

Findings-#2836



# Adaptive BPE Tokenization for Enhanced Vocabulary Adaptation in Finetuning Pretrained Language Models

**Gunjan Balde**\*, Soumyadeep Roy\*, Mainack Mondal, and Niloy Ganguly

Indian Institute of Technology Kharagpur

\*Equal Contribution

[balde.gunjan0812@kgpian.iitkgp.ac.in](mailto:balde.gunjan0812@kgpian.iitkgp.ac.in)

**EMNLP  
2024** 

# Byte-Pair Encoding (BPE<sup>[1]</sup>)

- Popular tokenization algorithm used in BART<sup>[2]</sup>, LLama<sup>[3]</sup>, and Mistral<sup>[4]</sup>
- Utilizes concept of ranked-based merge rules to tokenize a word



[1] Neural Machine Translation of Rare Words with Subword Units (Sennrich et al., ACL 2016)

[2] BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension (Lewis et al., ACL 2020)

[3] Llama 2: Open foundation and fine-tuned chat models. (Hugo, et al., 2023)

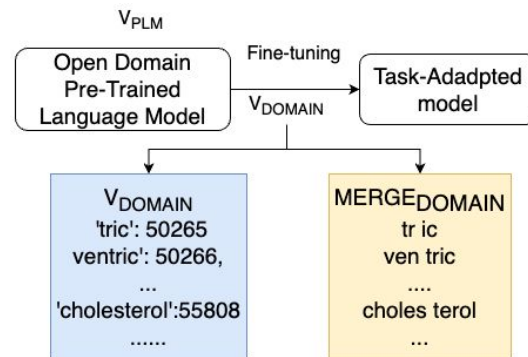
[4] Mistral 7B (Jiang et al., 2023)

# Vocabulary Adaptation for PLM using BPE

- Strategy to adapt PLM to a domain *during fine tuning*
- Works best when adapting to *expert domain*
- Build a domain-specific vocabulary ( $V_{DOMAIN}$ )
- And learn corresponding merge rules ( $MERGE_{DOMAIN}$ )
- Appends it at **the end** of existing PLM Vocabulary and Merges

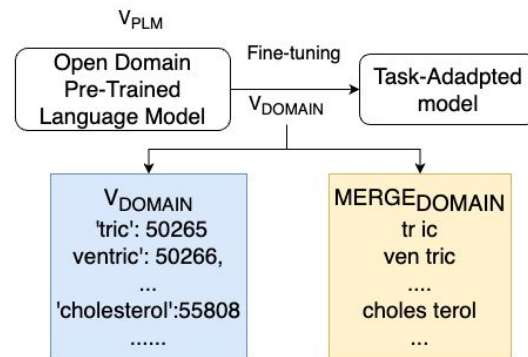
# Vocabulary Adaptation for PLM using BPE

- Strategy to adapt PLM to a domain *during fine tuning*
- Works best when adapting to *expert domain*
- Build a domain-specific vocabulary ( $V_{DOMAIN}$ )
- And learn corresponding merge rules ( $MERGE_{DOMAIN}$ )
- Appends it at **the end** of existing PLM Vocabulary and Merges



# Vocabulary Adaptation for PLM using BPE

- Strategy to adapt PLM to a domain *during fine tuning*
- Works best when adapting to *expert domain*
- Build a domain-specific vocabulary ( $V_{DOMAIN}$ )
- And learn corresponding merge rules ( $MERGE_{DOMAIN}$ )
- Appends it at **the end** of existing PLM Vocabulary and Merges



**Not all the added vocabulary tokens are utilized**

hypercholesterolemia: *hyper, ch, olester, ole, mia*

# How does BPE work

Split the input text down to character level

u-n-r-e-l-a-t-e-d

# How does BPE work

Split the input text down to character level

u-n-r-e-l-a-t-e-d

Find **applicable merges** on the list and get *their rank*

u-n: 148

r-e: 6

...

e-d: 22

# How does BPE work

Split the input text down to character level

u-n-r-e-l-a-t-e-d

Find **applicable merges** on the list and get *their rank*

u-n: 148

r-e: 6

...

e-d: 22

Apply the **top-ranked** merge rule

u-n-**re**-l-a-t-e-d



# How does BPE work

u-n-r-e-l-a-t-e-d    Merge **r-e**  
u-n-re-l-a-t-e-d

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at- <span style="border: 1px solid black; padding: 0 2px;">e-d</span>	Merge <b>e-d</b>

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at- <span style="border: 1px solid black;">e-d</span>	Merge <b>e-d</b>
u-n-re-l-at-ed	

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
<span style="border: 1px solid black;">u-n</span> -re-l-at-ed	Merge <b>u-n</b>

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
<span style="border: 1px solid black;">u-n</span> -re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>

Repeat merges till no merge applicable



# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>
un-re-l-ated	

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>
un-re-l-ated	Merge <b>re-l</b>

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>
un-re-l-ated	Merge <b>re-l</b>
un-rel-ated	

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>
un-re-l-ated	Merge <b>re-l</b>
un-rel-ated	Merge <b>rel-ated</b>

Repeat merges till no merge applicable

# How does BPE works

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>
un-re-l-ated	Merge <b>re-l</b>
un-rel-ated	Merge <b>rel-ated</b>
un-related	

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>
un-re-l-ated	Merge <b>re-l</b>
un-rel-ated	Merge <b>rel-ated</b>
<span>un-related</span>	Merge <b>un-related</b>

Repeat merges till no merge applicable

# How does BPE work

u-n-r-e-l-a-t-e-d	Merge <b>r-e</b>
u-n-re-l-a-t-e-d	Merge <b>a-t</b>
u-n-re-l-at-e-d	Merge <b>e-d</b>
u-n-re-l-at-ed	Merge <b>u-n</b>
un-re-l-at-ed	Merge <b>at-ed</b>
un-re-l-ated	Merge <b>re-l</b>
un-rel-ated	Merge <b>rel-ated</b>
un-related	Merge <b>un-related</b>
<u><b>unrelated</b></u>	

Repeat merges **till no merge applicable**

# Issue: Added Merges receive low priority!

- Added vocabulary **appended at the end** of existing PLM vocabulary
- Priority of  $MERGE_{DOMAIN} < MERGE_{PLM}$
- Added merge rules are **never utilized for multiple instances**
- **Resulting in ill tokenization of added vocabulary tokens**



# Issue: Added Merges receive low priority!

- Added vocabulary **appended at the end** of existing PLM vocabulary
- Priority of  $MERGE_{DOMAIN} < MERGE_{PLM}$
- Added merge rules are **never utilized for multiple instances**
- **Resulting in ill tokenization of added vocabulary tokens**

**Can we modify standard BPE to mitigate ill tokenization?**

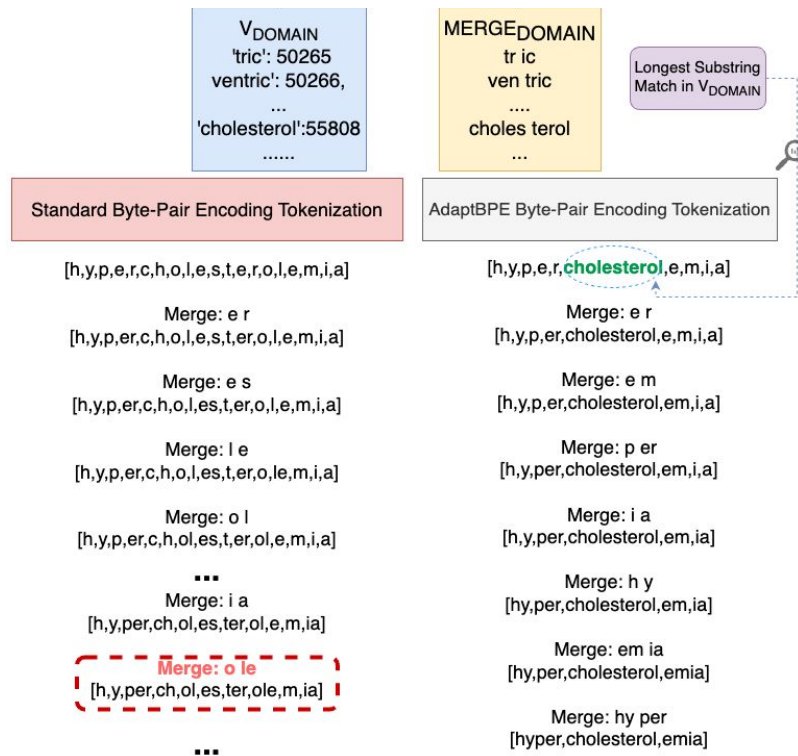
# ADAPTBPE: Mitigate Ill-tokenization

We propose a **fundamental change** in BPE

- BPE starts by splitting the text to character level
- **Instead:**
  - find **longest substring match** in  $V_{\text{DOMAIN}}$  iteratively
  - Preserve **the match** as is
- Run merge operations of BPE

# ADAPTBPE: Mitigate Ill-tokenization

In standard BPE because of merge rules from PLM merges, the word **cholesterol** could **never be formed** resulting in the ill-tokenization issue



However in ADAPTBPE, the subword **cholesterol** is **correctly captured in longest substring match phase** mitigating the ill-tokenization issue

# ADAPTBPE improves Classification (AVocaDo<sup>[1]</sup>)

- SoTA in vocabulary adaptation for classification tasks
- 4 classification datasets from different domain
- Overall improvement of **3.57%**

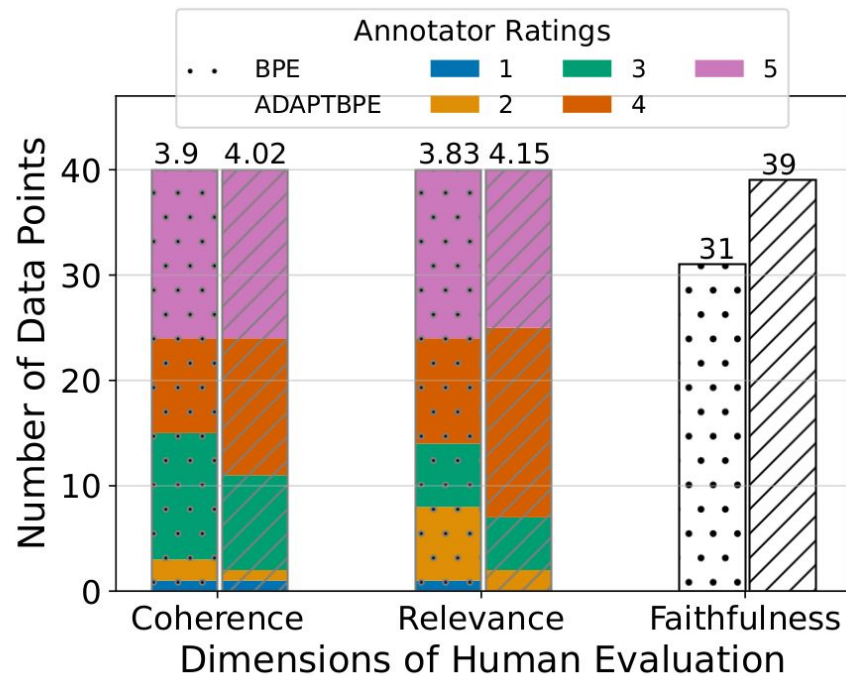
Dataset	Tokenizer	Accuracy	Macro-F1
CHEMPROT (BIOMED)	BPE	81.43 <sub>0.55</sub>	54.88 <sub>1.66</sub>
	ADAPTBPE	81.40 <sub>0.40</sub>	55.02 <sub>0.47</sub>
ACL-ARC (SCIENTIFIC)	BPE	69.03 <sub>5.05</sub>	55.04 <sub>8.24</sub>
	ADAPTBPE	73.02 <sub>4.21</sub>	62.00 <sub>4.95</sub>
HYP (NEWS)	BPE	77.84 <sub>5.20</sub>	74.23 <sub>7.01</sub>
	ADAPTBPE	82.16 <sub>2.50</sub>	80.64 <sub>3.03</sub>
AMAZON (REVIEWS)	BPE	83.13 <sub>3.64</sub>	68.34 <sub>0.47</sub>
	ADAPTBPE	86.26 <sub>0.53</sub>	69.90 <sub>0.29</sub>

# ADAPTBPE improves Summarization (MEDVOC<sup>[1]</sup>)

- SoTA in vocabulary adaptation for summarization
- 4 medical summarization datasets:
  - 2 query-focussed
  - 2 consumer health query
- Overall improvement of **1.87%**
- Better in high OOV concentration

Dataset	Tokenizer	R-L (All)	R-L (H-O)
EBM	BPE	20.65	19.23
	ADAPTBPE	20.73	21.43
BioASQ	BPE	48.02	39.23
	ADAPTBPE	47.72	42.95
MeQSum	BPE	55.88	75.56
	ADAPTBPE	58.00	82.64
CHQ	BPE	40.59	33.77
	ADAPTBPE	41.92	37.60

# Medical experts find summaries to be more relevant and faithful



# Conclusion and Takeaways

- First to show **incorrect BPE tokenization issue** for vocabulary adaptation
- ADAPTBPE is applicable to any vocabulary adaptation strategy
- ADAPTBPE improves over BPE in SoTA classification and summarization

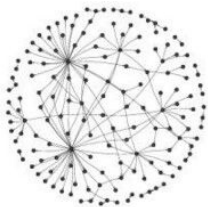


# Thank You!

**P M R F**  
Prime Minister's Research Fellows



Soumyadeep Roy   Mainack Mondal   Niloy Ganguly



**CNeRG**



**Codebase**



**Preprint**



# Any Questions?

Correspondence email:  
[balde.gunjan0812@kgpian.iitkgp.ac.in](mailto:balde.gunjan0812@kgpian.iitkgp.ac.in)