# Identifying Open Challenges in Language Identification

Rob van der Goot

# Language identification

▶ Task: Map text into language (ISO 639-3) codes
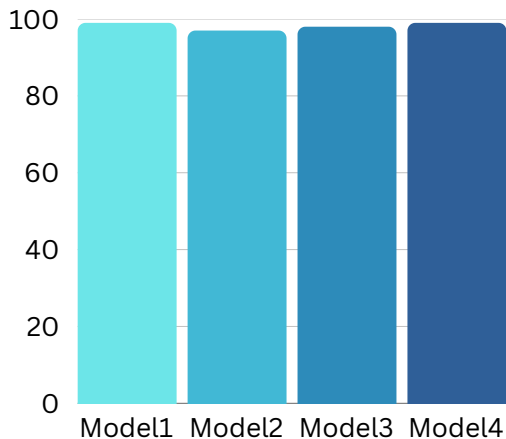
```
wêr jammerje jo no oer?
Skal vi snakke lidt?
```

# Language identification

▶ Task: Map text into language (ISO 639-3) codes

```
Wêr jammerje jo no oer?      FRY
Skal vi snakke lidt?         DAN
```

# Current state

# Data

| Dataset | langs | scripts | fams | domains |
|---|---|---|---|---|
| MIL-TALE | 2,110 | 47 | 139 | wiki, political, religious, grammar |
| UDHR | 397 | 38 | 61 | rights |
| OpenLID | 139 | 25 | 16 | literature, news, wiki, social, grammar, subtitles, spoken |
| MassiveSumm | 77 | 24 | 13 | news |
| TwitUser | 59 | 20 | 13 | social |
| UD | 54 | 11 | 17 | medical, news, academic, wiki, legal, nonfiction, learner-essays, fiction, social, grammar-examples, reviews, religious, spoken |
| Total | 2,176/ 7,850 | 51/ 163 | 145/ 298 | |

(note that glotlid (`https://huggingface.co/cis-lmu/glotlid`) has a similar setup/scope)

# Models

| Type | model | size | |
|------|-------|-----:|---|
| N-gram overlaps | Textcat | 40,000 | 🐜 |
| ML | NB | 100,000 | 🐸 |
| Embeddings | FastText | 4,434,860 | 🐧 |
| NN | LSTM | 15,158,772 | 🐊 |
| Transformer LM | GLOT500 | 395,687,155 | 🐳 |

# Size

- 10, 100, 1,000 languages
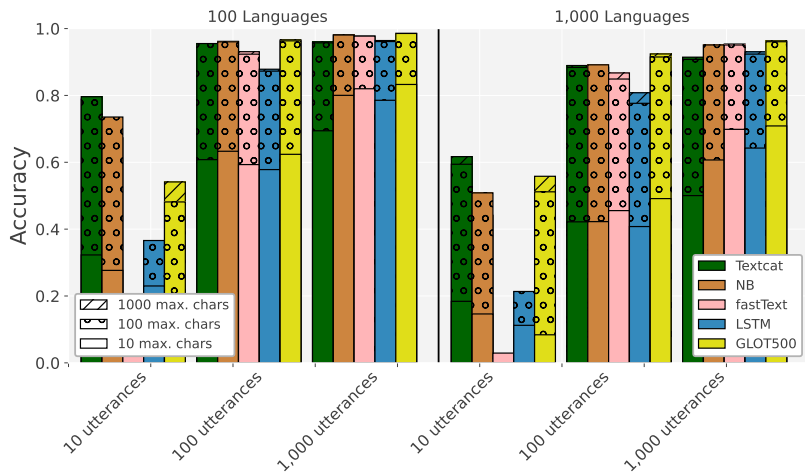- 10, 100, 1,000 utterances input
- 10, 100, 1,000 characters input

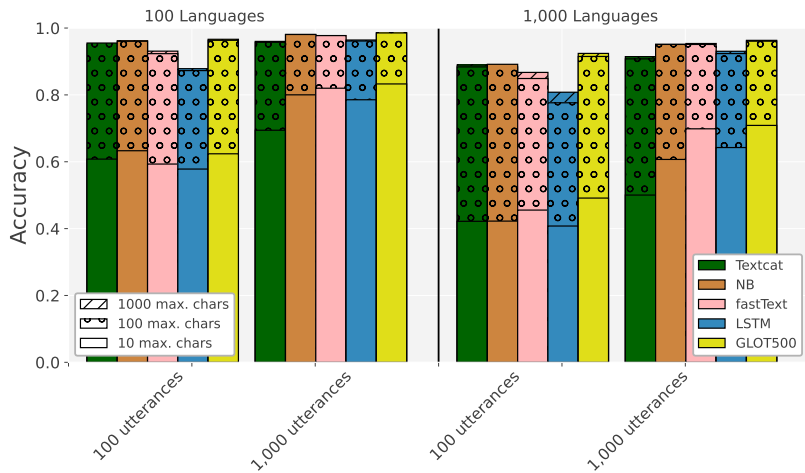1,000 utterances for testing

# Size

▶ Number of languages

# Size

▶ Number of utterances
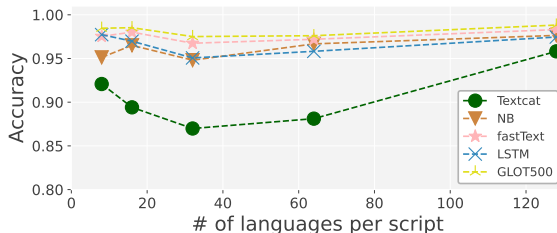
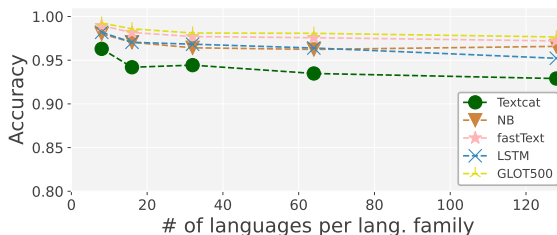# Size

- Number of characters
- Models

# Family/script

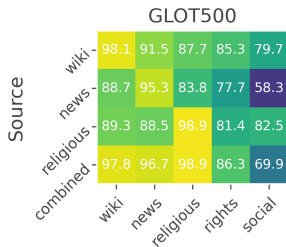- Make subsets with N languages per family/script

# Family/script

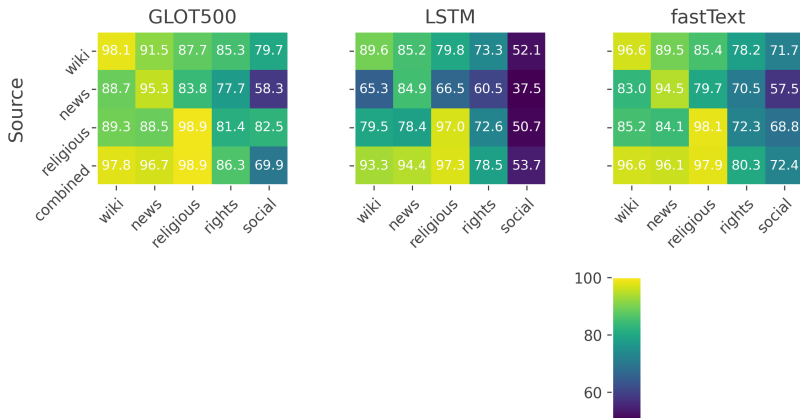▶ Make subsets with N languages per family/script

# Domains

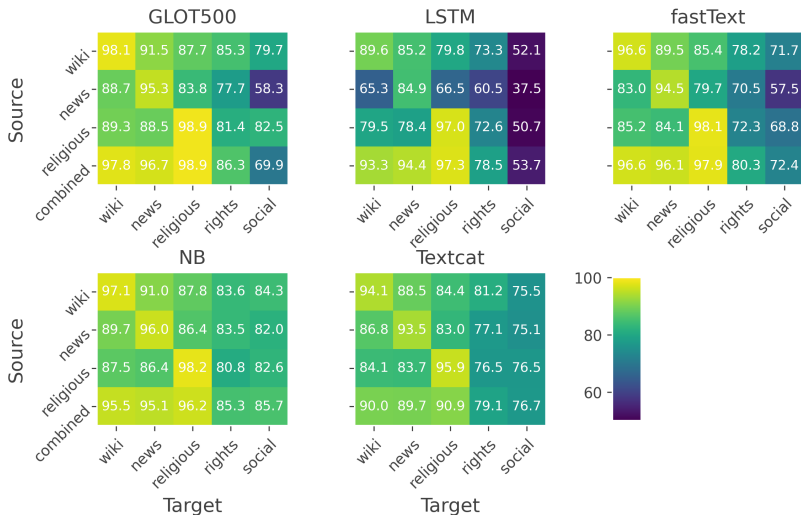- ▶ Evaluate across domains

# Domains

▶ Evaluate across domains

# Domains

▶ Evaluate across domains

# Conclusions

- ▶ Small number of utterances (1,000) is enough
- ▶ Family/script small effect
- ▶ Cross domain still challenging
- ▶ Larger models better in-dataset, small models more robust!
- ▶ Final models with $> 2,000$ languages released

# Conclusions

- ▶ Small number of utterances (1,000) is enough
- ▶ Family/script small effect
- ▶ Cross domain still challenging
- ▶ Larger models better in-dataset, small models more robust!
- ▶ Final models with $> 2,000$ languages released