# BYTE IS DEAD! LONG LIVE MYTE!

Tomasz Limisiewicz, Terra Blevins, Hila Gonen, Orevaoghene Ahia, Luke Zettlemoyer

## Morphological Bytes

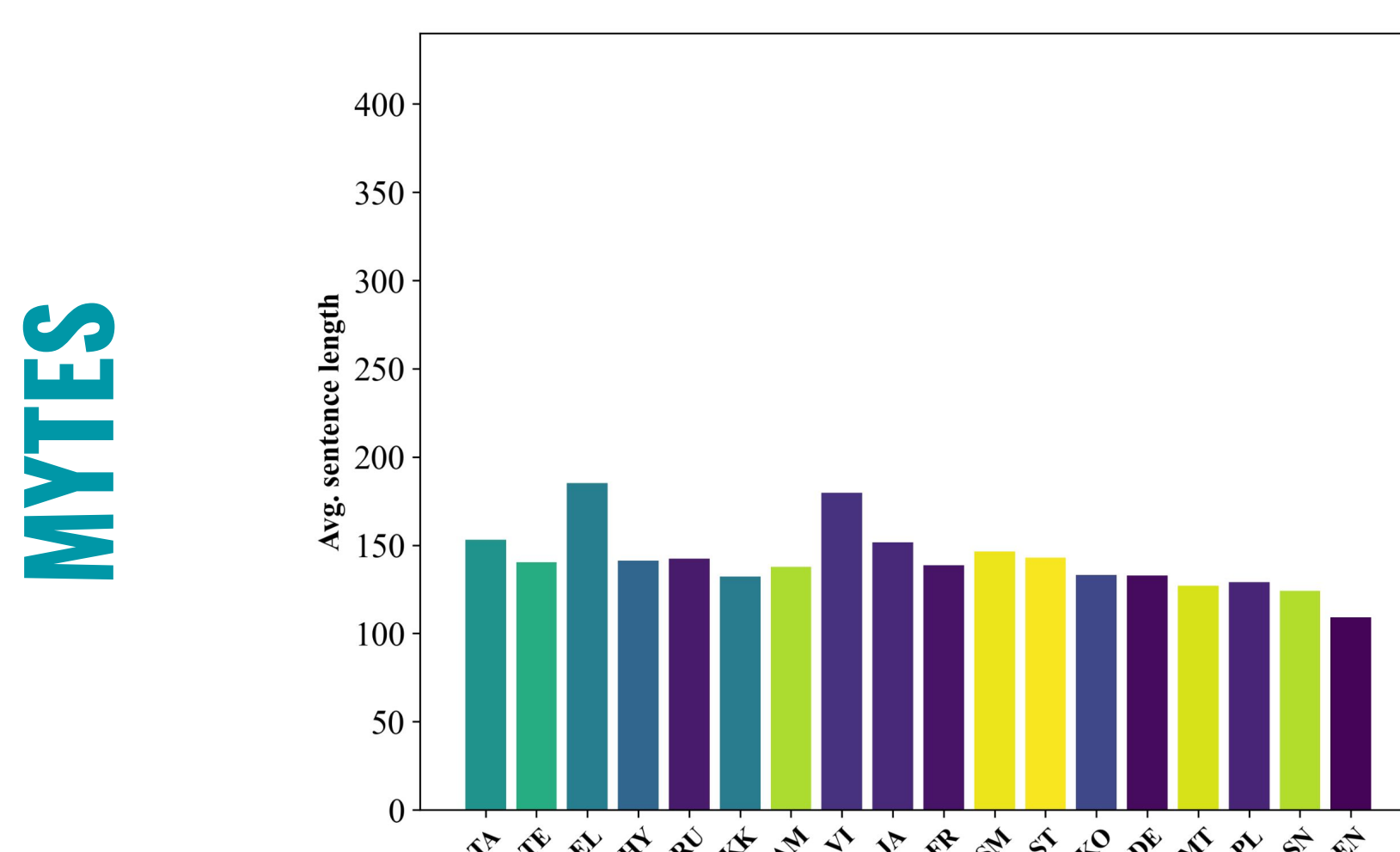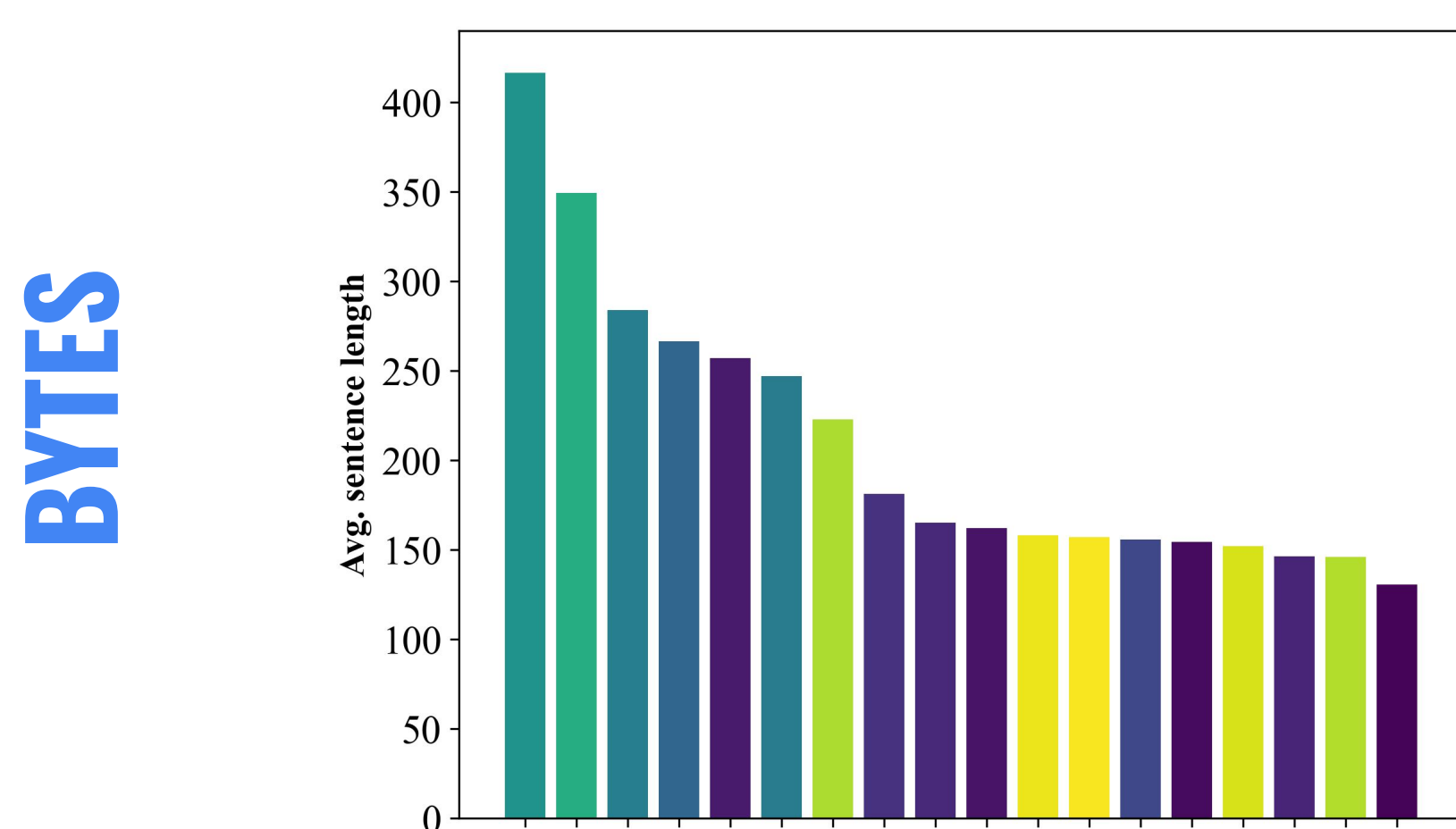We use morphological segmenter to redefine byte codespace.



## More Efficient

MYTE sequences are from 2 to 70% shorter than BYTEs ↘

## Fairer across Languages

Average length of parallel sentences FLORES 200 encoded in:

**BYTES**



**MYTES**



## MyT5 Models

We train T5 models in three sizes: small, base, and large and compare their performance with similar byte-level: ByT5 ➡

```
from transformers import T5ForConditionalGeneration
from transformers import MyT5Tokenizer

MODEL_SIZE = "large" # small, base, or large
MODEL = f"Tomlim/myt5_{MODEL_SIZE}"

model = T5ForConditionalGeneration.from_pretrained(
    MODEL, use_safetensors=True)

tokenizer = MyT5Tokenizer.from_pretrained(MODEL)
```

## Problem: Byte encoding long and suboptimal for many languages

We propose **M**orphological B**YTES** to improve over **UTF-8** encoding in:

- **Fairness**: comparable sequence length for the same information

- **Efficiency**: shorter sequence length



- **MYTE outperforms UTF-8 Bytes, efficiently representing texts in diverse languages, especially in non-Latin scripts.**
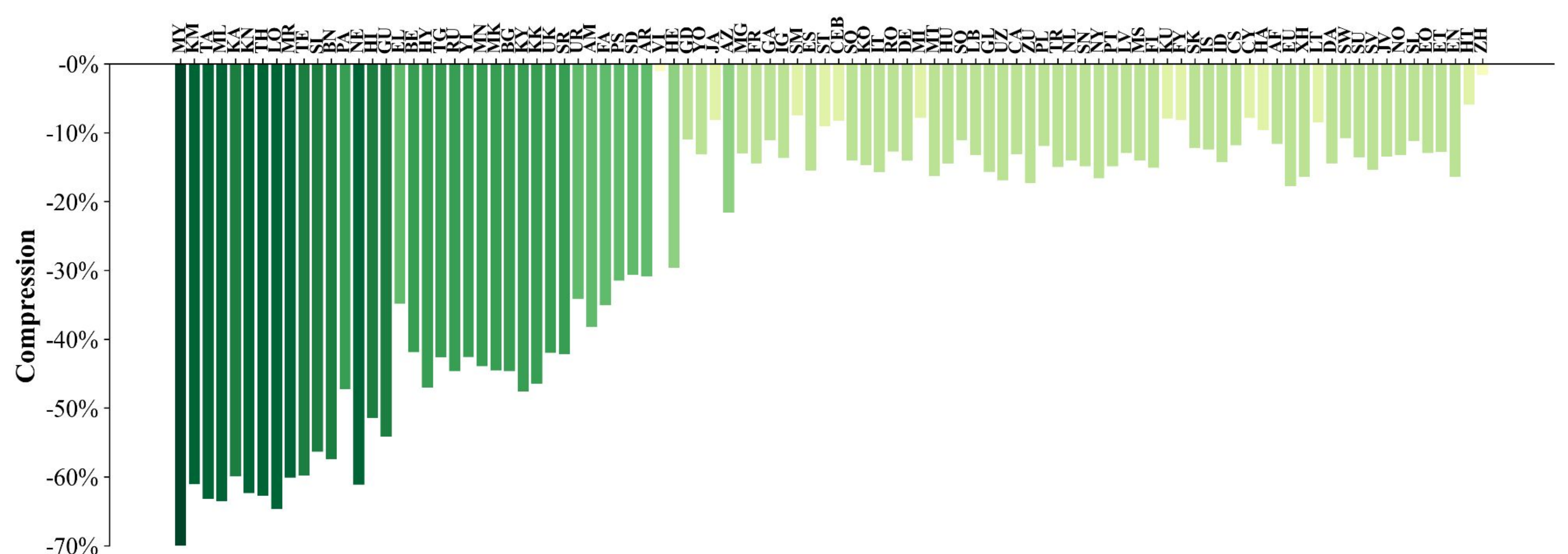
- **All 104 tested languages are encoded in less MYTEs than BYTEs**

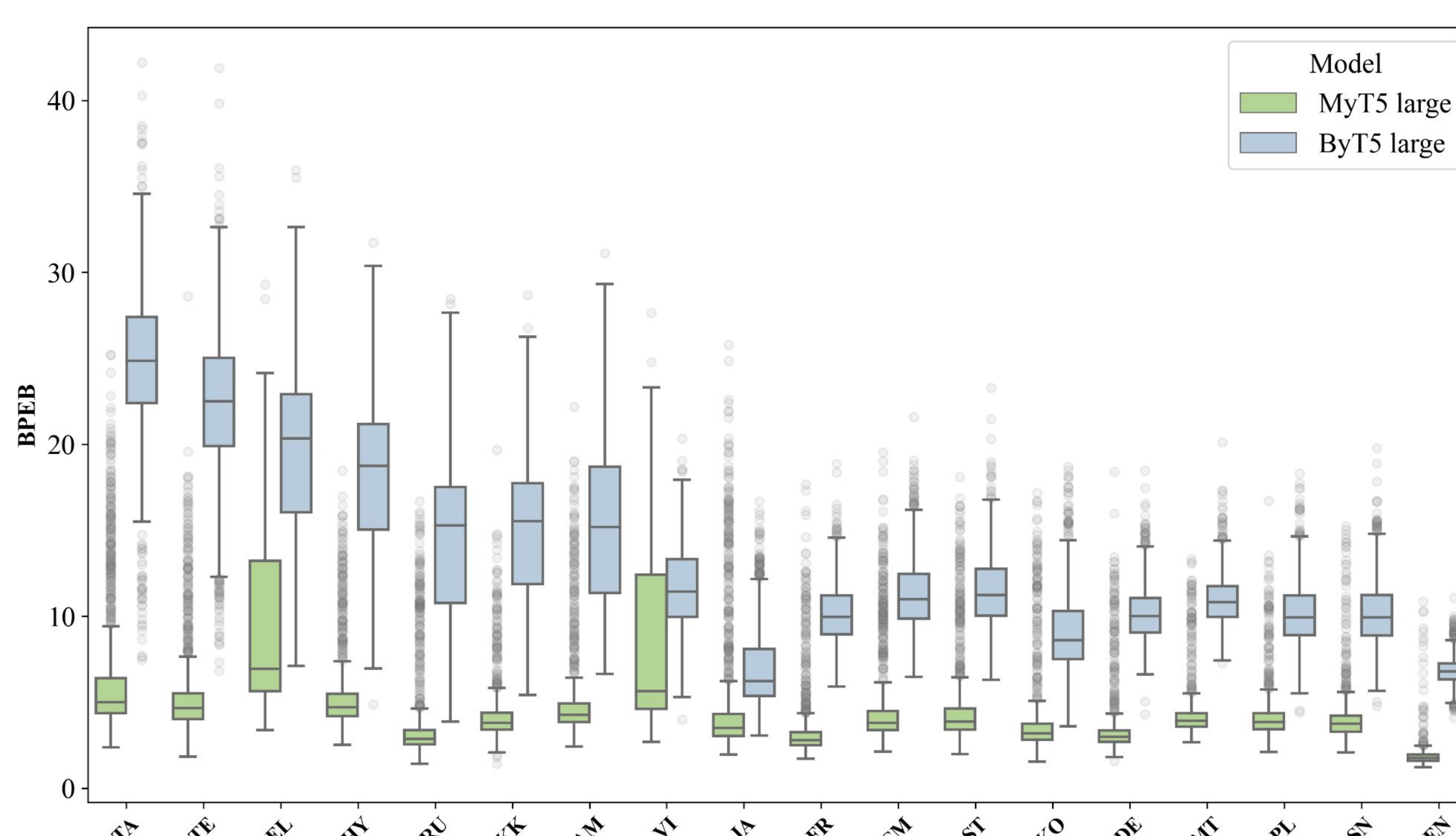- **MYTE speeds up and improves multilingual language modeling**

- **MYTE is more efficient for end-tasks with comparable results**

### Compression: Length in Mytes vs. UTF-8 Bytes



### LM Results: Perplexity of Parallel Sentences



### End-tasks

|  | ByT5 | | MyT5 | |
|---|---|---|---|---|
|  | score | time | score | time |
| QA | 73.2 | 36.2 | **75.3** | **35.6** |
| NER | **81.5** | 13.8 | 80.8 | **12.6** |
| SemP | **25.1** | 13.2 | 19.6 | **12.4** |
| MT | 20.1 | 15.9 | **20.4** | **12.6** |