmacists. ” :

D’ici la fin de l’année, nous aurons sept nouveaux pharmaciens.

# Using Multilingual LLMs

---

## 1. Fine-tune task-adapted models on L1 and process L2 with zero-shot model transfer

Only requires annotations in L1

## 2. Multi-purpose, multilingual text generation via prompting

Translate into English

“ By the end of the year, we will have seven new pharmacists. ” :

D’ici la fin de l’année, nous aurons sept nouveaux pharmaciens.

### **mLLMs are a blessing**

- hardly more difficult than mLLMs
- excel in multilingual tasks
- enable X-lingual transfer

# mLLMs need multilingual texts

---

\_Tous \_les \_être s \_humain s \_na issent \_libre s \_et \_ég aux \_en \_digni té \_et \_en  
\_droits . \_Ils \_sont \_do u és \_de \_raison \_et \_de \_conscience \_et \_doivent \_agir  
\_les \_uns \_en vers \_les \_autres \_dans \_un \_esprit \_de \_frater n ité .

\_Všichni \_lidé \_rod í \_se \_svobod ní \_a \_sobě \_rov ní \_co \_do \_d ů stoj nosti \_a  
\_práv . \_Jsou \_na dán i \_rozum em \_a \_s vědomí m \_a \_mají \_spolu \_jedna t \_v  
\_du chu \_brat r ství .

\_Tutti \_gli \_esse ri \_umani \_na scono \_liberi \_ed \_e gu ali \_in \_digni tà \_e \_diritti  
\_Es si \_sono \_do tati \_di \_ragione \_e \_di \_coscienza \_e \_devono \_agir e \_gli \_uni  
\_verso \_gli \_altri \_in \_spirito \_di \_fra tella nza .

# mLLMs need multilingual texts

---

\_Tous \_les \_être s \_humain s \_na issent \_libre s \_et \_ég aux \_en \_digni té \_et \_en  
\_droits . \_Ils \_sont \_do u és \_de \_raison \_et \_de \_conscience \_et \_doivent \_agir  
\_les \_uns \_en vers \_les \_autres \_dans \_un \_esprit \_de \_frater n ité .

\_Všichni \_lidé \_rod í \_se \_svobod ní \_a \_sobě \_rov ní \_co \_do \_d ů stoj nosti \_a  
\_práv . \_Jsou \_na dán i \_rozum em \_a \_s vědomí m \_a \_mají \_spolu \_jedna t \_v  
\_du chu \_brat r ství .

\_Tutti \_gli \_esse ri \_umani \_na scono \_liberi \_ed \_e gu ali \_in \_digni tà \_e \_diritti  
\_Es si \_sono \_do tati \_di \_ragione \_e \_di \_coscienza \_e \_devono \_agir e \_gli \_uni  
\_verso \_gli \_altri \_in \_spirito \_di \_fra tella nza .

# mLLMs need multilingual texts

---

\_Tous \_les \_être s \_humain s \_na issent \_libre s \_et \_ég aux \_en \_digni té \_et \_en  
\_droits . \_Ils \_sont \_do u és \_de \_raison \_et \_de \_conscience \_et \_doivent \_agir  
\_les \_uns \_en vers \_les \_autres \_dans \_un \_esprit \_de \_frater n ité .

\_Všichni \_lidé \_rod í \_se \_svobod ní \_a \_sobě \_rov ní \_co \_do \_d ů stoj nosti \_a  
\_práv . \_Jsou \_na dán i \_rozum em \_a \_s vědomí m \_a \_mají \_spolu \_jedna t \_v  
\_du chu \_brat r ství .

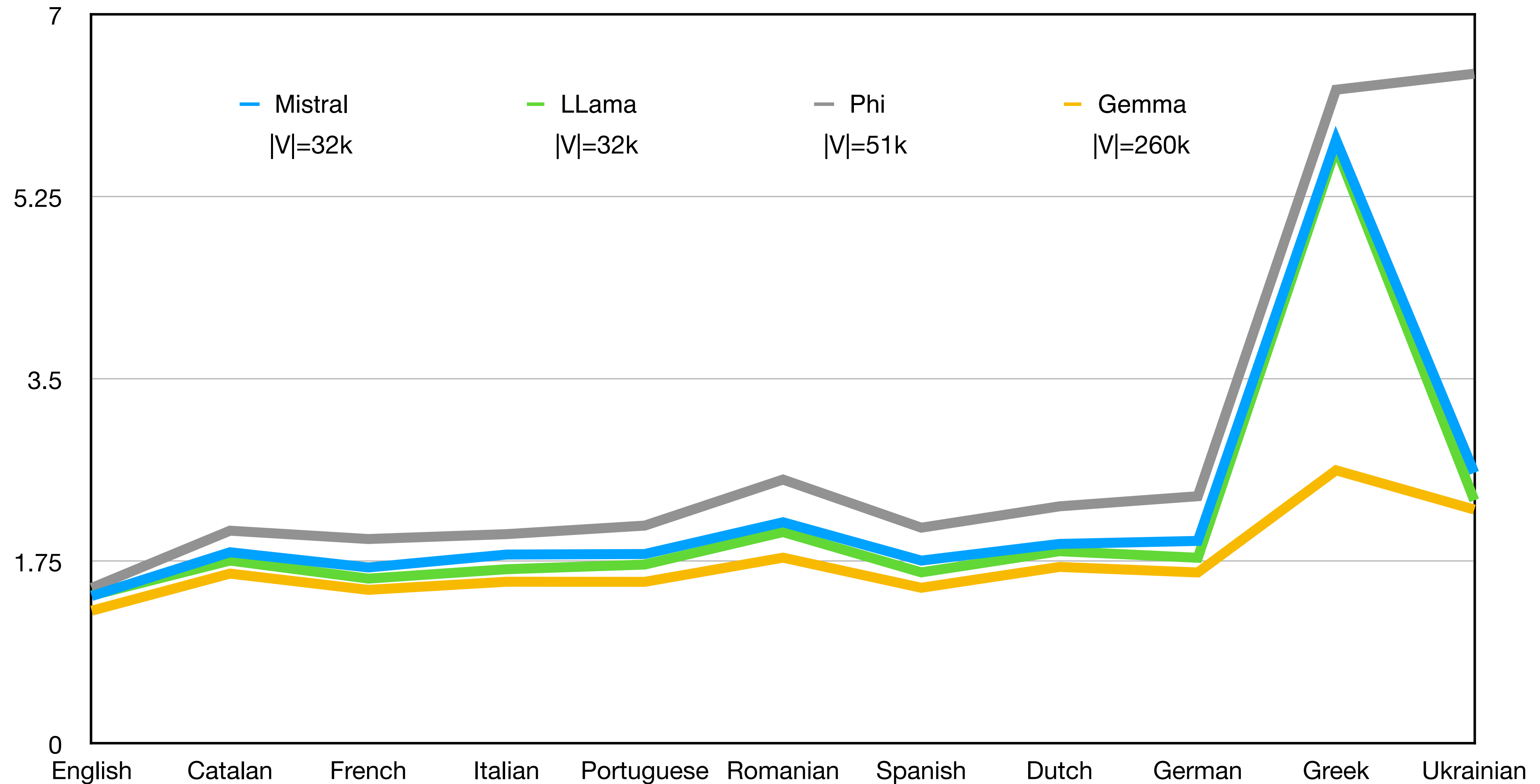
\_Tutti \_gli \_esse ri \_umani \_na scono \_liberi \_ed \_e gu ali \_in \_digni tà \_e \_diritti  
\_Es si \_sono \_do tati \_di \_ragione \_e \_di \_coscienza \_e \_devono \_agir e \_gli \_uni  
\_verso \_gli \_altri \_in \_spirito \_di \_fra tella nza .

## Subword tokenizers are trainable

- require **mixed-language, mixed-script** training corpora
- **parameter** sharing for same-script languages
- **larger language get more units**, are better segmented, perform better

# LMM tokenization: fertility

---



# Multilingualism in tokenizers

GlotScript / README.md

Preview Code Blame

Raw [copy] [download] [edit] [menu]

## Install

### from pip

```
pip3 install GlotScript
```

### from git

```
pip3 install GlotScript@git+https://github.com/cisnlp/GlotScript
```

## Usage

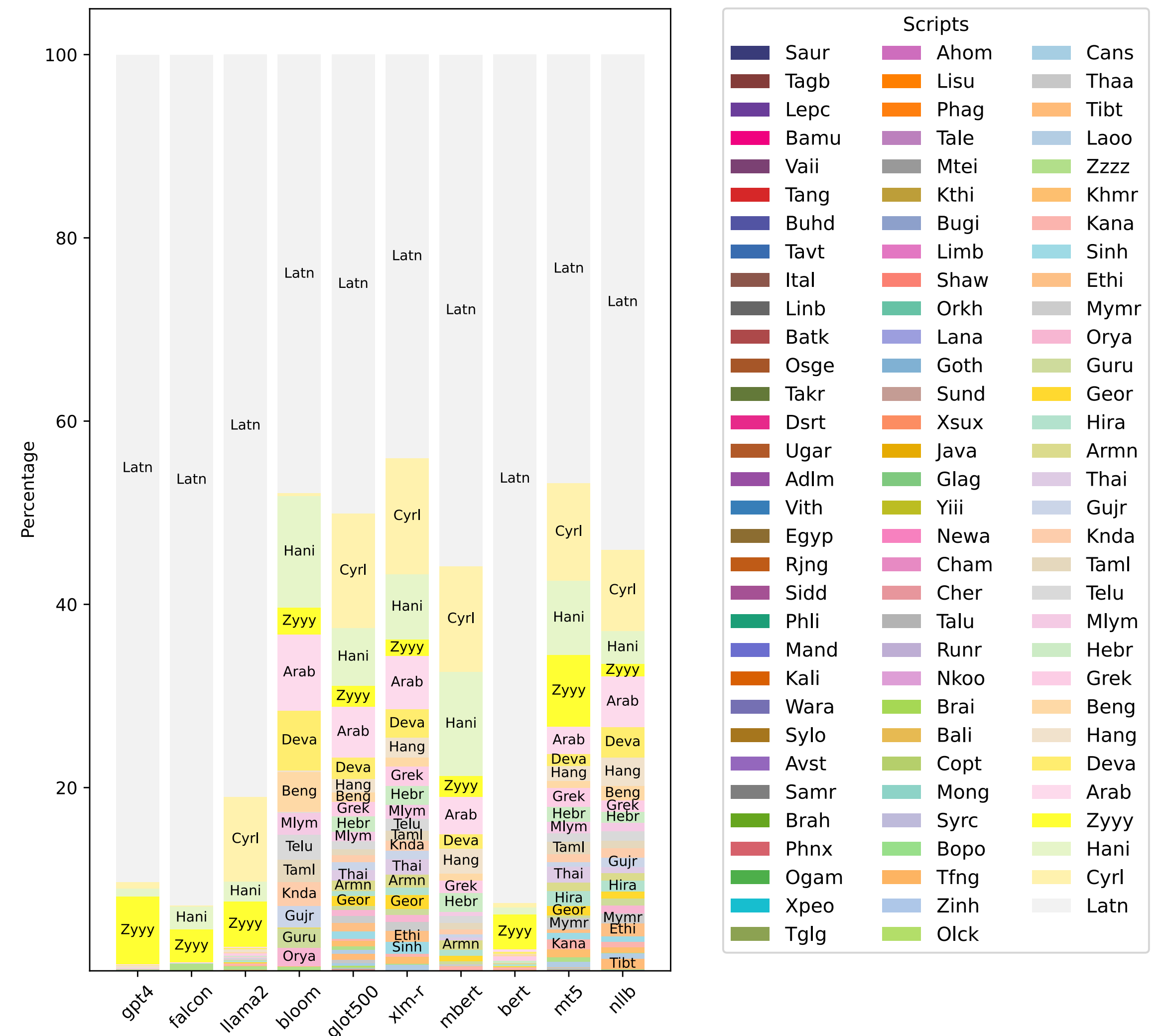
### Script Detection

```
from GlotScript import sp

sp('これは日本人です')
>> ('Hira', 0.625, {'details': {'Hira': 0.625, 'Hani': 0.375}, 'tie': False, 'interval': 0.25})

sp('This is Latin')[1]
>> ('Latn', 1.0)

sp('මෙහි සිංහල')[0]
>> 'Sinh'
```



GlotScript: A Resource and Tool for Low Resource Writing System Identification (Kargaran et al, proc LREC COLING 2024)



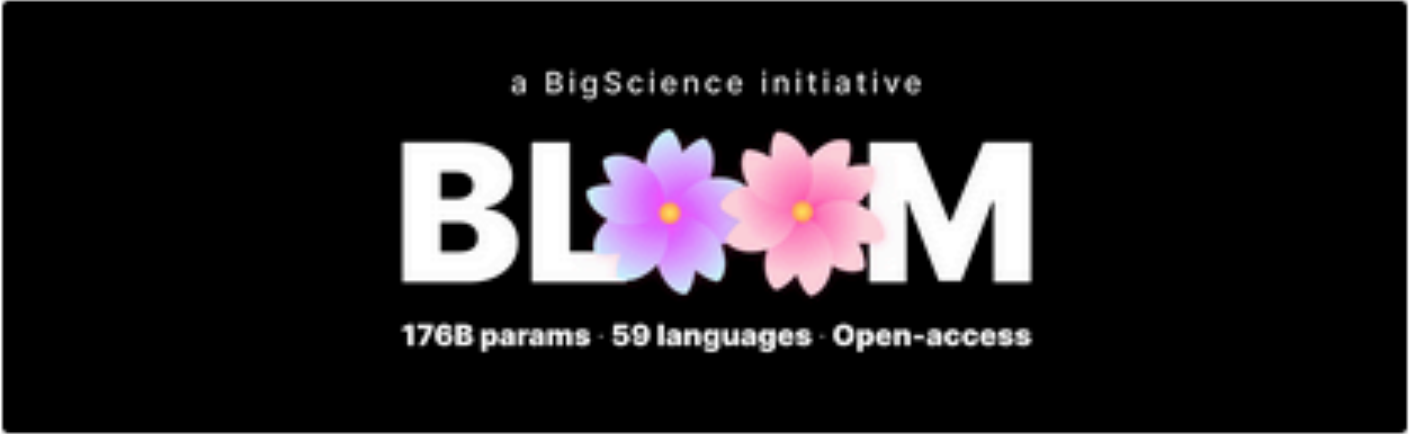
**mLMMs:  $m = ?$**

**Examining Model Cards**

Exa

Model card Files Metrics Community 276 Settings

Edit model card



BigScience Large Open-science Open-access Multilingual Language Model  
Version 1.3 / 6 July 2022

Current Checkpoint: **Training Iteration 95000**

Link to paper: [here](#)

Total seen tokens: **366B**

---

### Model Details

BLOOM is an autoregressive Large Language Model (LLM), trained to continue text from a prompt on vast amounts of text data using industrial-scale computational resources. As such, it is able to output coherent text **in 46 languages** and 13 programming languages that is hardly distinguishable from text written by humans. BLOOM can also be instructed to perform text tasks it hasn't been explicitly trained for, by casting them as text generation tasks.

ards

# Multilingualism in LMMs

---

*m* = 1 : monolingual LLMs

- ara: AraBERT, JASMINE, ...
- eng: BERT, ELMO, GPT-2, T5, ...
- fra: CamemBERT, BARThez, GPT-FR

# Multilingualism in LMMs

---

## $m = 1$ : monolingual LLMs

- ara: AraBERT, JASMINE, ...
- eng: BERT, ELMO, GPT-2, T5, ...
- fra: CamemBERT, BARThez, GPT-FR

## $m = 1 + \epsilon$ : 'almost' monolingual LLMs

- eng: OPT, Jurassic, Gopher, Chinchilla, Turing-NLG, Megatron-LM, Falcon, OLMO, Phi-2, ...

# Multilingualism in LLMs

---

## $m = 1$ : monolingual LLMs

- ara: AraBERT, JASMINE, ...
- eng: BERT, ELMO, GPT-2, T5, ...
- fra: CamemBERT, BARThez, GPT-FR

## $m = 1 + \epsilon$ : 'almost' monolingual LLMs

- eng: OPT, Jurassic, Gopher, Chinchilla, Turing-NLG, Megatron-LM, Falcon, OLMO, Phi-2, ...

## $m = 1 + 1$ : bilingual LLMs

- eng+rus: YaLM
- eng+zho: GLM-130, PanGu, Baichuan
- eng+fra: CroissantLLM
- eng+ara: JAIS

# Multilingualism in LLMs

---

## $m = 1$ : monolingual LLMs

- ara: AraBERT, JASMINE, ...
- eng: BERT, ELMO, GPT-2, T5, ...
- fra: CamemBERT, BARThez, GPT-FR

## $m = 1 + \epsilon$ : 'almost' monolingual LLMs

- eng: OPT, Jurassic, Gopher, Chinchilla, Turing-NLG, Megatron-LM, Falcon, OLMO, Phi-2, ...

## $m \in [5 : 25]$ : 'familial' LLMs

- Indic: IndicBERT (12), MuRIL (17)
- African: AfriBERTA (11), AfroXLM (17)
- Nordic: GPT-SW3 (5+1)
- European: Occiglot-v0 (5), Occiglot-v1 (24)

## $m = 1 + 1$ : bilingual LLMs

- eng+rus: YaLM
- eng+zho: GLM-130, PanGu, Baichuan
- eng+fra: CroissantLLM
- eng+ara: JAIS

# Multilingualism in LMMs

---

## $m \in [5 : 50]$ : 'opportunistic' LMMs

- Unicoder (15), VECO (50)
- PolyLM (20)
- mBART-25 (25)
- LLaMA, TowerLM (10), LLaMA-2 (20)
- BLOOM (46)
- Command+R (23)

## $m \in [5 : 25]$ : 'familial' LLMs

- Indic: IndicBERT (12), MuRIL (17)
- African: AfriBERTA (11), AfroXLM (17)
- Nordic: GPT-SW3 (5+1)
- European: Occiglot-v0 (5), Occiglot-v1 (24)

## $m = 1 + \epsilon$ : 'almost' monolingual LLMs

- eng: OPT, Jurassic, Gopher, Chinchilla, Turing-NLG, Megatron-LM, Falcon, OLMO, Phi-2, ...

## $m = 1 + 1$ : bilingual LLMs

- eng+rus: YaLM
- eng+zho: GLM-130, PanGu, Baichuan
- eng+fra: CroissantLLM
- eng+ara: JAIS

# Multilingualism in LMMs

---

## $m \in [5 : 50]$ : 'opportunistic' LMMs

- Unicoder (15), VECO (50)
- PolyLM (20)
- mBART-25 (25)
- LLaMA, TowerLM (10), LLaMA-2 (20)
- BLOOM (46)
- Command+R (23)

## $m > [50]$ : 'Massively Multilingual' LLMs

- mBERT (104), remBERT (104), XLM-R (100), VECO-2.0 (109)
- GPT-3 (94), XGLM (114)
- PALM (100+)
- mT5 (100)
- Aya (101)

## $m \in [5 : 25]$ : 'familial' LLMs

- Indic: IndicBERT (12), MuRIL (17)
- African: AfriBERTA (11), AfroXLM (17)
- Nordic: GPT-SW3 (5+1)
- European: Occiglot-v0 (5), Occiglot-v1 (24)

## $m = 1 + 1$ : bilingual LLMs

- eng+rus: YaLM
- eng+zho: GLM-130, PanGu, Baichuan
- eng+fra: CroissantLLM
- eng+ara: JAIS



# Multilingualism in LMMs

---

## $m \in [5 : 50]$ : 'opportunistic' LMMs

- Unicoder (15), VECO (50)
- PolyLM (20)
- mBART-25 (25)
- LLaMA, TowerLM (10), LLaMA-2 (20)
- BLOOM (46)
- Command+R (23)

## $m > [50]$ : 'Massively Multilingual' LLMs

- mBERT (104), remBERT (104), XLM-R (100), VECO-2.0 (109)
- GPT-3 (94), XGLM (114)
- PALM (100+)
- mT5 (100)
- Aya (101)

## $m \in [5 : 25]$ : 'familial' LLMs

- Indic: IndicBERT (12), MuRIL (17)
- African: AfriBERTA (11), AfroXLM (17)
- Nordic: GPT-SW3 (5+1)
- European: Occiglot-v0 (5), Occiglot-v1 (24)

## $m \gg 100$ : 'Extremely Multilingual' LMMs

- Serengeti (517)
- MadLad-400 (419)
- Glot500-m (511)
- MALA500 (534)

# Multilingualism in LMMs

---

## $m \in [5 : 50]$ : 'opportunistic' LMMs

- Unicoder (15), VECO (50)
- PolyLM (20)
- mBART-25 (25)
- LLaMA, TowerLM (10), LLaMA-2 (20)
- BLOOM (46)
- Command+R (23)

## $m > [50]$ : 'Massively Multilingual' LLMs

- mBERT (104), remBERT (104), XLM-R (100), VECO-2.0 (109)
- GPT-3 (94), XGLM (114)
- PALM (100+)
- mT5 (100)
- Aya (101)

## $m > 0$ : 'Vaguely documented' LMMs

- Mistral
- ChatGPT
- Claude
- ... and many more

## $m \gg 100$ : 'Extremely Multilingual' LMMs

- Serengeti (517)
- MadLad-400 (419)
- Glot500-m (511)
- MALA500 (534)

# Going Beyond Language counts

# Language distributions in mLLMs

---

**Comparing published statistics is \*hard\***

- variability in language / variety names (and scripts)
- errors in LID
- units of measures (tokens vs. bytes vs. #docs)
- corpus building details (genre matters)
- language sampling details

**When the information is even given...**

**Only correct solution: go back to data (if possible)... and count**

# Glott500: design choices

---

## 1. Select reliable sources

- curated multilingual corpora
- new data crawls
- multiple domains: web, news, science, religion, etc
- excludes toxicity by design

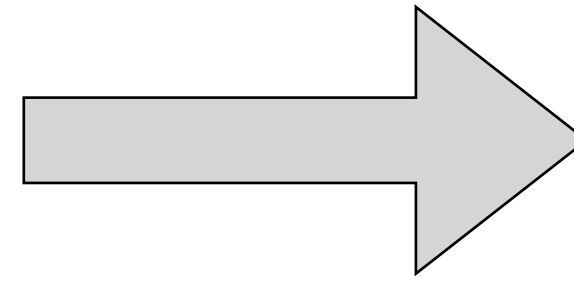


# Glott500: design choices

---

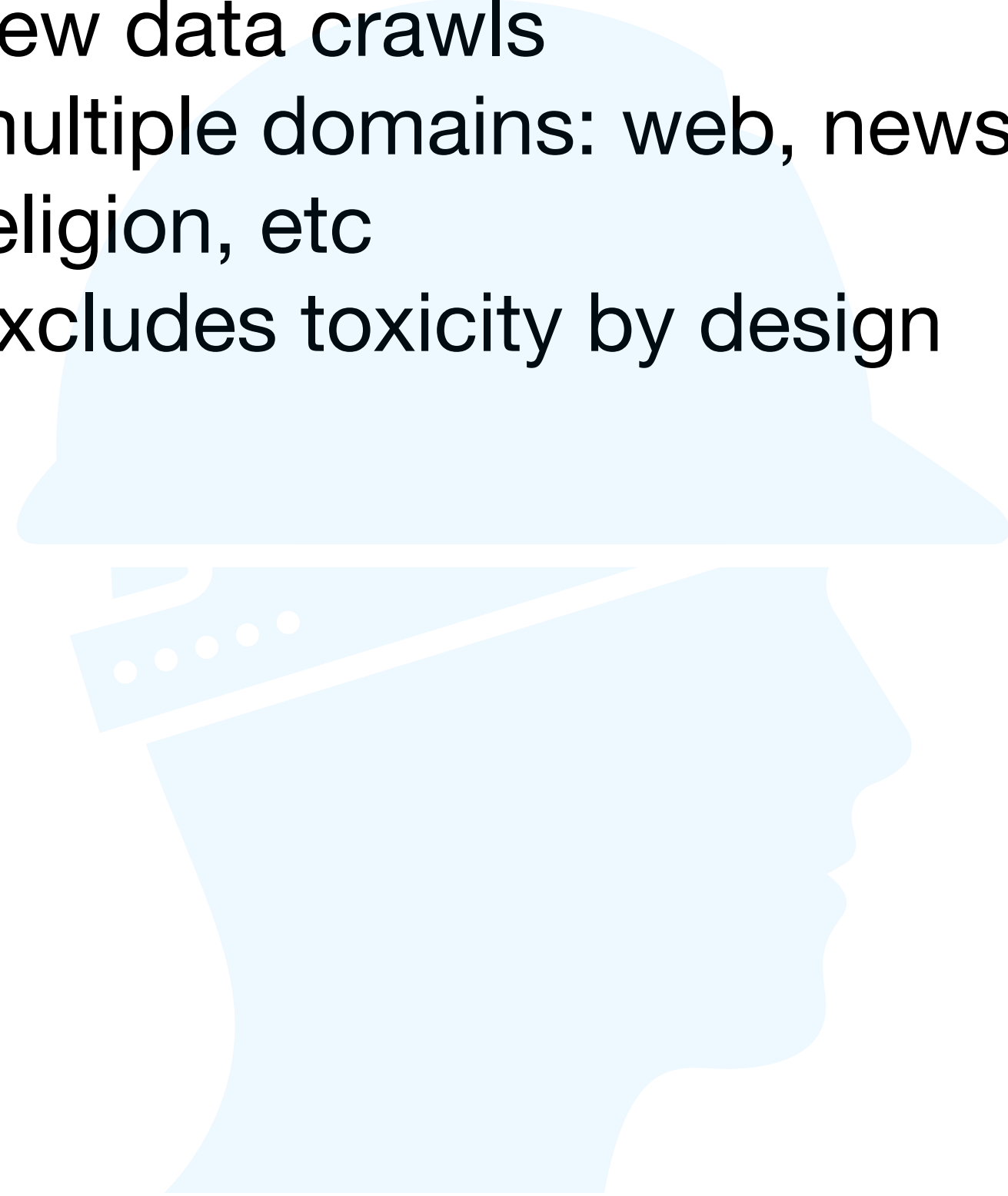
## 1. Select reliable sources

- curated multilingual corpora
- new data crawls
- multiple domains: web, news, science, religion, etc
- excludes toxicity by design



## 2. Language Identification

- per sentence LID
- joint detection of language + script

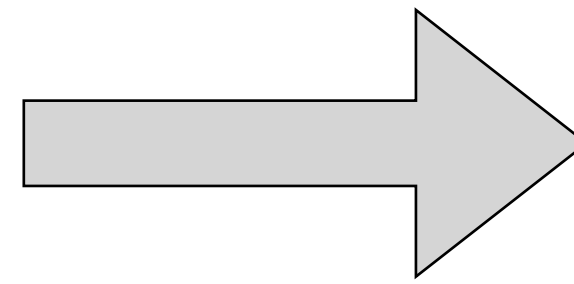


# Glott500: design choices

---

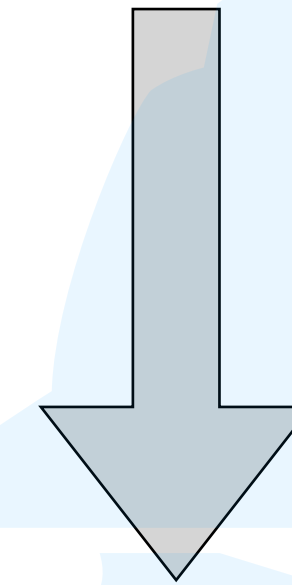
## 1. Select reliable sources

- curated multilingual corpora
- new data crawls
- multiple domains: web, news, science, religion, etc
- excludes toxicity by design



## 2. Language Identification

- per sentence LID
- joint detection of language + script



## 3. Chunks and sentence

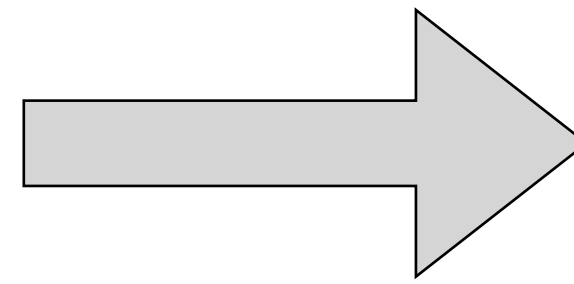
- high character repetitions
- normal / special char ratios
- insufficient number of words
- wrong language / script
- char LM filtering
- duplicate removal

# Glott500: design choices

---

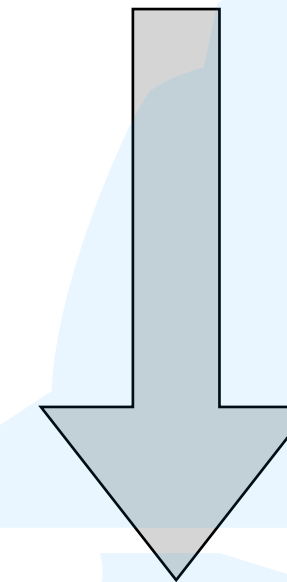
## 1. Select reliable sources

- curated multilingual corpora
- new data crawls
- multiple domains: web, news, science, religion, etc
- excludes toxicity by design



## 2. Language Identification

- per sentence LID
- joint detection of language + script

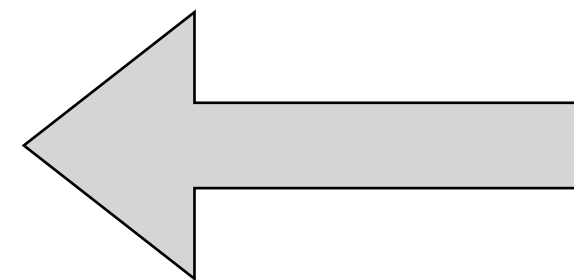


## 3. Chunks and sentence

- high character repetitions
- normal / special char ratios
- insufficient number of words
- wrong language / script
- char LM filtering
- duplicate removal

## 4. Final selection = Glott500-c

- > 30k sentences
- from 2000+ to 511 language / 34 scripts
- 610Gb of text 1.5b sentences
- **head** vs. **tail** languages
- 1000 dev+test sentences / language





# Glott500-c: m = 511 languages

100 “head”  
411 “tail”

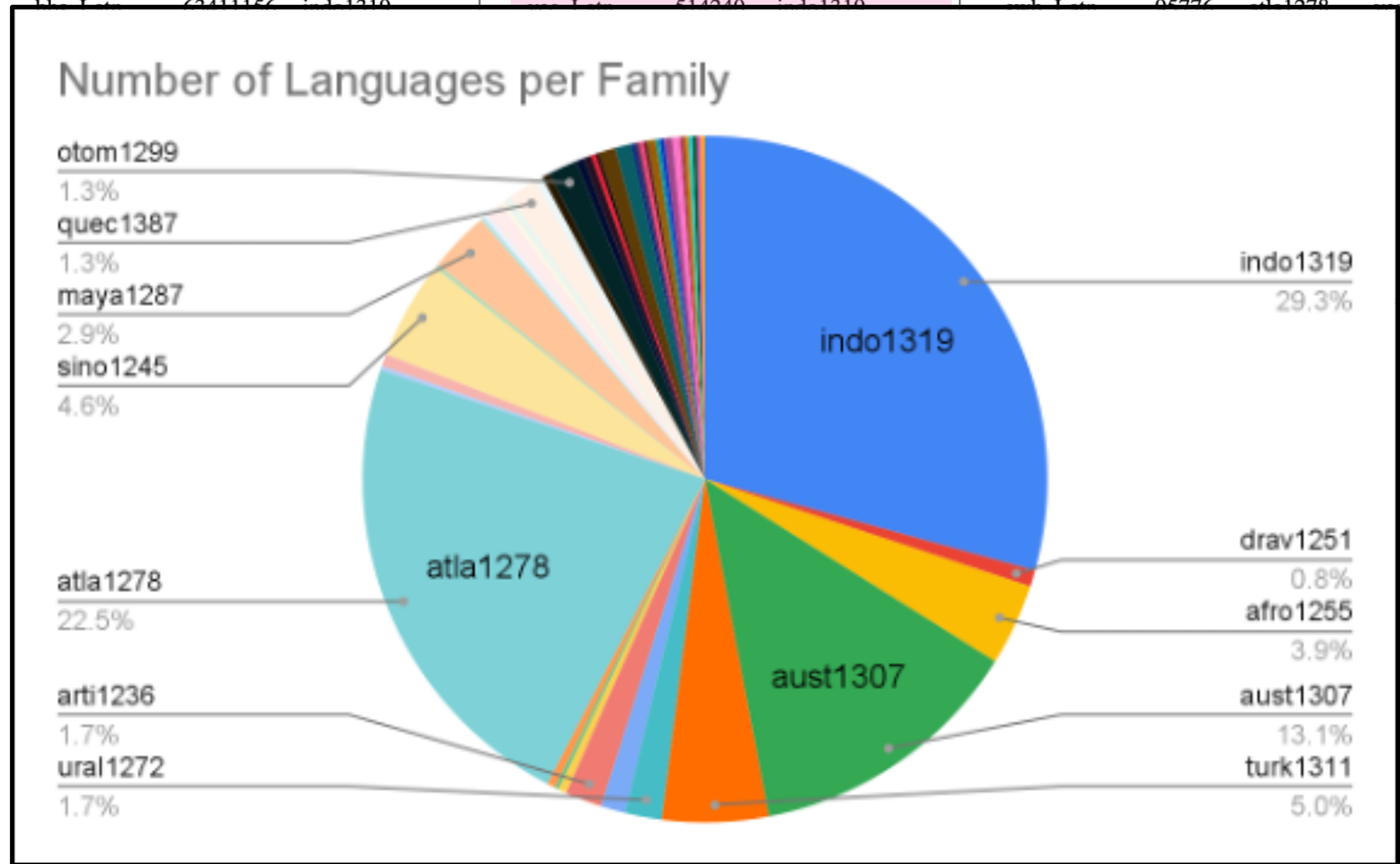
Language-Script	Sent	Family	Head	Language-Script	Sent	Family	Head	Language-Script	Sent	Family	Head
hbs_Latn	63411156	indo1319		vec_Latn	514240	indo1319		swh_Latn	95776	atla1278	yes
mal_Mlym	48098273	drav1251	yes	jpn_Jpan	510722	japo1237	yes	alt_Cyrl	95148	turk1311	
aze_Latn	46300705		yes	lus_Latn	509250	sino1245		rmn_Grek	94533	indo1319	
guj_Gujr	45738685	indo1319	yes	crs_Latn	508755	indo1319		miq_Latn	94343	misu1242	
ben_Beng	43514870	indo1319	yes	kqn_Latn	507913	atla1278		kaa_Cyrl	88815	turk1311	
kan_Knda	41836495	drav1251	yes	ndo_Latn	496613	atla1278		kos_Latn	88603	aust1307	
tel_Telu	41580525	drav1251	yes	snd_Arab	488730	indo1319	yes	grn_Latn	87568		
mlt_Latn	40654838	afro1255		yue_Hani	484700	sino1245		lhu_Latn	87255	sino1245	
fra_Latn	39197581	indo1319	yes	tiv_Latn	483064	atla1278		lzh_Hani	86035	sino1245	
spa_Latn	37286756	indo1319	yes	kua_Latn	473535	atla1278		ajp_Arab	83297	afro1255	
eng_Latn	36122761	indo1319	yes	kwy_Latn	473274	atla1278		cmn_Hani	80745	sino1245	yes
fil_Latn	33493255	aust1307	yes	hin_Latn	466175	indo1319		gcf_Latn	80737	indo1319	
nob_Latn	32869205	indo1319		iku_Cans	465011			rmn_Cyrl	79925	indo1319	
rus_Cyrl	31787973	indo1319	yes	kal_Latn	462430	eski1264		kjh_Cyrl	79262	turk1311	
deu_Latn	31015993	indo1319	yes	tdt_Latn	459818	aust1307		rng_Latn	78177	atla1278	
tur_Latn	29184662	turk1311	yes	gsw_Latn	449240	indo1319		mgh_Latn	78117	atla1278	
pan_Guru	29052537	indo1319	yes	mfe_Latn	447435	indo1319		xmv_Latn	77896	aust1307	
mar_Deva	28748897	indo1319	yes	swc_Latn	446378	atla1278		ige_Latn	77114	atla1278	
por_Latn	27824391	indo1319	yes	mon_Latn	437950	mong1349		rmy_Latn	76991	indo1319	
nld_Latn	25061426	indo1319	yes	mos_Latn	437666	atla1278		srm_Latn	76884	indo1319	
ara_Arab	24524122		yes	kik_Latn	437228	atla1278		bak_Latn	76809	turk1311	
zho_Hani	24143786		yes	cnh_Latn	436667	sino1245		gur_Latn	76151	atla1278	
ita_Latn	23539857	indo1319	yes	gil_Latn	434529	aust1307		idu_Latn	75106	atla1278	
ind_Latn	23018106	aust1307	yes	pon_Latn	434522	aust1307		yom_Latn	74818	atla1278	
ell_Grek	22033282	indo1319	yes	umb_Latn	431589	atla1278		tdx_Latn	74430	aust1307	
bul_Cyrl	21823004	indo1319	yes	lvs_Latn	422952	indo1319		mzn_Arab	73719	indo1319	
swe_Latn	20725883	indo1319	yes	sco_Latn	411591	indo1319		cfm_Latn	70227	sino1245	
ces_Latn	20376340	indo1319	yes	ori_Orya	410827		yes	zpa_Latn	69237	otom1299	
isl_Latn	19547941	indo1319	yes	arg_Latn	410683	indo1319		kbd_Cyrl	67914	abkh1242	
pol_Latn	19339945	indo1319	yes	kur_Latn	407169	indo1319	yes	lao_Lao	66966	taik1256	yes
ron_Latn	19190217	indo1319	yes	dhv_Latn	405711	aust1307		nap_Latn	65826	indo1319	
dan_Latn	19174573	indo1319	yes	luo_Latn	398974	nilo1247		qub_Latn	64973	quec1387	

Usual suspects  
Minority script  
Minority variety  
Minority language

# Glott500-c: m = 511 languages

100 "head"  
411 "tail"

Language-Script	Sent	Family	Head	Language-Script	Sent	Family	Head	Language-Script	Sent	Family	Head
...	62411156	indo1319	yes	...	514249	indo1319	yes	...	95776	indo1319	yes



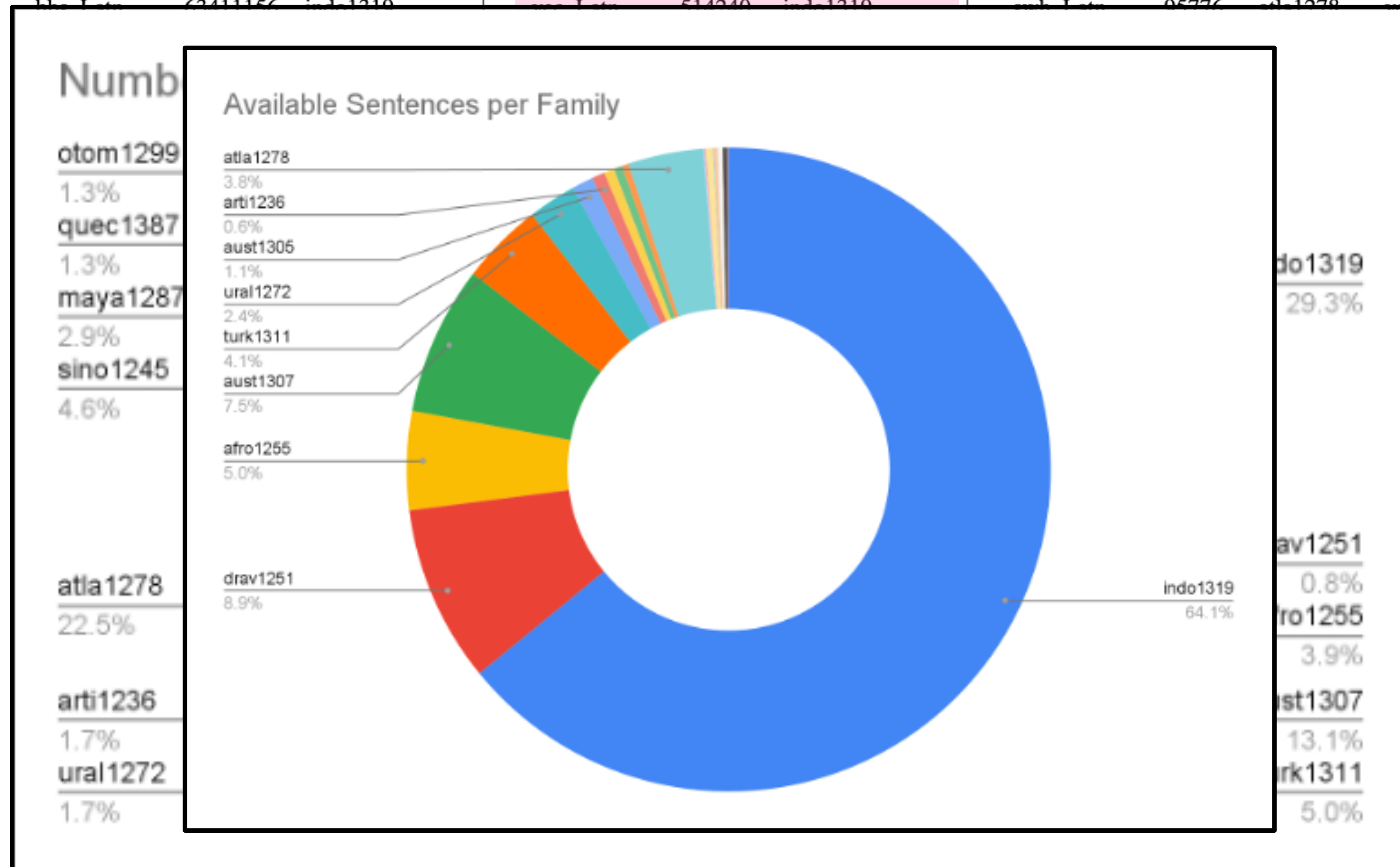
ron_Latn	19190217	indo1319	yes	dhv_Latn	405711	aust1307	no	nap_Latn	65826	indo1319	no
dan_Latn	19174573	indo1319	yes	luo_Latn	398974	nilo1247	no	qub_Latn	64973	quec1387	no

- Usual suspects
- Minority script
- Minority variety
- Minority language

# Glott500-c: m = 511 languages

100 "head"  
411 "tail"

Language-Script	Sent	Family	Head	Language-Script	Sent	Family	Head	Language-Script	Sent	Family	Head
...	62411156	indo1319	...	...	514249	indo1319	...	...	95776	indo1319	...



ron_Latn	19190217	indo1319	yes	dhv_Latn	405711	aust1307	nap_Latn	65826	indo1319
dan_Latn	19174573	indo1319	yes	luo_Latn	398974	nilo1247	qub_Latn	64973	quec1387

- Usual suspects
- Minority script
- Minority variety
- Minority language

# Language ID is key

---

## 2. Language Identification

- per sentence LID
- joint detection of language + script

# Language ID is key

---

## 2. Language Identification

- per sentence LID
- joint detection of language + script

**Is this reliable ? Will it scale ?  
Also domains ? genres ?**

# Auditing data with GlotScript

	Corpus Code: ISO 639-3	Scripts	ACC↑	ACC70↑	ACC50↑	
Highest ACC	mC4	st:sot (S Sotho)	<u>Latn</u> :1000	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
		fil:fil (Filipino)	<u>Latn</u> :998, <u>Cyrl</u> :1, <u>Hani</u> :1	0.998	0.999	<b>1.000</b>
		ro:ron (Romanian)	<u>Latn</u> :996, <u>Zyyy</u> :4, <u>Cyrl</u> :1	0.995	0.997	<b>1.000</b>
		id:ind (Indonesian)	<u>Latn</u> :995, <u>Zyyy</u> :3, <u>Hani</u> :1, <u>Hebr</u> :1	0.995	<b>1.000</b>	<b>1.000</b>
		sw:swa (Swahili)	<u>Latn</u> :995, <u>Zyyy</u> :5	0.995	<b>1.000</b>	<b>1.000</b>
Lowest ACC	mC4	ne:nep (Nepali)	<u>Deva</u> :609, <u>Hani</u> :219, <u>Latn</u> :88, <u>Hang</u> :44, <u>Thai</u> :12, <u>Laoo</u> :8, <u>Zyyy</u> :8, <u>Orya</u> :7, <u>Other</u> :5	0.609	0.730	<b>0.797</b>
		mn:mon (Mongolian)	<u>Cyrl</u> :502, <u>Hebr</u> :348, <u>Latn</u> :135, <u>Zyyy</u> :14, <u>Hani</u> :1	0.502	0.557	<b>0.570</b>
		cy:cym (Welsh)	<u>Gre</u> :603, <u>Latn</u> :367, <u>Zyyy</u> :11, <u>Hebr</u> :9, <u>Cyrl</u> :5, <u>Zzzz</u> :4, <u>Arab</u> :1	<b>0.367</b>	0.338	0.295
		sd:snd (Sindhi)	<u>Latn</u> :654, <u>Arab</u> :329, <u>Zyyy</u> :12, <u>Zzzz</u> :2, <u>Cyrl</u> :1, <u>Hang</u> :1, <u>Tel</u> :1	<b>0.329</b>	0.271	0.222
		mr:mar (Marathi)	<u>Hani</u> :454, <u>Thai</u> :252, <u>Latn</u> :119, <u>Deva</u> :116, <u>Zyyy</u> :34, <u>Guru</u> :10, <u>Beng</u> :4, <u>Khmr</u> :3, <u>Other</u> : 8	0.116	0.136	<b>0.141</b>
Highest ACC	OSCAR	id:ind (Indonesian)	<u>Latn</u> :998, <u>Zyyy</u> :2	0.998	<b>1.000</b>	<b>1.000</b>
		war:war (Waray)	<u>Latn</u> :997, <u>Zyyy</u> :3	<b>0.997</b>	<b>0.997</b>	0.996
		als:gsw (Swiss G)	<u>Latn</u> :996, <u>Zyyy</u> :3, <u>Cyrl</u> :1	0.996	0.996	<b>1.000</b>
		vo:vol (Volapük)	<u>Latn</u> :994, <u>Arab</u> :4, <u>Cyrl</u> :1	0.994	<b>1.000</b>	<b>1.000</b>
		nds:nds (Low G)	<u>Latn</u> :994, <u>Zyyy</u> :2, <u>Cyrl</u> :2, <u>Hang</u> :1, <u>Thaa</u> :1	0.994	<b>1.000</b>	<b>1.000</b>
Lowest ACC	OSCAR	am:amh (Amharic)	<u>Ethi</u> :822, <u>Latn</u> :164, <u>Zyyy</u> :12, <u>Hani</u> :1, <u>Arab</u> :1	0.822	0.883	<b>0.940</b>
		gu:guj (Gujarati)	<u>Gujr</u> :802, <u>Latn</u> :180, <u>Zyyy</u> :12, <u>Deva</u> :6	0.802	0.863	<b>0.883</b>
		si:sin (Sinhala)	<u>Sinh</u> :801, <u>Latn</u> :188, <u>Zyyy</u> :11	0.801	0.905	<b>0.948</b>
		th:tha (Thai)	<u>Thai</u> :800, <u>Latn</u> :181, <u>Zyyy</u> :18, <u>Hani</u> :1	0.800	0.883	<b>0.917</b>
		te:tel (Telugu)	<u>Tel</u> :799, <u>Latn</u> :188, <u>Zyyy</u> :9, <u>Deva</u> :3, <u>Cyrl</u> :1	0.799	0.880	<b>0.908</b>

Script-languages mismatches detection finds many errors

# Language Identification at scale ---

## Existing tools are limited

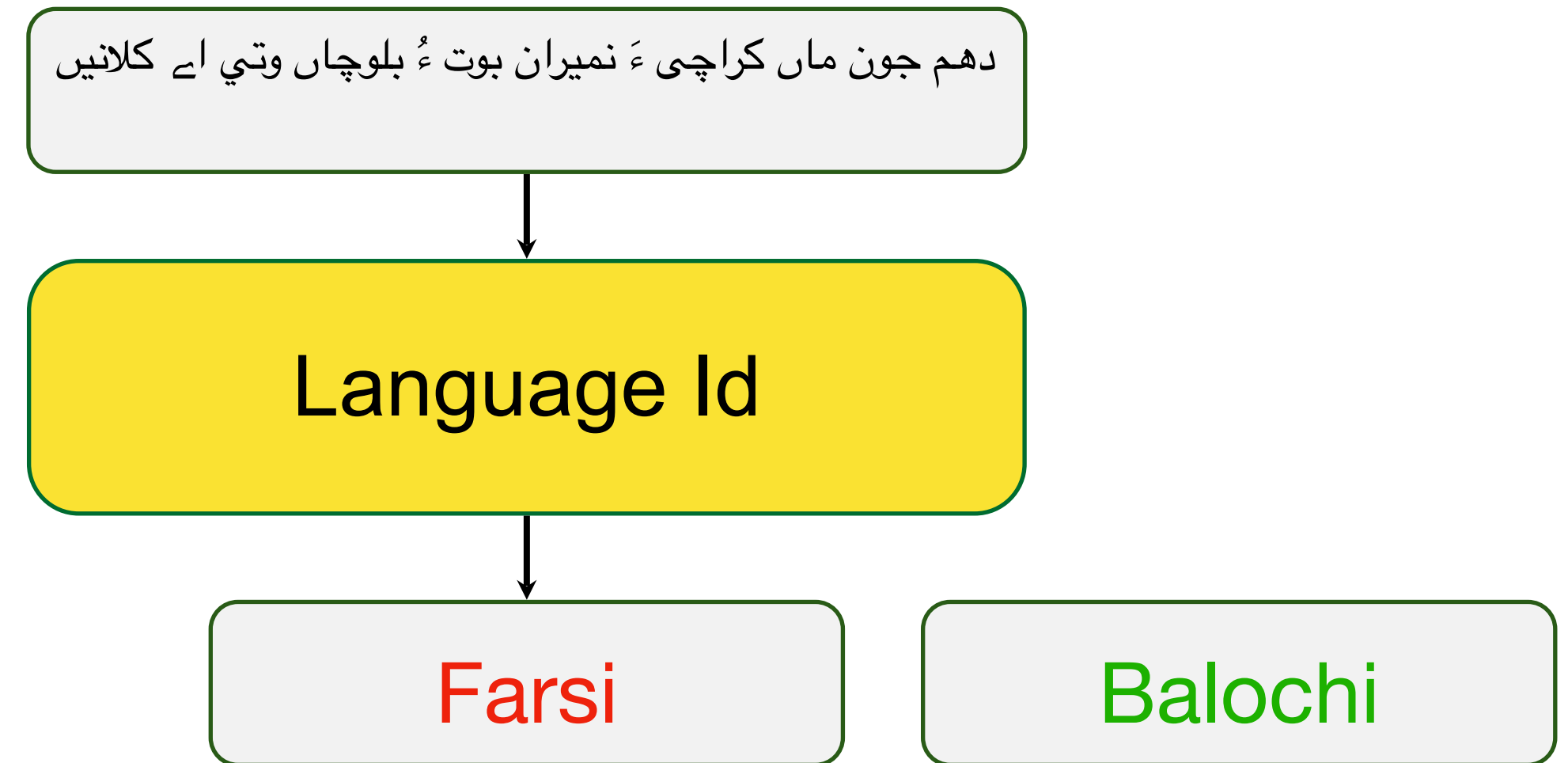
Fasttext LID	176
Compact Language Detector	107
whatlang	69
OpenLID	200
franc-s	82
franc-m	187
franc-l	417

[Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets \(Kreutzer et al, 2022\)](#),

# Language Identification at scale \_\_\_\_\_

## Existing tools are limited

Fasttext LID	176
Compact Language Detector	107
whatlang	69
OpenLID	200
franc-s	82
franc-m	187
franc-l	417



« out-of-model cousin error »

[Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets \(Kreutzer et al, 2022\),](#)



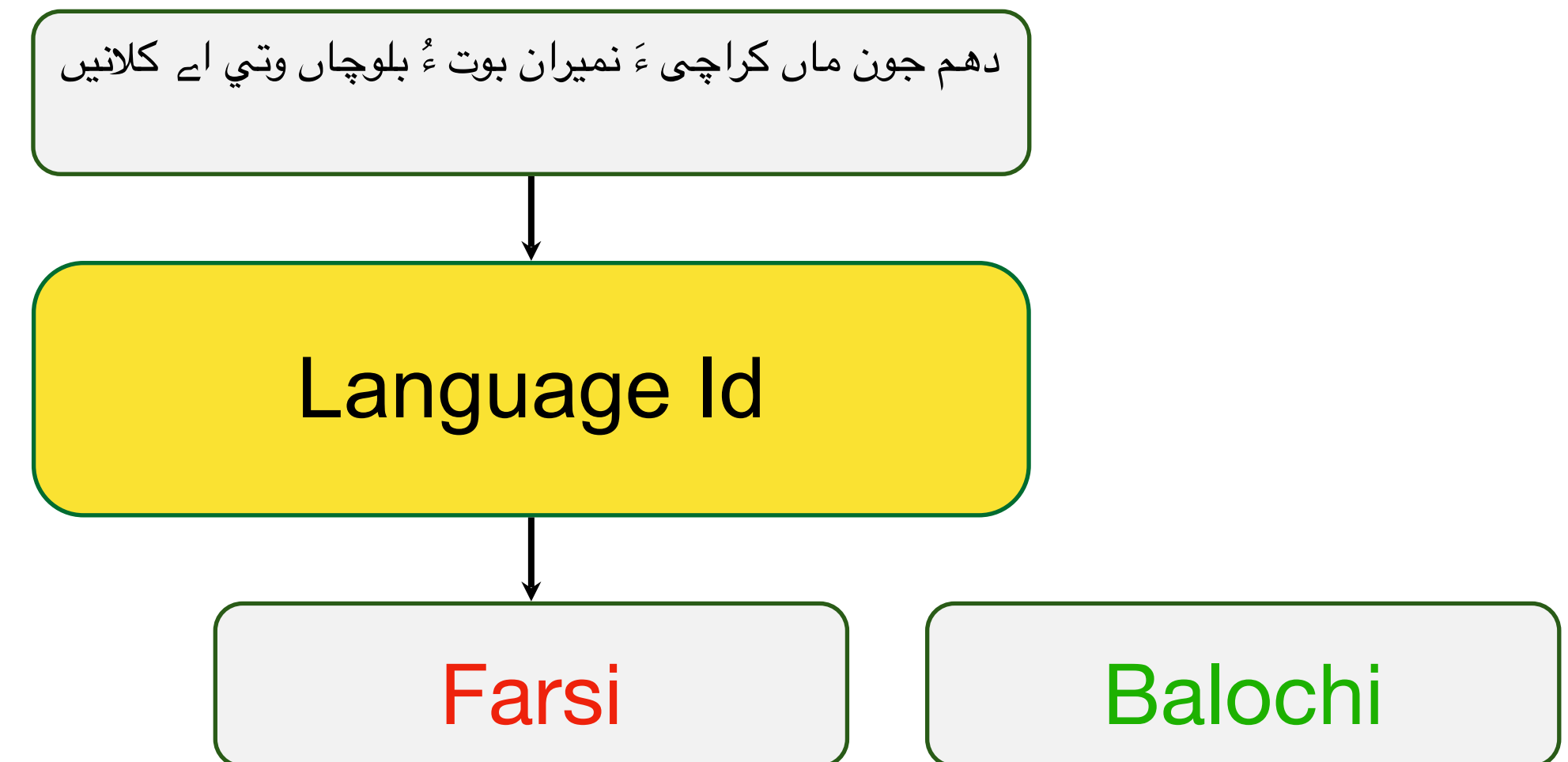
# Language Identification at scale

## Existing tools are limited

Fasttext LID	176
Compact Language Detector	107
whatlang	69
OpenLID	200
franc-s	82
franc-m	187
franc-l	417

## Other issues

- speed
- implementation / deployment
- lack of openness
- lack of documentation
- lack **confidence estimation**
- lack of **rejection model**
- **errors**



« out-of-model cousin error »

[Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets \(Kreutzer et al, 2022\)](#),

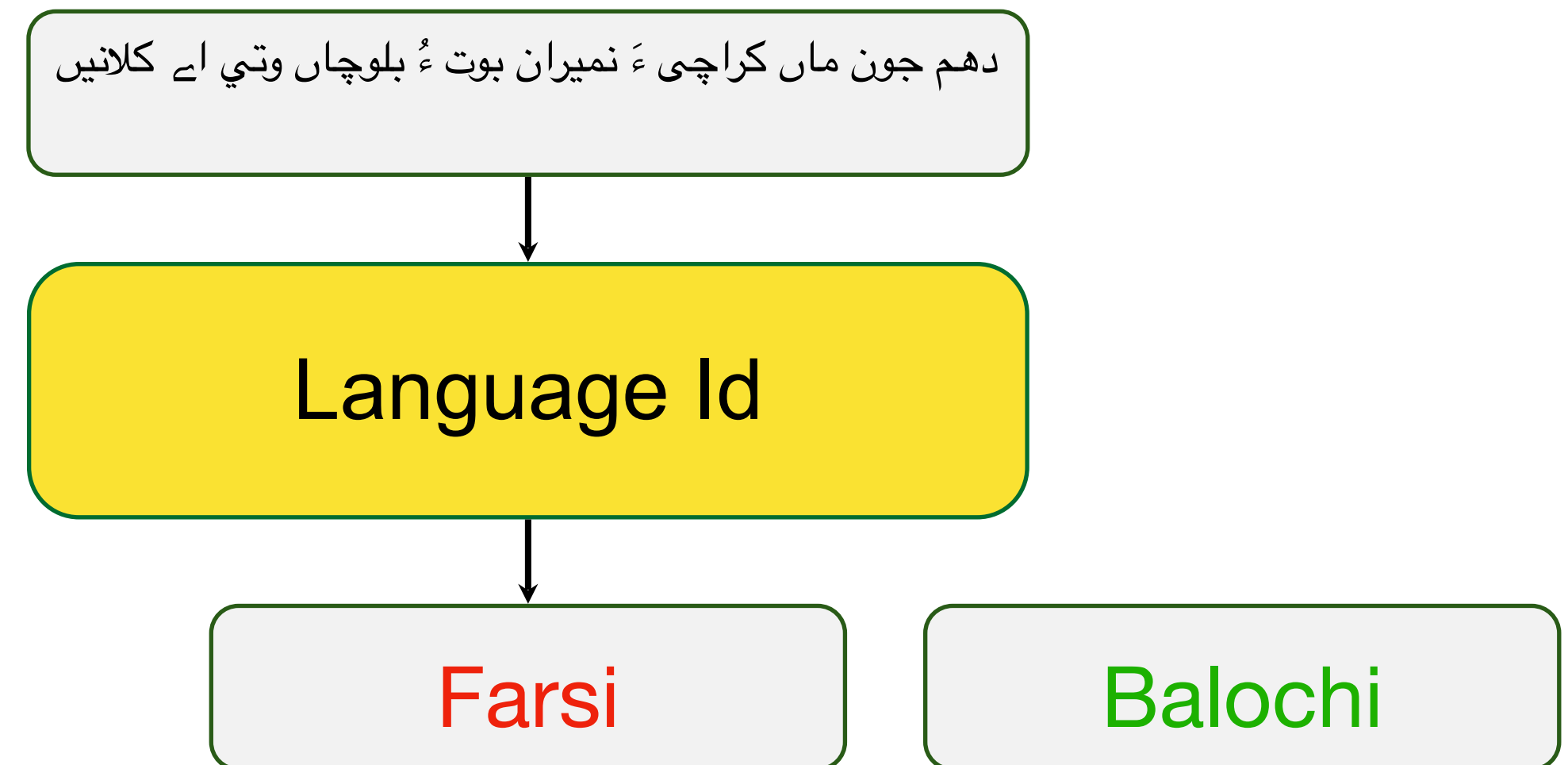
# Language Identification at scale

## Existing tools are limited

Fastext LID	176
Compact Language Detector	107
whatlang	69
OpenLID	200
franc-s	82
franc-m	187
franc-l	417

## Other issues

- speed
- implementation / deployment
- lack of openness
- lack of documentation
- lack **confidence estimation**
- lack of **rejection model**
- **errors**



« out-of-model cousin error »

[Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets \(Kreutzer et al, 2022\)](#),

## Designing your own

- selecting / naming languages
- curating corpora
- the language mix, again

# Training and testing GlotLID ---

## Train corpus (V1.0)

- >1800+ languages scripts
- multiple reliable sources
- mixture of domains:
  - Wikipedia
  - religious texts
  - collaborative translations
  - academia
  - storybooks
  - news sites

## Implementation: FastText

- linear classifier
- char & word n-gram features
- highly optimized for speed (training and inference)
- use with a pinch of salt

[Bag of Tricks for Efficient Text Classification](#) (Joulin et al., EACL 2017)

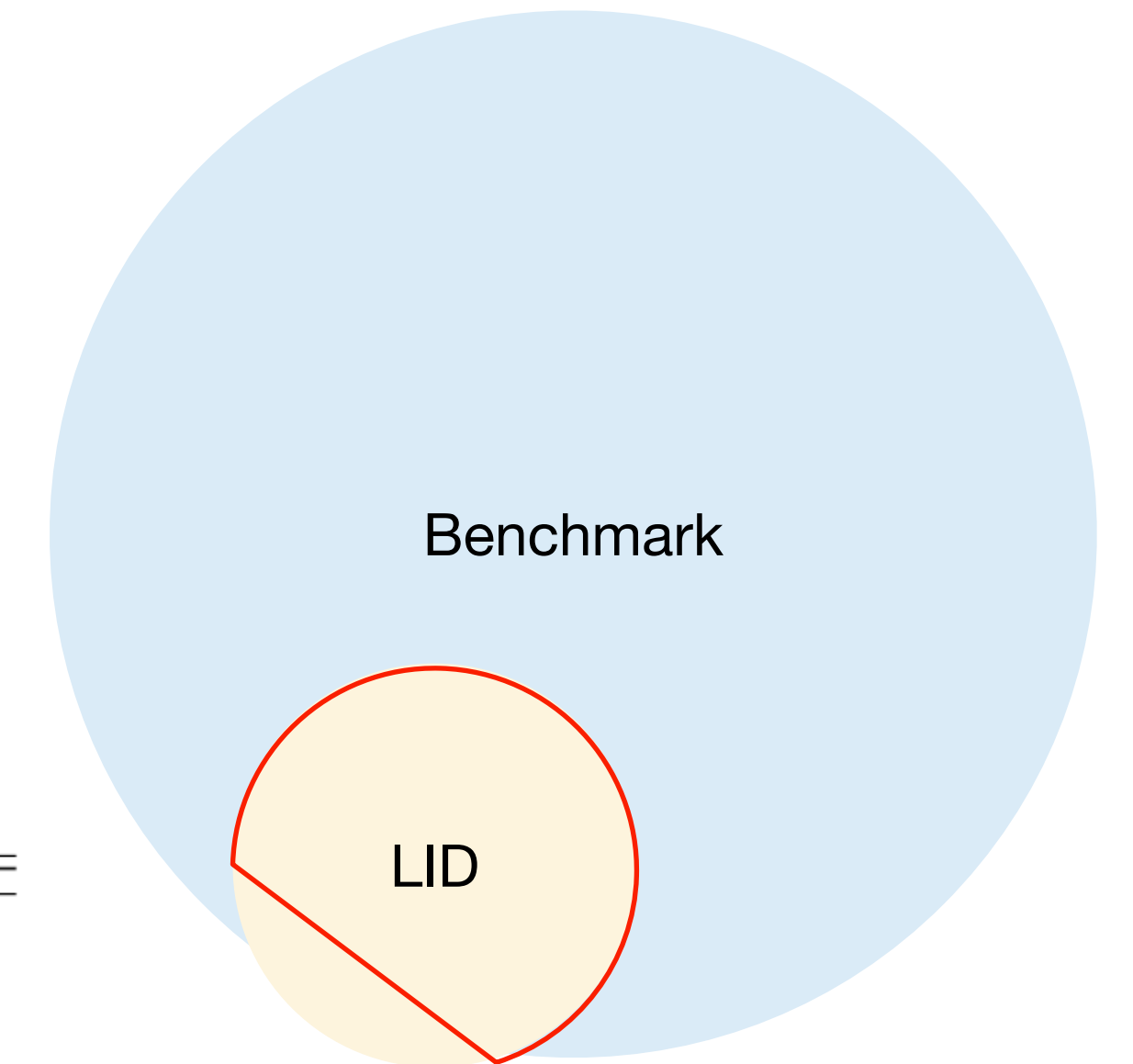
## Test data

- GlotLID-c (testset) (1800+)
- Universal Declaration of Human Rights (UDHR) (204)
- Flores 200

# Language verification

## Closed set assumption

- set of possible languages known (differs for each baseline)

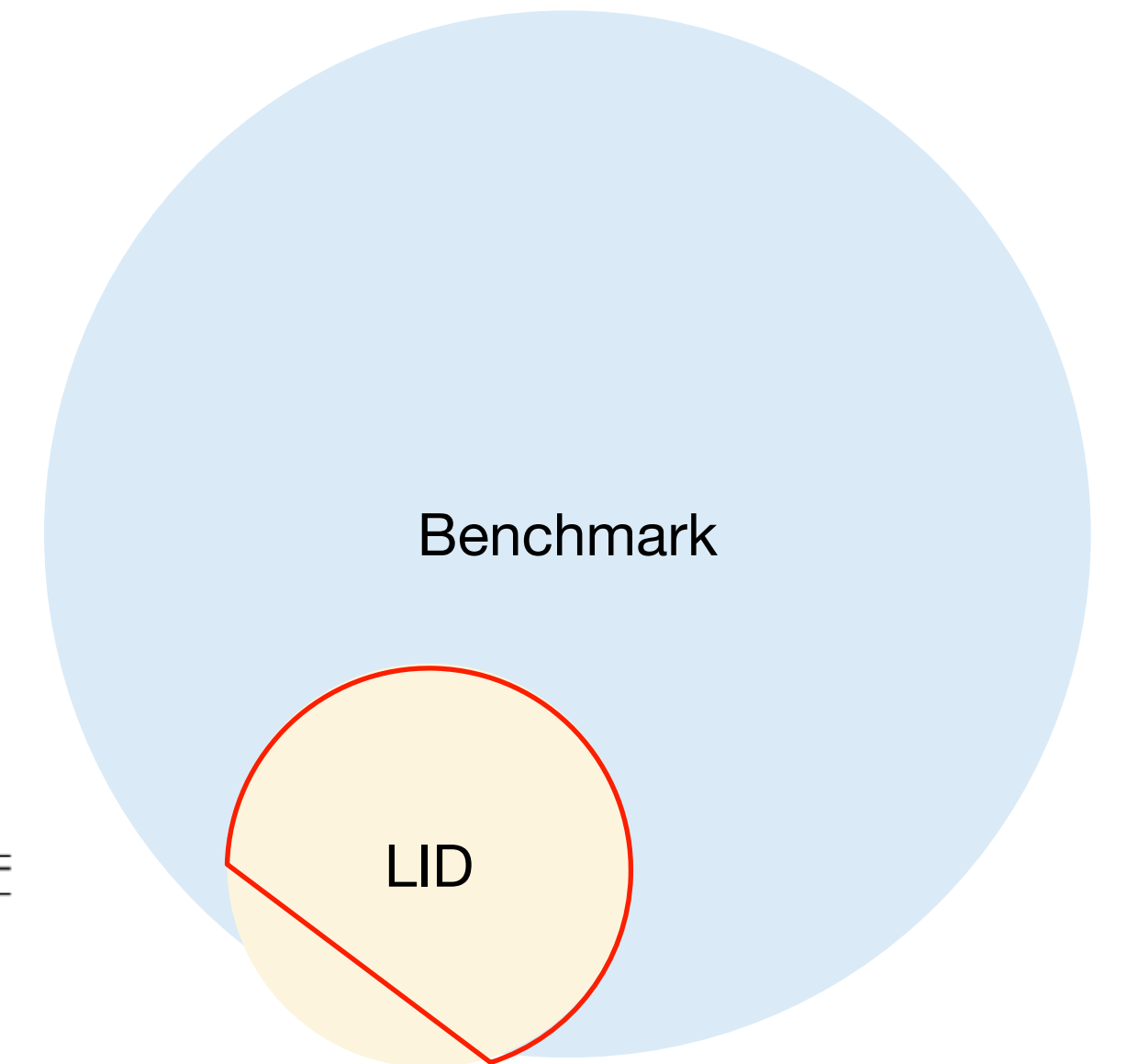


LID Model	$\theta$	FLORES-200								UDHR							
		CLD3		FT176		OpenLID		NLLB		CLD3		FT176		OpenLID		NLLB	
		$ L  = 96$	$ L  = 96$	$ L  = 108$	$ L  = 108$	$ L  = 195$	$ L  = 195$	$ L  = 188$	$ L  = 188$	$ L  = 100$	$ L  = 100$	$ L  = 124$	$ L  = 124$	$ L  = 159$	$ L  = 159$	$ L  = 172$	$ L  = 172$
		F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$
baselines	.0	<u>.952</u>	<b>.0104</b>	<u>.881</u>	<b>.0093</b>	<b>.923</b>	<b>.0051</b>	<u>.950</u>	<b>.0053</b>	<u>.922</u>	<u>.0101</u>	<u>.739</u>	<b>.0081</b>	<u>.881</u>	<b>.0063</b>	<u>.854</u>	<b>.0058</b>
GlottLID-M	.0	<b>.983</b>	<b>.0104</b>	<b>.991</b>	<b>.0093</b>	<u>.922</u>	<b>.0051</b>	<b>.954</b>	<b>.0053</b>	<b>.952</b>	<b>.0100</b>	<b>.927</b>	<b>.0081</b>	<b>.926</b>	<u>.0064</u>	<b>.925</b>	<u>.0060</u>

# Language verification

## Closed set assumption

- set of possible languages known (differs for each baseline)



LID Model	$\theta$	FLORES-200								UDHR							
		CLD3		FT176		OpenLID		NLLB		CLD3		FT176		OpenLID		NLLB	
		$ L  = 96$	$ L  = 108$	$ L  = 195$	$ L  = 188$	$ L  = 100$	$ L  = 124$	$ L  = 159$	$ L  = 172$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$
baselines	.0	<u>.952</u>	<b>.0104</b>	<u>.881</u>	<b>.0093</b>	<b>.923</b>	<b>.0051</b>	<u>.950</u>	<b>.0053</b>	<u>.922</u>	<u>.0101</u>	<u>.739</u>	<b>.0081</b>	<u>.881</u>	<b>.0063</b>	<u>.854</u>	<b>.0058</b>
GlottLID-M	.0	<b>.983</b>	<b>.0104</b>	<b>.991</b>	<b>.0093</b>	<u>.922</u>	<b>.0051</b>	<b>.954</b>	<b>.0053</b>	<b>.952</b>	<b>.0100</b>	<b>.927</b>	<b>.0081</b>	<b>.926</b>	<u>.0064</u>	<b>.925</b>	<u>.0060</u>

## Corpus building for low-resource languages

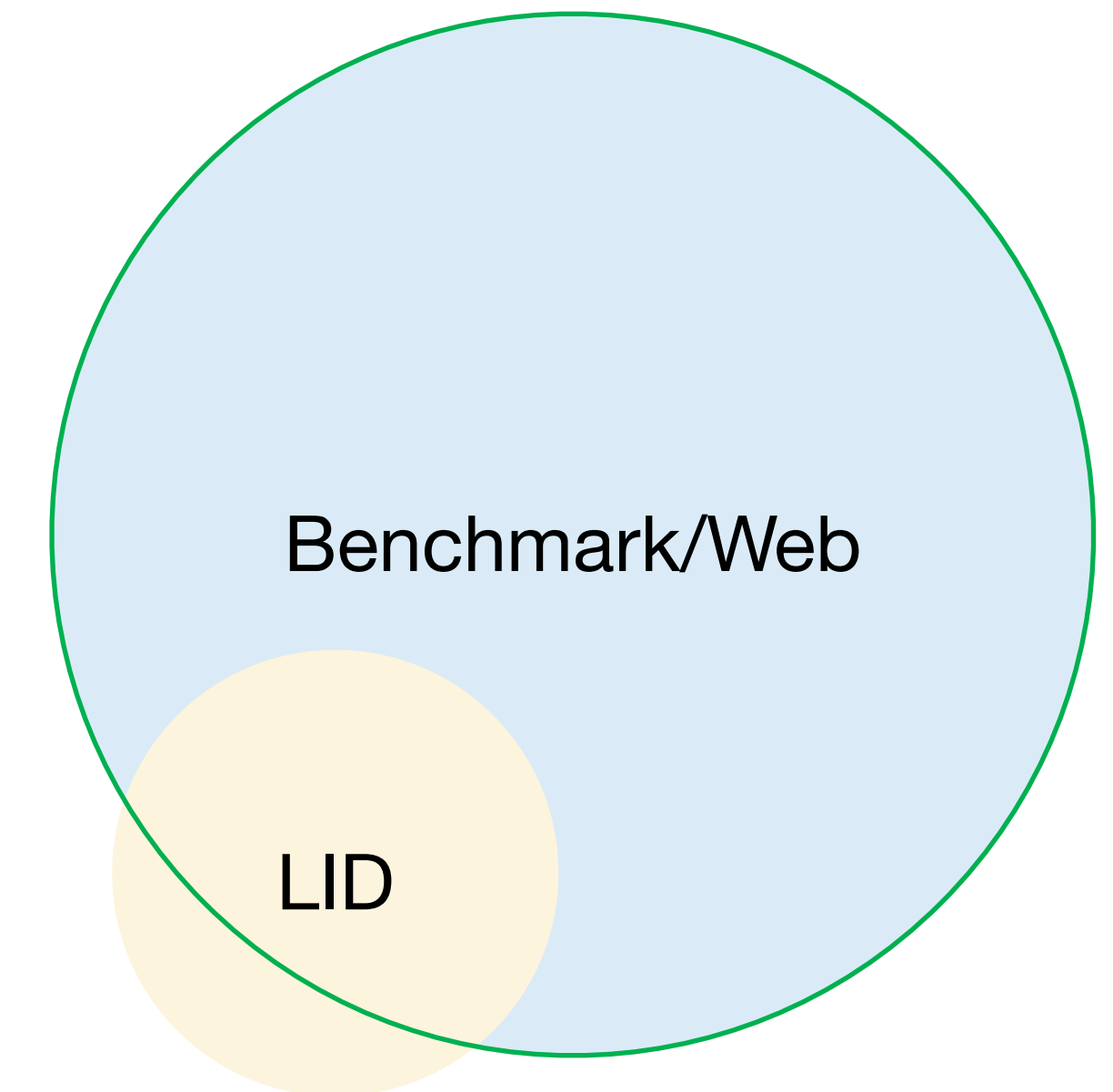
- Negatives >> Positives
- FPR matters
- F1 not important for HR

# Evaluating GlotLID

## Open set assumption : a realistic setting

- set of possible languages is unknown
- rejection matters

Benchmark		L	- rejection		+ rejection	
			GlottLID-M, $\theta=.0$		GlottLID-M, $\theta=.5$	
			F1 $\uparrow$	FPR $\downarrow$	F1 $\uparrow$	FPR $\downarrow$
GlottLID-C	all	1832	.940	.0005	.938	.0003
GlottLID-C	subset	1665	.977	.0003	.973	.0002
UDHR	all	374	.750	.0015	.734	.0007
UDHR	subset	342	.784	.0014	.770	.0006
FLORES-200	all	196	.917	.0042	.887	.0013
FLORES-200	subset	177	.957	.0029	.924	.0010



# Language Identification unsolved

	FLORES-200						UDHR					
	language	FP	cl	top FP source	#FP	%	language	FP	cl	top FP source	#FP	%
most errors	arb:St Arabic	3787	.18	ars:Najdi Arabi	829	.22	cmn:Mandarin Ch	596	.38	chr:Cherokee	81	.14
	arz:Egyptian Ar	1726	.32	apc:Levantine A	440	.25	qub:Huellaga Hu	247	.00	qvh:Huamalies-D	55	.22
	pes:Ir. Persian	1495	.40	prs:Dari	905	.61	fin:Finnish	224	.22	krl:Karelian	138	.62
	cmn:Mandarin Ch	1008	.00	yue:Yue Chinese	1008	.99	wuu:Wu Chinese	172	.24	hak:Hakka Chine	44	.26
	hin:Hindi	977	.51	awa:Awadhi	693	.71	rus:Russian	157	.28	niv:Gilyak	44	.28
most noisy	arb:St Arabic	3787	.18	ars:Najdi Arabi	829	.22	evn:Evenki	36	.23	oaa:Orok	19	.53
	arz:Egyptian Ar	1726	.32	apc:Levantine A	440	.25	quz:Cusco Quech	82	.40	qxu:Arequipa-La	61	.74
	prs:Dari	338	.24	pbt:S Pashto	310	.92	hrv:Croatian	84	.42	bos:Bosnian	39	.46
	dyu:Dyula	255	.25	bam:Bambara	255	.99	tzm:C Atlas Tam	52	.02	zgh:St Moroccan	52	.99
	apc:Levantine A	161	.42	ajp:S Levantine	70	.43	uzn:N Uzbek	72	.46	cbu:Candoshi-Sh	16	.22

Kargaran et al 2023: <https://arxiv.org/abs/2310.16248>

# Language Identification unsolved

	FLORES-200						UDHR					
	language	FP	cl	top FP source	#FP	%	language	FP	cl	top FP source	#FP	%
most errors	arb:St Arabic	3787	.18	ars:Najdi Arabi	829	.22	cmn:Mandarin Ch	596	.38	chr:Cherokee	81	.14
	arz:Egyptian Ar	1726	.32	apc:Levantine A	440	.25	qub:Huellaga Hu	247	.00	qvh:Huamalies-D	55	.22
	pes:Ir. Persian	1495	.40	prs:Dari	905	.61	fin:Finnish	224	.22	krl:Karelian	138	.62
	cmn:Mandarin Ch	1008	.00	yue:Yue Chinese	1008	.99	wuu:Wu Chinese	172	.24	hak:Hakka Chine	44	.26
	hin:Hindi	977	.51	awa:Awadhi	693	.71	rus:Russian	157	.28	niv:Gilyak	44	.28
most noisy	arb:St Arabic	3787	.18	ars:Najdi Arabi	829	.22	evn:Evenki	36	.23	oaa:Orok	19	.53
	arz:Egyptian Ar	1726	.32	apc:Levantine A	440	.25	quz:Cusco Quech	82	.40	qxu:Arequipa-La	61	.74
	prs:Dari	338	.24	pbt:S Pashto	310	.92	hrv:Croatian	84	.42	bos:Bosnian	39	.46
	dyu:Dyula	255	.25	bam:Bambara	255	.99	tzm:C Atlas Tam	52	.02	zgh:St Moroccan	52	.99
	apc:Levantine A	161	.42	ajp:S Levantine	70	.43	uzn:N Uzbek	72	.46	cbu:Candoshi-Sh	16	.22

Kargaran et al 2023: <https://arxiv.org/abs/2310.16248>

## New issues

- more realistic data (register, domains, etc)
- code-mixed inputs
- improved calibration



# Language Identification unsolved

	FLORES-200						UDHR					
	language	FP	cl	top FP source	#FP	%	language	FP	cl	top FP source	#FP	%
most errors	arb:St Arabic	3787	.18	ars:Najdi Arabi	829	.22	cmn:Mandarin Ch	596	.38	chr:Cherokee	81	.14
	arz:Egyptian Ar	1726	.32	apc:Levantine A	440	.25	qub:Huallaga Hu	247	.00	qvh:Huamalies-D	55	.22
	pes:Ir. Persian	1495	.40	prs:Dari	905	.61	fin:Finnish	224	.22	krl:Karelian	138	.62
	cmn:Mandarin Ch	1008	.00	yue:Yue Chinese	1008	.99	wuu:Wu Chinese	172	.24	hak:Hakka Chine	44	.26
	hin:Hindi	977	.51	awa:Awadhi	693	.71	rus:Russian	157	.28	niv:Gilyak	44	.28
most noisy	arb:St Arabic	3787	.18	ars:Najdi Arabi	829	.22	evn:Evenki	36	.23	oaa:Orok	19	.53
	arz:Egyptian Ar	1726	.32	apc:Levantine A	440	.25	quz:Cusco Quech	82	.40	qxu:Arequipa-La	61	.74
	prs:Dari	338	.24	pbt:S Pashto	310	.92	hrv:Croatian	84	.42	bos:Bosnian	39	.46
	dyu:Dyula	255	.25	bam:Bambara	255	.99	tzm:C Atlas Tam	52	.02	zgh:St Moroccan	52	.99
	apc:Levantine A	161	.42	ajp:S Levantine	70	.43	uzn:N Uzbek	72	.46	cbu:Candoshi-Sh	16	.22

Kargaran et al 2023: <https://arxiv.org/abs/2310.16248>

## New issues

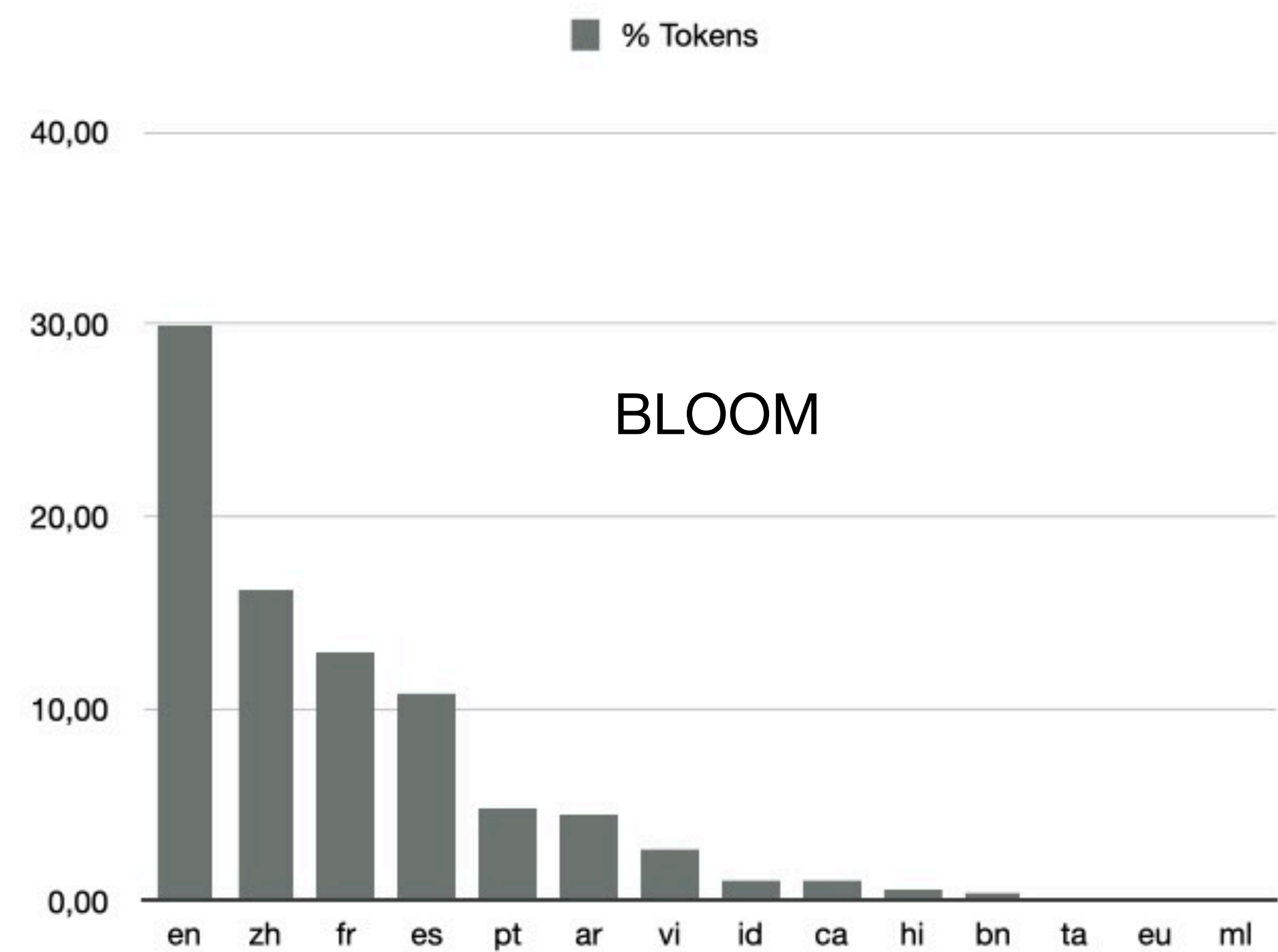
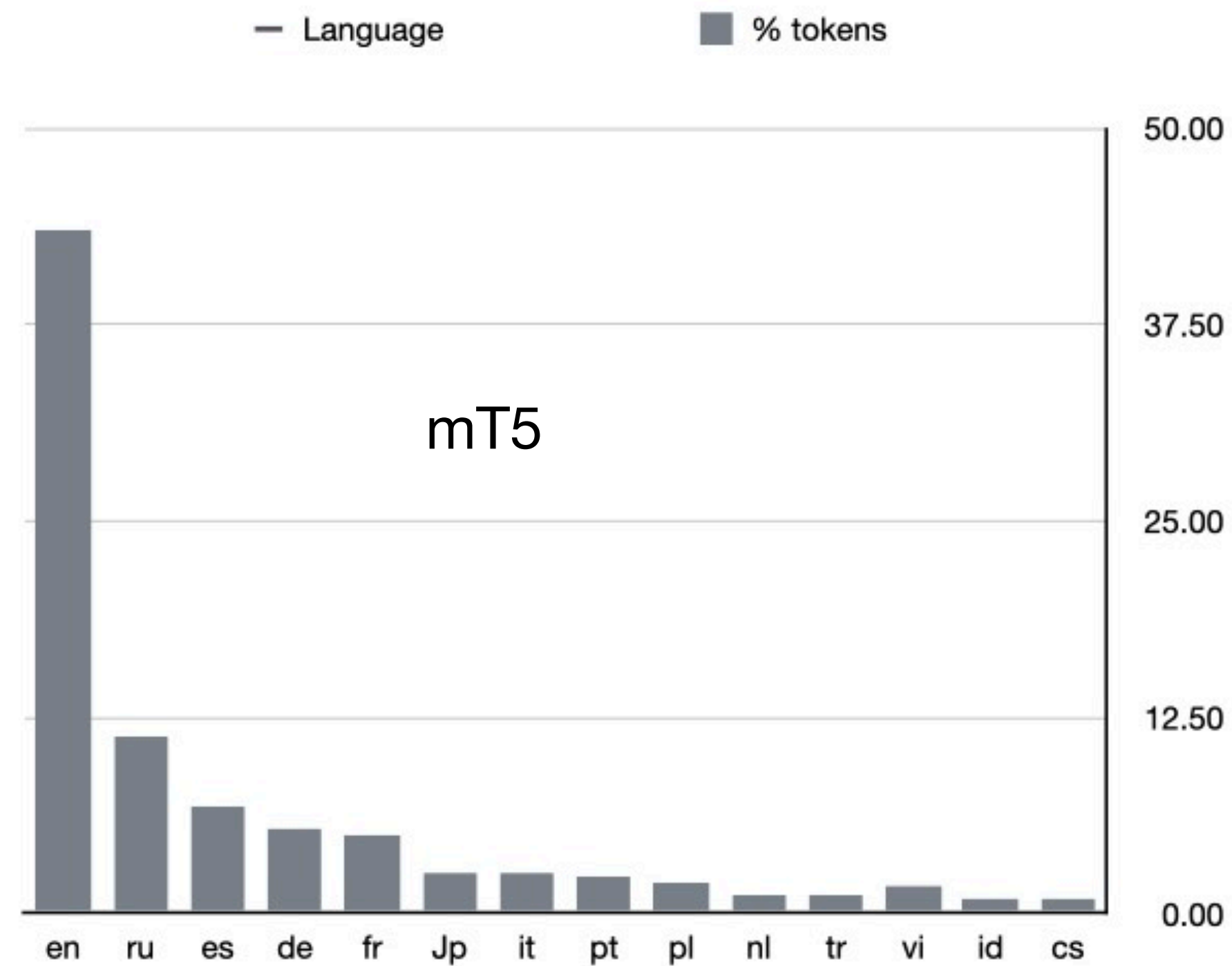
- more realistic data (register, domains, etc)
- code-mixed inputs
- improved calibration

The screenshot shows the GlotLID web interface. At the top, there is a logo for GlotLID with a bar chart and a globe icon. Below the logo, there are social media and repository links: Hugging Face Model, GitHub, license Apache-2.0, stars 41, and arXiv 2310.16248. The main text states: "GlotLID is an open-source language identification model with support for more than 1600 languages." There are two buttons: "Input a Sentence" and "Upload a File". Under "Choose model", there are two radio buttons: "v1" (selected) and "v2" (GlotLID version 2 (more data and languages)). The "Sentence:" field contains the text: "The work towards constructing the Latinxua Sinwenz (Chinese: 拉丁化新文字; pinyin: Lādīnghuà Xīn Wéi)". A "Submit" button is below the text. The result shows "Label: eng\_Latn, Language: English" with a horizontal bar chart showing a confidence score of approximately 0.9. At the bottom, there is a URL: <https://huggingface.co/spaces/cis-lmu/glotlid-space>.

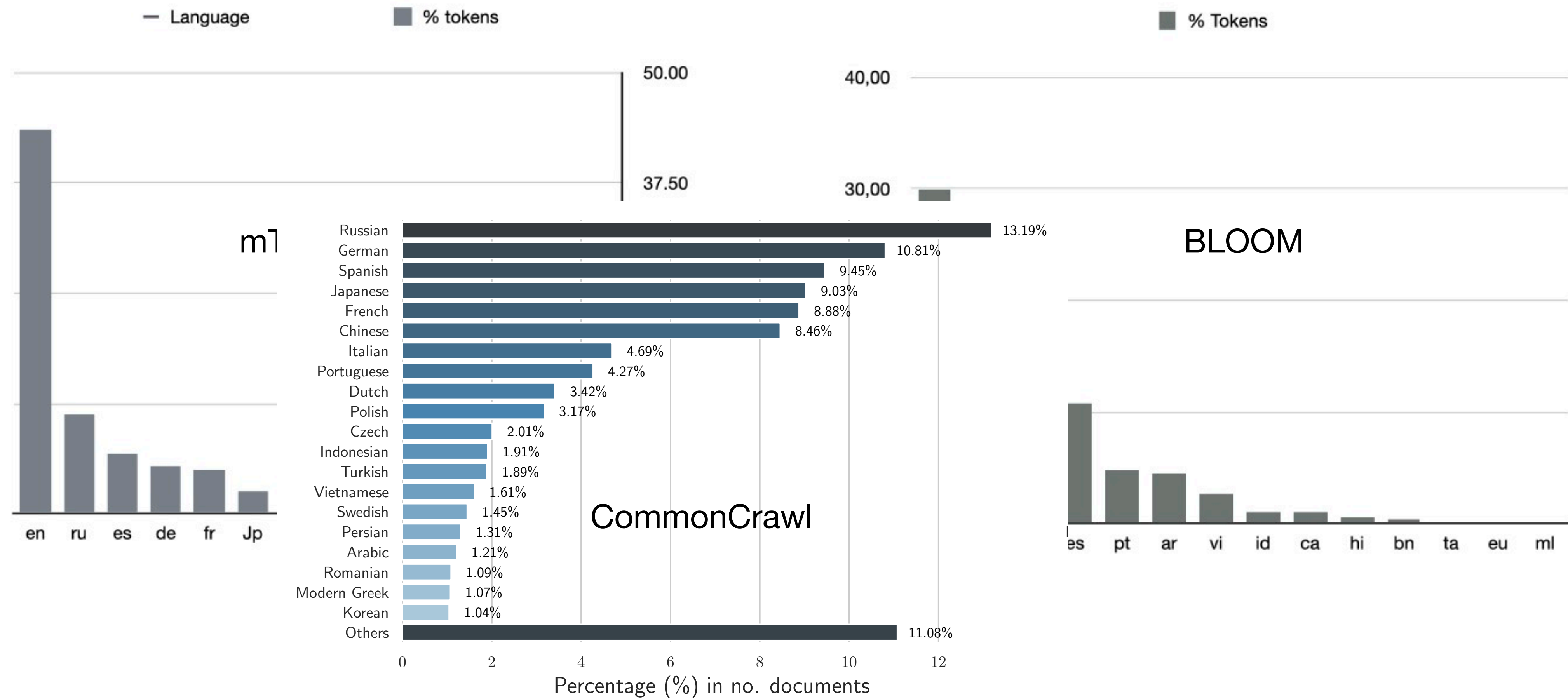
# Data unbalance in mLLMs ---

# Data unbalance in mLLMs

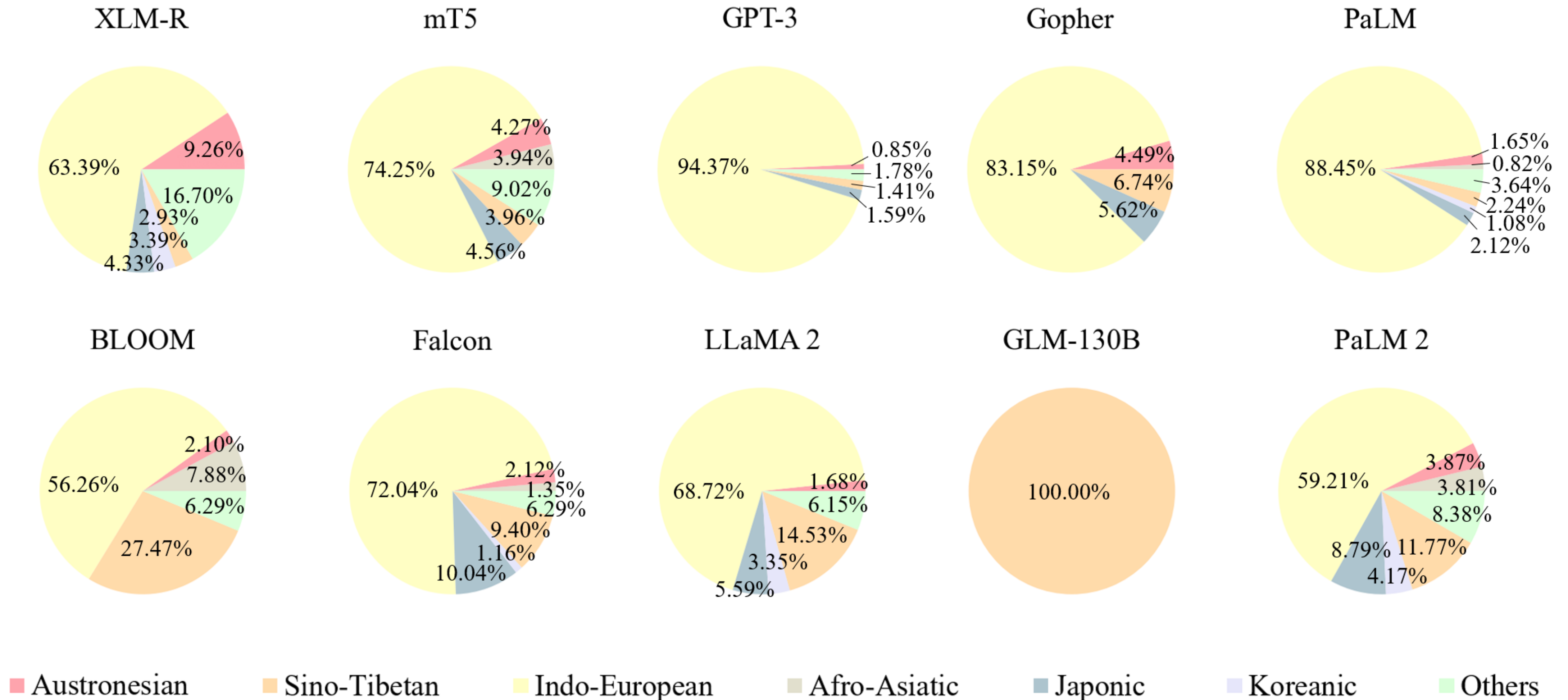
---



# Data unbalance in mLLMs



# Data unbalance in mLLMs



***“there is no multilingualism,  
only proofs of multilingualism”***

# Training Glot500-m

---

## Starting point: XLM-R-(B)

- 100 head languages
- « Pure encoder »
- Trained on CommonCrawl with MLM loss
- 250k vocab with *SentencePiece (SP)*
- base (270m) and large (550m) parameters

[Unsupervised Cross-lingual Representation Learning at Scale](#) (Conneau et al., ACL 2020)

**An “academic” configuration**

# Training Glot500-m

---

## Starting point: XLM-R-(B)

- 100 head languages
- « Pure encoder »
- Trained on CommonCrawl with MLM loss
- 250k vocab with *SentencePiece (SP)*
- base (270m) and large (550m) parameters

[Unsupervised Cross-lingual Representation Learning at Scale](#) (Conneau et al., ACL 2020)

## Extended subword vocabulary

- train SP model on Glot500-m (250k)
- temp = 0.3
- merge « old » and « new » types
- 401k vocabulary

**An “academic” configuration**



# Training Glot500-m

---

## Starting point: XLM-R-(B)

- 100 head languages
- « Pure encoder »
- Trained on CommonCrawl with MLM loss
- 250k vocab with *SentencePiece (SP)*
- base (270m) and large (550m) parameters

[Unsupervised Cross-lingual Representation Learning at Scale](#) (Conneau et al., ACL 2020)

## Extended subword vocabulary

- train SP model on Glot500-m (250k)
- temp = 0.3
- merge « old » and « new » types
- 401k vocabulary

## Training regime

- random language mixtures (temp = 0.3)
- MLM loss
- no change in model size
- two weeks of computation

**An “academic” configuration**

# Proofs of multilingualism

## mLLMs as a set of monolingual models

## mLLMs as representation learners

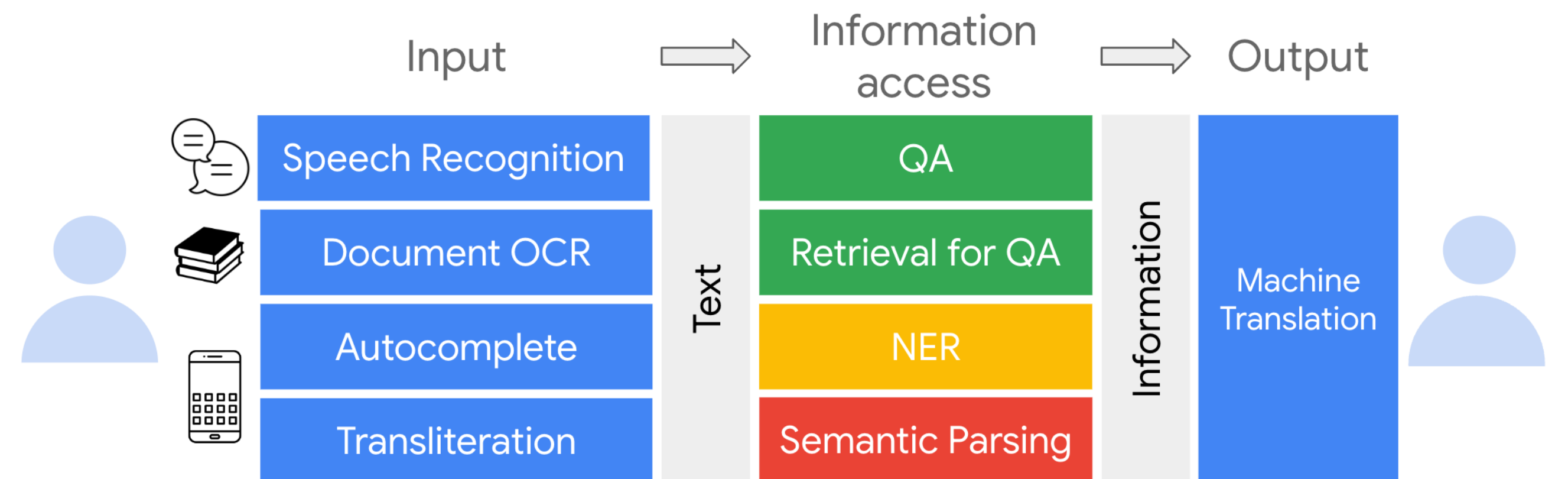
- recovering good bilingual associations
- yielding good (word, sentence) alignments
- encoding linguistic features
- effective cross-lingual performance

## A wealth of multilingual benchmarks

- XTreme, XTreme-R, XTreme-Up
- X-GLUE
- Mega, MegaVerse
- BUFFET

## mLLMs for text generation

- good models of multilingual texts
- good machine translation systems
- generating realistic mixed-language



[Xtremeup: A User-Centric Scarce-Data Benchmark for Under-Represented Languages \(Ruder et al, 2023\)](#)

# Proofs of multilingualism

## mLLMs as a set of monolingual models

## mLLMs as representation learners

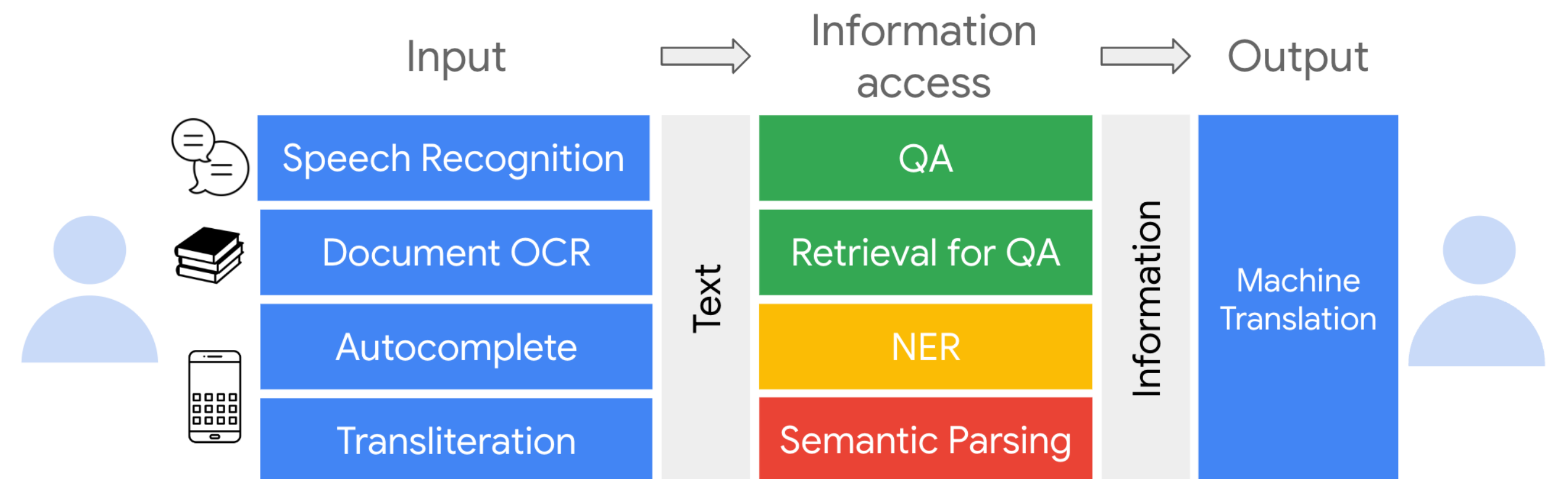
- recovering good bilingual associations
- yielding good (word, sentence) alignments
- encoding linguistic features
- effective cross-lingual performance

## A wealth of multilingual benchmarks

- XTreme, XTreme-R, XTreme-Up
- X-GLUE
- Mega, MegaVerse
- BUFFET

## mLLMs for text generation

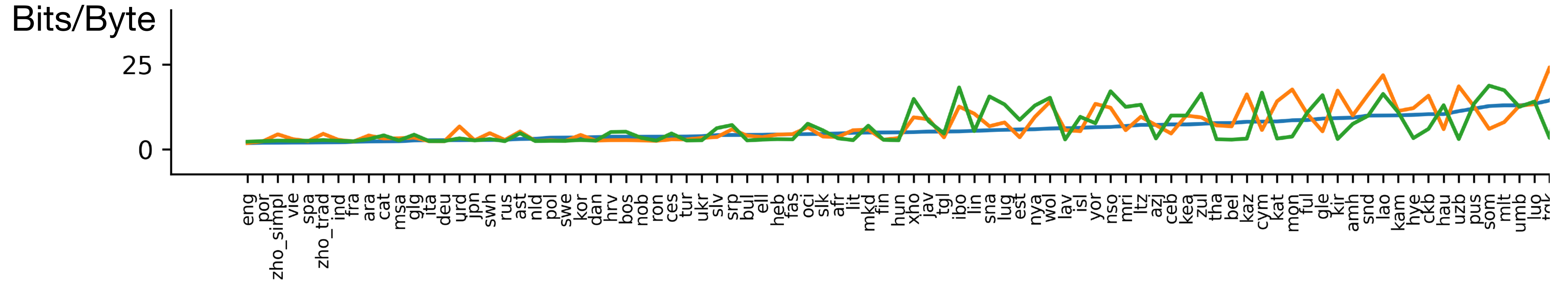
- good models of multilingual texts
- good machine translation systems
- generating realistic mixed-language



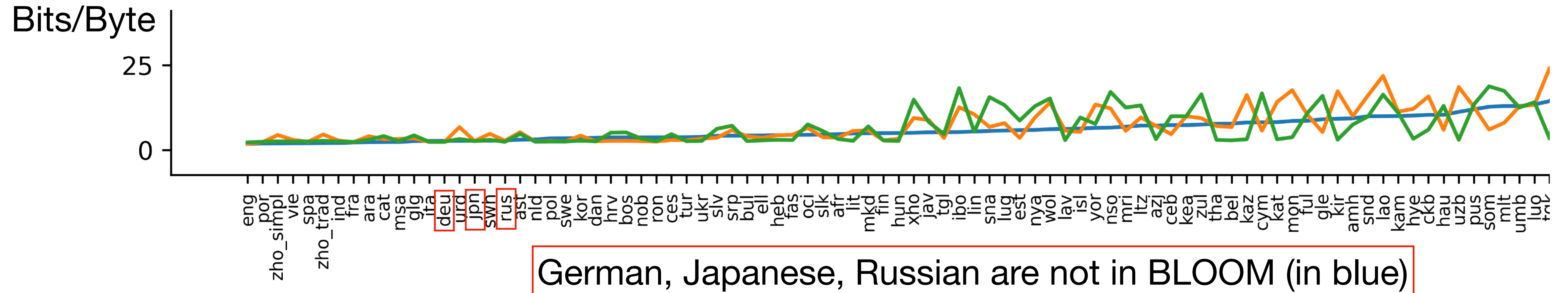
[Xtremeup: A User-Centric Scarce-Data Benchmark for Under-Represented Languages \(Ruder et al, 2023\)](#)

**Require annotated benchmarks - possibly via MT**

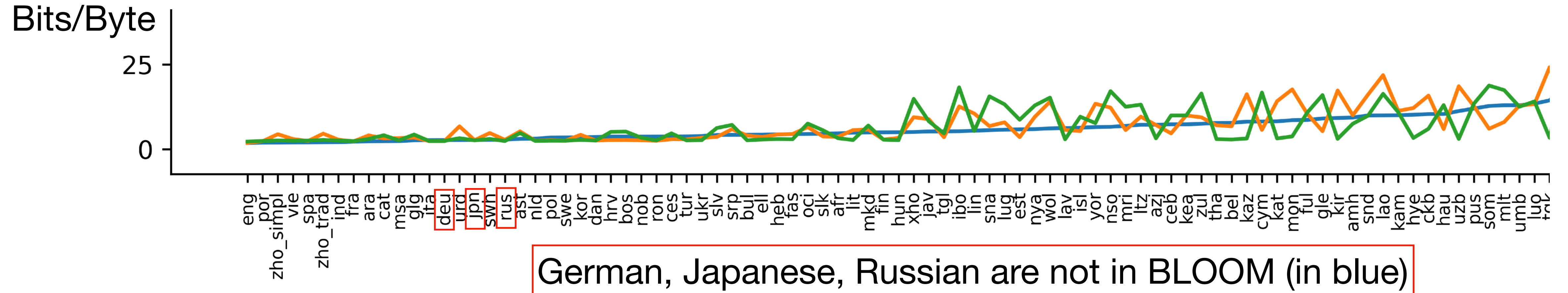
# PPL is a weak signal



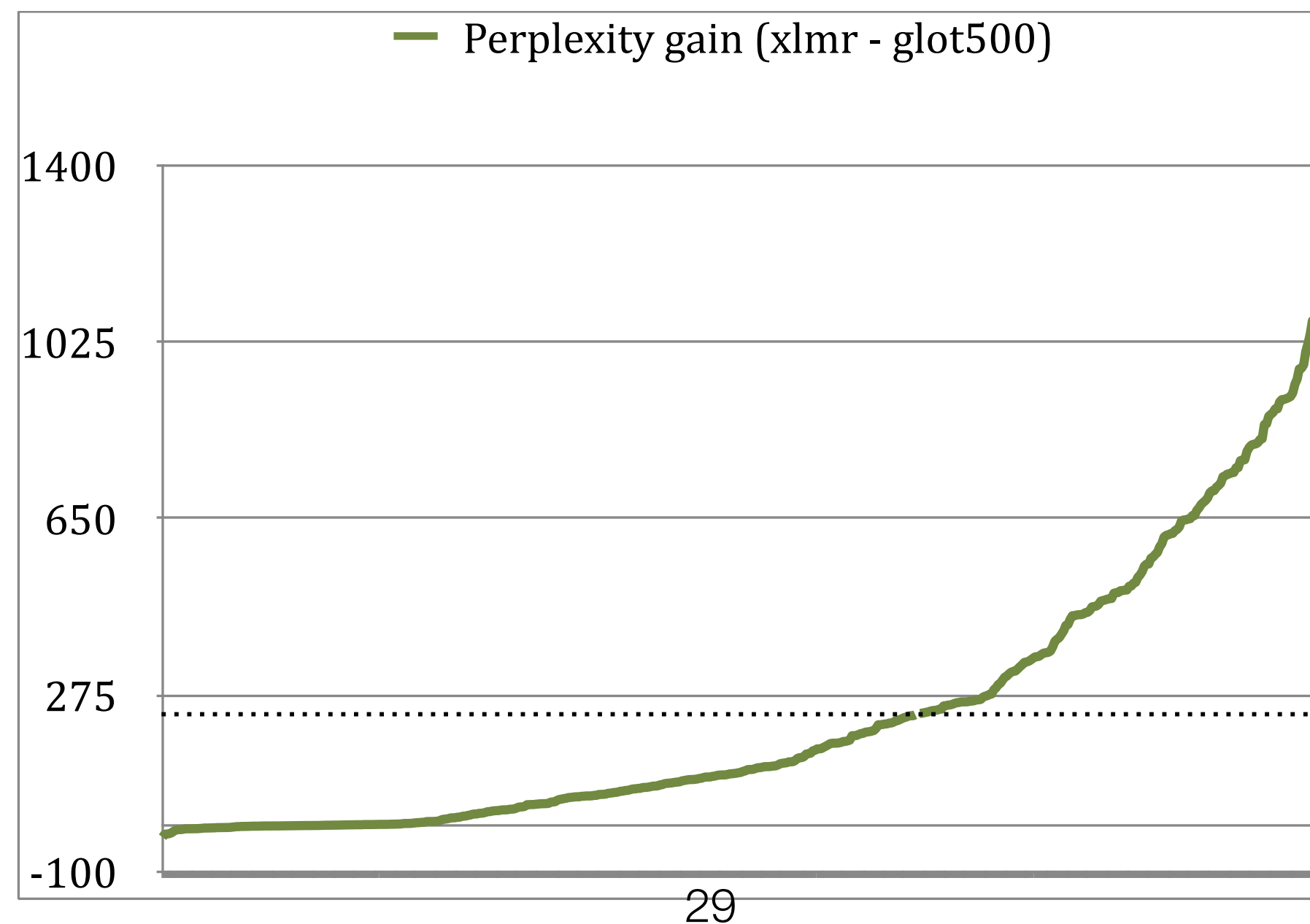
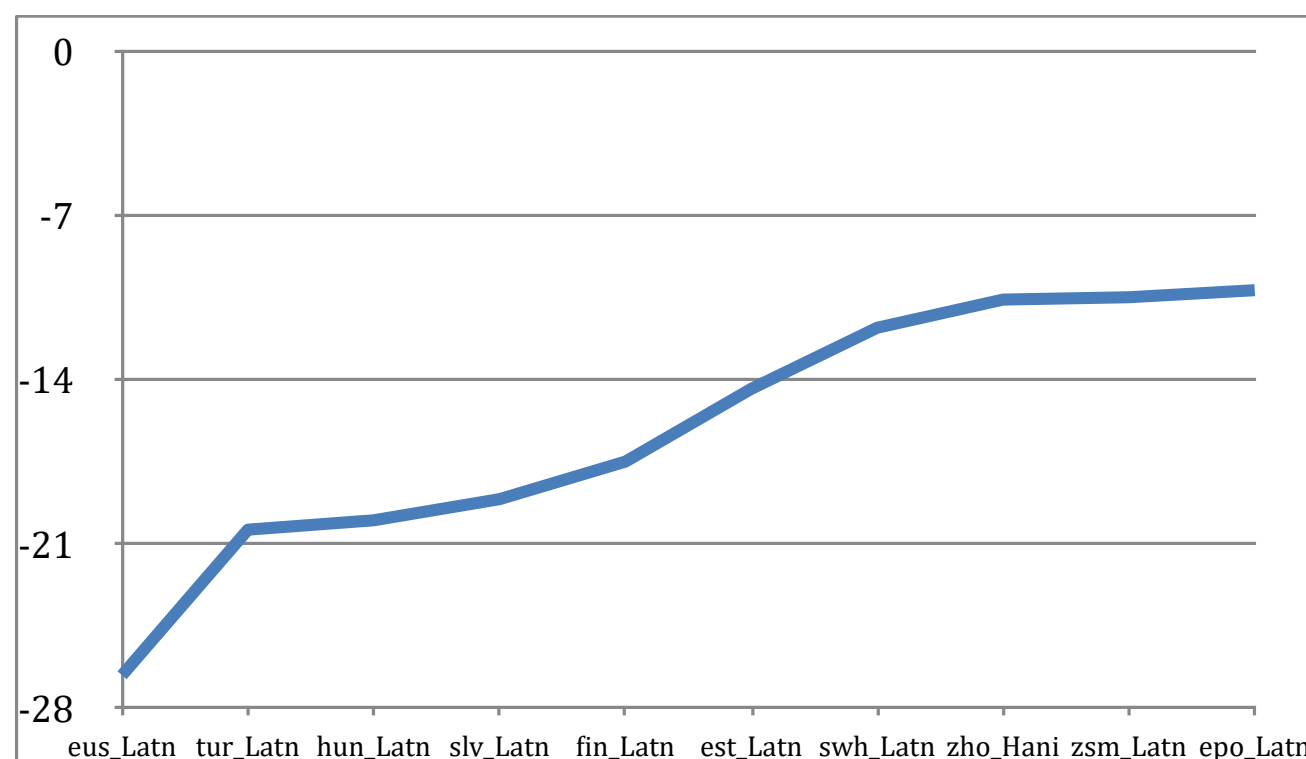
# PPL is a weak signal



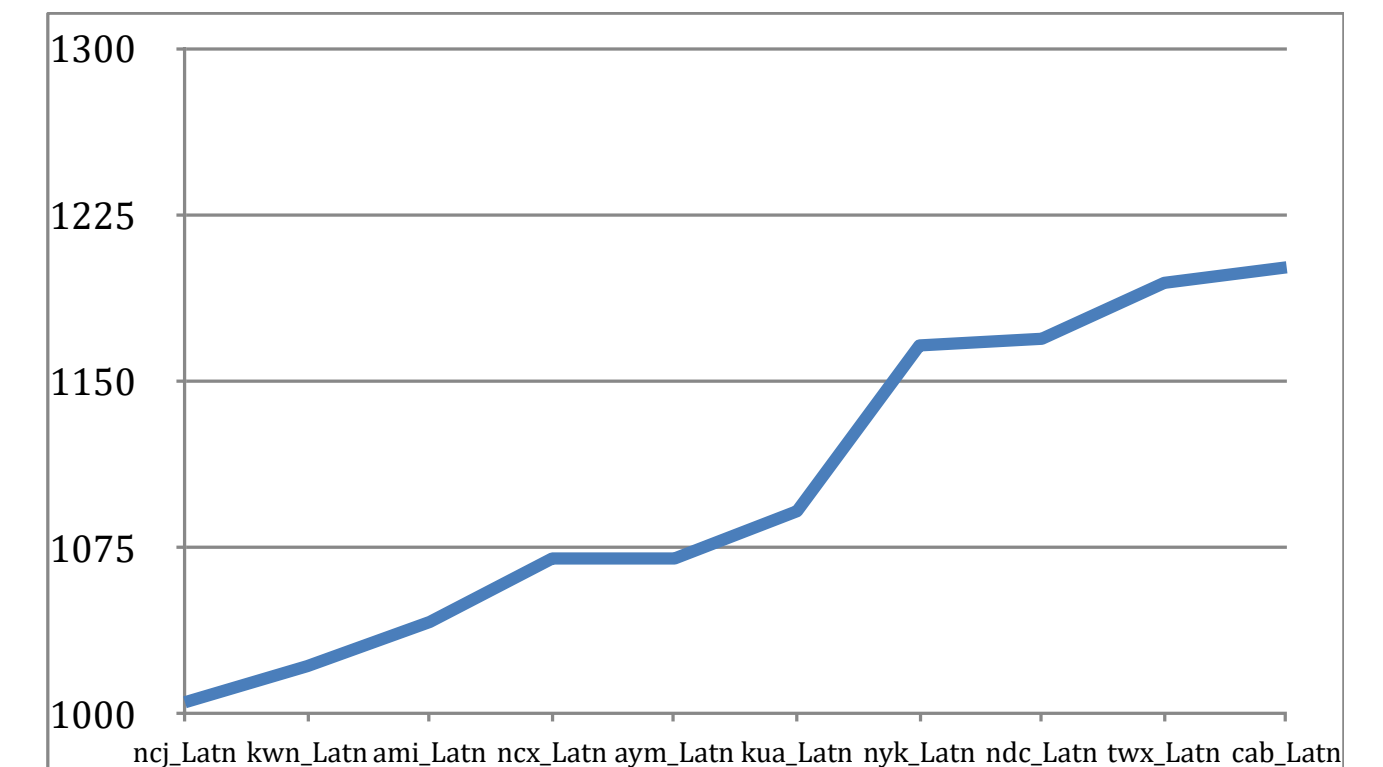
# PPL is a weak signal



“losers” (87)

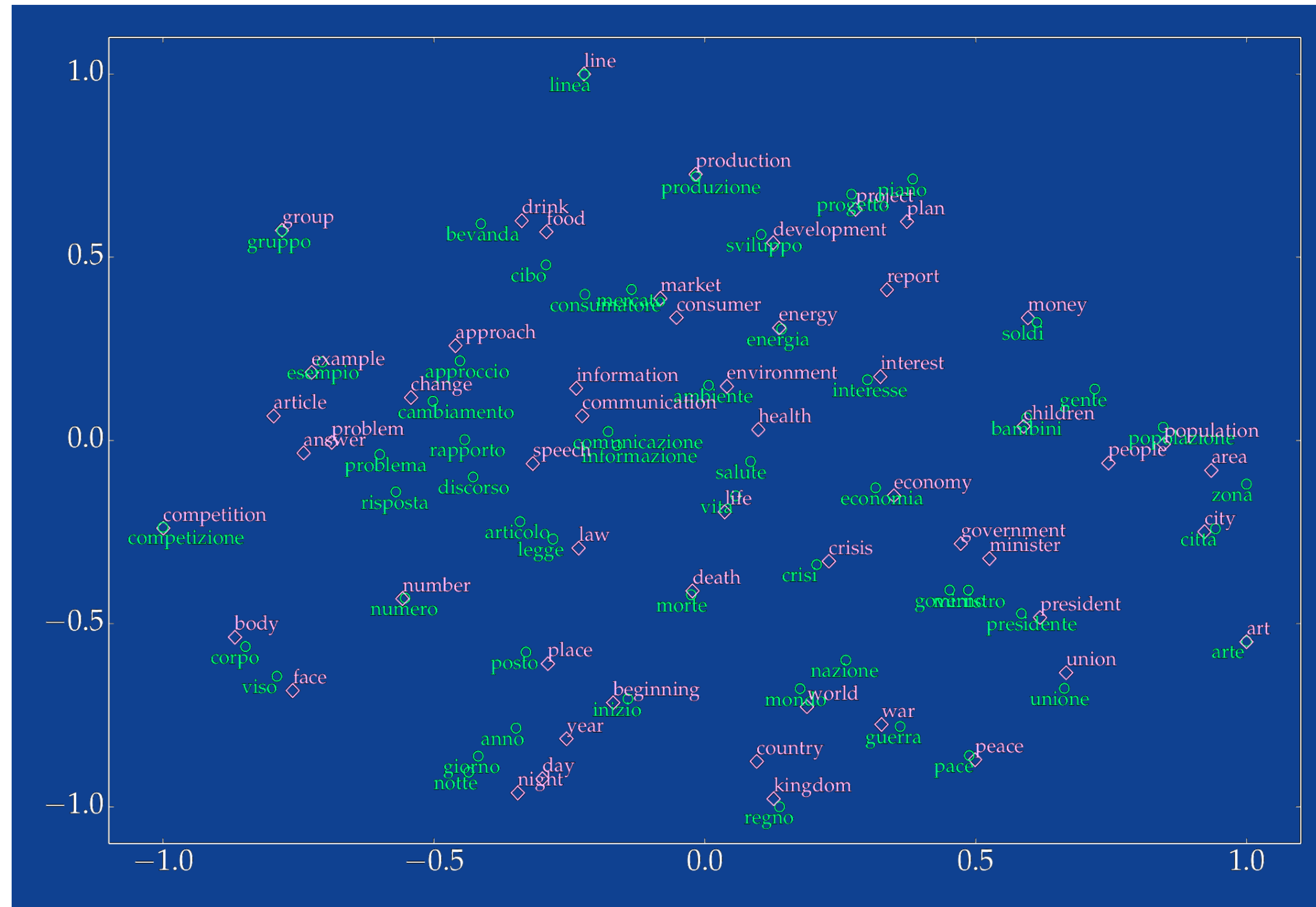


“winners” (424)



# mReps yield good alignments

---



[Unsupervised Cross-Lingual Representation Learning](#) (Ruder et al., ACL)

**Bilingual Lexicon Induction with context independent models**

# mReps yield good alignments ---

## **Similarity $\Rightarrow$ Alignment**

+ monotonicity, low distortion

+ symmetry, low fertility, “pigeon hole” principle



# mReps yield good alignments

## Sentence level

In the gayest and happiest spirits she set forward with her father;	Elle partit avec son père, le visage souriant;
not always listening, but always agreeing to what he said;	elle n' écoutait pas toujours, mais elle acquiesçait de confiance.
They arrived .	Ils arrivèrent .
It is Frank and Miss Fairfax, said Mrs. Weston .	– C'est Frank et Mlle Fairfax, dit aussitôt Mme Weston .
I was just going to tell you of our agreeable surprize in seeing him arrive this morning.	– J'allai justement vous faire part de l'agréable surprise que nous avons eue en le voyant arriver.
He stays till tomorrow, and Miss Fairfax has been persuaded to spend the day with us .	Il reste jusqu'à demain et Mlle Fairfax a bien voulu, sur notre demande , venir passer la journée.

## Similarity $\Rightarrow$ Alignment

+ monotonicity, low distortion

+ symmetry, low fertility, “pigeon hole” principle

# mReps yield good alignments

## Sentence level

In the gayest and happiest spirits she set forward with her father;	Elle partit avec son père, le visage souriant;
not always listening, but always agreeing to what he said;	elle n' écoutait pas toujours, mais elle acquiesçait de confiance.
They arrived .	Ils arrivèrent .
It is Frank and Miss Fairfax, said Mrs. Weston .	– C'est Frank et Mlle Fairfax, dit aussitôt Mme Weston .
I was just going to tell you of our agreeable surprize in seeing him arrive this morning.	– J'allai justement vous faire part de l'agréable surprise que nous avons eue en le voyant arriver.
He stays till tomorrow, and Miss Fairfax has been persuaded to spend the day with us .	Il reste jusqu'à demain et Mlle Fairfax a bien voulu, sur notre demande , venir passer la journée.

## Similarity $\Rightarrow$ Alignment

+ monotonicity, low distortion

+ symmetry, low fertility, “pigeon hole” principle

## Word level

Penguin Sir Nils Olav III. was knighted by the Norwegian king

Pinguin Nils Olav der Dritte wurde vom norwegischen König zum Ritter geschlagen

Sir Nils Olav III. です ペンギン knighted by el rey noruego

Nils Olav der Dritte is a penguin nominato cavaliere par un roi norvégien

# Assessing alignment

---

## Sentence Retrieval

- Tatoeba (68 head, 28 tail), acc@10
- Bible (94 head, 276 tail), acc@10

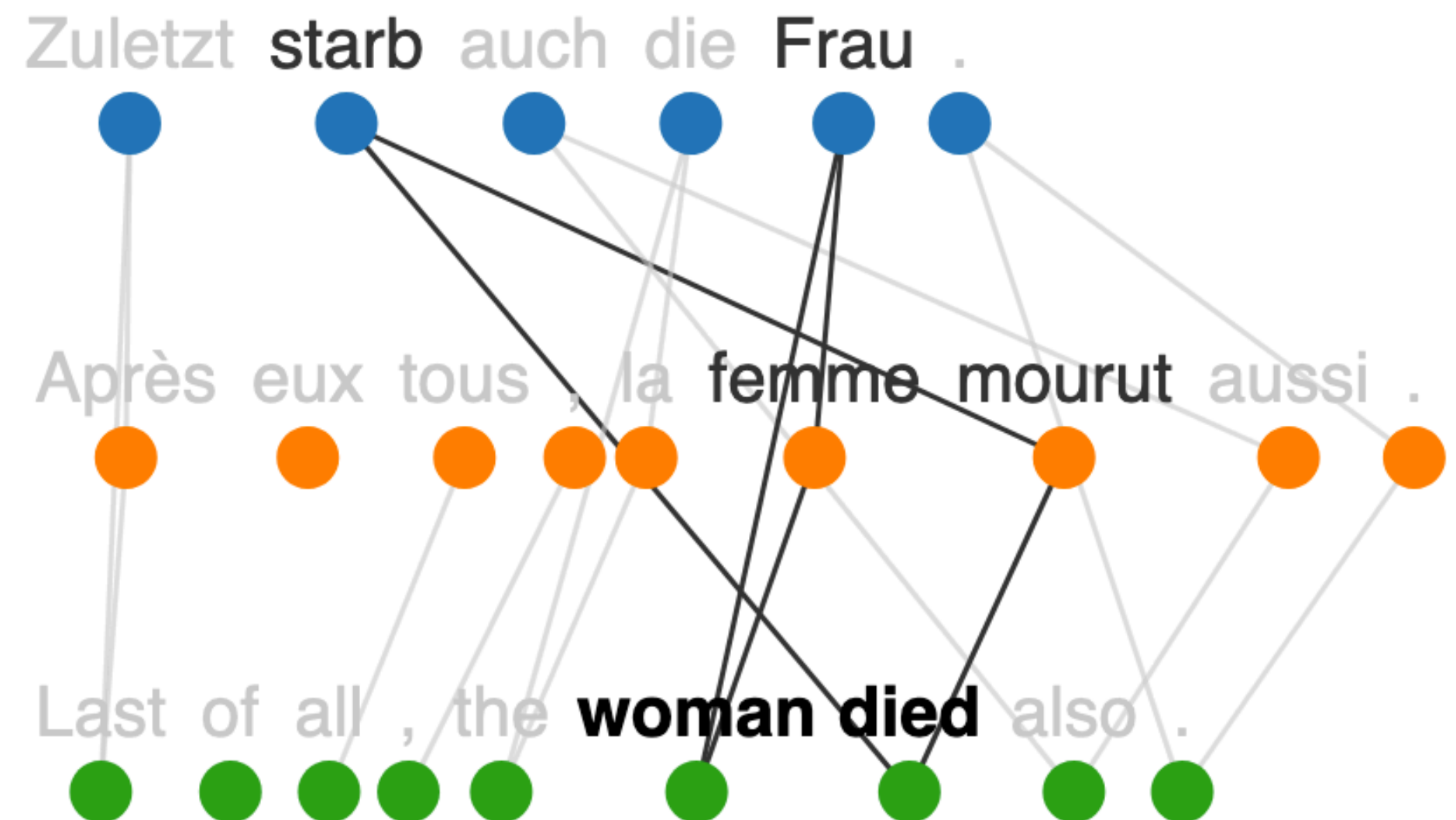
Find nearest foreign neighbor in multilingual space for 1000 (resp. 500) English sentences  
Requires parallel data

## Unsupervised multilingual evaluation

## Round trip Alignment

- Bible (95 head, 288 tail), acc@10

Word *SimAlign*  $L_1 \rightarrow L_2 \rightarrow L_3 \rightarrow L_4 \rightarrow L_1$   
Report exact matches, averaged over 5 runs  
Requires parallel data



# Assessing alignment

	<b>all</b>			<b>head</b>			<b>tail</b>		
	XLM-R-B	XLM-R-L	Glott500-m	XLM-R-B	XLM-R-L	Glott500-m	XLM-R-B	XLM-R-L	Glott500-m
SR [Tatoeba]	56.6	60.4	<b>70.7</b>	66.2	71.1	<b>75.0</b>	32.6	33.6	<b>59.8</b>
SR [Bible]	19.3	20.1	<b>47.3</b>	54.2	58.3	<b>59.0</b>	7.4	7.1	<b>43.2</b>
RTA	2.8	3.3	<b>4.7</b>	3.4	4.1	<b>5.5</b>	2.6	3.1	<b>4.5</b>

## Glott500-m vs XLM-R

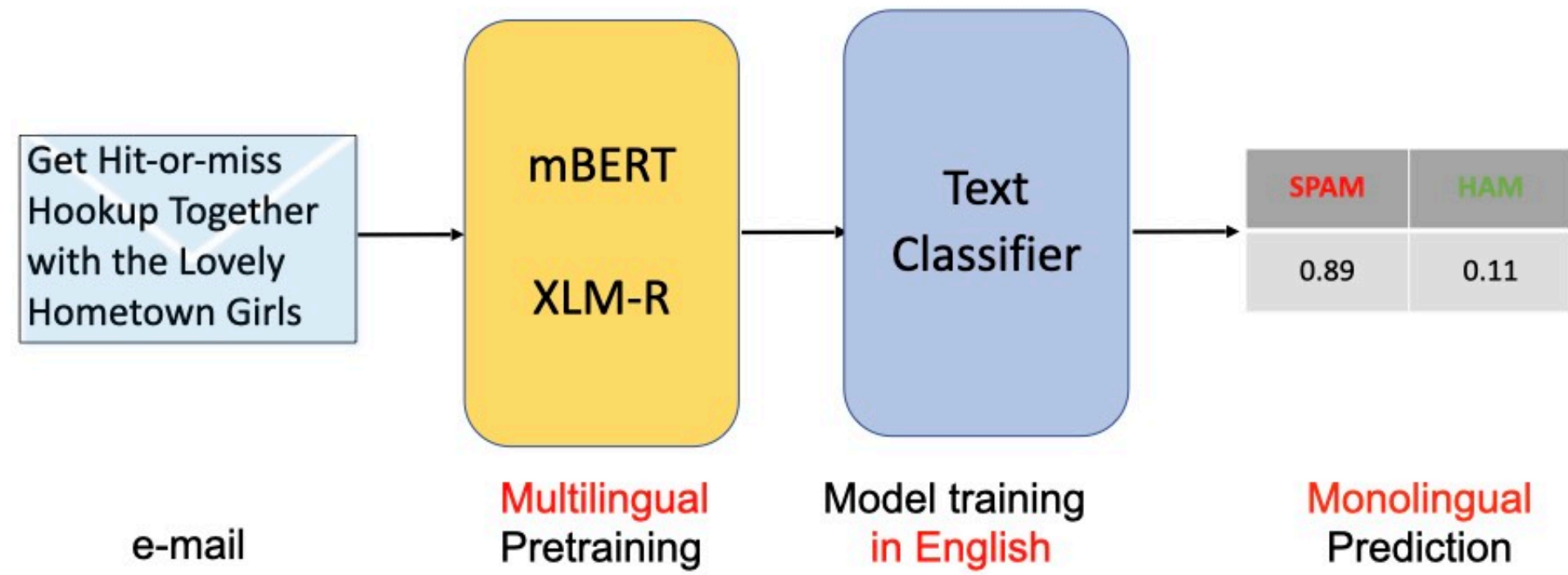
- outperforms all models on average
- better than XLM-R-B for head languages
- much better than XLM-R-\* for tail languages

## Caveats

- Only ~ 100 Tatoeba languages
- Tatoeba & Bible are very peculiar
- RTT is very hard

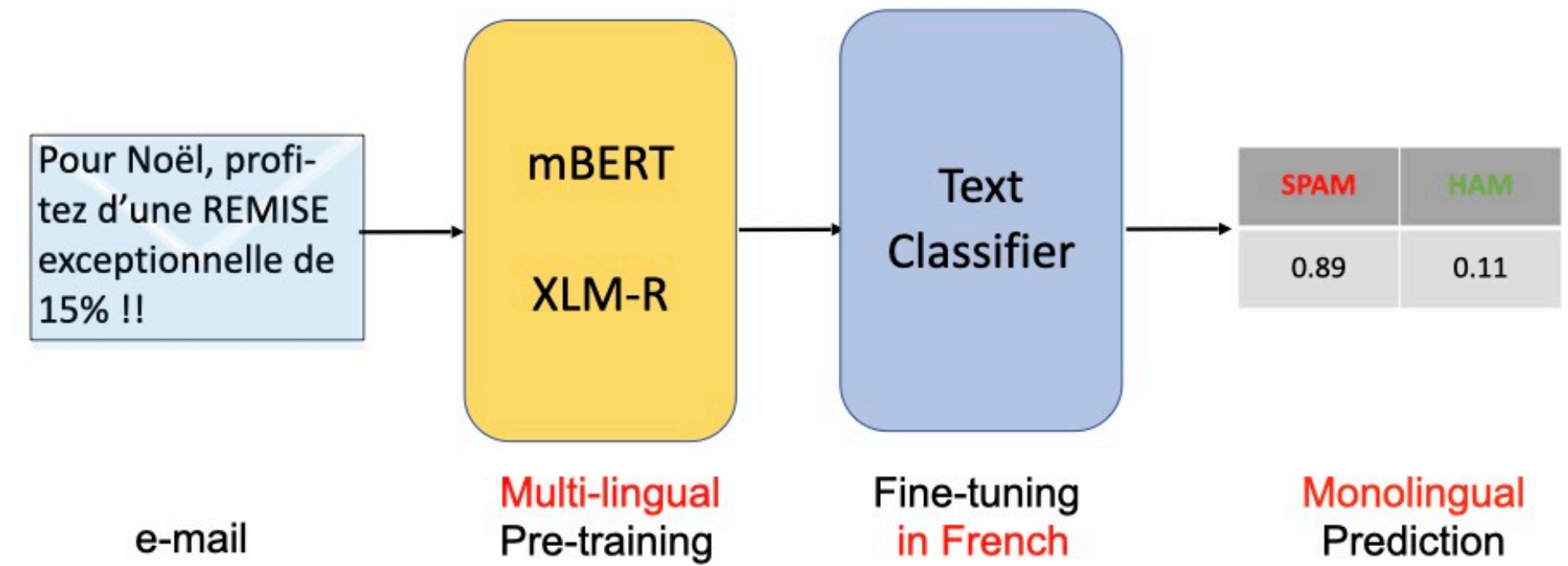
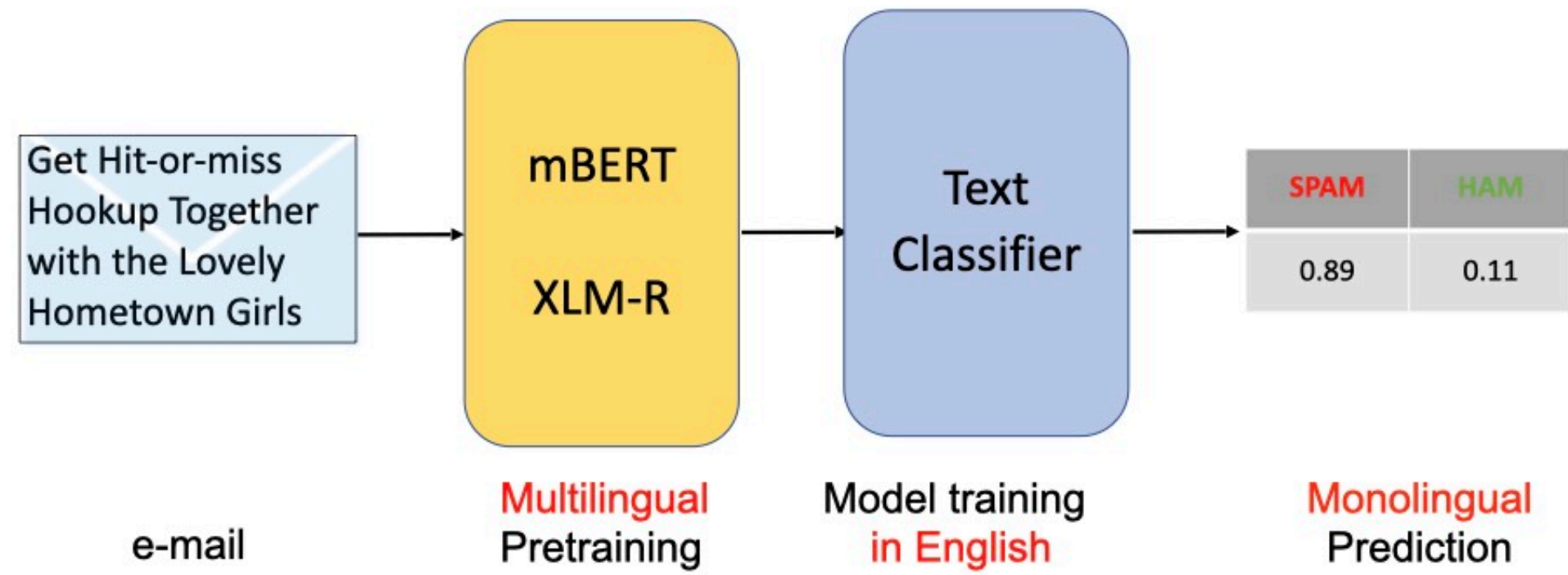
# Zero-shot X-lingual transfer

---

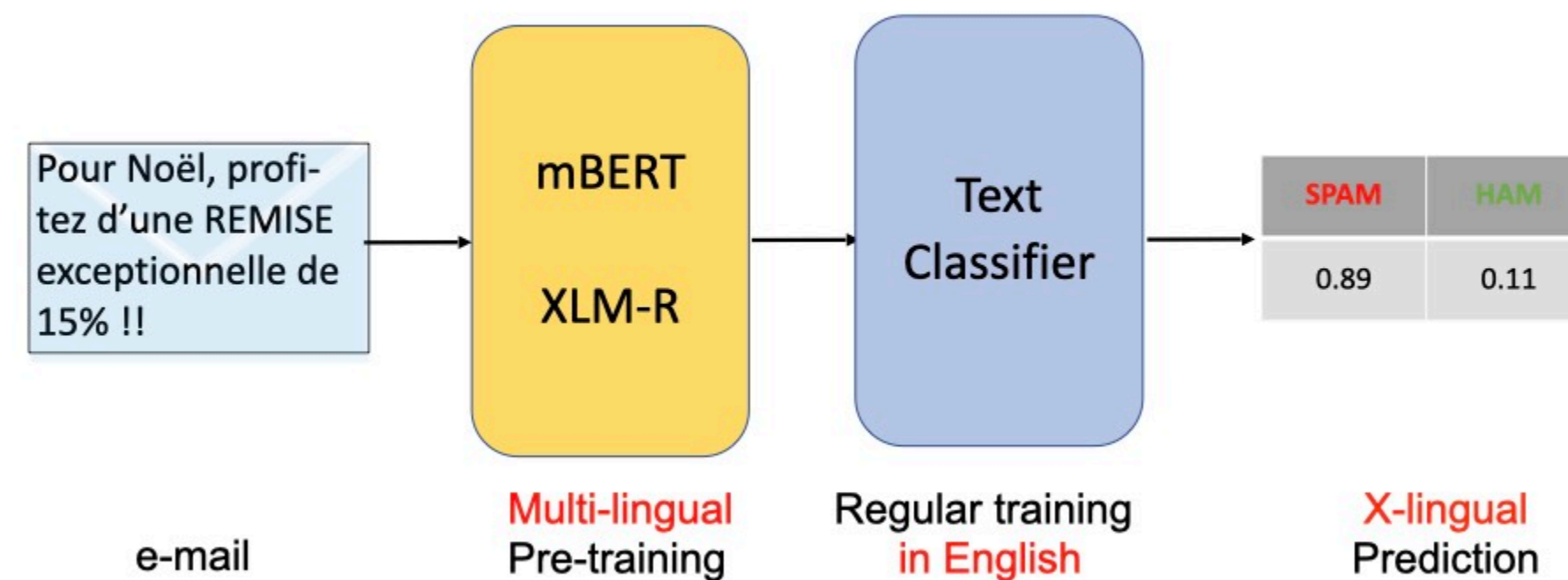
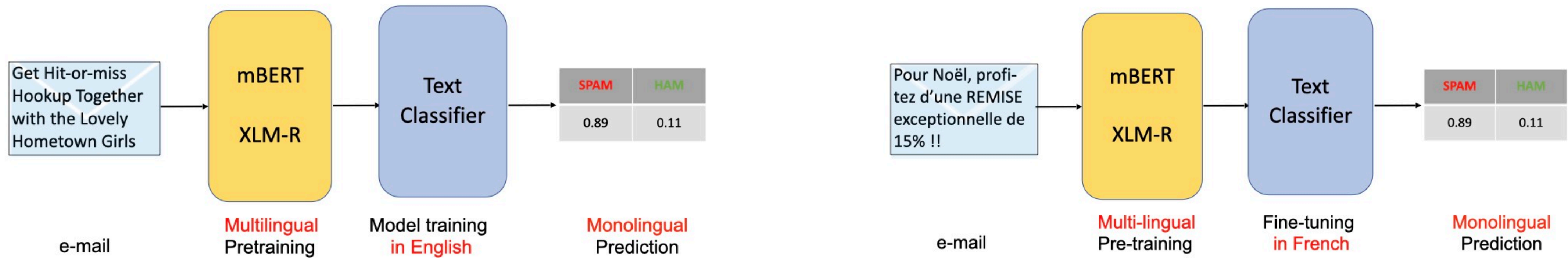


# Zero-shot X-lingual transfer

---



# Zero-shot X-lingual transfer



**Similar representations**



**Identical labels**

# Evaluating zero-shot transfer

---

## Text classification

- Sentence Classification (90 head, 284 tail) F1 using Taxi1500

Zero-shot transfer from English  
6-way classification for test data (Bible)

[Taxi1500: A Multilingual Dataset for Text Classification in 1500 Languages](#)  
C Ma et al (2023) arXiv preprint arXiv:2305.08487,

## Sentence labeling

- NER (89 head, 75 tail), F1 using wikiAnn
- POS (63 head, 28 tail), F1 using UD

Zero-shot transfer from English  
Requires gold labels

**Standard benchmarks with fine-grained annotations**  
**Mostly head languages**



# Improving tail languages

	all			head			tail		
	XLM-R-B	XLM-R-L	Glott500-m	XLM-R-B	XLM-R-L	Glott500-m	XLM-R-B	XLM-R-L	Glott500-m
TextClass	23.3	25.8	<b>48.7</b>	51.3	<b>60.5</b>	54.7	13.7	13.9	<b>46.6</b>
NER	55.3	59.5	<b>62.4</b>	61.8	<b>66.0</b>	63.9	47.5	51.8	<b>60.7</b>
POS	65.8	67.7	<b>71.8</b>	76.4	<b>78.4</b>	76.0	41.7	43.5	<b>62.3</b>

## Glott500-m vs XLM-R

- outperforms all models on average
- better than XLM-R-B for head languages
- much better than XLM-R-\* for tail languages

## Also

- are averages that informative ?
- tail language scores remain poor

# Where Glot500-m helps (or doesn't)

		language-script	XLMR	Glot500	gain		language-script	XLMR	Glot500	gain
high end	SentRetr Tatoeba	tat C Tatar	10.3	70.3	60.0	SentRetr Bible	uzn C Northern Uzbek	5.4	87.0	81.6
		nds L Low German	28.8	77.1	48.3		crs L Seselwa Creole	7.4	80.6	73.2
		tuk L Turkmen	16.3	63.5	47.3		srn L Sranan Tongo	6.8	79.8	73.0
		ile L Interlingue	34.6	75.6	41.0		uzb C Uzbek	6.2	78.8	72.6
		uzb C Uzbek	25.2	64.5	39.3		bcl L Central Bikol	10.2	79.8	69.6
low end	SentRetr Bible	dtp L Kadazan Dusun	5.6	21.1	15.5	xav L Xavánte	2.2	5.0	2.8	
		kab L Kabyle	3.7	16.4	12.7	mau L Huautla Mazatec	2.4	3.6	1.2	
		pam L Pampanga	4.8	11.0	6.2	ahk L Akha	3.0	3.2	0.2	
		lvs L Standard Latvian	73.4	76.9	3.5	aln L Gheg Albanian	67.8	67.6	-0.2	
		nob L Bokmål	93.5	95.7	2.2	nob L Bokmål	82.8	79.2	-3.6	
high end	NER	div T Dhivehi	0.0	50.9	50.9	POS	mlt L Maltese	21.3	80.3	59.0
		che C Chechen	15.3	61.2	45.9		sah C Yakut	21.9	76.9	55.0
		mri L Maori	16.0	58.9	42.9		sme L Northern Sami	29.6	73.6	44.1
		nan L Min Nan	42.3	84.9	42.6		yor L Yoruba	22.8	64.2	41.4
		tgk C Tajik	26.3	66.4	40.0		quc L K'iche'	28.5	64.1	35.6
low end	NER	zea L Zeeuws	68.1	67.3	-0.8	lzh H Literary Chinese	11.7	18.4	6.7	
		vol L Volapük	60.0	59.0	-1.0	nap L Neapolitan	47.1	50.0	2.9	
		min L Minangkabau	42.3	40.4	-1.8	hyw A Western Armenian	79.1	81.1	2.0	
		wuu H Wu Chinese	28.9	23.9	-5.0	kmr L Northern Kurdish	73.5	75.2	1.7	
		lzh H Literary Chinese	15.7	10.3	-5.4	aln L Gheg Albanian	54.7	51.2	-3.5	

## “Winners”:

- Languages not in XLM-R + large training data
- Scripts not in XLM-R
- “Cluster” effects

## “Losers”:

- Neighbors or superclass in XLM-R
- Small training data

# Glott500-m: complementary results

## Sentence Retrieval Bible

family	$ L_G $	$ L_X $	XLM-R-B	Glott500-m	gain
indo1319	91	50	41.5	61.4	19.9
atla1278	69	2	5.5	45.2	39.6
aust1307	53	6	13.7	47.0	33.2
turk1311	22	7	20.1	62.9	42.8
sino1245	22	2	7.6	38.9	31.3
maya1287	15	0	3.8	20.3	16.4
afro1255	12	5	13.0	34.3	21.4

Analysis per language family

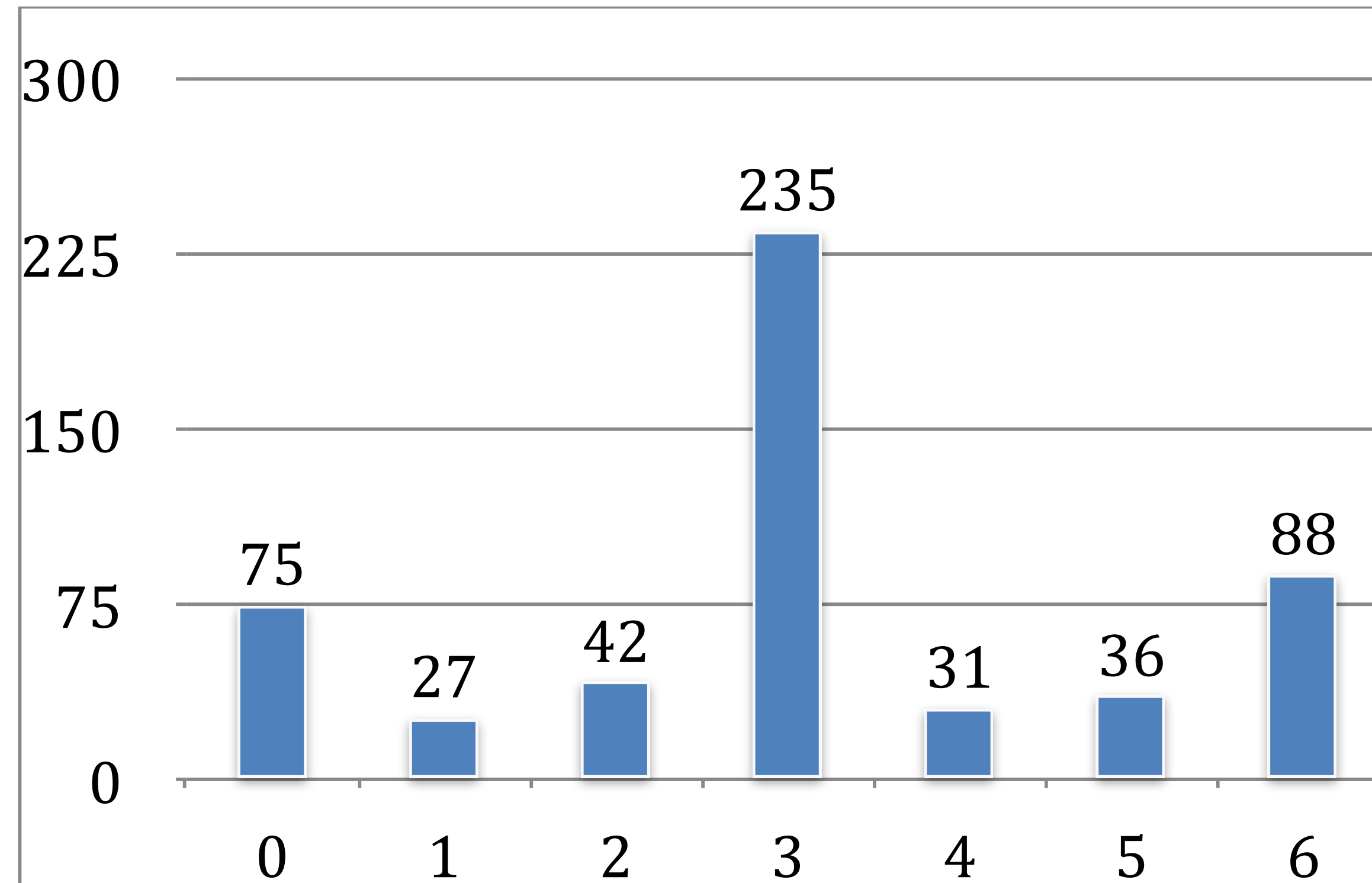
lang-script		XLM-R-B	Glott500	gain
uig_Arab	head	0.458	0.562	0.104
uig_Latn	tail	0.098	0.628	0.530
hin_Deva	head	0.670	0.766	0.096
hin_Latn	tail	0.136	0.432	0.296
uzb_Latn	head	0.548	0.676	0.128
uzb_Cyrl	tail	0.062	0.788	0.726
kaa_Cyrl	tail	0.176	0.738	0.562
kaa_Latn	tail	0.092	0.434	0.342
kmr_Cyrl	tail	0.040	0.424	0.384
kmr_Latn	tail	0.358	0.630	0.272
tuk_Cyrl	tail	0.136	0.650	0.514
tuk_Latn	tail	0.096	0.662	0.566

One language, two scripts

**Also: the “curse of multilinguality”**

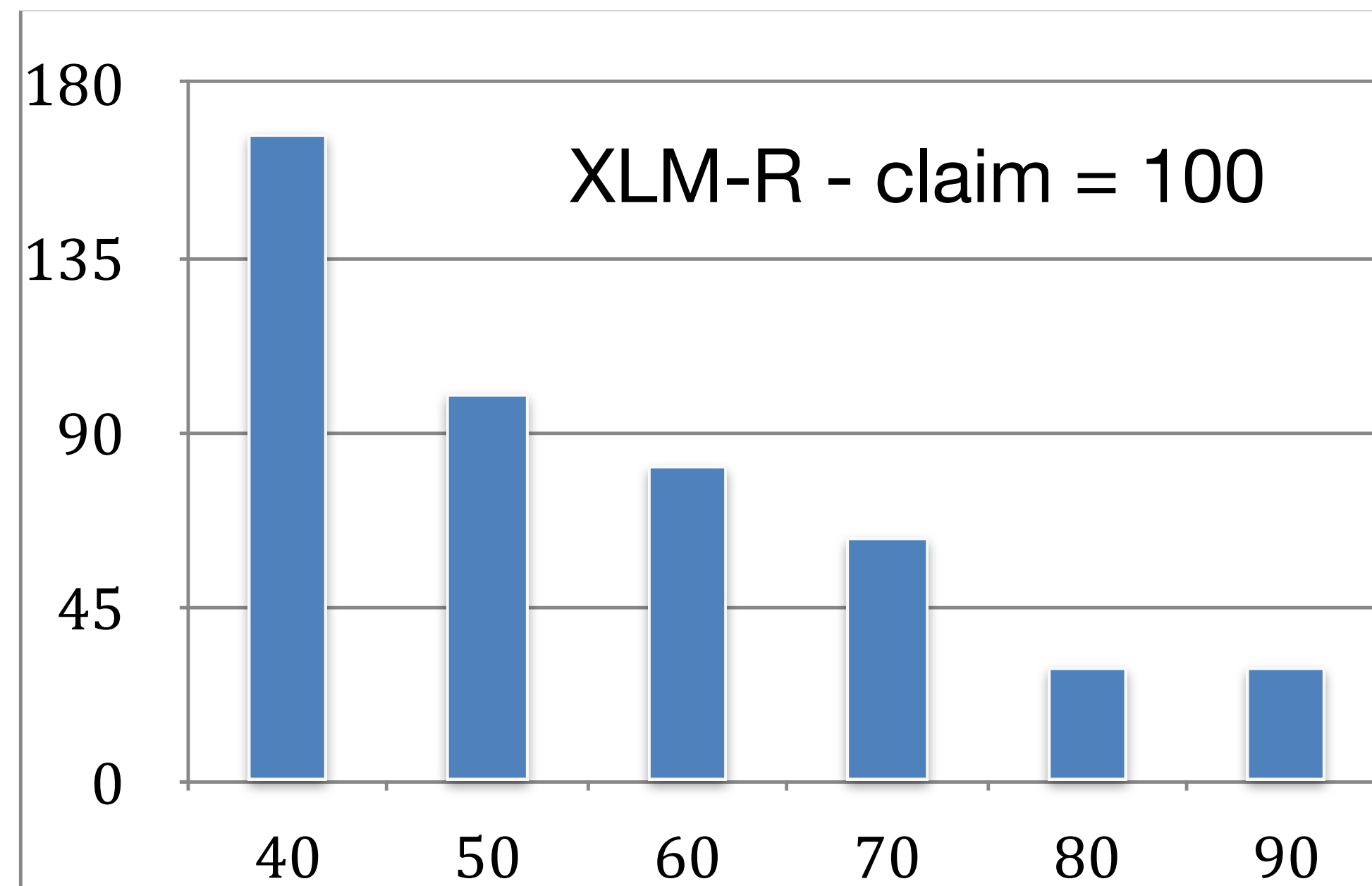
**Check paper for complete results / language**

# Glott500-m: what coverage? ---



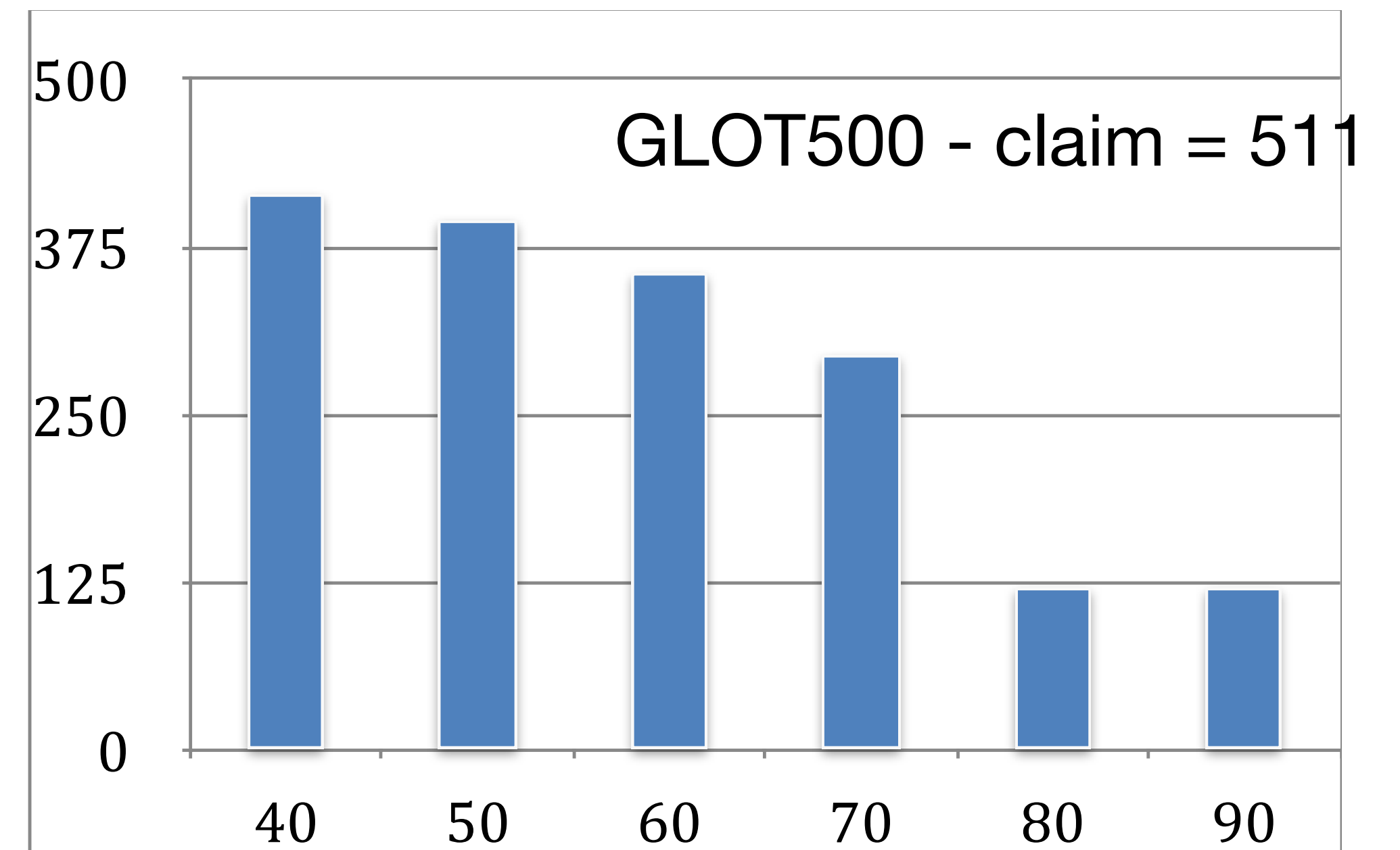
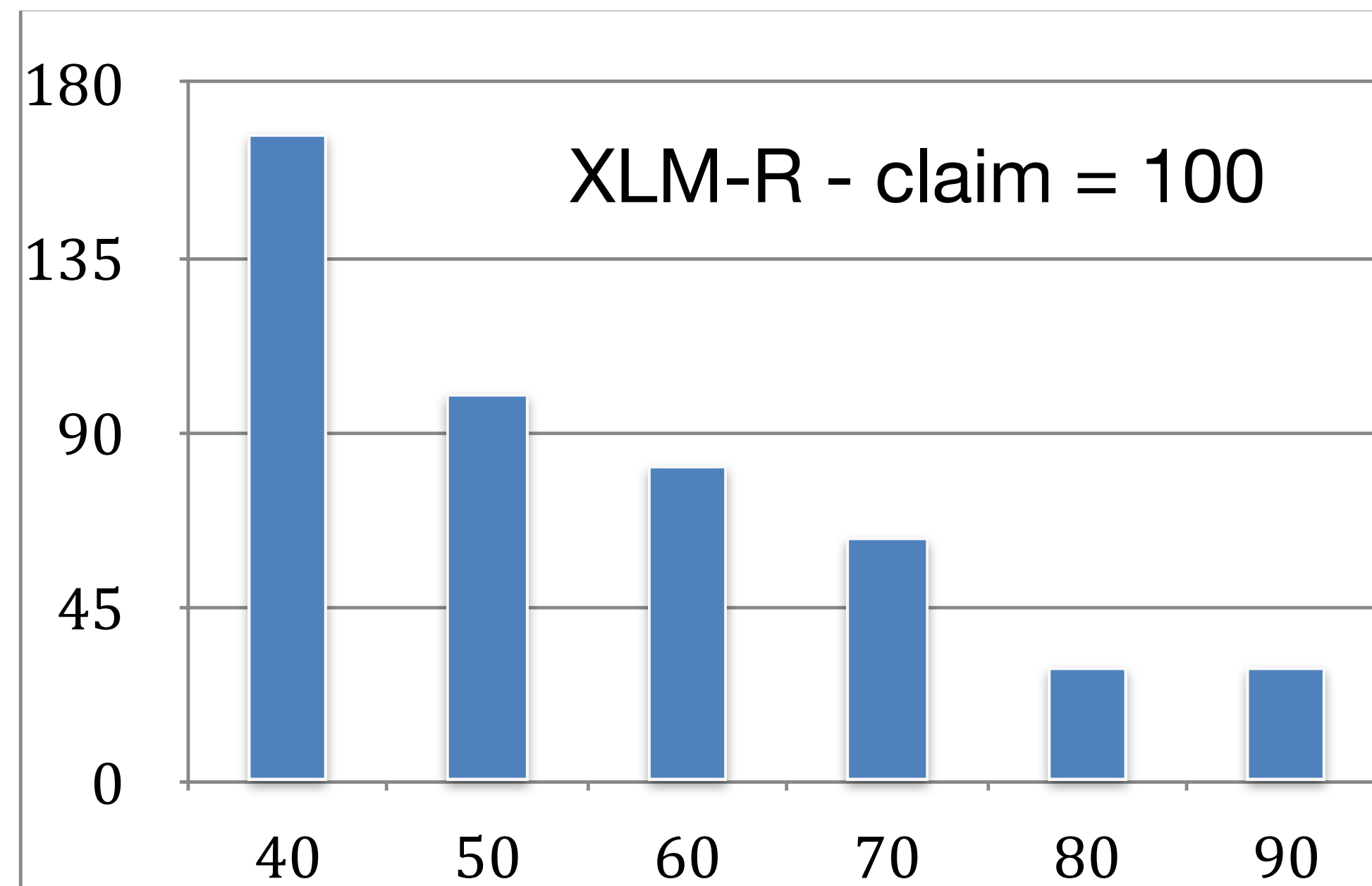
**Number of tasks evaluated**

# Glott500-m: what coverage? \_\_\_\_\_



**Performance wrt French**

# Glott500-m: what coverage? ---



**Performance wrt French**

# Glott500, some lessons learned

## A useful artefact

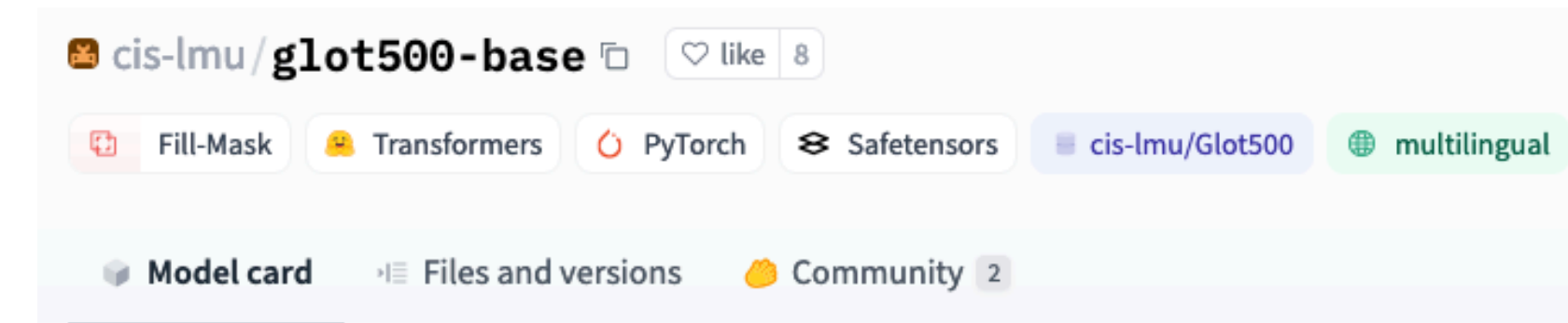
- extends the scope of “modern” NLP for tail languages
- open-source, documented corpus collection

## To be continued

- analysis of results
- analysis by language families / scripts
- linguistic analysis of representations
- evaluation of text generation abilities

## Limitations

- corpus selection and curation
- subword vocab size and training
- language choices
- filtering strategies



### Glott500 (base-sized model)

Glott500 model (Glott500-m) pre-trained on 500+ languages using a masked language modeling (MLM) objective. It was introduced in [this paper](#) (ACL 2023) and first released in [this repository](#).

### Usage

You can use this model directly with a pipeline for masked language modeling:

```
>>> from transformers import pipeline
>>> unmasker = pipeline('fill-mask', model='cis-lmu/glott500-base')
>>> unmasker("Hello I'm a <mask> model.")
```

Here is how to use this model to get the features of a given text in PyTorch:

```
>>> from transformers import AutoTokenizer, AutoModelForMaskedLM

>>> tokenizer = AutoTokenizer.from_pretrained('cis-lmu/glott500-base')
>>> model = AutoModelForMaskedLM.from_pretrained("cis-lmu/glott500-base")

>>> # prepare input
>>> text = "Replace me by any text you'd like."
>>> encoded_input = tokenizer(text, return_tensors='pt')

>>> # forward pass
>>> output = model(**encoded_input)
```

# mLLMs are still a wonder ---

## Open question and issues

mLLMs: more than a collection of monolingual models ?

- which linguistic properties help / break transfer?
- how critical is **alignment? parallel data?**
- how to measure positive / negative **interference?**
- how to measure language **coverage?**
- how to measure **non-linguistic biases** ?
  - racial, social, cultural, etc



# mLLMs are still a wonder ---

## Open question and issues

mLLMs: more than a collection of monolingual models ?

- which linguistic properties help / break transfer?
- how critical is **alignment? parallel data?**
- how to measure positive / negative **interference?**
- how to measure language **coverage?**
- how to measure **non-linguistic biases** ?
  - racial, social, cultural, etc

How about the “**curse of multilinguality**” ?

- impact of language distributions ?
- impact of model size?
- impact of vocabulary size?
  
- how to achieve **fairness** in mLLMs design?

# mLLMs are still a wonder

---

## Open question and issues

mLLMs: more than a collection of monolingual models ?

- which linguistic properties help / break transfer?
- how critical is **alignment? parallel data?**
- how to measure positive / negative **interference?**
- how to measure language **coverage?**
- how to measure **non-linguistic biases** ?
  - racial, social, cultural, etc

How about the “**curse of multilinguality**” ?

- impact of language distributions ?
- impact of model size?
- impact of vocabulary size?
  
- how to achieve **fairness** in mLLMs design?

How about actual **multilingual tasks**?

- machine translation
- generating code-switched language
- summarization from multilingual texts

What happens with **multi-step training**?

- Effect of **multilingual finetuning**
- Effect of **multilingual instruction tuning**
- Effect of **multilingual alignment**

# mLLMs are still a wonder

---

## Open question and issues

mLLMs: more than a collection of monolingual models ?

- which linguistic properties help / break transfer?
- how critical is **alignment? parallel data?**
- how to measure positive / negative **interference?**
- how to measure language **coverage?**
- how to measure **non-linguistic biases** ?
  - racial, social, cultural, etc

How about the “**curse of multilinguality**” ?

- impact of language distributions ?
- impact of model size?
- impact of vocabulary size?
- how to achieve **fairness** in mLLMs design?

How about actual **multilingual tasks**?

- machine translation
- generating code-switched language
- summarization from multilingual texts

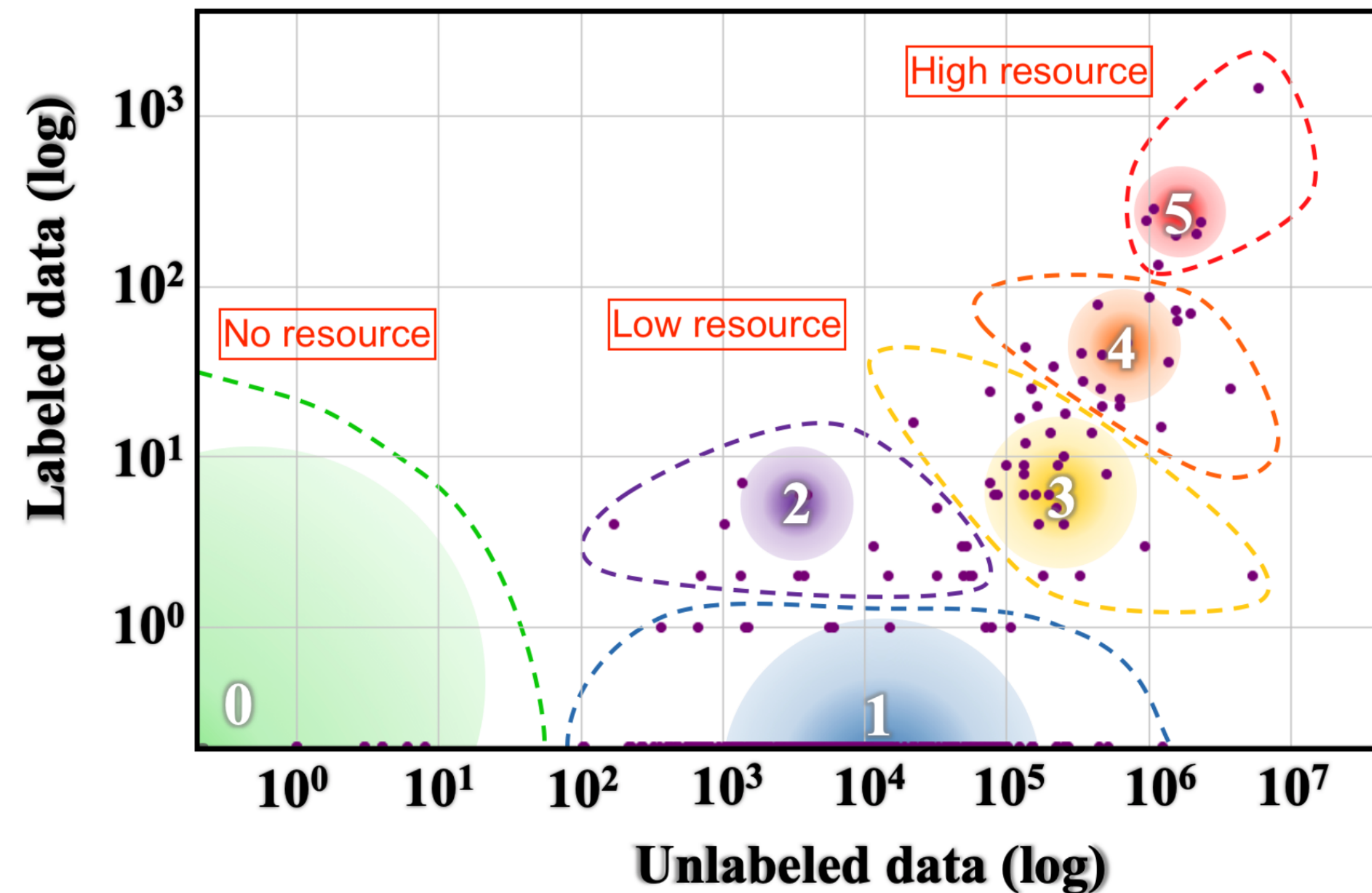
What happens with **multi-step training**?

- Effect of **multilingual finetuning**
- Effect of **multilingual instruction tuning**
- Effect of **multilingual alignment**

## What we need to do

- more tools (lid, cs id, variety id, etc)
- massively multilingual benchmarks
- better models of heterogeneous data
- more diverse samples of language use
- **rethink data collection, evaluation and scores**

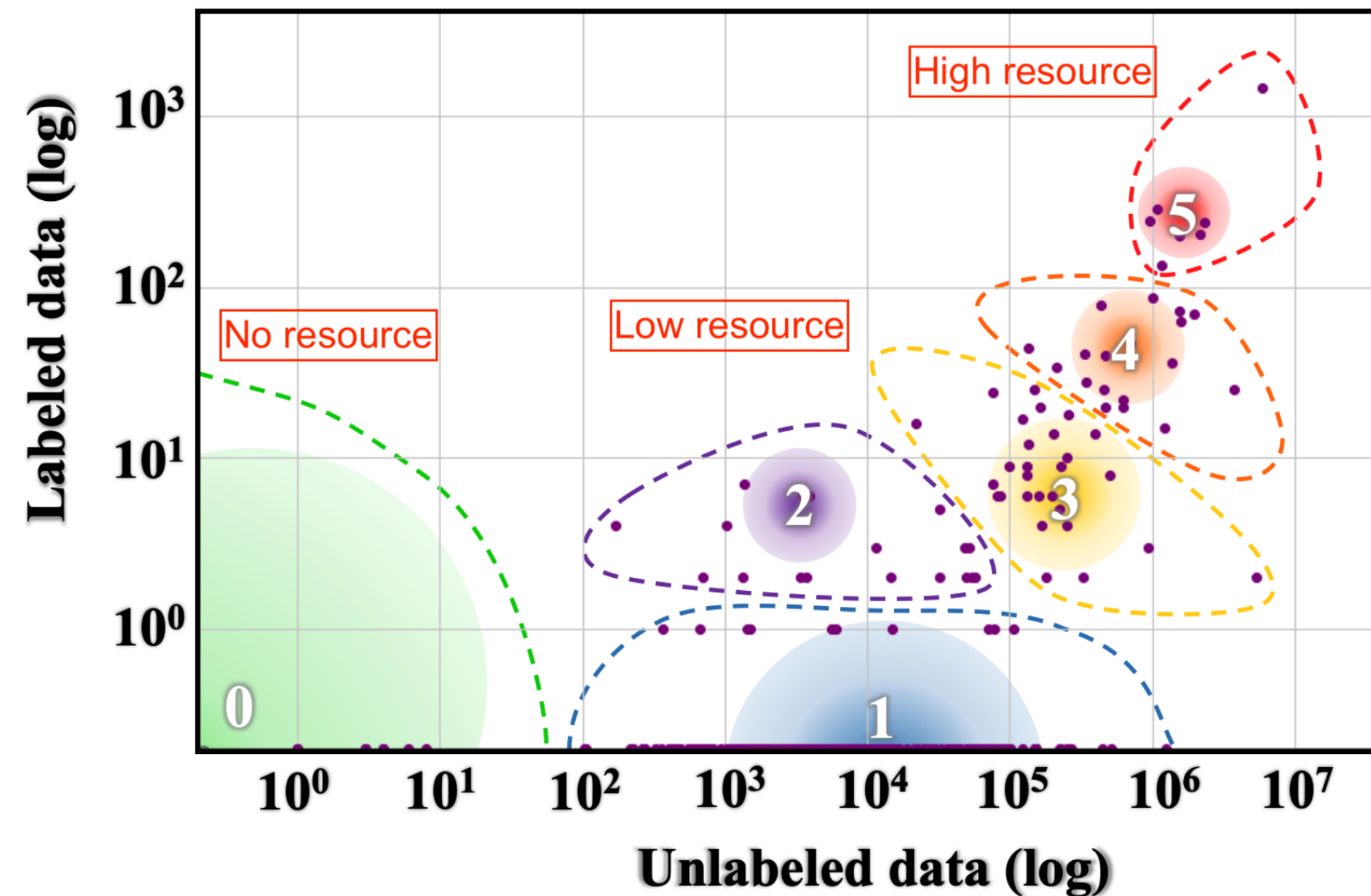
# A final word of caution



Towards “the next 1000 languages”, really ?

Class	5 Example Languages	#Langs	#Speakers	% of Total Langs
0	Dahalo, Warlpiri, Popoloca, Wallisian, Bora	2191	1.2B	88.38%
1	Cherokee, Fijian, Greenlandic, Bhojpuri, Navajo	222	30M	5.49%
2	Zulu, Konkani, Lao, Maltese, Irish	19	5.7M	0.36%
3	Indonesian, Ukranian, Cebuano, Afrikaans, Hebrew	28	1.8B	4.42%
4	Russian, Hungarian, Vietnamese, Dutch, Korean	18	2.2B	1.07%
5	English, Spanish, German, Japanese, French	7	2.5B	0.28%

# A final word of caution



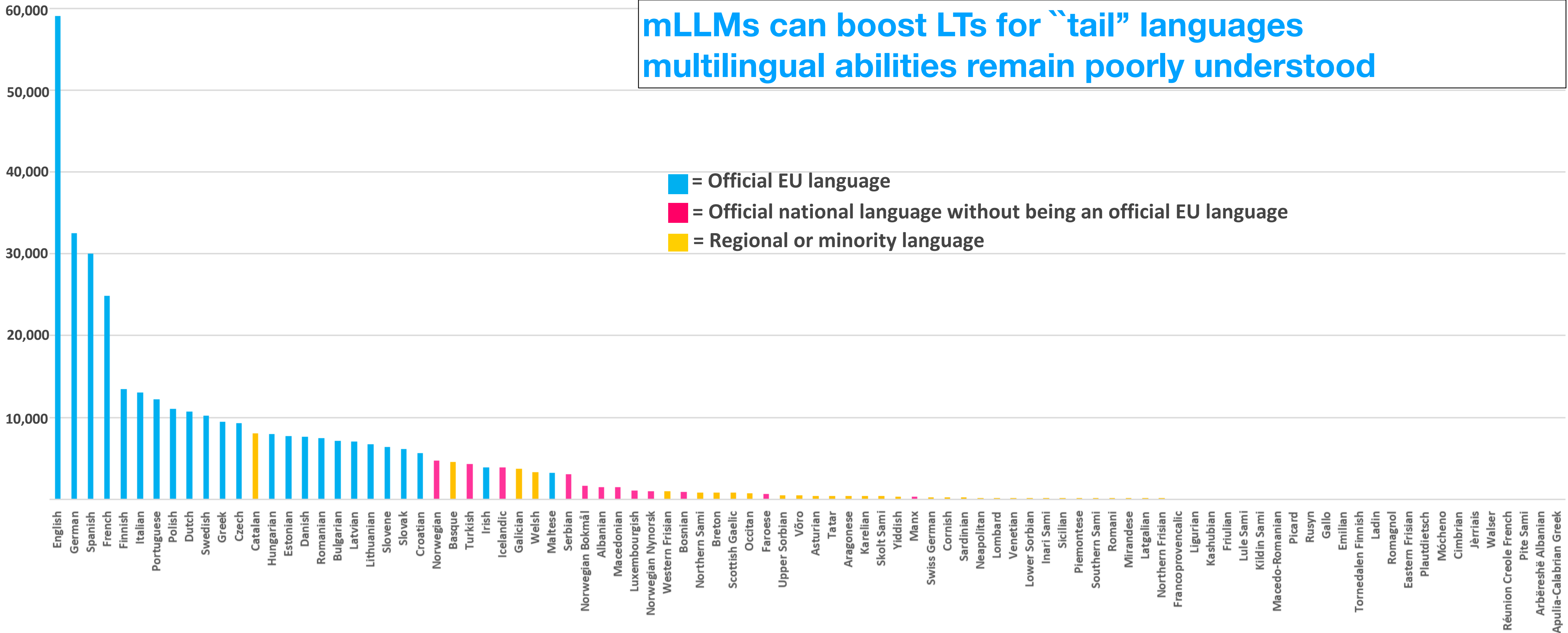
Towards “the next 1000 languages”, really ?

Class	5 Example Languages	#Langs	#Speakers	% of Total Langs
0	Dahalo, Warlpiri, Popoloca, Wallisian, Bora	2191	1.2B	88.38%
1	Cherokee, Fijian, Greenlandic, Bhojpuri, Navajo	222	30M	5.49%
2	Zulu, Konkani, Lao, Maltese, Irish	19	5.7M	0.36%
3	Indonesian, Ukranian, Cebuano, Afrikaans, Hebrew	28	1.8B	4.42%
4	Russian, Hungarian, Vietnamese, Dutch, Korean	18	2.2B	1.07%
5	English, Spanish, German, Japanese, French	7	2.5B	0.28%



# Conclusions & Take aways

mLLMs can boost LTs for “tail” languages  
 multilingual abilities remain poorly understood



# Special thanks to the LMU/CIS team

Ayyoob ImaniGooghari  
Masoud Jalili Sabet  
Amir Hossein Kargaran  
Nora Kassner  
Lütfi Kerem Şenel  
Peiqin Lin



Chunlan Ma  
André Martins  
Silvia Severini  
Helmut Schmid  
Hinrich Schütze

and to the BigScience “Evaluation” WG