

# بابه تی 4: له تکردنی تۆکن خولی سه رهیللی زمانه وانیی کۆمپیوتهری

سینا ئه حمه دی

[sinaahmadi.github.io/KurdishCL](https://sinaahmadi.github.io/KurdishCL)

1 له تکردنی تۆکن

2 له تکردنی تۆکن بۆ کوردی



# لهتکردنی تۆکن



# چۆن دەق لە سەر کۆمپیوتەر نیشان دەدەین؟

لە زۆر بابەتی پڕۆسەس کردنی زماندا پیۆیستمان بە داشکاندنی دەق (segmentation) بۆ دەرھێنانی ییکھاتە زمانەوانییەکان یان ئاسانکاریی تەکنیکی ھەیە.

- ھەنگاوێکی پێش-پڕۆسەس کردنی زۆرینە ی ئەرکەکانی پڕۆسەس کردنی زمانە
- لەتەکان یەکە ی جیاوازی دەقن (پیت، وشە، ھتد)
- ئاستی زبر وەک داشکاندنی دەق بە سەر پاراگراف و رستەکاندا

مافەکانی مەژۆف کۆمەڵێک بنچینە یان رێسای رەوشتین کە لەسەر پێوانەی ئاکاری مەژۆفەکان ساز کراون و بەگۆڕە ی یاسای ناوخرۆیی و نۆدەولەتی پارێزگاریان لێ دمرکێت. بەگشتی دەتوانین ئەم بنەمایانە بەم شێوەیە ناو بنین کە مافی رەوای ھەموو مەژۆفێک و شیاوی گۆڕین و دەستکاریی نین و بەیەکسانی و بەھۆی مەژۆفبوونی کەسەکەو پێی رەوا ببنراو. بێ گۆیدانە تەمەن، نەتەو، ئایین، زمان، رەنگ و، رەگەز ھەمووان وەکو مەژۆف ھاوتان. ئەم مافانە لە ھەموو شوێنێک و لە ھەموو کاتێکدا جێبەجێ دەبن بەسەر مەژۆفەکاندا، لە ھەر کۆیەک بن لە جیھاندا، بێ جیاوازی. لە سایە ی سەرورە ی یاسا و سۆزی مەژۆفانەو چاومەوان دمرکێت کە ھەموو پێرە ی رێزگرتنی مافەکانی مەژۆفی ھەبێت بۆ مەژۆفەکانی تر، ھێچ کات ئەو مافانە لە مەژۆف ناترازی، مافی مەژۆف بۆ ھەموو کەسێک مافیکی رەسەنە و پێو ی لکاو، مەگەر لە دۆخیکی یاسایی تایبەتی و بە ھۆکاری سەپاندنی لە رەوتیکی یاسایدا.

وەرگیراو لە وتاری «مافەکانی مەژۆف»، ویکیپیدیای کوردی



- له تکردن یان که رتکردنی دهق چه شنیکه له داشکاندنی وشه و بهو له تانه دهلیین «توکن» (token)
- بۆ یه کهم جار له سالی 1992 دا ئاماژه بهو له تانه کراوه [1]
- بنه مای زۆر ئهرکی پرۆسه سکردنی زمان له سهر توکنه کانه

## توکن

- «پریز بهندی یهک له دوا ی یهک و نابه تالی گرافیم یان فونیمه کان له ده قدا» [2]
- «له تکردنی دهق به سهر توکندا بریتییه له دۆزینه وهی سنووری وشه کان» [3]
- «توکن خه یالیکی پروپووچی ئەنداز یاریکی کو مپیوتره بۆ ئاسانکردنی کاری خو ی» [4]

# له تکر دنی توکن: سنووری وشه

پیناسه ی توکن به پیی دۆزینه وهی سنووری وشه

Dr. O'Neil presents state-of-the-art results. •

• « دواکهوتنی شیوازهکانی به رهه مهینان »

• دواکهوتنی شیوازهکانی به رهه مهینان

• دوا کهوتنی شیوازهکانی به رهه مهینان

• دواکهوتن ی شیواز هکان ی به رهه مهینان

• دوا کهوتن ی شیواز هکان ی به رهه مهینان

• سنووری وشه له Scriptio continua هکاندا

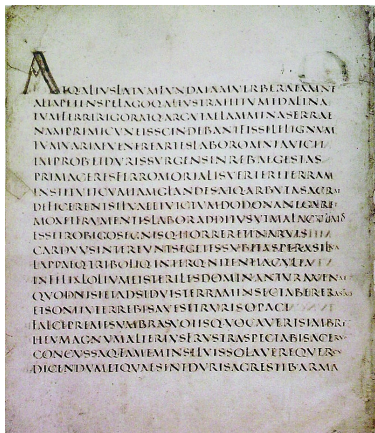
• سنووری وشه له زمانهکانی دیکه دا

• 我是中国人

• つ目の入力

• Sprachwissenschaft

پیناسه کردنی توکن ئه رکیکی ئاسان نییه!



# لهتکردنی تۆکن: روانگهی رینووس

روانگهی رینووس: پشت به شیوازی نووسینی دهقه که بیهسته

• داشکاندن به پپی هیماکانی خالبهندی وهک «!?!:;»,» بۆشایی، هیلی نوێ و ته ب

• دواکهوتنی شیوازهکانی بهرهمهینان

• باشی +: ئاسانه

• خراپی -

• داهینانی ههموو یاساکان زهمهته

• پیویستی به زانیاری له سهر زمانهکان و رینووسهکانیان ههیه

• ناتوانی چارهسهری سنووره دیارینهکراوهکانی وشه بکا

• ههلهی رینووسی سازکردنی وشه دانه که تیگ دهدا



# له تکر دنی تۆکن: روانگهی زمانه وانی

روانگهی زمانه وانی: پۆلی زمانه وانی پیکهاته کان وهک پۆلی وشه کان له رسته سازیدا

- له تکر دن به پیی پیکهاته کان له وشه دا، وهک مۆرفیمه کان
- زۆر جار، تۆکنه کان به پیی سه رچاوه یکی دیکه ده رده هیئدرین، وهک فرههنگ

دواکه وتن ی شیواز هکان ی به رهه مهینان .

دوا که وتن ی شیواز هکان ی به رهه مهینان .

• باشی +

- ئاسانه به لام نه به قهت روانگهی رینووس
- به سووده لهو ئه رکانه ی که پیوستان به بۆچوونیکی زمانه وانانه ههیه وهک لیکدانه وهی رسته

• خراپی -

- پیوستانی به زانیاری زمانه وانی له سه ر زمانه کان و رینووسه کانیا ن ههیه
- یاسا زمانه وانیه کان ده توانن ناته واو بن
- هه له ی رینووسی سازکردنی وشه دانه که تیگ دهدا
- ئه گه ر وشه یه ک له فرهه نگدا نه بی و یاسا کان کاریان لی نه کا؟ ← وشه ی نه ناسراو



# له تکر دنی تۆکن: وشه نه ناسراوه کان

زمان پره له وشه و هه موویان له فرههنگدا ده رناکهون، وهک «تورته مالی»، «کوڤید» یان «لاتک»

- ئه و تۆکنانه ی له وشه دان یان فرههنگی زمانه که دا نه هاتوون به  $\langle unk \rangle$  دیاری ده کرین

- زۆر وشه ی ده گمه ن یان نیوی که س به وشه ی نه ناسراو ده رده چن

دوا کهوتن ی شیواز هکان ی  $\langle unk \rangle$  هینان .

- نه ناسینه وه ی تۆکنیک ده بیته هو ی له ده ست چوونی زانیاری

- چۆن ئه م کیشه یه چاره سه ر بکه ین به بی ئه وه ی پیویستمان به زانیاری له سه ر هه موو زمانه کدا هه بی؟



# له تکر دنی توکن: روانگهی پیدراوه

روانگهی پیدراوه: چ پیکهاته گه لیک باوتره، با به توکنیان دانیین

- ئەم بۆچوونه چاره سه ریكه بۆ كێشه كانی بۆچوونه پیشووه كان

- پشت به فراوانیی ئەو پیتانه ده به ستی كه وێكرا دهرده كه ون و پێیان ده لێن وشه له (subword)

- Byte pair encoding (BPE): ناسراوترین و باوترین ئەلگوریتمه (algorithm) كه پشت به م

بۆچوونه وه ده به ستی [5]

- «پالنه ری سه ره کیی ئەم بۆچوونه ئەوهیه كه وه رگی پانی هه ندیک وشه پروونه، به و مانایه ی كه ده کریت

له لایه ن وه رگی پانی لی ها تو وه وه وه ربگی پ درین ته نانه ت ئەگه ر بۆ ئەو نویش بن، به پشت به ستن به وه رگی پانی

یه كه كانی ژیر وشه ی ناسرا وه ك مۆرفیم یان فۆنیم.»

دوا كه وتن ی شیواز هكان ی #به ر هه م هیان .

Neural Machine Translation of Rare Words with Subword Units

Rico Sennrich and Barry Haddow and Alexandra Birch

School of Informatics, University of Edinburgh

{rico.sennrich,a.birch}@ed.ac.uk, bhaddow@inf.ed.ac.uk



## ههنگاهه کانی راهینان:

- 1 سازکردنی وشه دان: دهست پیده کات به دابه شکردنی وشه کان بو تاکه پیته کان (هه ریه کیکیان هیمایه که له وشه دانی کو تاییدا)
- 2 دۆزینه وهی جووته هیمای فراوان له وشه دانی ئیستادا. بهم ئهرکه دهلین یه که خستن (merge)  
بو نموونه ئه گهر پیتی «د» به فراوانییه کی زۆره وه له گهل «ه»  
دهرکه وئی، «ده» وه کوو یه که هیما له وشه دانه که زیاد ده کا  
زیاد کردنی هیما یه که خراوه کان له وشه دانه که
- 3 دووباره کردنه وهی ههنگاهه کانی (2) و (3) هه تا به ژماره ی دیاریکراوی یه که خران ده گهین یان هیچ  
پیکهاته یه کی نویی هیماکان به فراوانی پیویسته وه نه میئیت

# لهتکردنی تۆکن: ریکاری BPE - نمونه

low low low low low lowest lowest newer newer  
newer newer newer newer wider wider wider new new

1 سازکردنی وشه دان: دابه شکردنی وشه کان بۆ تاکه پیته کان. دهقه که به پیی بۆشایی داشکینه و کۆتایی وشه کان به هیمایه کی تایبته دیاری که وهک «</w>» یان «\_»

کۆرپس

'l o w </w>': 5  
'l o w e s t </w>': 2  
'n e w e r </w>': 6  
'w i d e r </w>': 3  
'n e w </w>': 2

وشه دان

d, e, i, l, n, o, r, s, t, w, </w>

# لهتکردنی تۆکن: ریکاری BPE - نموونه

## فراوانی جووتهکان

('e', 'r'): 9  
( 'r', '</w>' ): 9  
( 'w', 'e' ): 8  
( 'n', 'e' ): 8  
( 'e', 'w' ): 8  
( 'l', 'o' ): 7  
( 'o', 'w' ): 7  
( 'w', '</w>' ): 7  
( 'w', 'i' ): 3  
( 'i', 'd' ): 3  
( 'd', 'e' ): 3  
( 'e', 's' ): 2  
( 's', 't' ): 2  
( 't', '</w>' ): 2

## 1 سازکردنی وشه‌دان

کۆرپس  
'l o w </w>': 5  
'l o w e s t </w>': 2  
'n e w e r </w>': 6  
'w i d e r </w>': 3  
'n e w </w>': 2

## وشه‌دان

d, e, i, l, n, o, r, s, t, w, </w>

## 2 یه‌کخستنی جووته فراوانه‌کان

'e', 'r' ده‌بن به یه‌ک وشه‌له  $\Leftarrow$  er

# لهتکردنی تۆکن: ریکاری BPE - نموونه

## فراوانیی جووتهکان

('er', '</w>'): 9  
( 'n', 'e'): 8  
( 'e', 'w'): 8  
( 'l', 'o'): 7  
( 'o', 'w'): 7  
( 'w', '</w>'): 7  
( 'w', 'er'): 6  
( 'w', 'i'): 3  
( 'i', 'd'): 3  
( 'd', 'er'): 3  
( 'w', 'e'): 2  
( 'e', 's'): 2  
( 's', 't'): 2  
( 't', '</w>'): 2

1 سازکردنی وشه دان

2 یهکخستنی جووته فراوانهکان

3 زیادکردنی هیما یهکخراوهکان له وشه دانه که

کۆریس  
'l o w </w>': 5  
'l o w e s t </w>': 2  
'n e w e r </w>': 6  
'w i d e r </w>': 3  
'n e w </w>': 2

وشه دان  
d, e, i, l, n, o, r, s, t, w, </w>, er

4 ههنگاوی 2 و 3 دووپات بکهوه

# لهتکردنی تۆکن: ریکاری BPE - نمونه

ههنگاوی 2 و 3 دووپات بکهوه

فراوانیی جووتهکان

('n', 'e'): 8  
( 'e', 'w'): 8  
( 'l', 'o'): 7  
( 'o', 'w'): 7  
( 'w', '</w>'): 7  
( 'w', 'er</w>'): 6  
( 'w', 'i'): 3  
( 'i', 'd'): 3  
( 'd', 'er</w>'): 3  
( 'w', 'e'): 2  
( 'e', 's'): 2  
( 's', 't'): 2  
( 't', '</w>'): 2

$er</w> \Leftarrow 'er', '</w>'$

کۆرپس

'l o w </w>': 5  
'l o w e s t </w>': 2  
'n e w er</w>': 6  
'w i d er</w>': 3  
'n e w </w>': 2

وشه‌دان

d, e, i, l, n, o, r, s, t,  
w, </w>, er, er</w>

# لهتکردنی تۆکن: ریکاری BPE - نموونه

ههنگاوی 2 و 3 دووپات بکهوه

## فراوانیی جووتهکان

('ne', 'w'): 8  
( 'l', 'o': 7  
( 'o', 'w': 7  
( 'w', '</w>': 7  
( 'w', 'er</w>': 6  
( 'w', 'i': 3  
( 'i', 'd': 3  
( 'd', 'er</w>': 3  
( 'w', 'e': 2  
( 'e', 's': 2  
( 's', 't': 2  
( 't', '</w>': 2

ne  $\Leftarrow$  'n', 'e'

کۆرپس

'l o w </w>': 5  
'l o w e s t </w>': 2  
'ne w er</w>': 6  
'w i d er</w>': 3  
'ne w </w>': 2

وشه‌دان

d, e, i, l, n, o, r, s, t,  
w, </w>, er</w>, er, ne



# له تکر دنی توکن: ریکاری BPE - نمونه

ریکاری BPE له پاش 10 یه کخستن دهگا بهمانه:

## فراوانیی جووتهکان

('wid', 'er</w>'): 3  
( 'low', 'e'): 2  
( 'e', 's'): 2  
( 's', 't'): 2  
( 't', '</w>'): 2  
( 'new', '</w>'): 2

## کورپس

'low</w>': 5  
'low e s t </w>': 2  
'newer</w>': 6  
'wider</w>': 3  
'new </w>': 2

د، e، i، l، n، o، r، s، t، w، </w> وشهدان  
er، er</w>، ne، new، lo،  
low، newer</w>، low</w>، wi، wid



# له تکر دنی توکن: ریکاری BPE

- له کاردا، زۆر بهی مۆدیلەکان و ئامیڕەکان 32 بۆ 64 ههزار وشە له یان له وشە داندا ههیه بۆ نموونه، 50257 GPT2 وشە لهی ههیه
- وشە لهکان دهتوانن مۆرفیمهکان بن وهک «یک» یان وشە باوهکان وهک «ههیه» یان ههر پیکهاتهیهکی باو
- له پاش ئەژماریکی دیاریکراوی یه کخستن له سهر کۆرپسیکی گهوره، وشه دانکه بۆ له تکر دنی توکنهکان به کار دی
- له تکر دنی توکنهکان بۆ دهقیکی نوێ به پێی ئەو وشە لانه دهکری که فیریان بووین کهوا بی، گرینگ نییه دهقه نوییه که چی تیدایه. وشه دانکه تهنیا پشت به کۆرپسه که ده بهستی.
- ئەم ریکاره له زۆر مۆدیل و ئامیڕدا به کار دی، به تایبەت له مۆدیلەکانی زمان (language models)
- BPE تهنیا ریکار له سهر وشه لهکان نییه [6] Unigram و [7] WordPiece



# لهتکردنی تۆکن و دهربرینی فرهوشه

- دهربرینی فرهوشه (multiword expression) پیکهاتوو له چهند بهش که ویکرا واتایهکی جودایان ههیه

• بپاریان دا ← بپار یان دا

• خشکه و پشکه ← خشکه و پشکه

• گوئی له مست ← گوئی له مست

• وهره وهر ← وهره وهر

• bi-can-û-bên ← bi can û bên

- چۆن ئەم دهربرینانه لهت بکرین؟

- بهکارهینانی فرههنگی تایبته به دهربرینه فرهوشهیهکان
- بهکارهینانی پوانگهی پیدراوه به رهچاوکردنی بۆشایی؟



# لهتکردنی تۆکن: چه مکی ئیستا

- مۆدیله نۆرالهکان (neural models) ی زمان سیسته مگه لی ته نیشته به ته نیشته (end-to-end) ن له پرۆسه س کردنی زماندا و جیگه ی پرۆه وه نه ریتییه کانیان گرتووه ته وه.
- چه مکی ئیستای لهتکردنی تۆکن: پشتبهستن به پیدراوه له باتی یاسای زمانهوانی
  - دابهشکردن بۆ یه که ی بچووکترو ناتایپوگرافی وهک وشه لهکان
  - تۆکنهکان له گه ل رینووس یان یه که ی زمانهوانیدا یهک ناگره وه  $\Leftarrow$  تۆکن  $\neq$  وشه
  - لهتکردنی تۆکنهکان به پپی رینووس یان زمانهوانی  $\Leftarrow$  «لهتکردنی پیشهکی» (pre-tokenization)
  - تۆکنه نه ریتییهکان  $\Leftarrow$  «پیش-تۆکن» (pre-token)



# لهتکردنی تۆکن بۆ کوردی



له نووسینی کوردیدا، بۆشایی و نیشانهی نیو بۆشایی (ZWNJ) بۆ جوداکردنهوهی وشهکان بهکار دێن، بهلام کیشهی زۆریان ههیه:

- یهکلانهبوونهوهی رینووس:

- 18ê, 18-ê, 18'ê

- بهکار دێ، بهکار دێ، بهکار دێ، بهکار دێ

- نووساندنی زۆری وشهکان:

- لهویشدايه ← له وی ش دا یه

- وشه لیكدراوهکان:

مردوو مرو

خواخراو بۆکردگ

«له کوردستانهوه کووتایی به پێوهندییهکه هیئا.» ← له ... هوه؟ کووتایی هیئا؟



# له تکرڊنی تۆکن بۆ کوردی له KLPT دا

- له تکرڊنی تۆکن له KLPT دا پشت به فهرههنگیکی وشه و شیکاریی مۆرفۆلۆژی ده به ستنی
  - گونجاوه بۆ له تکرڊنی زمانه وانی
  - یارمه تی دۆزینه وهی سنووری وشه ده دا
  - ریخۆشکهره بۆ ئهرکی دیکه له زمانه وانی کۆمپیوته ریدا
- نمونه ی له تکرڊن:

دهق: به رهه مه یانی شیوازه کانی دواکه وتنی  
دوای له تکرڊن: دوا-که وتن-ی شیوازه کان-ی به رهه م-هیئان-ی

- له م وتاره دا باسی کراوه: [10]

KLPT: <https://github.com/sinaahmadi/klpt>





## نموونه‌ی کوردیی سه‌روو

"bi-can-û-bên":

"bicanûbên"

"bi canûbên"

"bican ûbên"

"bi can ûbên"

"bicanû bên"

"bi canû bên"

## نموونه‌ی کوردیی ناوه‌ندی

"ئاخر-و-ئۆخر":

"ئاخرو ئۆخر"

"ئاخرووئۆخر"

"ئاخر و ئۆخر"

"ئاخروئۆخر"

"ئاخر وئۆخر"

1 کۆ کردنه‌وه‌ی وشه

2 وشه لێکدراوه‌کان بدۆژه‌وه  
کوردیی ناوه‌ندی: 8180 وشه (1513 لێکدراو)  
کوردیی سه‌روو: 9970 وشه (1507 لێکدراو)

3 هه‌موو فۆرمه ئالۆز و نا‌ئالۆزه‌کانیان ساز بکه

- لیستی پێشگر، پاشگر و کلیتی که کان:  
پێشگر: له، وه، ده، ره، به  
پاشگر: هکه، هکان، ان، گهل  
جیناوه لکاوهکان: م، ت، ی، مان، تان، یان
- نموونهی شیکاریی مۆرفۆلۆژی:

بهرزترهکه ← بهرز + تر + هکه

- ژمارهی مۆرفیمهکان:  
سۆرانی: 161 پاشگر و 11 پێشگر  
کورمانجی: 46 پاشگر و 17 پێشگر

## نموونه

«٢٥ کهسیان پێی کهوتن»

1. 25 کهسیان پێی کهوتن

2. 25 کهسیان —پێی— کهوتن—

3. 25 —کهسیان— —پێی— کهوتن—

4. 25 —کهس— —یان— —پێی— کهوتن—

1 پێش- پرۆسهس کردن:

یه کهخستنی نووسینی پیتەکان

زیادکردنی بۆشایی له دهوری خالبهندی و ژمارهکان

2 لهتکردنی وشه لیكدراوهکان:

دۆزینهوهی وشه لیكدراوهکان له فهرههنگدا

جوداکردنهویان به نیشانهی ژیرهیڵ —

3 لهتکردنی وشه:

جیاکردنهوه به بۆشایی

دۆزینهوهی وشهکان له فهرههنگدا

4 شیکاریی مۆرفۆلۆژی:

دۆزینهوهی پێشگر و پاشگرهکان

جیاکردنهوهی مۆرفیمه لیكدراوهکان

- قەبارە ی بچووکی فەرھەنگ:
- کاریگەری لە سەر کوالیتی لەتکردن
- ئالۆز بوونی واتای پیکهاتهکان:
- بۆ نموونه وشە ی «لاوین» که دهکری «لاو» + «ین» یان وهک یهک وشه بی، بهلام سیستمه که به  
'\_لاوین' له تی دهکا
- پیکهاته پەرژ و بلاوهکان:
- وهک له «بریاریکی گهوره یان دا» دا، هه رچه ند «بریار دان» کرداریکی لیکدراوه بهلام وشه ی دیکه یان  
که وتوووته نیوان
- کیشە ی کارایی:
- گه ران له فەرھەنگ کاتگیره
- روانگە ی پیدراوه؟

# هه‌سه‌نگاندنی له‌تکردنی تۆکن

«دواکه‌وتنی شیواز هکانی به‌ره‌مه‌هینان له‌م ئابووریانه‌دا ده‌گه‌رێته‌وه‌ بۆ: نه‌بوونی هۆیه‌کانی ته‌کنیکی تازهی هاورده‌ تا به‌ره‌مه‌هینه‌کان به‌کاری بێن.»

(4000) BPE

دواکه‌وت نی شیواز هکانی به‌ره‌مه‌هینان  
له‌م ئابووریانه‌دا ده‌گه‌رێته‌وه‌ بۆ: نه‌بوونی  
هۆیه‌کانی ته‌کنیکی تازهی هاورده  
ه‌ تا به‌ره‌مه‌هین هکان به‌کاری بێن.

KLPT

دواکه‌وتنی شیواز هکانی به‌ره‌مه‌هینان  
له‌م ئابووریانه‌دا ده‌گه‌رێته‌وه‌ بۆ: نه‌بوونی  
هۆیه‌کانی ته‌کنیکی تازهی هاورده  
تا به‌ره‌مه‌هینه‌کان به‌کاری بێن.

هه‌سه‌نگه‌ندی له‌تکردنی تۆکن پێوه‌ندی به‌ ئه‌رکه‌وه‌ هه‌یه‌.



سەرچاوەی ھەموو وێنەکان ویکیمیדיا: [https://commons.wikimedia.org/wiki/Main\\_Page](https://commons.wikimedia.org/wiki/Main_Page)

- [1] Jonathan J. Webster and Chunyu Kit.  
Tokenization as the initial phase in NLP.  
In *COLING 1992 Volume 4: The 14th International Conference on Computational Linguistics*, 1992.
- [2] Eric de la Clergerie and Lionel Clément.  
MAF: a morphosyntactic annotation framework.  
*Actes de LTC*, pages 90–94, 2005.
- [3] David D Palmer.  
Tokenisation and sentence segmentation.  
*Handbook of natural language processing*, pages 11–35, 2000.
- [4] Sabrina J Mielke, Zaid Alyafeai, Elizabeth Salesky, Colin Raffel, Manan Dey, Matthias Gallé, Arun Raja, Chenglei Si, Wilson Y Lee, Benoît Sagot, et al.  
Between words and characters: a brief history of open-vocabulary modeling and tokenization in NLP.  
*arXiv preprint arXiv:2112.10508*, 2021.
- [5] Rico Sennrich, Barry Haddow, and Alexandra Birch.  
Neural machine translation of rare words with subword units.  
In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August 2016.  
Association for Computational Linguistics.

- [6] Taku Kudo.  
Subword regularization: Improving neural network translation models with multiple subword candidates.  
*arXiv preprint arXiv:1804.10959*, 2018.
- [7] Mike Schuster and Kaisuke Nakajima.  
Japanese and Korean voice search.  
In *2012 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5149–5152. IEEE, 2012.
- [8] Kris Cao and Laura Rimell.  
You should evaluate your language model on marginal likelihood over tokenisations.  
*arXiv preprint arXiv:2109.02550*, 2021.
- [9] Aleksandar Petrov, Emanuele La Malfa, Philip HS Torr, and Adel Bibi.  
Language model tokenizers introduce unfairness between languages.  
*arXiv preprint arXiv:2305.15425*, 2023.
- [10] Sina Ahmadi.  
A tokenization system for the Kurdish language.  
In *Proceedings of the 7th Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 114–127, 2020.

