# Language Modelling with Pixels

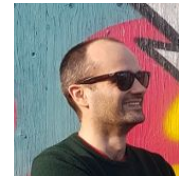Phillip Rust[1]    Jonas F. Lotz[1]    Emanuele Bugliarello[1]    Elizabeth Salesky[3]    Miryam de Lhoneux[1,2,4]    Desmond Elliott[1]

[1]University of Copenhagen    [2]KU Leuven    [3]Johns Hopkins University    [4]Uppsala University

# Summary

We train a **pixel-based encoder of language (PIXEL)**,
a **language model** trained solely on **images** of **rendered text.**

## Some of PIXEL's strengths are

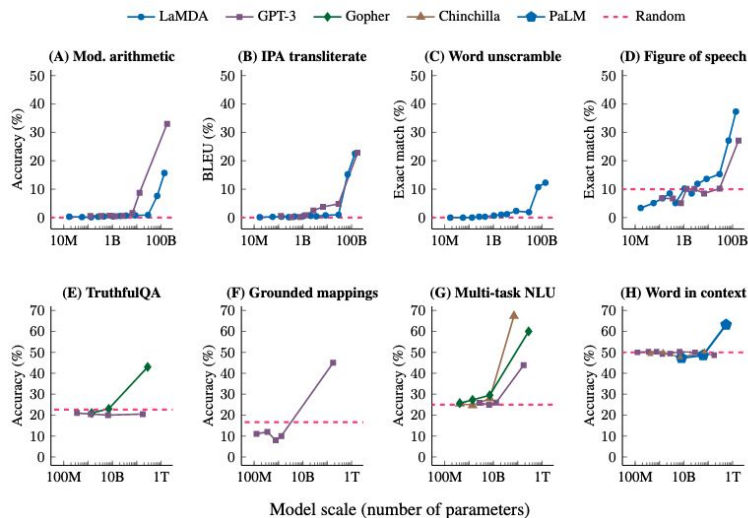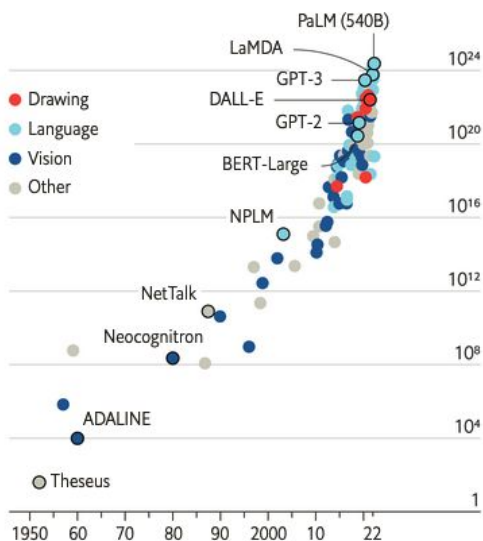Out-of-the-box **transfer** to **unseen languages** and **scripts**

**Robustness** to orthographic attacks & code-switching**\***

**\*** See our paper for the code-switching results

# NLP in the Era of Scale



## The blessings of scale
AI training runs, estimated computing resources used
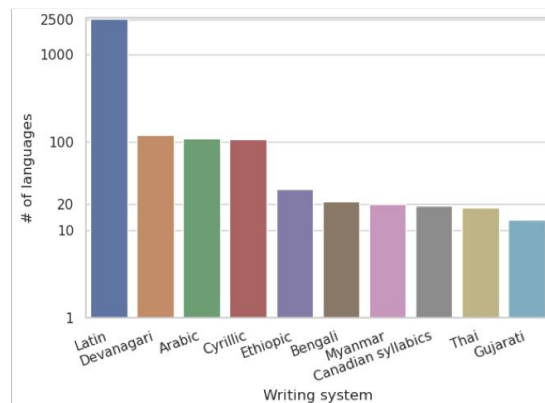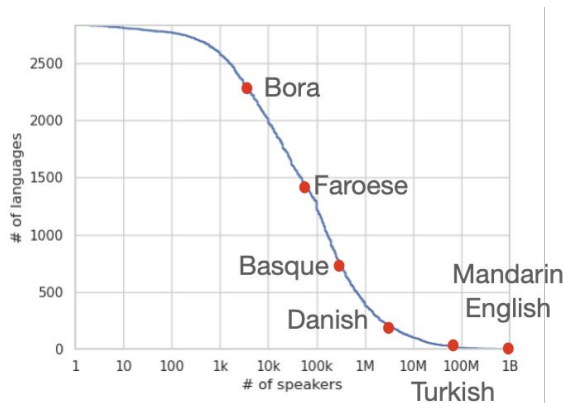Floating-point operations, selected systems, by type, log scale

Emergent Abilities of Large Language Models
*(Wei+ TMLR'22)*

# NLP for all written languages?

There are **~7000 spoken languages**, of which **~3000** are **written** and at least **400** have **>1M speakers**

**Most NLP** only covers **100 languages** *(van Esch+ LREC'22)*
→ **Lack** of **technological inclusion** for **billions of people**
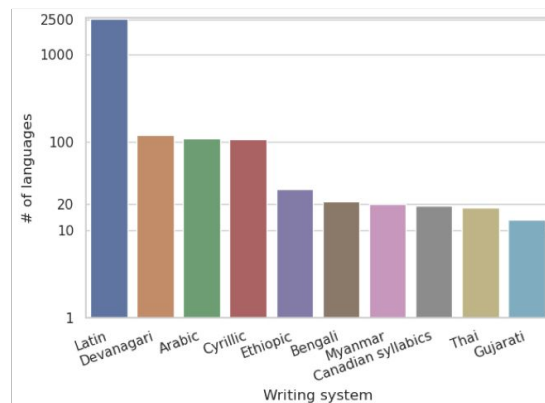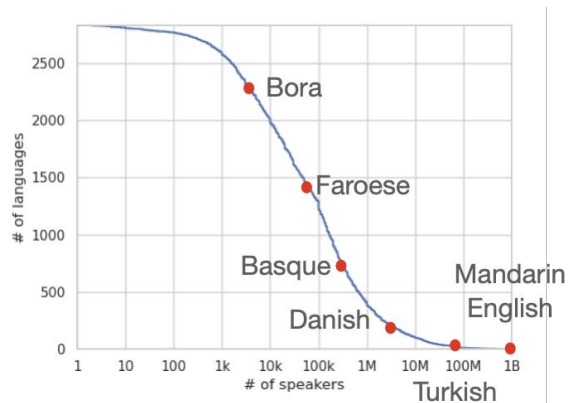


*Slide credit: Sebastian Ruder*

# What's left? NLP for all written languages

There are **~7000 spoken languages**, of which **~3000** are **written** and at least **400** have **>1M speakers**

**Most NLP** only covers **100 languages** *(van Esch+ LREC'22)*
→ **Lack** of **technological inclusion** for **billions of people**



*Slide credit: Sebastian Ruder*

# Question: What's stopping us?

NLP is an **open vocabulary problem**.

A language model's ability to **process unseen words** is **determined by its vocabulary**:

**1. "Trained" over a corpus:** Byte-Pair Encoding *(Sennrich+ ACL'16)*
→ Unseen tokens not in the vocabulary (unless w/ byte-level fallback)

**2. Corpus-independent**: characters *(Clark+ TACL'22)* / bytes *(Xue+ ACL'22)*
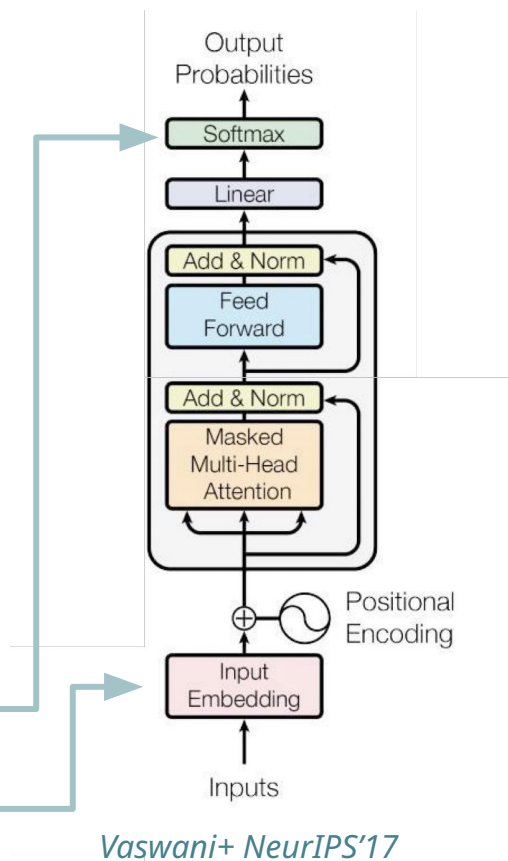→ Need to deal with longer sequence lengths

# Answer: The *Vocabulary Bottleneck*

**Language models** have **discrete** input and output **vocabularies** expressed over a **finite inventory** of tokens, characters, words, sub-words, etc.

→ **This creates a bottleneck in two places**

*Computational bottleneck* in the output layer

*Representational bottleneck* in the embedding layer

*Vaswani+ NeurIPS'17*

# TL;DR of our paper
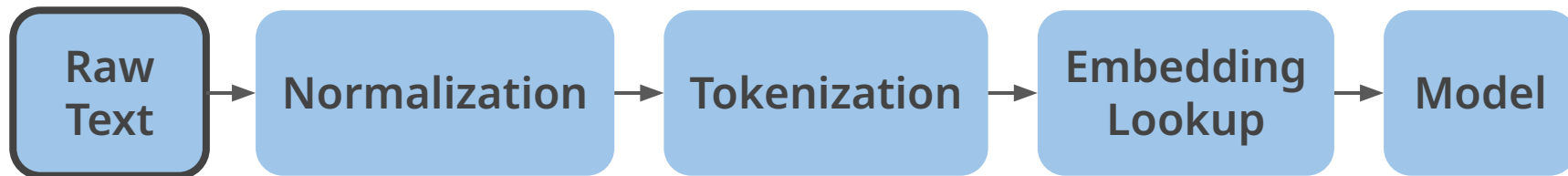
We attempt to crack the *vocabulary bottleneck* with pixels.
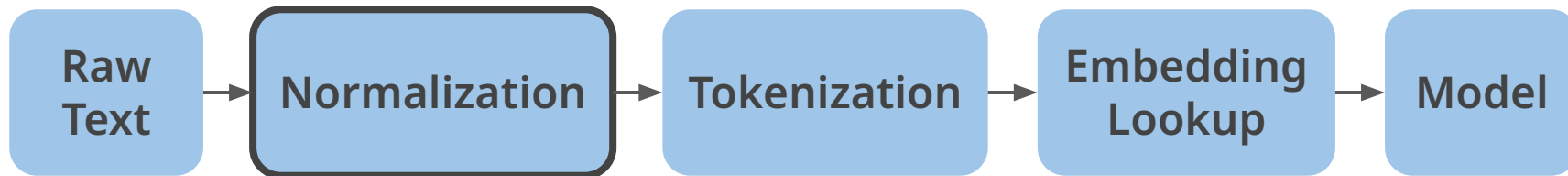
But what does that mean?

# The NLP pipeline simplified
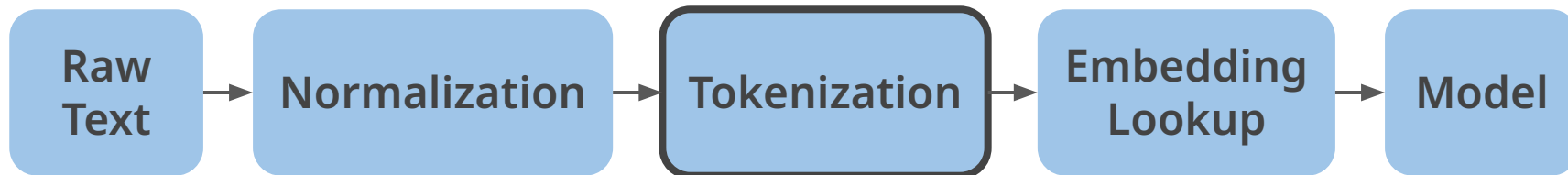
# The NLP pipeline simplified

Raw Text → Normalization → Tokenization → Embedding Lookup → Model

*My cat, Dr. Beans II., sleeps 22h a day.*

# The NLP pipeline simplified

Raw Text → Normalization → Tokenization → Embedding Lookup → Model

*My cat , Dr . Beans II . , sleeps 22h a day .*

# The NLP pipeline simplified

Raw Text → Normalization → Tokenization → Embedding Lookup → Model

$[CLS]^{101}$ $My^{1422}$ $cat^{5855}$ $,^{117}$ $Dr^{1987}$ $.^{119}$ $Bean^{21561}$ $\#\#s^{1116}$ $II^{1563}$ $.^{119}$ $,^{117}$
$sleep^{2946}$ $\#\#s^{1116}$ $22^{1659}$ $\#\#h^{1324}$ $a^{170}$ $day^{1285}$ $.^{119}$ $[SEP]^{102}$

12

# The NLP pipeline simplified

| Raw Text | → | Normalization | → | Tokenization | → | Embedding Lookup | → | Model |
|---|---|---|---|---|---|---|---|---|

| | | | |
|---|---|---|---|---|
| 4.4960e-01<br>9.7664e-02<br>-2.0737e-01<br>...<br>-9.5125e-02 | -8.0239e-02<br>9.0034e-01<br>3.5708e-01<br>...<br>1.1732e-01 | ... | -5.0748e-01<br>6.4085e-01<br>9.4414e-01<br>...<br>1.2228e-01 | -8.9100e-02<br>-3.8999e-01<br>4.0695e-01<br>...<br>6.2151e-01 |

# The NLP pipeline simplified

| Raw Text | → | Normalization | → | Tokenization | → | Embedding Lookup | → | Model |

| 6.8479e-01<br>7.4406e-02<br>-5.7043e-02<br>...<br>-5.6763e-02 | 5.5677e-01<br>-3.8491e-01<br>5.5702e-01<br>...<br>2.2008e-01 | ... | 5.0547e-01<br>9.9931e-02<br>1.7320e-01<br>...<br>-1.8281e-01 | 1.7858e+00<br>2.3194e-01<br>1.4286e-01<br>...<br>-1.3138e-01 |

# Cracking the *vocabulary bottleneck* with pixels
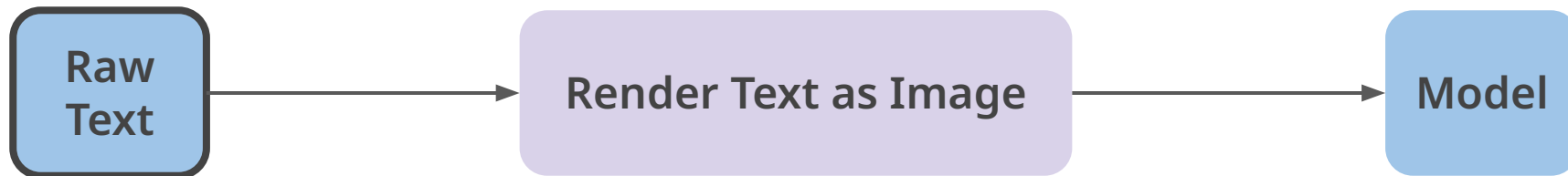
Treat **language processing** as **visual processing**

Raw Text → Normalization → Tokenization → Embedding Lookup → Model

# Cracking the *vocabulary bottleneck* with pixels

Treat **language processing** as **visual processing**

Raw Text → Render Text as Image → Model

# Cracking the *vocabulary bottleneck* with pixels

Treat **language processing** as **visual processing**

Raw Text → Render Text as Image → Model

*My cat, Dr. Beans II., sleeps 22h a day.*

# Cracking the *vocabulary bottleneck* with pixels

Treat **language processing** as **visual processing**

**Raw Text** → **Render Text as Image** → **Model**

My cat, Dr. Beans II., sleeps 22h a day. ▮

# Cracking the *vocabulary bottleneck* with pixels

Treat **language processing** as **visual processing**

**Raw Text** → **Render Text as Image** → **Model**

| -4.1020e-01<br>2.7750e-01<br>3.3202e-01<br>...<br>2.9194e-01 | -3.6107e-01<br>2.0695e-01<br>1.7878e+00<br>...<br>9.4824e-02 | ... | 3.2538e-01<br>1.2356e+00<br>-1.0839e+00<br>...<br>6.8341e-01 | 1.7513e-01<br>1.1834e+00<br>-4.8054e-01<br>...<br>6.8465e-01 |

# Inspiration



**Robust Open-Vocabulary Translation from Visual Text Representations**
*(Salesky+ EMNLP'21)*



**Masked Autoencoders are Scalable Visual Learners**
*(He+ CVPR'22)*

# Pixel-based Encoder of Language (PIXEL)



$$\text{MSE} = \frac{1}{m}\frac{1}{n}\sum_{i=1}^{m}\sum_{j=1}^{n}(Y_j^i - \hat{Y}_j^i)^2$$

Decoder — 8 Layers

Encoder — 12 Layers

CLS

③ CLS Embedding & Span Mask $m$ patches

② Projection + Position Embedding

My cat ᒼᑫᘮ enjoys eating warm oatmeal for lunch and dinner.

16x16 patch resolution

Google Noto Fonts

PyGame / PangoCairo

① Render Text

My cat ᒼᑫᘮ enjoys eating warm oatmeal for lunch and dinner.

# PIXEL learns to reconstruct text



**100k steps**

# Downstream Task Fine-Tuning

# Flexible Text Renderer

My cat 🐈 loves pancakes 🥞 and grapes 🍇.

一隻貓正在吃碗中的猫粮   قط جاثم على غصن شجرة

አቢሲኒያውያን በጣም ጠንካራ ድመቶች ናቸው ፡፡

24

# The Benefits of Pixels

PIXEL can **process anything that can be rendered**
→ **Open vocabulary** which is easily extensible to **unseen text**
→ Support all written languages

**Complete parameter sharing** from the input representation
(unlike separate-but-related subwords in an embedding matrix)

**Nothing language-specific** in the **input / output**
→ **Greater flexibility** to process written language in **different forms**
(PDFs, scanned newspapers, etc.)

# **Experiments**

# Pretraining

| | |
|---|---|
| **Dataset** | |
| **Masking** | |
| **Max. Seq. Length** | |
| **Compute** | |
| **Parameters** | |

There is only **0.05% non-English** text in our **pretraining data** (estimated by Blevins and Zettlemoyer 2022)

The **Great Wall of China** (traditional Chinese: 萬里長城; simplified Chinese: 万里长城; pinyin: *Wànlǐ Chángchéng*)
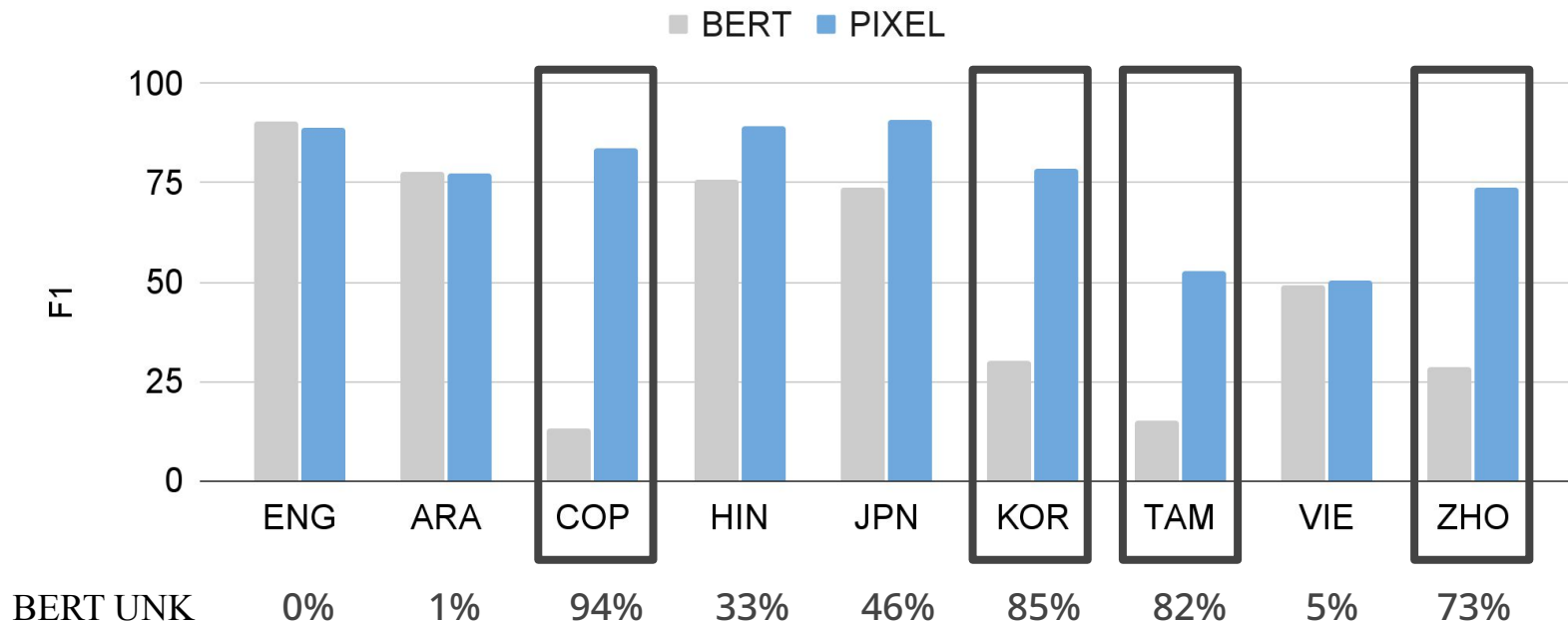
# Finetuning Experiments

| Datasets | |
|----------|----------|
| | . |

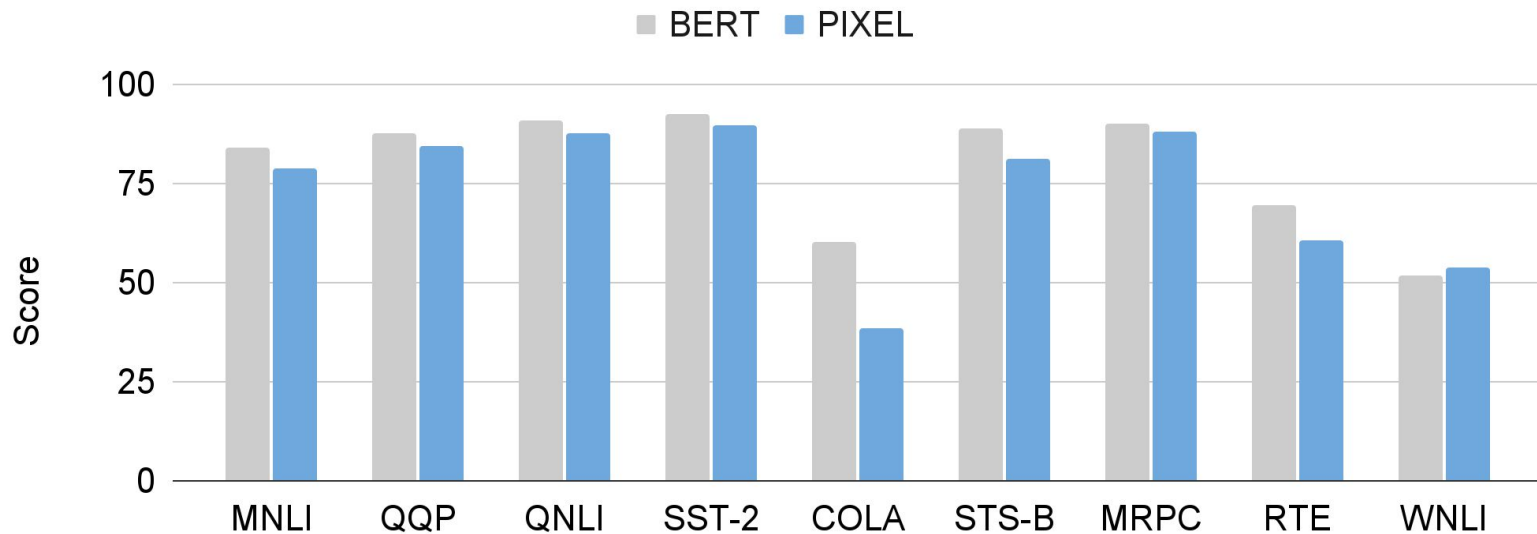# Finetuning Experiments

| Datasets | |
|---|---|
| | . |

# Dependency Parsing Results



PIXEL (vastly) outperforms BERT on unseen scripts

# GLUE Results
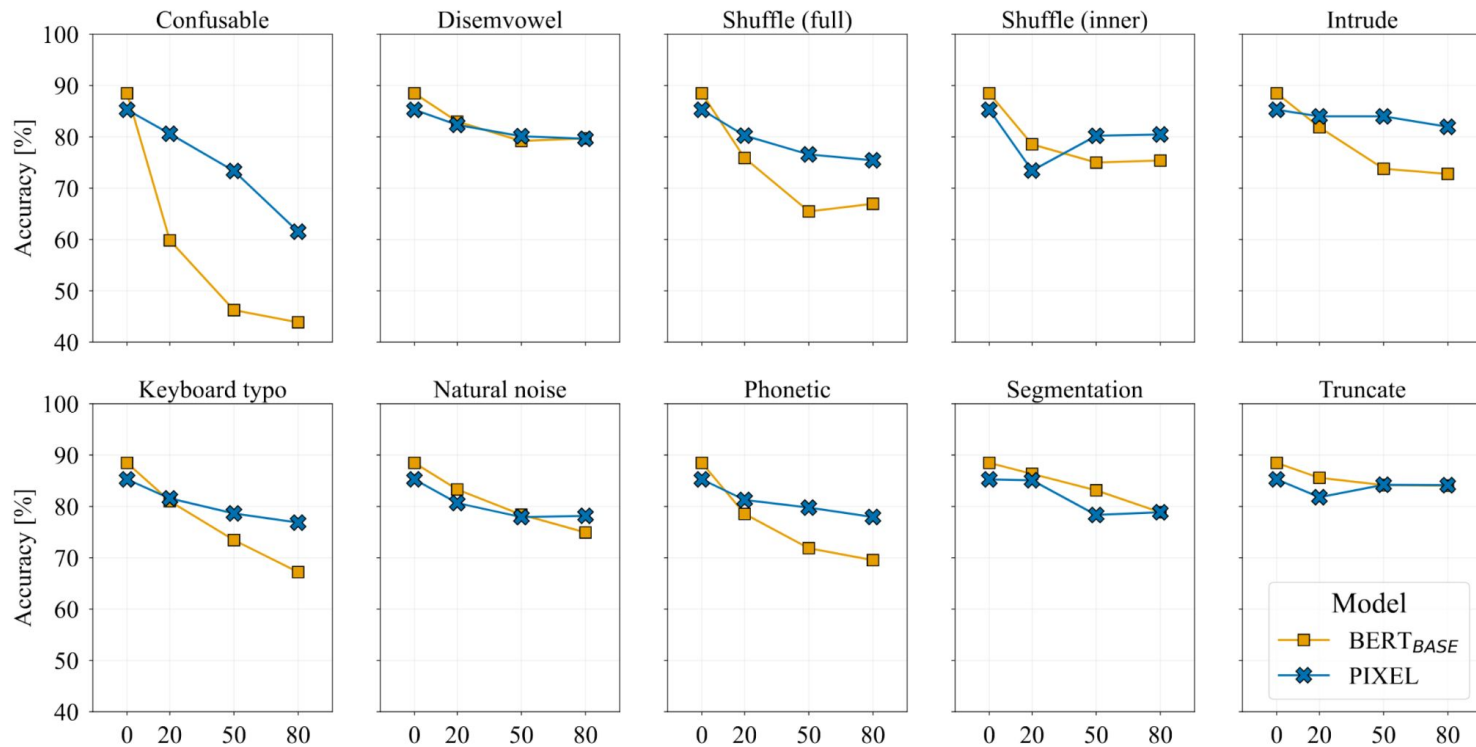


BERT outperforms PIXEL on English sentence-level tasks

# Robustness against orthographic attacks (Zeroé)

| Attack | Sentence |
|--------|----------|
| NONE | Penguins are designed to be streamlined |

# Robustness against orthographic attacks (Zeroé)

| Attack | Sentence |
|--------|----------|
| NONE | Penguins are designed to be streamlined |

# PIXEL is more robust than BERT

# Conclusions

**PIXEL** is a **new type of language model** that renders **text as images** instead of splitting text into a finite set of tokens.

**Rendered text** makes it possible to achieve **high-quality transfer** to **unseen scripts** in syntactic and semantic tasks.

Pixel-based learning could be a **promising research direction** to **make NLP technology accessible** to more people.

# PIXEL Resources

https://github.com/xplip/pixel

https://huggingface.co/Team-PIXEL