

Token(s) of Appreciation for BPE

Vilém Zouhar et al, September 2024



Token(s) of Appreciation for BPE

Vilém Zouhar et al, September 2024

Token(s) of Appreciation for BPE

Vilém Zouhar et al, September 2024

A Formal Perspective on Byte-Pair Encoding
formalization & bounds..

Token(s) of Appreciation for BPE

Vilém Zouhar et al, September 2024

A Formal Perspective on Byte-Pair Encoding
formalization & bounds..

Tokenization and the Noiseless Channel
how to apriori choose the best tokenization and how we published a wrong paper..

Token(s) of Appreciation for BPE

Vilém Zouhar et al, September 2024

A Formal Perspective on Byte-Pair Encoding
formalization & bounds..

Tokenization and the Noiseless Channel
how to apriori choose the best tokenization and how we published a wrong paper..

Distributional Properties of Subword Regularization
and we use tokenization incorrectly..

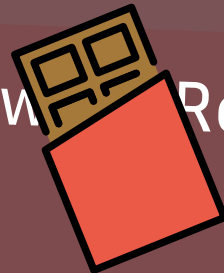
Token(s) of Appreciation for BPE

Vilém Zouhar et al, September 2024

A Formal Perspective on Byte-Pair Encoding
formalization & bounds..

Tokenization and the Noiseless Channel
how to apriori choose the best tokenization and how we published a wrong paper..

Distributional Properties of Subword Re
and we use tokenization incorrectly..



Trick question for a
chocolate at the end.

tokenization

Scholar About 22'600 results (0.07 sec) YEAR

Since 2020

byte-pair encoding

Scholar About 8'100 results (0.09 sec) YEAR

Since 2020

```
vilda@karotte:~$ grep -i "token\|BPE\|byte.pair" acl-anthology.bib | wc -l
3109
vilda@karotte:~$
```

tokenization

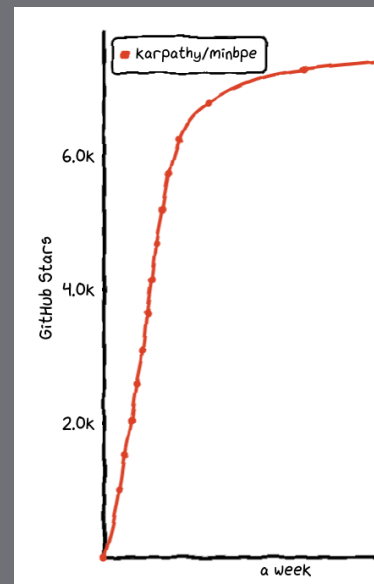
Scholar About 22'600 results (0.07 sec) YEAR

Since 2020

byte-pair encoding

Scholar About 8'100 results (0.09 sec) YEAR

Since 2020



```
vilda@karotte:~$ grep -i "token\|BPE\|byte.pair" acl-anthology.bib | wc -l
3109
vilda@karotte:~$
```


How to represent text?

“he loved pickled pickles”

How to represent text?

“he loved pickled pickles”

[h, e, _, l, o, v, e, d, _, p, i, c, k, l, e, d, _, p, i, c, k, l, e, s]

[104, 101, 32, 108, 111, 118, 101, 100, 32, 112, 105, 99, 107, 108, 101, 100, 32, 112, 105, 99, 107, 108, 101, 115]

How to represent text?

“he loved pickled pickles”

[h, e, _, l, o, v, e, d, _, p, i,

meaningless items, long sequences

[104, 101, 32, 108, 111, 118, 101, 100, 32, 112, 105, 99, 107, 108, 101, 100, 32, 112, 105, 99, 107, 108, 101, 115]

How to represent text?

“he loved pickled pickles”

[h, e, _, l, o, v, e, d, _, p, i, **meaningless items, long sequences**

[104, 101, 32, 108, 111, 118, 101, 100, 32, 112, 105, 99, 107, 108, 101, 100, 32, 112, 105, 99, 107, 108, 101, 115]

[he, _, loved, _, pickled, _, pickles]

[1045, 0, 30123, 0, 6232, 0, 72057]

How to represent text?

“he loved pickled pickles”

[h, e, _, l, o, v, e, d, _, p, i,

meaningless items, long sequences

[104, 101, 32, 108, 111, 118, 101, 100, 32, 112, 105, 99, 107, 108, 101, 100, 32, 112, 105, 99, 107, 108, 101, 115]

[he, _, loved, _, pickled, _

high dimensionality, rare & unknown words

[1045, 0, 30123, 0, 6232, 0, 72057]

How to represent text?

“he loved pickled pickles”

[h, e, _, l, o, v, e, d, _, p, i, **meaningless items, long sequences**

[104, 101, 32, 108, 111, 118, 101, 100, 32, 112, 105, 99, 107, 108, 101, 100, 32, 112, 105, 99, 107, 108, 101, 115]

[he, _, loved, _, pickled, _ **high dimensionality, rare & unknown words**

[1045, 0, 30123, 0, 6232, 0, 72057]

[he, _, lov, ed, _, pickl, ed, _, pickl, es]

[532, 0, 20, 952, 0, 1911, 952, 0, 1911, 12]

How to represent text?

“he loved pickled pickles”

[h, e, _, l, o, v, e, d, _, p, i,

meaningless items, long sequences

[104, 101, 32, 108, 111, 118, 101, 100, 32, 112, 105, 99, 107, 108, 101, 100, 32, 112, 105, 99, 107, 108, 101, 115]

[he, _, loved, _, pickled, _

high dimensionality, rare & unknown words

[1045, 0, 30123, 0, 6232, 0, 72057]

[he, _, lov, ed, _, pickl, e

how to find those?

[532, 0, 20, 952, 0, 1911, 952, 0, 1911, 12]

Finding meaningful *subwords*

Morphology segmentation:

- [he, _, lov, ed, _, pickl, ed, _, pickl, es]

Finding meaningful *subwords*

Morphology segmentation:

- [he, _, lov, ed, _, pickl, ed, _, pickl, es]
- requires language knowledge
- what if we desire higher/lower granularity?

Finding meaningful *subwords*

Morphology segmentation:

- [he, _, lov, ed, _, pickl, ed, _, pickl, es]
- requires language knowledge
- what if we desire higher/lower granularity?

Dedicated algorithm (wish):

- unsupervised
- granularity/vocabulary size hyperparameter
 - [he, _, loved, _, pickl, ed, _, pickl, es]
 - [he, _, lov, ed, _, pic, kl, ed, _, pic, kl, es]

Finding meaningful *subwords*

Morphology segmentation:

- [he, _, lov, ed, _, pickl, ed, _, pickl, es]
- requires language knowledge
- what if we desire higher/lower granularity?

Byte-Pair Encoding:

- unsupervised
- granularity/vocabulary size hyperparameter
 - [he, _, loved, _, pickl, ed, _, pickl, es]
 - [he, _, lov, ed, _, pic, kl, ed, _, pic, kl, es]

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge p i → pi

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge c k → ck

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

merge pi ck → pick

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k → pick

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k → pick

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k l → pickl

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k → pick

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k l → pickl

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k → pick

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k l → pickl

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Unsupervised ?

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k → pick

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k l → pickl

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Unsupervised ✓ Granularity ?

Byte-Pair Encoding

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i → pi

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge c k → ck

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k → pick

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s
merge p i c k l → pickl

t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

Unsupervised ✓ Granularity ✓

Byte-Pair Encoding

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

Byte-Pair Encoding

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

Byte-Pair Encoding

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

Output example:

$\bar{x} = \text{t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s}$

$\bar{\mu} = \langle \text{p i} \rightarrow \text{pi}, \text{ c k} \rightarrow \text{ck}, \text{ pi ck} \rightarrow \text{pick}, \text{ pick l} \rightarrow \text{pickl} \rangle$

Byte-Pair Encoding

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

Runtime:

$O(M) \times ?$

Output example:

$\bar{x} = \text{t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s}$

$\bar{\mu} = \langle \text{p i} \rightarrow \text{pi}, \text{ c k} \rightarrow \text{ck}, \text{ pi ck} \rightarrow \text{pick}, \text{ pick l} \rightarrow \text{pickl} \rangle$

Byte-Pair Encoding

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

Runtime:

$O(M) \times (\operatorname{argmax} \operatorname{Freq} + \operatorname{ApplyMerge})$

Output example:

$\bar{x} = \text{t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s}$

$\bar{\mu} = \langle \text{p i} \rightarrow \text{pi}, \text{ c k} \rightarrow \text{ck}, \text{ pi ck} \rightarrow \text{pick}, \text{ pick l} \rightarrow \text{pickl} \rangle$

Byte-Pair Encoding

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

Runtime:

$O(M \times |\bar{x}|)$

Output example:

$\bar{x} = \text{t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s}$

$\bar{\mu} = \langle \text{p i} \rightarrow \text{pi}, \text{ c k} \rightarrow \text{ck}, \text{ pi ck} \rightarrow \text{pick}, \text{ pick l} \rightarrow \text{pickl} \rangle$

Byte-Pair Encoding

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

Runtime:

$O(M \times |\bar{x}|)$

Not good, will fix later

Output example:

$\bar{x} = \text{t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s}$

$\bar{\mu} = \langle \text{p i} \rightarrow \text{pi}, \text{ c k} \rightarrow \text{ck}, \text{ pi ck} \rightarrow \text{pick}, \text{ pick l} \rightarrow \text{pickl} \rangle$

BPE is Compression

$s_1 =$ t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

$$|s_1| = 27$$

$$|V_1| = 12$$

BPE is Compression

$s_1 =$ t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

$$|s_1| = 27$$

$$|V_1| = 12$$

$$\text{cost} = |s_1| \cdot \lceil \log_2 |V_1| \rceil = 22 \cdot 4 = 108 \text{ bits}$$

BPE is Compression

$s_1 = \text{t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s}$

$$|s_1| = 27$$

$$|V_1| = 12$$

$$\text{cost} = |s_1| \cdot \lceil \log_2 |V_1| \rceil = 27 \cdot 4 = 108 \text{ bits}$$

$s_2 = \text{t h e y _ p i c k e d _ _ p i c k l e d _ _ p i c k l e s}$

$$|s_2| = 14$$

$$|V_2| = |V_1| + |\{\text{pi, ck, pick, pickl, ed}\}| = 17$$

BPE is Compression

$s_1 =$ t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

$$|s_1| = 27$$

$$|V_1| = 12$$

$$\text{cost} = |s_1| \cdot \lceil \log_2 |V_1| \rceil = 27 \cdot 4 = 108 \text{ bits}$$

$s_2 =$ t h e y _ p i c k e d _ p i c k l e d _ p i c k l e s

$$|s_2| = 14$$

$$|V_2| = |V_1| + |\{\text{pi}, \text{ck}, \text{pick}, \text{pickl}, \text{ed}\}| = 17$$

$$\text{cost} = |s_2| \cdot \lceil \log_2 |V_2| \rceil = 14 \cdot 5 = 70 \text{ bits}$$

BPE is Greedy

BPE is Greedy

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

BPE is Greedy

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

BPE is Greedy

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

Decision example:

$\bar{x} = a \ b \ a \ a \ b \ b \ a \ a$

$\operatorname{Freq}(x, (\mu', \mu'')) = \langle a \ b \rightarrow 2, \quad b \ a \rightarrow 2, \quad a \ a \rightarrow 2, \quad b \ b \rightarrow 1 \rangle$

BPE is *not* Optimal

$$\begin{aligned}\bar{x} &= a \ b \ a \ a \ b \ b \ a \ a \\ &= ab \ a \ ab \ b \ a \ a \\ &= aba \ ab \ b \ a \ a\end{aligned}$$

BPE is *not* Optimal

$$\begin{aligned}\bar{x} &= a \ b \ a \ a \ b \ b \ a \ a \\ &= ab \ a \ ab \ b \ a \ a \\ &= aba \ ab \ b \ a \ a\end{aligned}$$

$$\bar{\mu}^\triangle = \langle a \ b \rightarrow ab, \ ab \ a \rightarrow aba \rangle$$

$$|\text{ApplyMerges}(\bar{\mu}^\triangle, \bar{x})| = 5$$

BPE is *not* Optimal

$$\begin{aligned}\bar{x} &= a b a a b b a a \\ &= ab a ab b a a \\ &= aba ab b a a\end{aligned}$$

$$\bar{\mu}^\blacktriangle = \langle a b \rightarrow ab, ab a \rightarrow aba \rangle \quad |\text{ApplyMerges}(\bar{\mu}^\blacktriangle, \bar{x})| = 5$$

$$\begin{aligned}\bar{x} &= a b a a b b a a \\ &= a ba a b ba a \\ &= a baa b baa\end{aligned}$$

$$\bar{\mu}^\star = \langle b a \rightarrow ba, ba a \rightarrow baa \rangle$$

$$\text{ApplyMerges}(\bar{\mu}^\star, \bar{x}) = a baa b baa \quad |\text{ApplyMerges}(\bar{\mu}^\star, \bar{x})| = 4$$

BPE is *Approximately*-Optimal

Given any $\bar{x} \in \Sigma^*$

c^\wedge = compression length with BPE

$|\text{BPE}(\bar{x})_0|$

c^\star = best compression length

$\min_{\bar{\mu}, |\bar{\mu}| \leq M} |\text{ApplyMerges}(\bar{\mu}, \bar{x})|$

$|\bar{x}| - c^\wedge$ = how many characters were saved

(higher is better)

BPE is *Approximately*-Optimal

Given any $\bar{x} \in \Sigma^*$

c^\blacktriangle = compression length with BPE $|\text{BPE}(\bar{x})|$

c^\star = best compression length $\min_{\bar{\mu}, |\bar{\mu}| \leq M} |\text{ApplyMerges}(\bar{\mu}, \bar{x})|$

$|\bar{x}| - c^\blacktriangle$ = how many characters were saved (higher is better)

$$\frac{|\bar{x}| - c^\blacktriangle}{|\bar{x}| - c^\star} \geq \frac{1 - e^{-\sigma(\bar{\mu}^\star)}}{\sigma(\bar{\mu}^\star)} \approx 37\%$$

BPE is *Approximately*-Optimal

Given any $\bar{x} \in \Sigma^*$

c^\blacktriangle = compression length with BPE

$|\text{BPE}(\bar{x})|$

c^\star = best compression length

$\min_{\bar{\mu}, |\bar{\mu}| \leq M} |\text{ApplyMerges}(\bar{\mu}, \bar{x})|$

$|\bar{x}| - c^\blacktriangle$ = how many characters were saved

(higher is better)

$$\frac{|\bar{x}| - c^\blacktriangle}{|\bar{x}| - c^\star} \geq \frac{1 - e^{-\sigma(\bar{\mu}^\star)}}{\sigma(\bar{\mu}^\star)} \approx 37\%$$

BPE is *Approximately*-Optimal: Proof

Let $\kappa(\bar{\mu}) = |x| - c^\wedge$

We get $\frac{|x| - c^\wedge}{|x| - c^\star} \geq \frac{1 - e^{-\sigma(\bar{\mu}^\star)}}{\sigma(\bar{\mu}^\star)}$ if $\kappa(\bar{\mu})$ is:

BPE is *Approximately*-Optimal: Proof

Let $\kappa(\bar{\mu}) = |x| - c^\wedge$

We get $\frac{|x| - c^\wedge}{|x| - c^\star} \geq \frac{1 - e^{-\sigma(\bar{\mu}^\star)}}{\sigma(\bar{\mu}^\star)}$ if $\kappa(\bar{\mu})$ is:

- (1) monotone non-decreasing, and
- (2) submodular, and
- (3) hierarchical sequence-submodular.

BPE is *Approximately*-Optimal: Proof

Let $\kappa(\bar{\mu}) = |x| - c^\wedge$

We get $\frac{|\bar{x}| - c^\wedge}{|\bar{x}| - c^\star} \geq \frac{1 - e^{-\sigma(\bar{\mu}^\star)}}{\sigma(\bar{\mu}^\star)}$ if $\kappa(\bar{\mu})$ is:

- (1) monotone non-decreasing, and
- (2) submodular, and
- (3) hierarchical sequence-submodular.

(1) monotone non-decreasing:

$$\kappa(\mu\nu) \geq \kappa(\mu)$$

BPE is *Approximately*-Optimal: Proof

Let $\kappa(\bar{\mu}) = |x| - c^\Delta$

We get $\frac{|\bar{x}| - c^\Delta}{|\bar{x}| - c^\star} \geq \frac{1 - e^{-\sigma(\bar{\mu}^\star)}}{\sigma(\bar{\mu}^\star)}$ if $\kappa(\bar{\mu})$ is:

- (1) monotone non-decreasing, and
- (2) submodular, and
- (3) hierarchical sequence-submodular.

(1) monotone non-decreasing:

$$\kappa(\mu\nu) \geq \kappa(\mu)$$

(2) submodular: $\kappa(\nu \mid \mu) \geq \kappa(\nu \mid \mu\omega) \Leftrightarrow \kappa(\mu\nu) - \kappa(\mu) \geq \kappa(\mu\omega\nu) - \kappa(\mu\omega)$

BPE is *Approximately*-Optimal: Proof

Let $\kappa(\bar{\mu}) = |x| - c^\wedge$

We get $\frac{|\bar{x}| - c^\wedge}{|\bar{x}| - c^\star} \geq \frac{1 - e^{-\sigma(\bar{\mu}^\star)}}{\sigma(\bar{\mu}^\star)}$ if $\kappa(\bar{\mu})$ is:

- (1) monotone non-decreasing, and
- (2) submodular, and
- (3) hierarchical sequence-submodular.

(1) monotone non-decreasing:

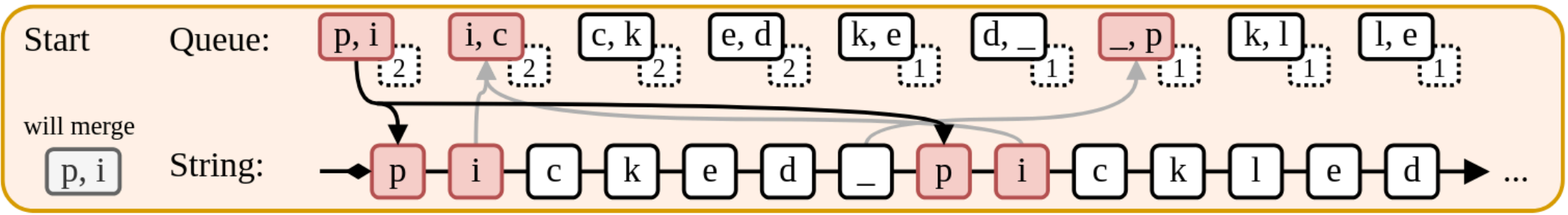
$$\kappa(\mu\nu) \geq \kappa(\mu)$$

(2) submodular: $\kappa(\nu \mid \mu) \geq \kappa(\nu \mid \mu\omega) \Leftrightarrow \kappa(\mu\nu) - \kappa(\mu) \geq \kappa(\mu\omega\nu) - \kappa(\mu\omega)$

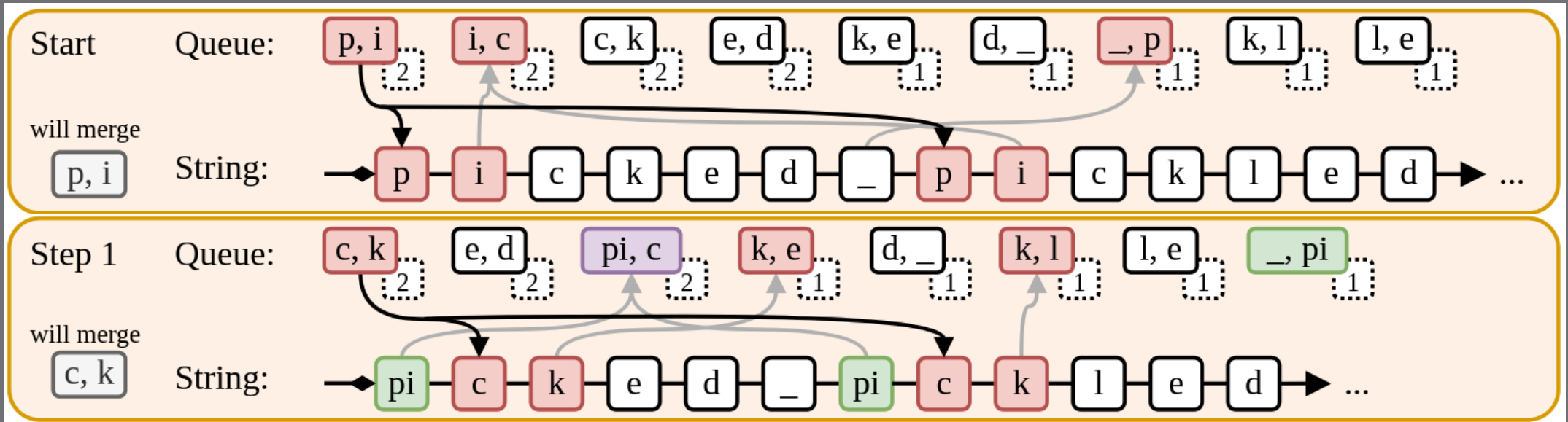
(2) hierarchical sequence-submodular:

$$\kappa(\nu' \mid \mu') \geq \kappa(\nu \mid \mu' \nu' \mu)$$

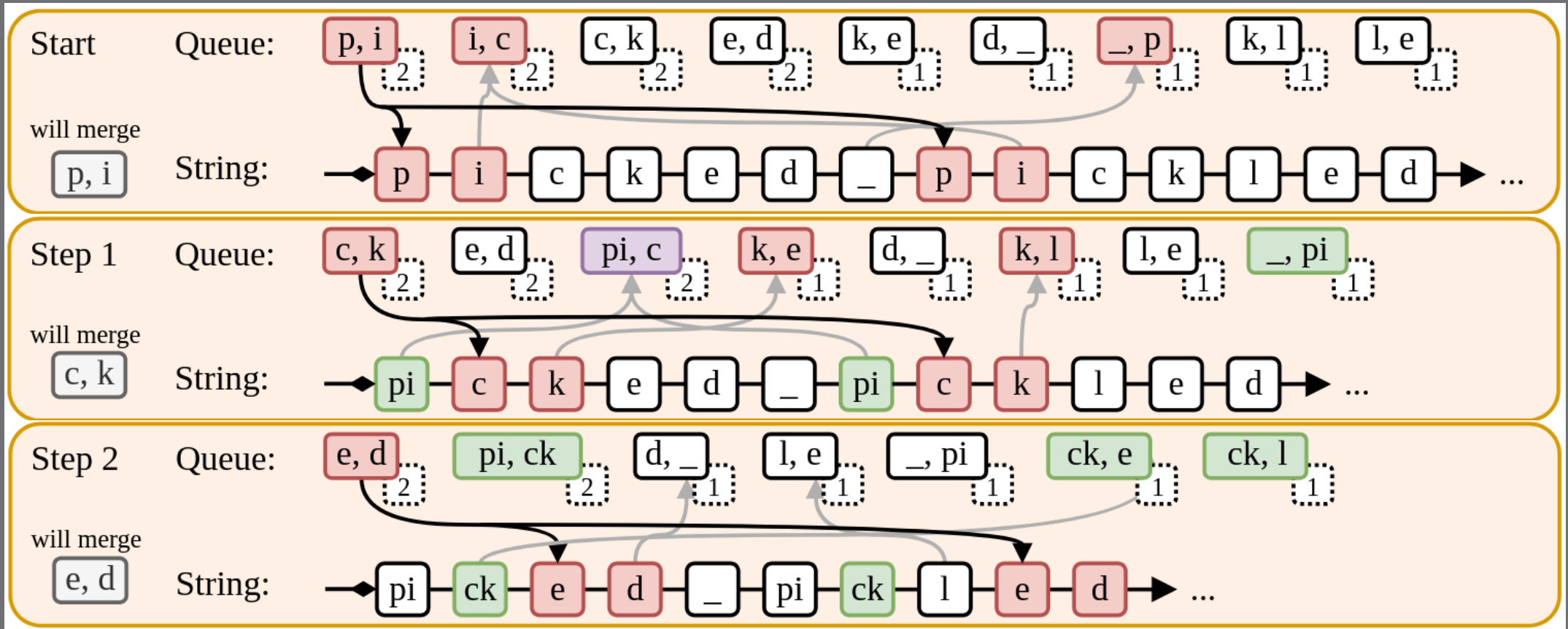
BPE can be faster



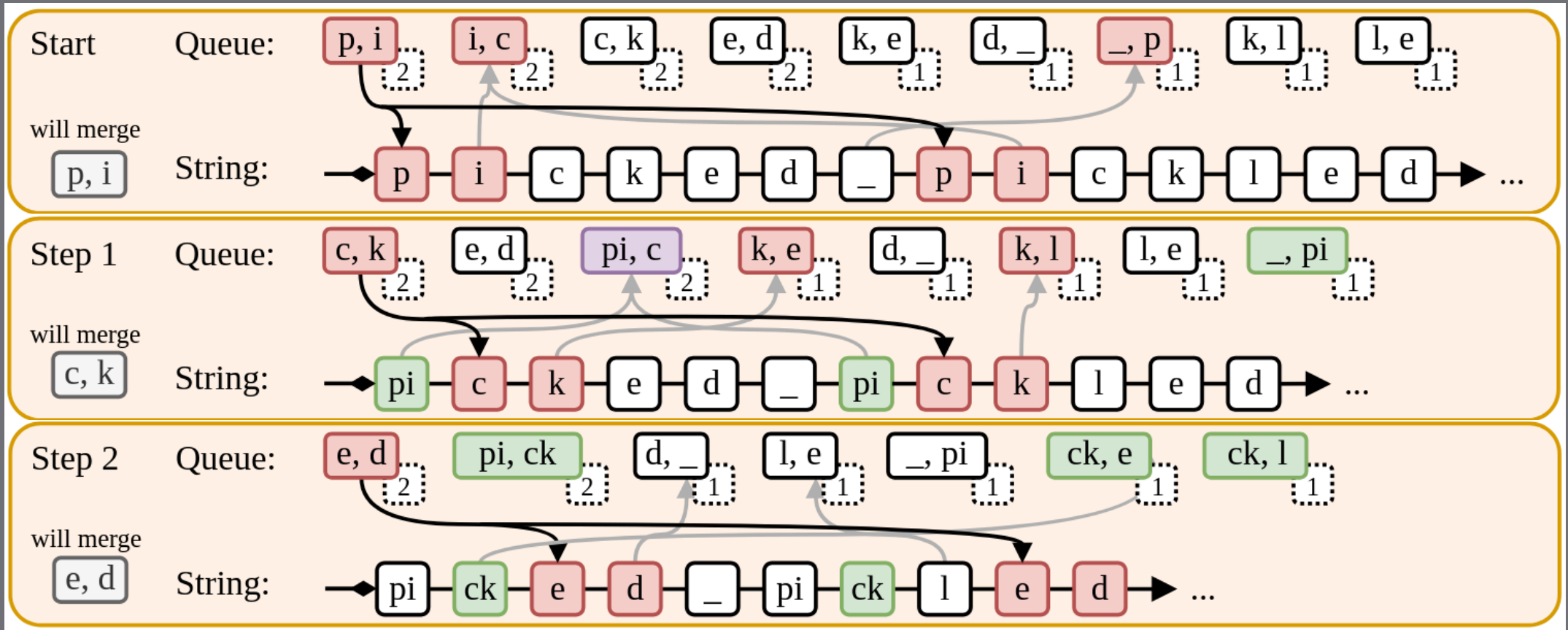
BPE can be faster



BPE can be faster



BPE can be faster



$$\text{Runtime } O\left(\sum_{t=0}^M R_t \log M\right) = O(N \log M)$$

(at most M merges)

BPE can be better: beam search

$\bar{x} = a \ b \ a \ a \ b \ b \ a \ a$

$\text{Freq}(x, (\mu', \mu'')) = \langle a \ b \rightarrow 2, \ b \ a \rightarrow 2, \ a \ a \rightarrow 2, \ b \ b \rightarrow 1 \rangle$

BPE can be better: beam search

$\bar{x} = a b a a b b a a$

$\text{Freq}(x, (\mu', \mu'')) = \langle a b \rightarrow 2, b a \rightarrow 2, a a \rightarrow 2, b b \rightarrow 1 \rangle$

• Step 1:

- ▶ Beam 1: $\langle a b \rightarrow ab \rangle$ $\langle ab a \rightarrow 1, a ab \rightarrow 1, ab b \rightarrow 1, b a \rightarrow 1, a a \rightarrow 1, \rangle$
- ▶ Beam 2: $\langle b a \rightarrow b a \rangle$ $\langle ba a \rightarrow 2, a ba \rightarrow 1, b ba \rightarrow 1, a b \rightarrow 1, \rangle$

BPE can be better: beam search

$\bar{x} = a b a a b b a a$

$\text{Freq}(x, (\mu', \mu'')) = \langle a b \rightarrow 2, b a \rightarrow 2, a a \rightarrow 2, b b \rightarrow 1 \rangle$

• Step 1:

▶ Beam 1: $\langle a b \rightarrow ab \rangle$ $\langle ab a \rightarrow 1, a ab \rightarrow 1, ab b \rightarrow 1, b a \rightarrow 1, a a \rightarrow 1, \rangle$

▶ Beam 2: $\langle b a \rightarrow b a \rangle$ $\langle ba a \rightarrow 2, a ba \rightarrow 1, b ba \rightarrow 1, a b \rightarrow 1, \rangle$

• Step 1:

▶ Beam 1: $\langle a b \rightarrow ab, ab a \rightarrow aba \rangle$ $aba ab b a a$

▶ Beam 2: $\langle b a \rightarrow b a, ba a \rightarrow baa \rangle$ $a baa b baa$

BPE can be better: optimal

Input: \bar{x} (text to compress)

Output: \bar{x}^* (optimally compressed text)

$\bar{\mu}^*$ (compression merges)

BPE can be better: optimal

1. $q \leftarrow \text{Stack}()$
2. $q.\text{push}(\langle \langle \rangle, \bar{x} \rangle)$
3. $\bar{\mu}^*, \bar{x}^* \leftarrow \langle \rangle, \bar{x}$
4. **while** q not empty **do**
5. $\bar{\mu}, \bar{x} \leftarrow q.\text{pop}()$
6. **if** $|\bar{\mu}| = M$ **then continue**
7. **for** $\mu \in \text{Pairs}(\bar{x})$ **do**
8. $\bar{x}' \leftarrow \text{Apply}(\bar{x}, \mu)$
9. $\bar{\mu}' \leftarrow \bar{\mu} \circ \mu$
10. **if** $|\bar{x}'| < |\bar{x}^*|$ **then** $\bar{\mu}^*, \bar{x}^* \leftarrow \bar{\mu}', \bar{x}'$
11. $q.\text{push}(\bar{\mu}', \bar{x}')$
12. **end for**
13. **end while**

Input: \bar{x} (text to compress)

Output: \bar{x}^* (optimally compressed text)

$\bar{\mu}^*$ (compression merges)

BPE can be better: optimal

1. $q \leftarrow \text{Stack}()$
2. $q.\text{push}(\langle \langle \rangle, \bar{x} \rangle)$
3. $\bar{\mu}^*, \bar{x}^* \leftarrow \langle \rangle, \bar{x}$
4. **while** q not empty **do**
5. $\bar{\mu}, \bar{x} \leftarrow q.\text{pop}()$
6. **if** $|\bar{\mu}| = M$ **then continue**
7. **for** $\mu \in \text{Pairs}(\bar{x})$ **do**
8. $\bar{x}' \leftarrow \text{Apply}(\bar{x}, \mu)$
9. $\bar{\mu}' \leftarrow \bar{\mu} \circ \mu$
10. **if** $|\bar{x}'| < |\bar{x}^*|$ **then** $\bar{\mu}^*, \bar{x}^* \leftarrow \bar{\mu}', \bar{x}'$
11. $q.\text{push}(\bar{\mu}', \bar{x}')$
12. **end for**
13. **end while**

Input: \bar{x} (text to compress)

Output: \bar{x}^* (optimally compressed text)

$\bar{\mu}^*$ (compression merges)

Runtime: $O(N^M)$

BPE can be better: optimal

Input: \bar{x} (text to compress)

Output: \bar{x}^* (optimally compressed text)

$\bar{\mu}^*$ (compression merges)

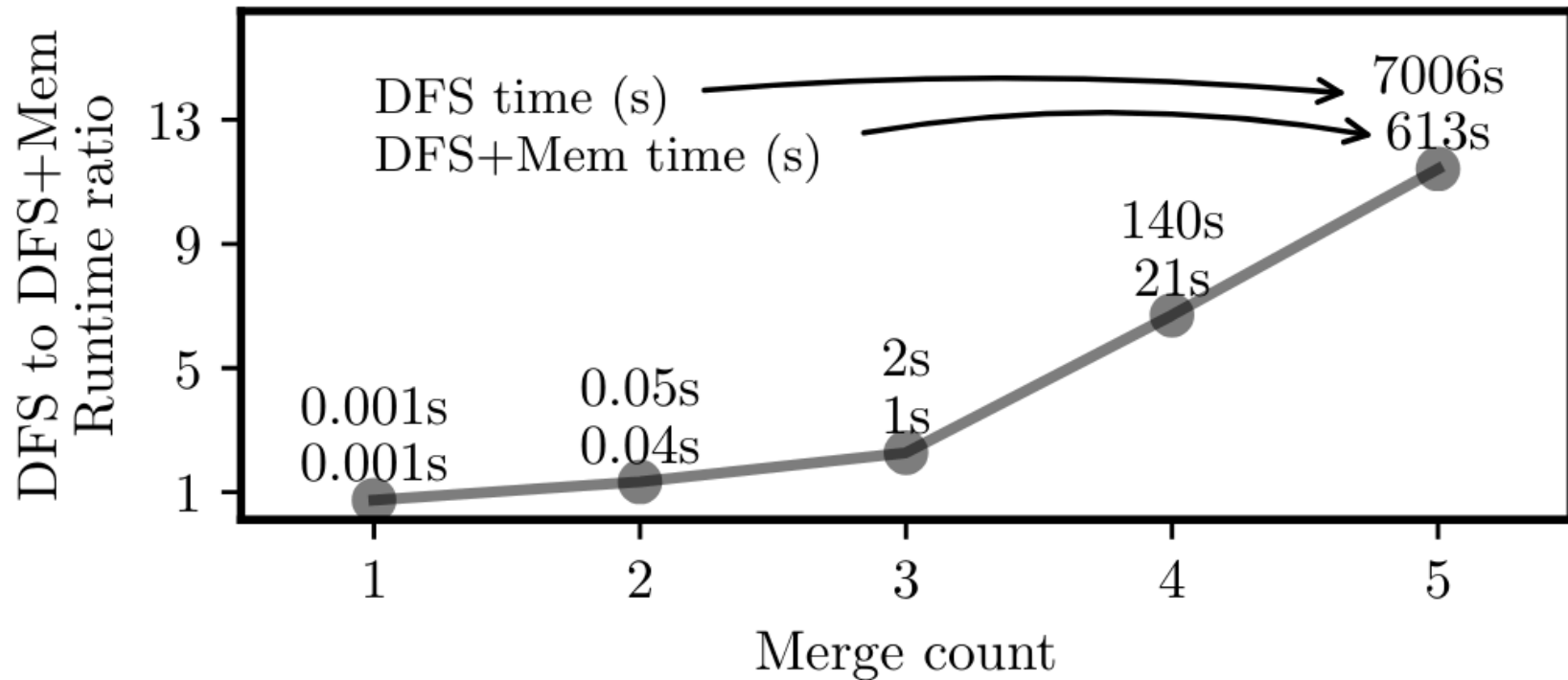
BPE can be better: optimal

1. $q \leftarrow \text{Stack}()$
2. $q.\text{push}(\langle \langle \rangle, \bar{x} \rangle)$
3. $\bar{\mu}^*, \bar{x}^* \leftarrow \langle \rangle, \bar{x}$
4. **while** q not empty **do**
5. $\bar{\mu}, \bar{x} \leftarrow q.\text{pop}()$
6. **if** $|\bar{\mu}| = M$ **then continue**
7. **for** $\mu \in \text{Pairs}(\bar{x})$ **do**
8. * **if** $|\bar{\mu}| \neq 0 \wedge \neg \mu \succ \bar{\mu}_{-1}$ **then continue**
9. $\bar{x}' \leftarrow \text{Apply}(\bar{x}, \mu)$
10. $\bar{\mu}' \leftarrow \bar{\mu} \circ \mu$
11. **if** $|\bar{x}'| < |\bar{x}^*|$ **then** $\bar{\mu}^*, \bar{x}^* \leftarrow \bar{\mu}', \bar{x}'$
12. $q.\text{push}(\bar{\mu}', \bar{x}')$
13. **end for**
14. **end while**

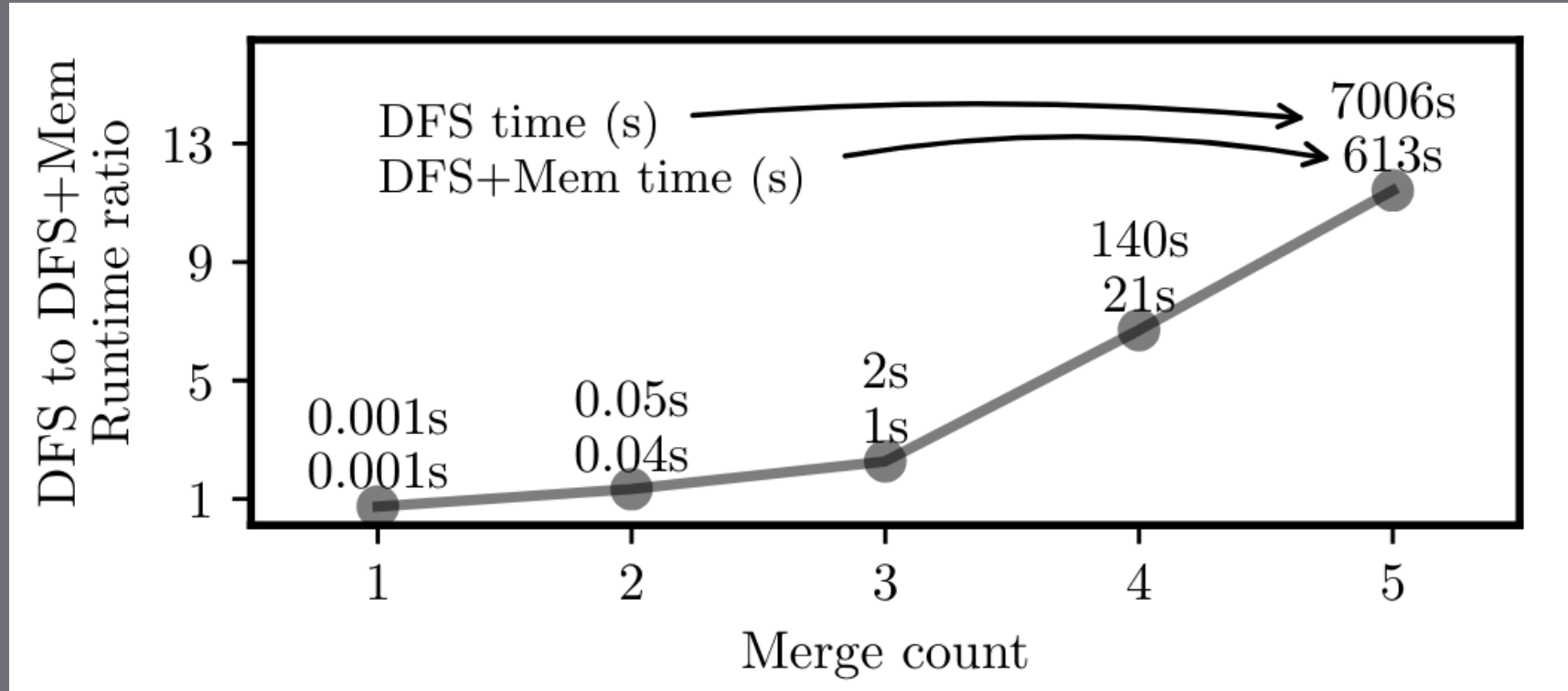
Input: \bar{x} (text to compress)

Output: \bar{x}^* (optimally compressed text)
 $\bar{\mu}^*$ (compression merges)

BPE on natural languages



BPE on natural languages



- Greedy BPE likely close to optimal anyway (proof/evidence missing)

Choosing the right Tokenization

Choosing the right Tokenization

before expensive model training

Choosing the right Tokenization

before expensive model training

th ey _ pi ck ed_ pi ck l ed _ pi ck l es

or

t h e y _ pick e d_ pickl e d _ pickl e s

The best tokenization..

“he loved pickled pickles”

..compresses text the most; shortest sequences

The best tokenization..

“he loved pickled pickles”

..compresses text the most; shortest sequences

each word is in the vocabulary; no segmentation
[he, _, loved, _, pickled, _, pickles]

The best tokenization..

“he loved pickled pickles”

..compresses text the most; shortest sequences

each word is in the vocabulary; no segmentation
[he, _, loved, _, pickled, _, pickles]

..compresses text the most; no redundancy;
no high-/low-frequency tokens

The best tokenization..

“he loved pickled pickles”

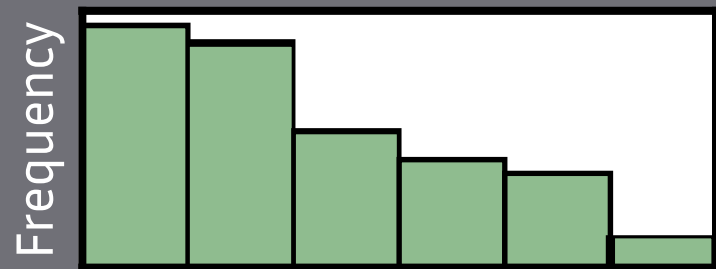
..compresses text the most; shortest sequences

each word is in the vocabulary; no segmentation
[he, _, loved, _, pickled, _, pickles]

..compresses text the most; no redundancy;
no high-/low-frequency tokens

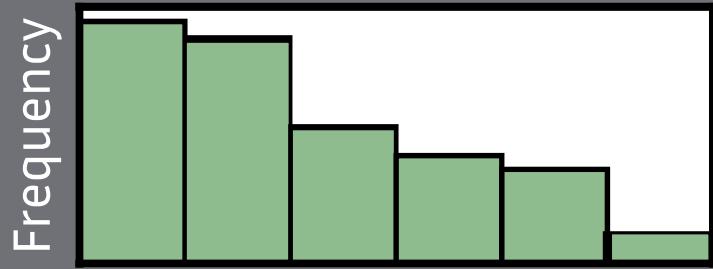
entropy: $\sum_{x \in V} p(x) \log p(x)$
[he, _, lov, ed, _, pickl, ed, _, pickl, es]

Why entropy is a good tokenization metric



p i c k l e d _ p i c k l e s

Why entropy is a good tokenization metric

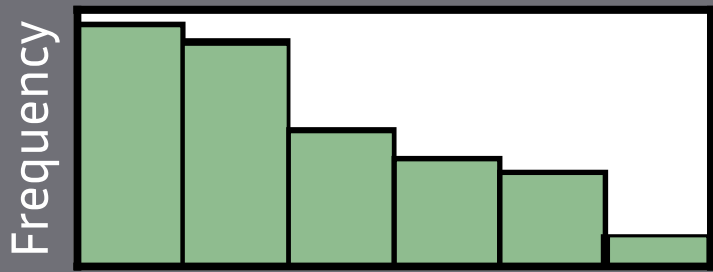


p i c k l e d _ p i c k l e s

high- & low-frequency units

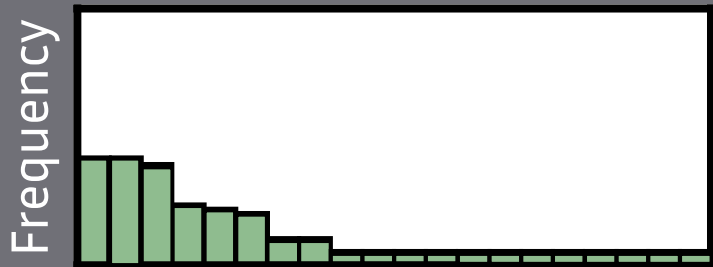
no discernable meaning of each token

Why entropy is a good tokenization metric



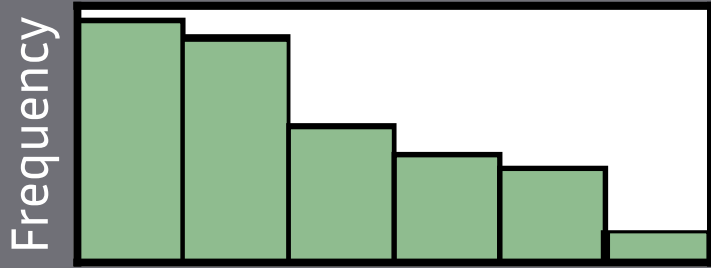
`p i c k l e d _ p i c k l e s`

high- & low-frequency units
no discernable meaning of each token



`pickled _ pickles`

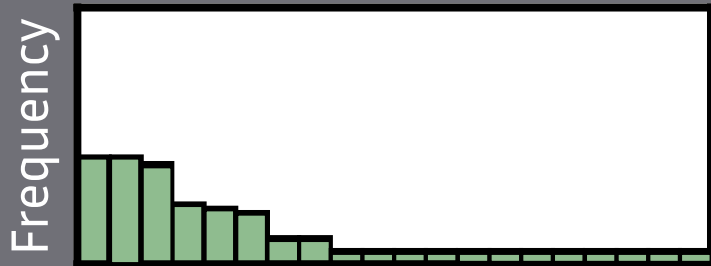
Why entropy is a good tokenization metric



`p i c k l e d _ p i c k l e s`

high- & low-frequency units

no discernable meaning of each token

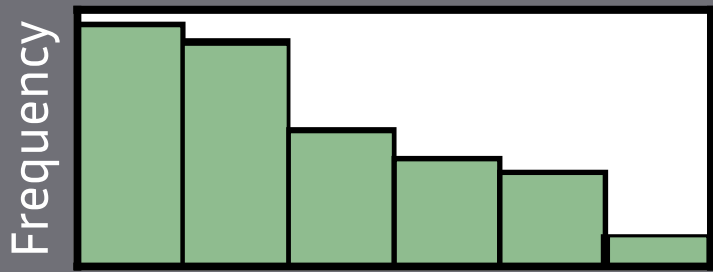


`pickled _ pickles`

low-frequency units

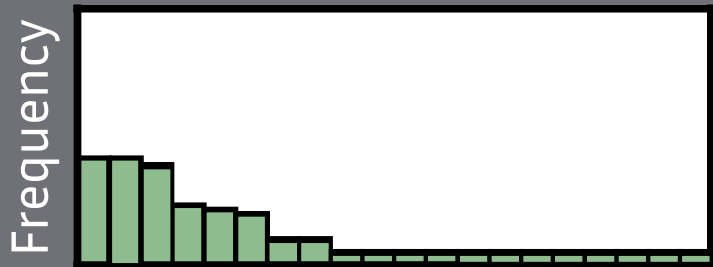
no overlap between words

Why entropy is a good tokenization metric



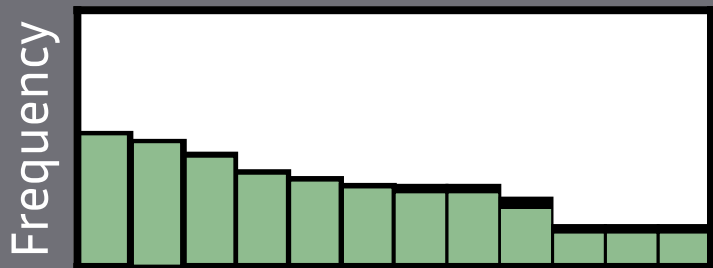
p i c k l e d _ p i c k l e s

high- & low-frequency units
no discernable meaning of each token



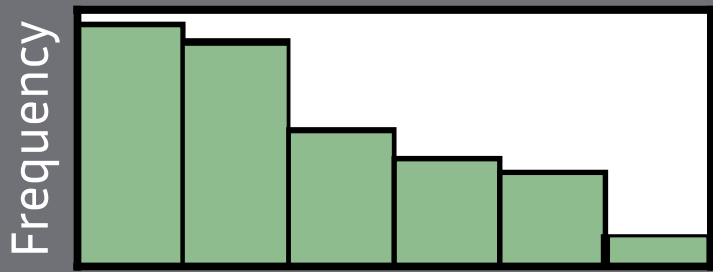
pickled _ pickles

low-frequency units
no overlap between words



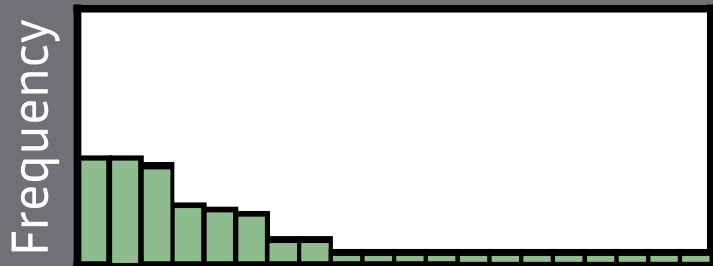
pi ck l ed _ pi ck l es

Why entropy is a good tokenization metric



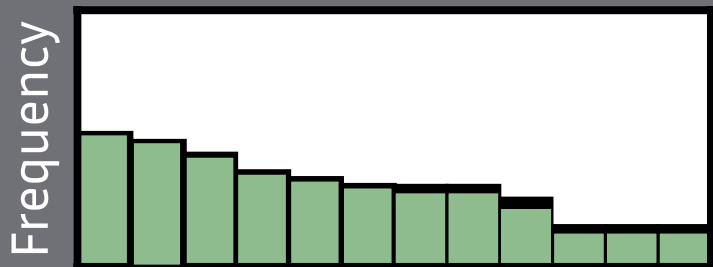
`p i c k l e d _ p i c k l e s`

high- & low-frequency units
no discernable meaning of each token



`pickled _ pickles`

low-frequency units
no overlap between words

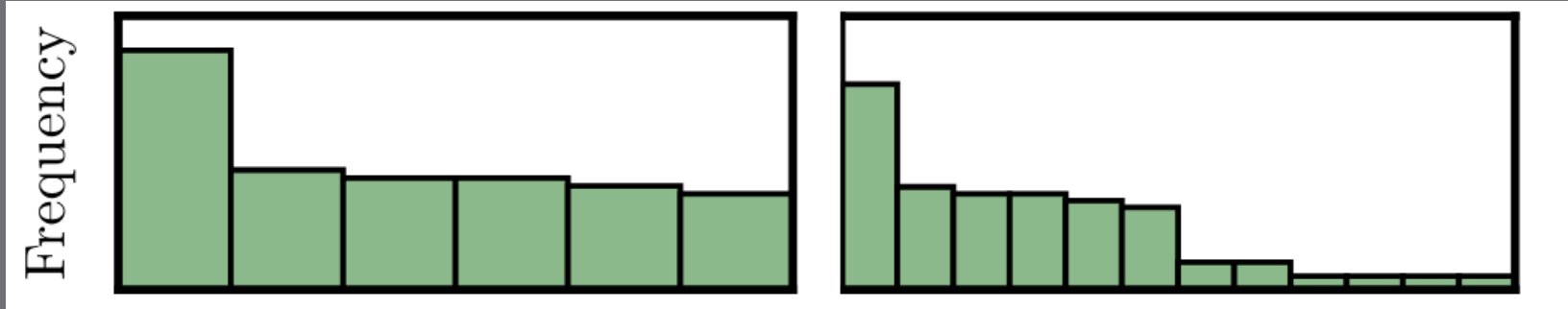


`pi ck l ed _ pi ck l es`

balanced, meaningful units

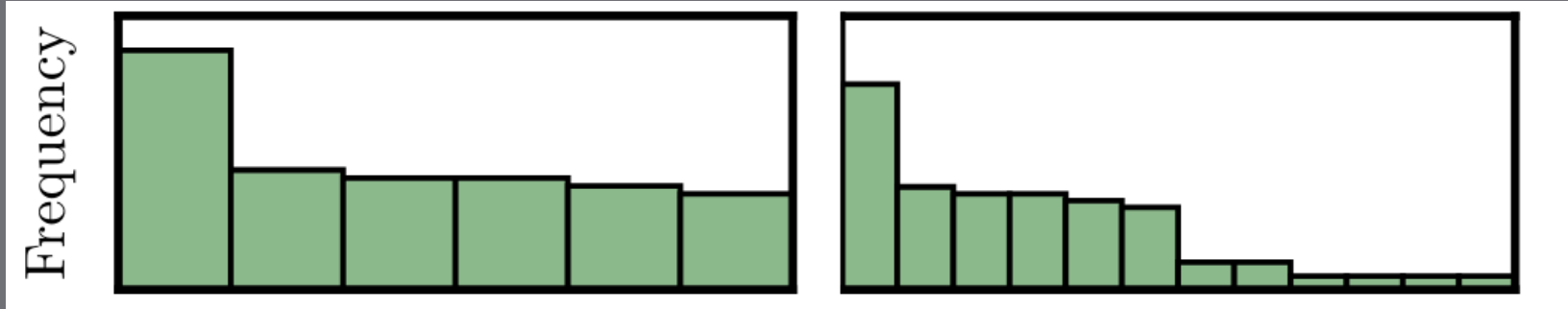
Entropies with varying vocabularies

Entropy bounds depend on support set size $H \in [0, \log|V|]$



Entropies with varying vocabularies

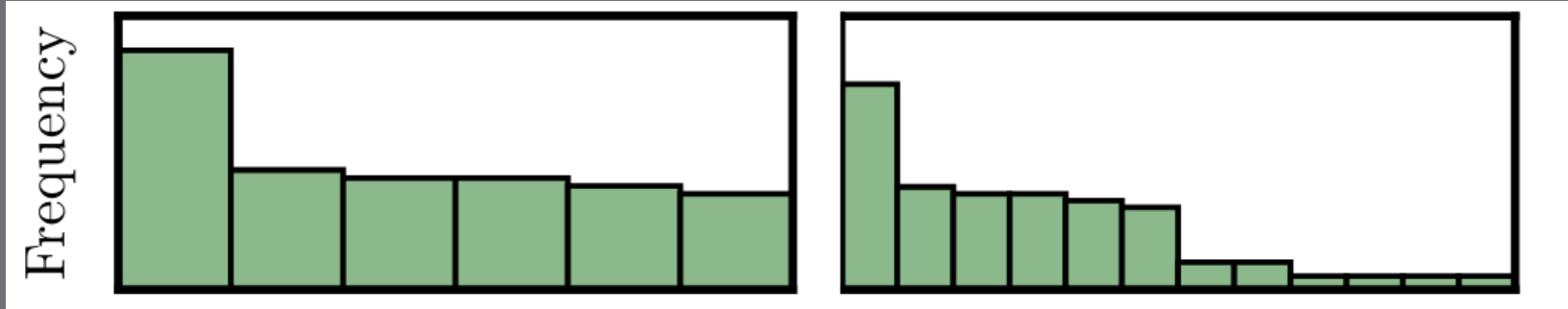
Entropy bounds depend on support set size $H \in [0, \log|V|]$



$$|V_{\text{left}}| = 6 \quad |V_{\text{right}}| = 12$$

Entropies with varying vocabularies

Entropy bounds depend on support set size $H \in [0, \log|V|]$

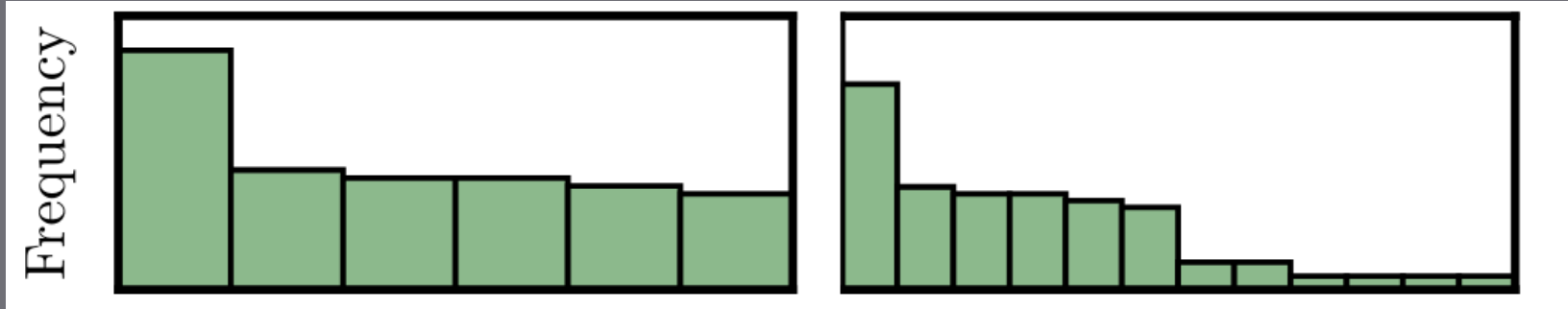


$$|V_{\text{left}}| = 6 \quad |V_{\text{right}}| = 12$$

$$H_{\text{left}} = 2.50 \quad H_{\text{right}} = 3.08$$

Entropies with varying vocabularies

Entropy bounds depend on support set size $H \in [0, \log|V|]$



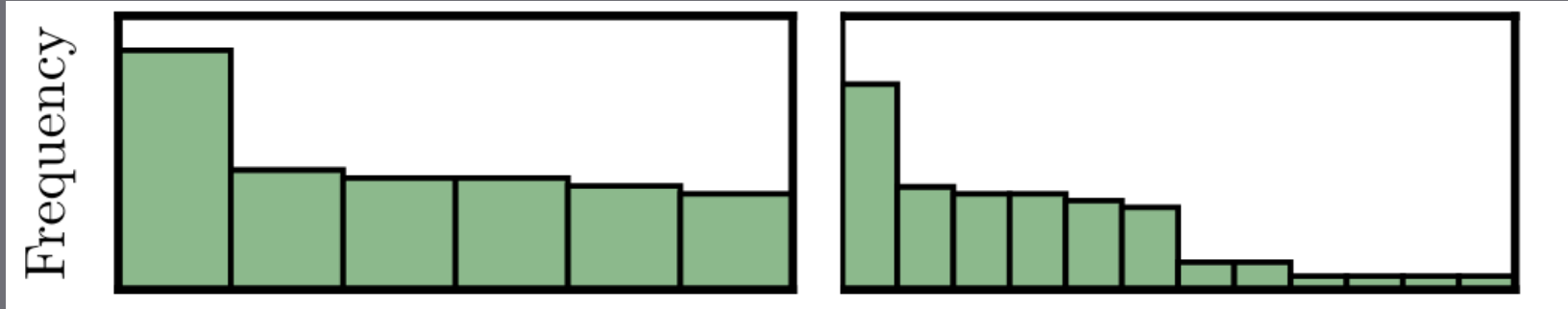
$$|V_{\text{left}}| = 6 \quad |V_{\text{right}}| = 12$$

$$H_{\text{left}} = 2.50 \quad H_{\text{right}} = 3.08$$

proportion to maximal possible entropy given $|V|$

Entropies with varying vocabularies

Entropy bounds depend on support set size $H \in [0, \log|V|]$



$$|V_{\text{left}}| = 6 \quad |V_{\text{right}}| = 12$$

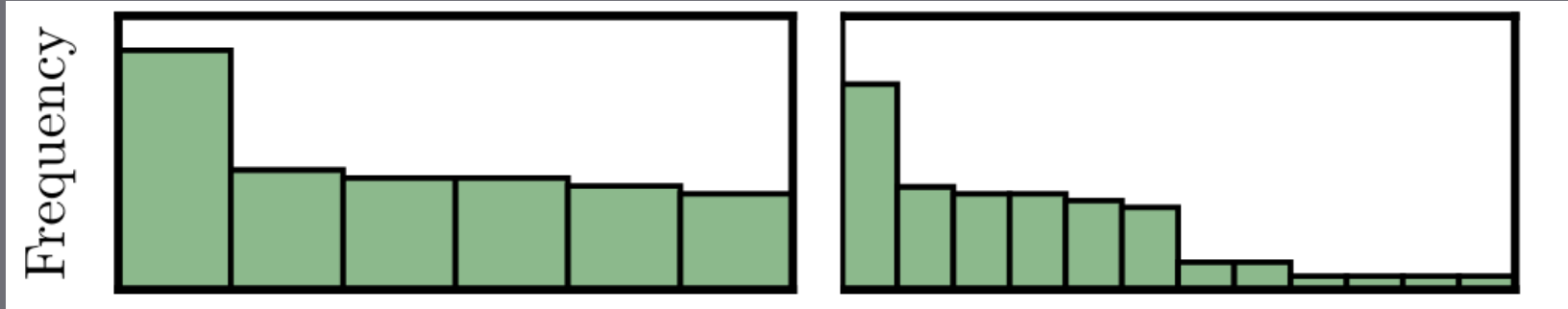
$$H_{\text{left}} = 2.50 \quad H_{\text{right}} = 3.08$$

proportion to maximal possible entropy given $|V|$

$$\text{efficiency: } \frac{\text{optimal encoding of } p}{\text{uninformed encoding of } p} = \frac{H_p}{H_U} = \frac{H_p}{\log|V|}$$

Entropies with varying vocabularies

Entropy bounds depend on support set size $H \in [0, \log|V|]$



$$|V_{\text{left}}| = 6 \quad |V_{\text{right}}| = 12$$

$$H_{\text{left}} = 2.50 \quad H_{\text{right}} = 3.08$$

proportion to maximal possible entropy given $|V|$

$$\text{efficiency: } \frac{\text{optimal encoding of } p}{\text{uninformed encoding of } p} = \frac{H_p}{H_U} = \frac{H_p}{\log|V|}$$

$$\text{Eff}_{\text{left}} = 97\% \quad \text{Eff}_{\text{right}} = 86\%$$

Entropies with varying penalization

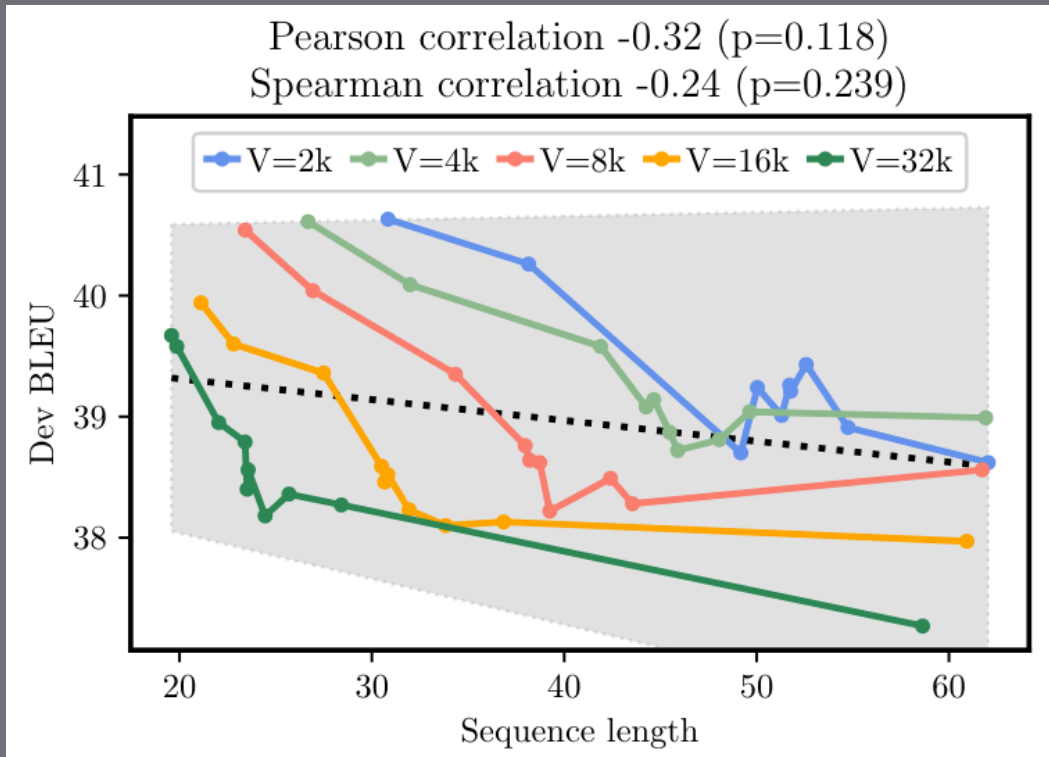
Shannon Entropy: $-\sum p \log p$ \approx expected optimal encoding length
low frequency p gets long codes

Entropies with varying penalization

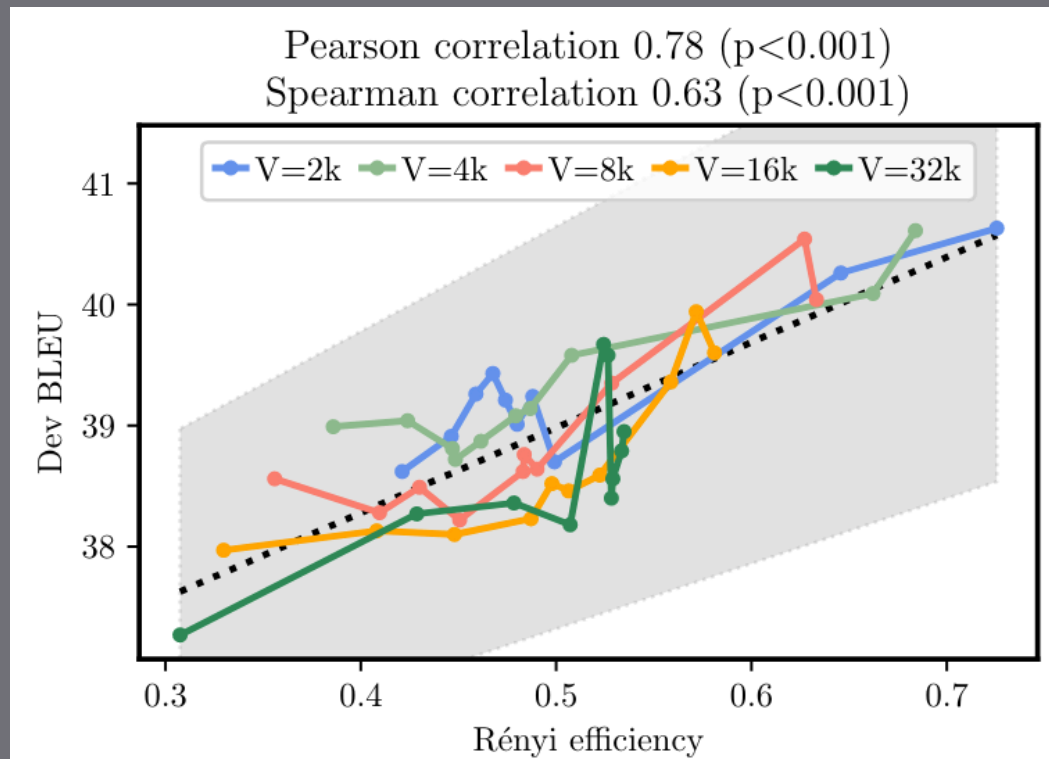
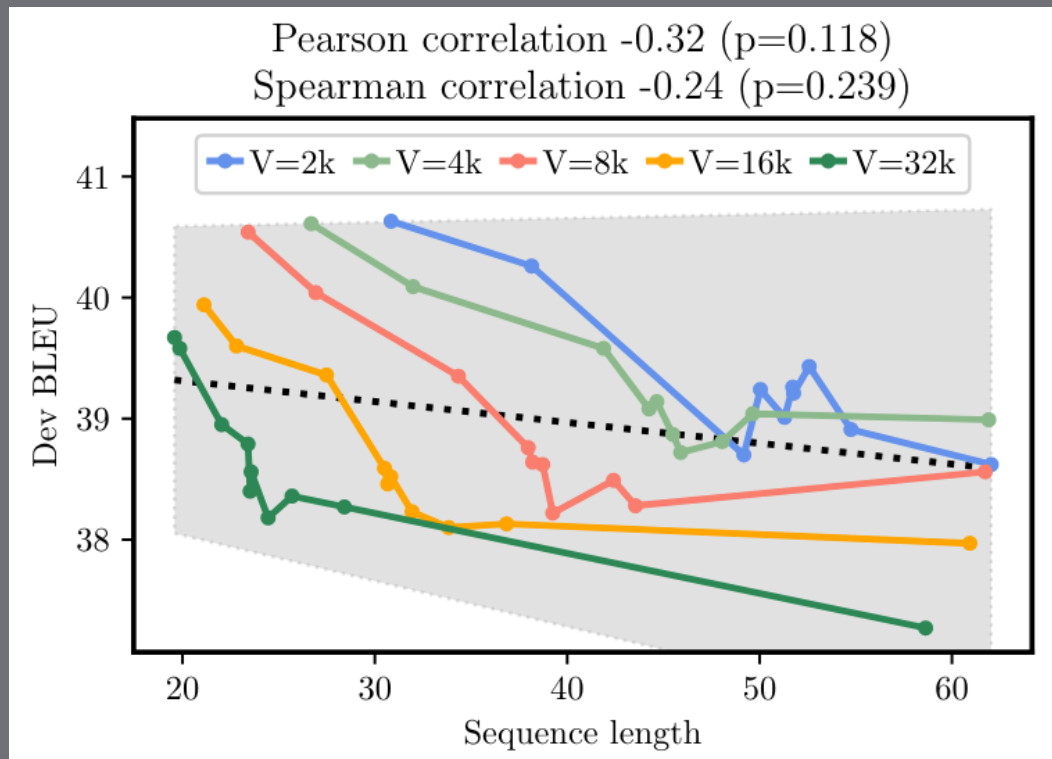
Shannon Entropy: $-\sum p \log p$ \approx expected optimal encoding length
low frequency p gets long codes

Rényi Entropy: $\frac{1}{1-\alpha} \log(\sum p^\alpha)$ \approx expected optimal encoding length
low frequency p gets long codes
too long/short codes (depending on α) are penalized

Choosing good tokenization



Choosing good tokenization



Choosing good tokenization

| Predictor | Pearson | Spearman | ρ^2 |
|------------------|----------------|----------------|----------|
| Sequence len. | -0.32 (=0.118) | -0.24 (=0.239) | 10% |
| Percentile freq. | 0.76 (<0.001) | 0.63 (<0.001) | 58% |
| Entropy | 0.22 (=0.281) | 0.12 (=0.578) | 5% |
| Entropy eff. | 0.56 (=0.004) | 0.38 (=0.006) | 31% |
| Rényi entropy | 0.49 (=0.001) | 0.38 (=0.006) | 24% |
| Rényi eff. | 0.78 (<0.001) | 0.66 (<0.001) | 61% |

tokenization-scorer

CLI

```
$ head -n 1 data1.txt
pick @@ed pick @@l @@ed pick @@les
$ head -n 1 data2.txt
pick @@e @@d pick @@l @@e @@d pick @@l
@@e @@s

$ pip install tokenization-scorer

$ tokenization-scorer -i data1.txt
0.8031528501359657
$ tokenization-scorer -i data2.txt
0.9105681923824472
```


tokenization-scorer

CLI

```
$ head -n 1 data1.txt
pick @@ed pick @@l @@ed pick @@les
$ head -n 1 data2.txt
pick @@e @@d pick @@l @@e @@d pick @@l
@@e @@s

$ pip install tokenization-scorer

$ tokenization-scorer -i data1.txt
0.8031528501359657
$ tokenization-scorer -i data2.txt
0.9105681923824472
```

Python

```
import tokenization_scorer
data1 = "pick @@ed pick @@l @@ed pick @@les"
tokenization_scorer.score(
    text1, metric="renyi", power=2.5
)
> 0.8031528501359657

data2 = "pick @@e @@d pick @@l @@e @@d pick @@l @@e @@s"
tokenization_scorer.score(
    text2, metric="renyi", power=2.5
)
> 0.9105681923824472
```

tokenization-scorer

CLI

```
$ head -n 1 data1.txt
pick @@ed pick @@l @@ed pick @@les
$ head -n 1 data2.txt
pick @@e @@d pick @@l @@e @@d pick @@l
@@e @@s

$ pip install tokenization-scorer

$ tokenization-scorer -i data1.txt
0.8031528501359657
$ tokenization-scorer -i data2.txt
0.9105681923824472
```

Python

```
import tokenization_scorer
data1 = "pick @@ed pick @@l @@ed pick @@les"
tokenization_scorer.score(
    text1, metric="renyi", power=2.5
)
> 0.8031528501359657

data2 = "pick @@e @@d pick @@l @@e @@d pick @@l @@e @@s"
tokenization_scorer.score(
    text2, metric="renyi", power=2.5
)
> 0.9105681923824472
```

tokenization-scorer

demo

A few questions about Tokenization and the Noiseless Channel



Marco <cognetta.marco@gmail.com>

25 Jul 2023, 04:17



to vzouhar 

Hi Vilém,

A few questions about Tokenization and the Noiseless Channel



Marco <cognetta.marco@gmail.com>

25 Jul 2023, 04:17




to vzouhar 

Hi Vilém,

I think your paper is incorrect

A few questions about Tokenization and the Noiseless Channel



Marco <cognetta.marco@gmail.com>
to vzouhar 

25 Jul 2023, 04:17



Hi Vilém,

I think your paper is incorrect

and I can prove it

A few questions about Tokenization and the Noiseless Channel ➤



Marco <cognetta.marco@gmail.com>

25 Jul 2023, 04:17



to vzouhar ▼

Hi Vilém,

I think your paper is incorrect

and I can prove it

Best, Marco

(untrue dramatic reenactment)

Counterexample?

Hypothesis: higher efficiency \rightarrow higher performance

Counterexample?

Hypothesis: higher efficiency \rightarrow higher performance

Counterexample: increase efficiency *but* lower performance

Counterexample?

Hypothesis: higher efficiency \rightarrow higher performance

Counterexample: increase efficiency *but* lower performance

```
she pick @@ed pick @@l @@ed pick @@les  
he pick @@ed pick @@l @@ed pick @@les  
they pick @@ed pick @@l @@ed pick @@les
```

\rightarrow

duplicate `pick` into `pick{1}` and `pick{2}`
 \rightarrow (swap randomly)

```
she pick{1} @@ed pick{1} @@l @@ed pick{1} @@les  
he pick{2} @@ed pick{1} @@l @@ed pick{2} @@les  
they pick{2} @@ed pick{2} @@l @@ed pick{1} @@les
```

Counterexample?

- the: 5%
- a: 4%
- -es: 1%
- -ing: 1%
- ...

Counterexample?

- the: 5%
 - a: 4%
 - -es: 1%
 - -ing: 1%
 - ...
-
- the{1}: 1%
 - the{2}: 1%
 - the{3}: 1%
 - the{4}: 1%
 - the{5}: 1%
 - a: 4%
 - -es: 1%
 - -ing: 1%
 - ...

Counterexample?

- the: 5%
- a: 4%
- -es: 1%
- -ing: 1%
- ...

→

- the{1}: 1%
- the{2}: 1%
- the{3}: 1%
- the{4}: 1%
- the{5}: 1%
- a: 4%
- -es: 1%
- -ing: 1%
- ...

→

increase entropy (vocabulary only in log)

→

$$\frac{H(p)}{\log V} \text{ goes up}$$

→

just noise, so lower performance

Counterexample?

| | Original BPE | Duplicated BPE |
|------------|--------------|----------------|
| Efficiency | 40% | 50% |

Counterexample?

| | Original BPE | Duplicated BPE |
|------------|--------------|----------------|
| Efficiency | 40% | 50% |
| BLEU | 33.59 | 32.49 |

Counterexample?

| | Original BPE | Duplicated BPE |
|------------|--------------|----------------|
| Efficiency | 40% | 50% |
| BLEU | 33.59 | 32.49 |

All is lost?

Counterexample?

| | Original BPE | Duplicated BPE |
|------------|--------------|----------------|
| Efficiency | 40% | 50% |
| BLEU | 33.59 | 32.49 |

All is lost?

- non-natural adversarial examples

Counterexample?

| | Original BPE | Duplicated BPE |
|------------|--------------|----------------|
| Efficiency | 40% | 50% |
| BLEU | 33.59 | 32.49 |

All is lost?

- non-natural adversarial examples
- invitation to find a better metric

Stochastic Tokenization

· So far `tokenize(pickles) = <pick les>`

(deterministic)

Stochastic Tokenization

- So far $\text{tokenize}(\text{pickles}) = \langle \text{pick le s} \rangle$ (deterministic)
- Vocabulary $V = \Sigma \cup \{\text{pi}, \text{ck}, \text{pick}, \text{le}\}$
- $\{\text{pick le s}, \text{pi ck les}, \text{pick le s}, \text{p i c k les}, \dots\} \subset \Sigma^*$

Stochastic Tokenization

- So far $\text{tokenize}(\text{pickles}) = \langle \text{pick le s} \rangle$ (deterministic)
- Vocabulary $V = \Sigma \cup \{\text{pi, ck, pick, le}\}$
- $\{\text{pick le s, pi ck les, pick le s, p i c k les, ...}\} \subset \Sigma^*$
- ~~$\text{tokenize}(\text{pickles, pickles, pickles, pickles, pickles})$
 $= \langle \text{pick le s, pick le s, pick le s, pick le s, pick le s} \rangle$~~
- $\text{tokenize}'(\text{pickles, pickles, pickles, pickles, pickles})$
 $= \langle \text{pi ck les, pick le s, pick les, p i ck le s, pi ck les} \rangle$

Stochastic Tokenization

- So far $\text{tokenize}(\text{pickles}) = \langle \text{pick le s} \rangle$ (deterministic)
- Vocabulary $V = \Sigma \cup \{\text{pi}, \text{ck}, \text{pick}, \text{le}\}$
- $\{\text{pick le s}, \text{pi ck le s}, \text{pick le s}, \text{p i c k le s}, \dots\} \subset \Sigma^*$
- ~~$\text{tokenize}(\text{pickles}, \text{pickles}, \text{pickles}, \text{pickles}, \text{pickles})$
 $= \langle \text{pick le s}, \text{pick le s}, \text{pick le s}, \text{pick le s}, \text{pick le s} \rangle$~~
- $\text{tokenize}'(\text{pickles}, \text{pickles}, \text{pickles}, \text{pickles}, \text{pickles})$
 $= \langle \text{pi ck le s}, \text{pick le s}, \text{pick le s}, \text{p i ck le s}, \text{pi ck le s} \rangle$
- Deterministic: $\text{tokenize} : \Sigma^* \rightarrow V^*$
- Stochastic: $\text{tokenize} : \Sigma^* \rightarrow \underbrace{[X \rightarrow V^*]}_{\text{distribution}}$ • Increases robustness

Stochastic Tokenization

· tokenize(pickles)

= {80% : pick les, 10% : pi ck les, 5% : pick le s, ...}

Stochastic Tokenization

- `tokenize(pickles)`
= {80% : pick les, 10% : pi ck les, 5% : pick le s, ...}
- BPE Dropout (modified from Provilkov et al., 2020)

Input: \bar{x} (word), merges $\bar{\mu}$

Output: \bar{x} (tokenized word)

1. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
2. **while** $\bar{\mu}' \neq \emptyset$ **do**
3. $(a, b) \leftarrow \bar{\mu}'_1$
4. $\bar{x} \leftarrow \text{ApplyMerge}(a\ b \rightarrow ab, \bar{x})$
5. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
6. **end while**
7. **return** \bar{x}

Stochastic Tokenization: BPE-Dropout

Input: \bar{x} (word), merges $\bar{\mu}$

Output: \bar{x} (tokenized word)

1. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
2. **while** $\bar{\mu}' \neq \emptyset$ **do**
3. $(a, b) \leftarrow \bar{\mu}'_1$
4. $\bar{x} \leftarrow \text{ApplyMerge}(a \ b \rightarrow ab, \bar{x})$
5. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
6. **end while**
7. **return** \bar{x}

Stochastic Tokenization: BPE-Dropout

Input: \bar{x} (word), merges $\bar{\mu}$

Output: \bar{x} (tokenized word)

1. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
2. **while** $\bar{\mu}' \neq \emptyset$ **do**
3. $(a, b) \leftarrow \bar{\mu}'_1$
4. $\bar{x} \leftarrow \text{ApplyMerge}(a\ b \rightarrow ab, \bar{x})$
5. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
6. **end while**
7. **return** \bar{x}

~~$\bar{\mu} = \langle p\ i \rightarrow pi, \ c\ k \rightarrow ck, \ pi\ ck \rightarrow pick, \ pick\ l \rightarrow pickl \rangle$~~

$\bar{\mu} = \langle p\ i \rightarrow pi, \ c\ k \rightarrow ck, \ pick\ l \rightarrow pickl \rangle$

Stochastic Tokenization: BPE-Dropout

Input: \bar{x} (word), merges $\bar{\mu}$

Output: \bar{x} (tokenized word)

1. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
2. **while** $\bar{\mu}' \neq \emptyset$ **do**
3. $(a, b) \leftarrow \bar{\mu}'_1$
4. $\bar{x} \leftarrow \text{ApplyMerge}(a \ b \rightarrow ab, \bar{x})$
5. $\bar{\mu}' \leftarrow \langle (\bar{x}_i, \bar{x}_{i+1}) \mid (\bar{x}_i, \bar{x}_{i+1}) \in \bar{\mu} \wedge \text{Rand}() < p \rangle$
6. **end while**
7. **return** \bar{x}

~~$\bar{\mu} = \langle p \ i \rightarrow pi, \ c \ k \rightarrow ck, \ pi \ ck \rightarrow pick, \ pick \ l \rightarrow pickl \rangle$~~

$\bar{\mu} = \langle p \ i \rightarrow pi, \ c \ k \rightarrow ck, \ pick \ l \rightarrow pickl \rangle$

~~$T(\text{pickles}) = pickl \ e \ s$~~ $T(\text{pickles}) = pi \ ck \ l \ e \ s$

Stochastic Tokenization is not Uniform

BPE-Dropout $p = 0.1$

| | |
|------------------|---------|
| to ken ization | 97.77% |
| to ke n ization | 1.89% |
| to k en ization | 0.25% |
| to ken iz ation | 0.04% |
| t oken ization | 0.03% |
| to k en iz ation | 0.01% |
| to ke n iz ation | 0.01% |
| to ken i z ation | < 0.01% |

Stochastic Tokenization can be Uniform

Input: \bar{x} (word), tokenizer T

Output: \bar{x} (tokenized word)

1. **if** $\text{Rand}() < p$ **then**
2. **return** $\text{RandomTokenization}(\bar{x})$
3. **else**
4. **return** $T(\bar{x})$
5. **end if**

Stochastic Tokenization can be Uniform

Input: \bar{x} (word), tokenizer T

Output: \bar{x} (tokenized word)

1. **if** $\text{Rand}() < p$ **then**
2. **return** $\text{RandomTokenization}(\bar{x})$
3. **else**
4. **return** $T(\bar{x})$
5. **end if**

How to get $\text{RandomTokenization}(\bar{x})$?

Modelling All Possible Tokenizations

- Idea 1: run BPE-Dropout for each word 10^7 times and store all unique tokenizations.

Modelling All Possible Tokenizations

- Idea 1: run BPE-Dropout for each word 10^7 times and store all unique tokenizations.
 - Very slow.

Modelling All Possible Tokenizations

- Idea 1: run BPE-Dropout for each word 10^7 times and store all unique tokenizations.
 - Very slow.
- Idea 2: construct the set for each word apriori.

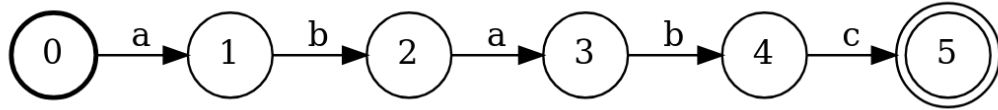
Modelling All Possible Tokenizations

- Idea 1: run BPE-Dropout for each word 10^7 times and store all unique tokenizations.
 - Very slow.
- Idea 2: construct the set for each word apriori.
 - Very big; unclear how.
 - For word of length $|\bar{x}|$ the tight upper bound is $2^{|\bar{x}|-1}$ (yes/no space).

Modelling All Possible Tokenizations

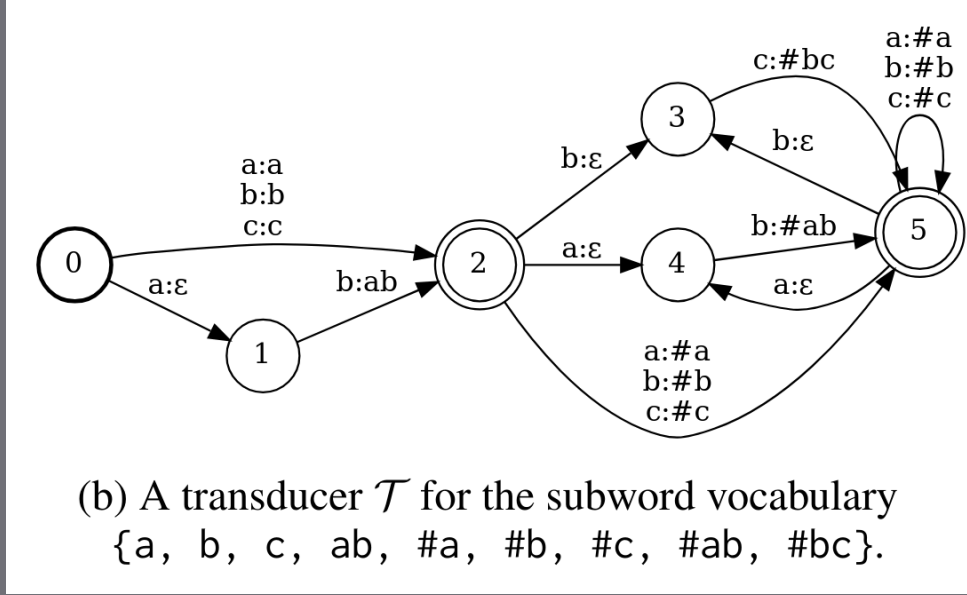
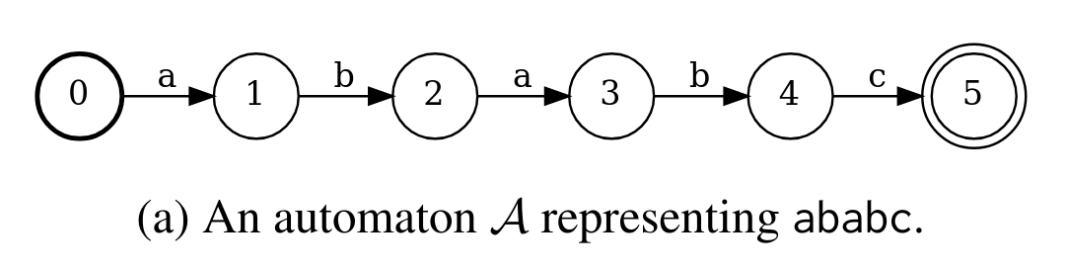
- Idea 1: run BPE-Dropout for each word 10^7 times and store all unique tokenizations.
 - Very slow.
- Idea 2: construct the set for each word apriori.
 - Very big; unclear how.
 - For word of length $|\bar{x}|$ the tight upper bound is $2^{|\bar{x}|-1}$ (yes/no space).
- Idea 3: function that can yield random tokenization.

Modelling All Possible Tokenizations

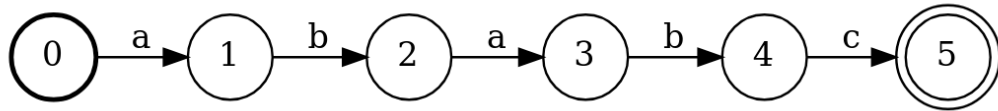


(a) An automaton \mathcal{A} representing $ababc$.

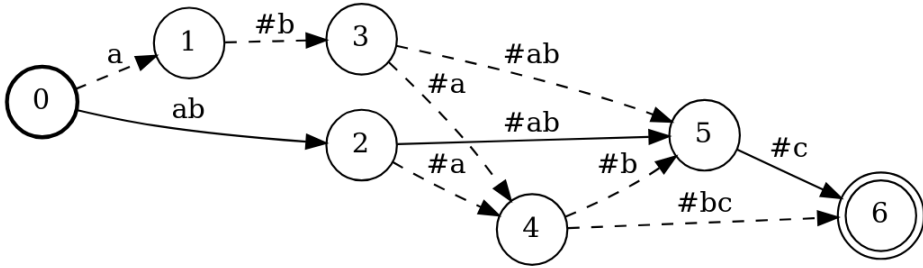
Modelling All Possible Tokenizations



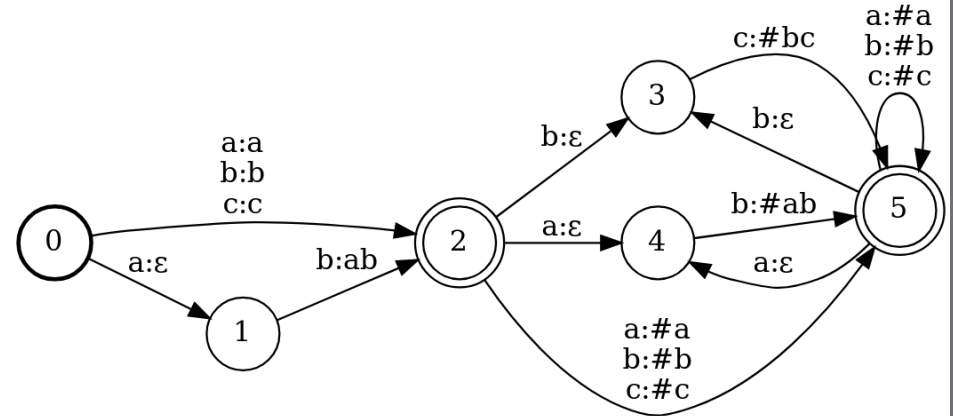
Modelling All Possible Tokenizations



(a) An automaton \mathcal{A} representing ababc.

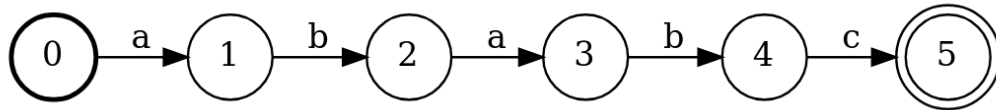


(c) A lattice, $\mathcal{A} \circ \mathcal{T}$, of all possible tokenizations of ababc.

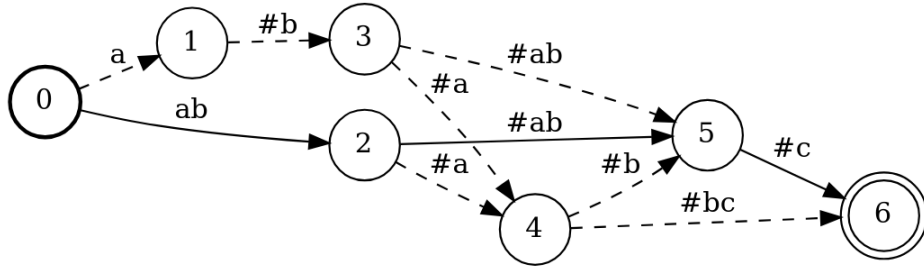


(b) A transducer \mathcal{T} for the subword vocabulary $\{a, b, c, ab, \#a, \#b, \#c, \#ab, \#bc\}$.

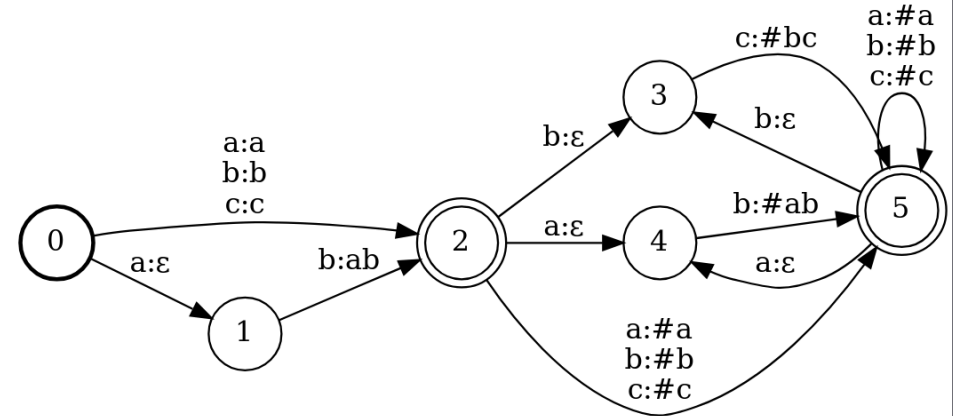
Modelling All Possible Tokenizations



(a) An automaton \mathcal{A} representing ababc.



(c) A lattice, $\mathcal{A} \circ \mathcal{T}$, of all possible tokenizations of ababc.



(b) A transducer \mathcal{T} for the subword vocabulary $\{a, b, c, ab, \#a, \#b, \#c, \#ab, \#bc\}$.

- Walk through lattice is tokenization
- Can be sampled randomly in linear time

Uniform Sampling Increases Diversity

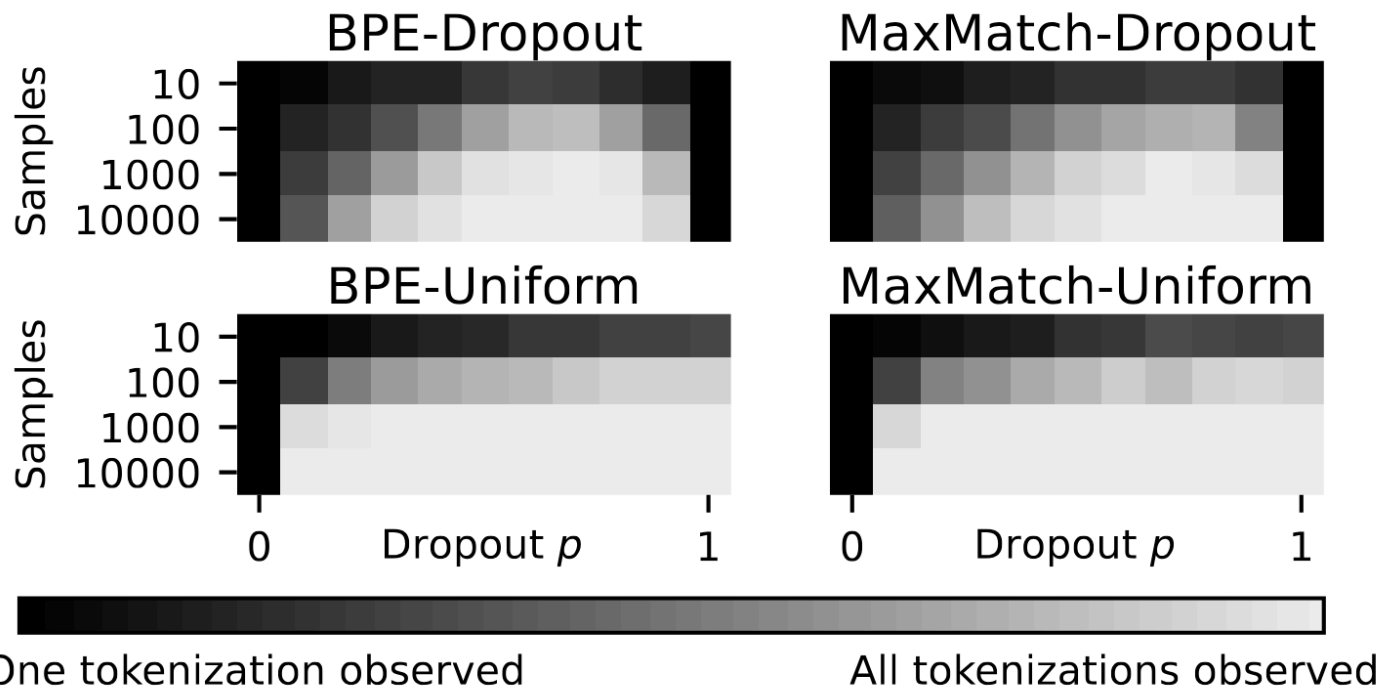


Figure 2: The number of unique, observed tokenizations of a word with N samples and dropout p .

Uniform Sampling Increases Efficiency (and Performance?)

| Tokenizer | Efficiency | BLEU | CHRF | COMET |
|---------------------------------|------------|--------------|--------------|--------------|
| BPE | 0.4524 | 23.56 | 53.20 | 81.03 |
| BPE + Dropout ($p=0.1$) | 0.4614 | 23.98 | 53.70 | 81.90 |
| BPE + Uniform ($p=0.1$) | 0.4594 | 23.83 | 53.67 | 82.00 |
| BPE + Uniform ($p=0.25$) | 0.4647 | 24.13 | 53.73 | 82.20 |
| MaxMatch | 0.4476 | 23.52 | 53.23 | 81.17 |
| MaxMatch + Dropout ($p=0.3$) | 0.4578 | 23.95 | 53.70 | 81.98 |
| MaxMatch + Uniform ($p=0.1$) | 0.4528 | 24.32 | 53.90 | 82.11 |
| MaxMatch + Uniform ($p=0.25$) | 0.4563 | 24.10 | 53.87 | 82.06 |
| Unigram ($\alpha=1$) | 0.4338 | 23.68 | 53.37 | 81.28 |
| Unigram ($\alpha=0.3$) | 0.4284 | 24.17 | 53.87 | 82.00 |

Miscellaneous Tricks

```
$ from transformers import AutoTokenizer  
$ vocab = AutoTokenizer.from_pretrained("gpt2").vocab  
$ [x for x in vocab if len(x) > 50]
```


Miscellaneous Tricks

```
$ from transformers import AutoTokenizer  
$ tokenizer = AutoTokenizer.from_pretrained("gpt2")
```

Miscellaneous Tricks

```
$ from transformers import AutoTokenizer  
$ tokenizer = AutoTokenizer.from_pretrained("gpt2")
```

```
$ tokenizer.tokenize("1945")
```

```
1945
```

```
$ tokenizer.tokenize("1969")
```

```
1969
```

```
$ tokenizer.tokenize("1961")
```

```
19 61
```

Miscellaneous Tricks

```
$ from transformers import AutoTokenizer
```

```
$ tokenizer = AutoTokenizer.from_pretrained("gpt2")
```

```
$ tokenizer.tokenize("66666666")  
66666666
```

```
$ tokenizer.tokenize("1945")  
1945
```

```
$ tokenizer.tokenize("6666666")  
6666 666
```

```
$ tokenizer.tokenize("1969")  
1969
```

```
$ tokenizer.tokenize("00200000")  
00200000
```

```
$ tokenizer.tokenize("1961")  
19 61
```

```
$ tokenizer.tokenize("00100000")  
00 100 000
```

Miscellaneous Tricks

Don't include numbers in subwords.

```
$ from transformers import AutoTokenizer
```

```
$ tokenizer = AutoTokenizer.from_pretrained("gpt2")
```

```
$ tokenizer.tokenize("66666666")
```

```
66666666
```

```
$ tokenizer.tokenize("1945")
```

```
1945
```

```
$ tokenizer.tokenize("6666666")
```

```
6666 666
```

```
$ tokenizer.tokenize("1969")
```

```
1969
```

```
$ tokenizer.tokenize("00200000")
```

```
00200000
```

```
$ tokenizer.tokenize("1961")
```

```
19 61
```

```
$ tokenizer.tokenize("00100000")
```

```
00 100 000
```


Miscellaneous Tricks

- Multilingual LM:
 - 1B English words
 - 10M Swiss German words
- Train BPE on 1010M words.

Miscellaneous Tricks

- Multilingual LM:
 - 1B English words
 - 10M Swiss German words
- Train BPE on 1010M words.

```
$ tokenizer("I see something that you don't.")  
I see something that you don 't .
```

```
$ tokenizer("Ich sehe öpis, das de nöd siehst.")  
I ch se he ö p is , das de n ö d s i e h st .
```

Miscellaneous Tricks

- Multilingual LM:
 - ▶ 1B English words
 - ▶ 10M Swiss German words
- Train BPE on 1010M words.

\$ tokenizer("I see something
I see something that you do

\$ tokenizer("Ich sehe öpis
I ch se he ö p is , das de

Why Romanization?

- Significant part of web data is code mixed and Romanized
- Major LLMs have seen such data
 - Chance to transfer from English
- Byte based BPE oversegments

But, Romanization doesn't

| Language | N | R | R-IndicNLP |
|-----------|-------|-------------|------------|
| Gujarati | 18.44 | 3.39 | 4.16 |
| Hindi | 7.36 | 2.98 | 3.98 |
| Malayalam | 12.85 | 5.04 | 5.56 |
| Marathi | 8.91 | 3.64 | 4.84 |
| Tamil | 12.11 | 4.89 | 5.35 |

Participants 2 Chat React Share Host tools Apps Pause/stop recording More

57

Subwords per token(?)

Miscellaneous Tricks

Supersample low-frequency data
(if that's what you want).

- Multilingual LM:
 - ▶ 1B English words
 - ▶ 10M Swiss German words
- Train BPE on 1010M words.

\$ tokenizer("I see something
I see something that you do

\$ tokenizer("Ich sehe öpis
I ch se he ö p is , das de

Why Romanization?

- Significant part of web data is code mixed and Romanized
- Major LLMs have seen such data
 - Chance to transfer from English
- Byte based BPE oversegments

But, Romanization doesn't

| Language | N | R | R-IndicNLP |
|-----------|-------|-------------|------------|
| Gujarati | 18.44 | 3.39 | 4.16 |
| Hindi | 7.36 | 2.98 | 3.98 |
| Malayalam | 12.85 | 5.04 | 5.56 |
| Marathi | 8.91 | 3.64 | 4.84 |
| Tamil | 12.11 | 4.89 | 5.35 |

Participants 2 Chat React Share Host tools Apps Pause/stop recording More

57

Subwords per token(?)

This code is wrong

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

This code is wrong

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

This code is wrong

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

$\bar{x} = a a a a a b a b a b a$

This code is wrong

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

$\bar{x} = a a a a a b a b a b a$

$\operatorname{Freq}(x, (\mu', \mu'')) = \langle a a \rightarrow 4, b a \rightarrow 3, a b \rightarrow 1, \rangle$

This code is wrong

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

$\bar{x} = a a a a a b a b a b a$

$\operatorname{Freq}(x, (\mu', \mu'')) = \langle a a \rightarrow 4, b a \rightarrow 3, a b \rightarrow 1, \rangle$

$\operatorname{ApplyMerge}(a a \rightarrow aa, \bar{x}) = aa aa a b a b a b a$

This code is wrong

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

$\bar{x} = a a a a a b a b a b a$

$\operatorname{Freq}(x, (\mu', \mu'')) = \langle a a \rightarrow 4, b a \rightarrow 3, a b \rightarrow 1, \rangle$

$\operatorname{ApplyMerge}(a a \rightarrow aa, \bar{x}) = aa aa a b a b a b a$

$|\bar{x}| - |\operatorname{ApplyMerge}(a a \rightarrow aa, \bar{x})| = 2 \neq 4$

This code is wrong

Fix Freq computation.

Input: \bar{x} (text to compress)

Output: \bar{x} (compressed text) $\bar{\mu}$ (compression merges)

1. $\bar{\mu} \leftarrow \langle \rangle$
2. **for** $i \in \{1, \dots, M\}$ **do**
3. $\mu \leftarrow \operatorname{argmax}_{\mu', \mu'' \in \bar{x}} \operatorname{Freq}(x, (\mu', \mu''))$
4. $\bar{x} \leftarrow \operatorname{ApplyMerge}(\mu, \bar{x})$
5. $\bar{\mu} \leftarrow \bar{\mu} \circ \langle \mu \rangle$
6. **end for**

$\bar{x} = a a a a a b a b a b a$

$\operatorname{Freq}(x, (\mu', \mu'')) = \langle a a \rightarrow 4, b a \rightarrow 3, a b \rightarrow 1, \rangle$

$\operatorname{ApplyMerge}(a a \rightarrow aa, \bar{x}) = aa aa a b a b a b a$

$|\bar{x}| - |\operatorname{ApplyMerge}(a a \rightarrow aa, \bar{x})| = 2 \neq 4$

A Formal Perspective on Byte-Pair Encoding

ACL 2023

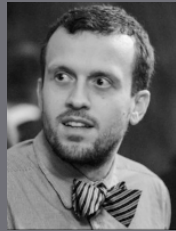
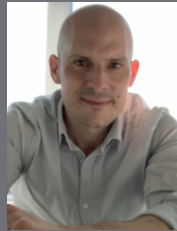
Tokenization and the Noiseless Channel

ACL 2023

Two Counterexamples to Tokenization and the Noiseless Channel Cognetta et al.

LREC-COLING 2024

Distributional Properties of Subword Tokenization Cognetta et al. preprint 2024



Marco
Cognetta

Mrinmaya
Sachan

Ryan
Cotterell

Juan
Gastaldi

Clara
Meister

Tim
Vieira

Sangwhan
Moon

Naoaki
Okazaki

Leo
Du

..and
others

- BPE is compression
- BPE is greedy & not optimal
- BPE is approximately optimal

- BPE is compression
- BPE is greedy & not optimal
- BPE is approximately optimal
- Think about distributions
in stochastic tokenizations

- BPE is compression
- BPE is greedy & not optimal
- BPE is approximately optimal
- Think about distributions in stochastic tokenizations
- High-entropy tokenization is good
- `pip install tokenization-scorer`
`tokenization-scorer -i data1.txt`
`tokenization-scorer -i data2.txt`

- BPE is compression
- BPE is greedy & not optimal
- BPE is approximately optimal
- Think about distributions in stochastic tokenizations
- High-entropy tokenization is good
- `pip install tokenization-scoring`
`tokenization-scoring -i data1.txt`
`tokenization-scoring -i data2.txt`

What makes some tokenizations better than others?

Future

- BPE is compression
- BPE is greedy & not optimal
- BPE is approximately optimal
- Think about distributions in stochastic tokenizations
- High-entropy tokenization is good
- `pip install tokenization-scorer`
`tokenization-scorer -i data1.txt`
`tokenization-scorer -i data2.txt`

What makes some tokenizations better than others?

Multilingual tokenization?

Future

- BPE is compression
- BPE is greedy & not optimal
- BPE is approximately optimal
- Think about distributions in stochastic tokenizations
- High-entropy tokenization is good
- `pip install tokenization-scorer`
`tokenization-scorer -i data1.txt`
`tokenization-scorer -i data2.txt`

What makes some tokenizations better than others?

Multilingual tokenization?

Future

Thank you

Rényi Entropy

$$H_{\alpha}(X) = \frac{1}{1 - \alpha} \log \left(\sum_{i=1}^n p_i^{\alpha} \right)$$