



# Byte-pair encoding

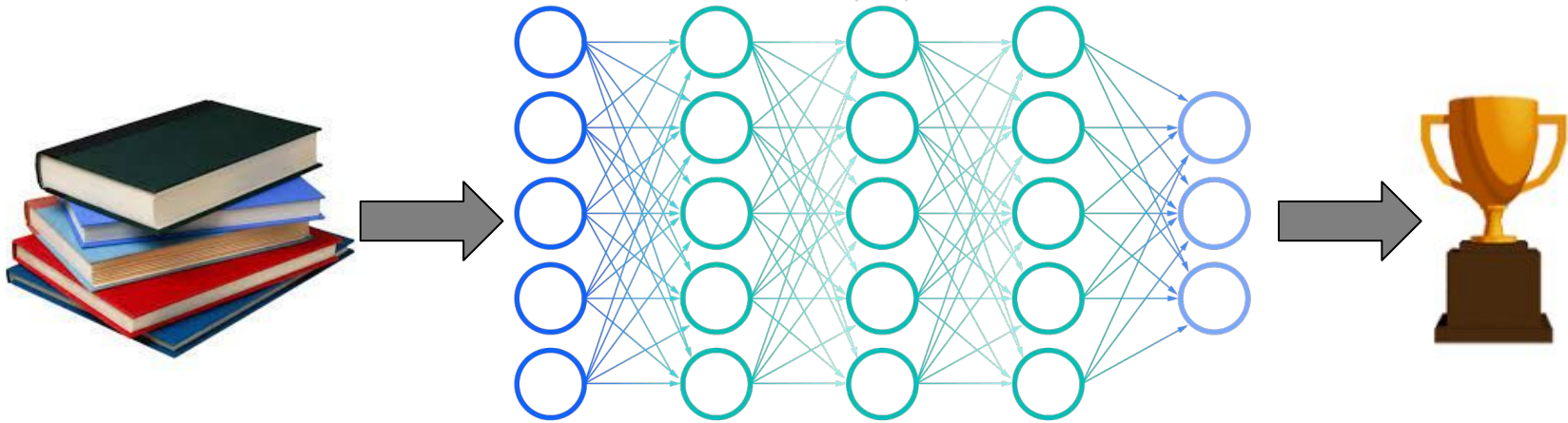
July 16, 2021



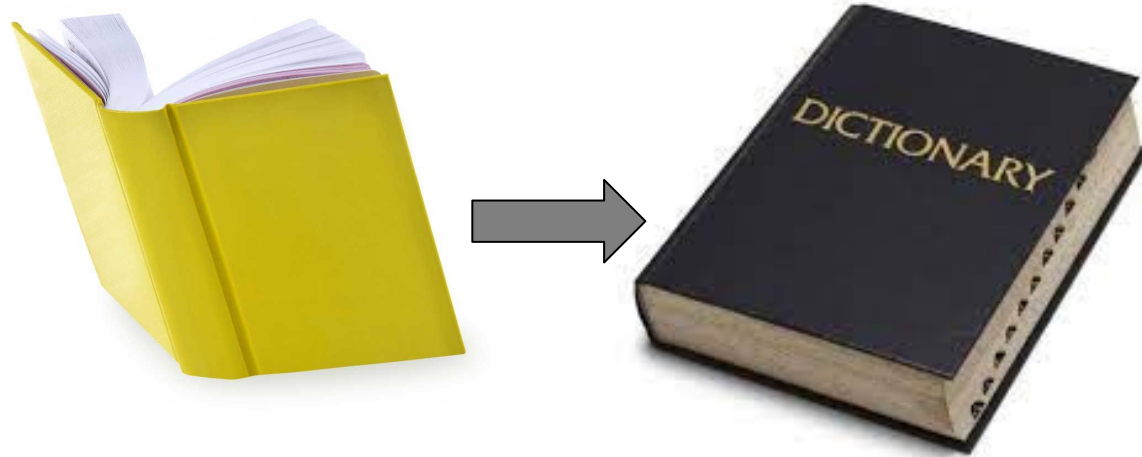
# Outline

- How do neural models process text?
- Why is BPE necessary?
- Walkthrough of BPE generation
- Alternatives: SentencePiece & WordPiece, Sequitur, Morfessor

# How do neural models read input?

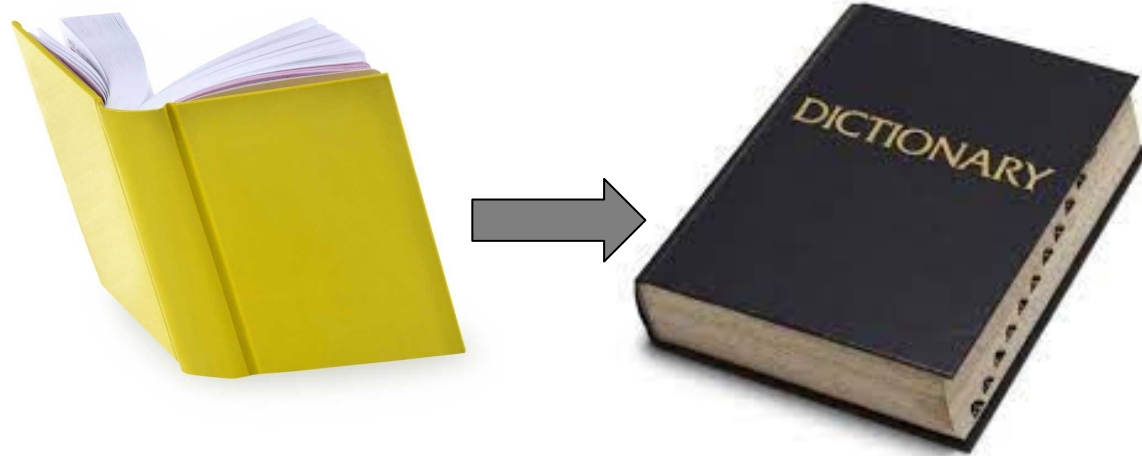


# Keep track of all of the words during training



Type	ID
sink	0
representative	1
exclusive	2
shortage	3
confront	4
root	5
technology	6
table	7
forest	8

# Keep track of all of the words during training



The birds in the forest heard a sound.  
52 1032 19 52 8 27 1941 36 2592

Type	ID
sink	0
representative	1
exclusive	2
shortage	3
confront	4
root	5
technology	6
table	7
forest	8

# Abstraction:

The birds in the forest heard a sound.

0	0
1	0
2	0
...	0
51	0
52	1
53	0
...	0
15000	0

0	0
1	0
2	0
...	0
1031	0
1032	1
1033	0
...	0
15000	0

0	0
1	0
2	0
...	0
18	0
19	1
20	0
...	0
15000	0

0	0
1	0
2	0
...	0
51	0
52	1
53	0
...	0
15000	0

...

# Characters or words?

	Words	Chars
Pros		
Cons		

# Characters or words?

	Words	Chars
Pros		
Cons	Less flexible - easier to get OOV	



# Characters or words?

	Words	Chars
Pros		Only need a vocabulary of ~50.
Cons	Less flexible - easier to get OOV	

# Characters or words?

	Words	Chars
Pros		Only need a vocabulary of ~50.
Cons	Less flexible - easier to get OOV	Letters don't always represent the same features; context is much larger.

# Characters or words?

	Words	Chars
Pros	If seen frequently enough, we can learn very precise embeddings.	Only need a vocabulary of ~50.
Cons	Less flexible - easier to get OOV	Letters don't always represent the same features; context is much larger.

# Subword processing

- Compromise between word-based and character-based models.
- Preserves frequent sequences, but allows some flexibility

The birds in the forest heard a sound: *The bird -s in the forest hear -d a sound*

# BPE

- Find common adjacent sequences, build from bottom up.

On the 24th of February, 1815, the look-out at Notre-Dame de la Garde signalled the three-master, the Pharaon from Smyrna, Trieste, and Naples.

As usual, a pilot put off immediately, and rounding the Château d'If, got on board the vessel between Cape Morgiou and Rion island.

Immediately, and according to custom, the ramparts of Fort Saint-Jean were covered with spectators; it is always an event at Marseilles for a ship to come into port, especially when this ship, like the Pharaon, has been built, rigged, and laden at the old Phocée docks, and belongs to an owner of the city.

# BPE

- Find common adjacent sequences, build from bottom up.

on the 24th of february , 1815 , the look - out at notre - dame de la garde signalled the three - master , the pharaon from smyrna , trieste , and naples . as usual , a pilot put off immediately , and rounding the château d' if , got on board the vessel between cape morgiou and rion island . immediately , and according to custom , the ramparts of fort saint - jean were covered with spectators ; it is always an event at marseilles for a ship to come into port , especially when this ship , like the pharaon , has been built , rigged , and laden at the old phocee docks , and belongs to an owner of the city .

# BPE

- Find common adjacent sequences, build from bottom up.

on\_the\_24th\_of\_february\_,\_1815\_,\_the\_look\_-\_out\_at  
\_notre\_-\_dame\_de\_la\_garde\_signalled\_the\_three\_-\_m  
aster\_,\_the\_pharaon\_from\_smyrna\_,\_trieste\_,\_and\_na  
ples\_.

as\_usual\_,\_a\_pilot\_put\_off\_immediately\_,\_and\_round  
ing\_the\_château\_d'\_if\_,\_got\_on\_board\_the\_vessel\_be  
tween\_cape\_morgiou\_and\_rion\_island\_.immediately\_  
\_,\_and\_according\_to\_custom\_,\_the\_ramparts\_of\_fort\_s  
aint\_-\_jean\_were\_covered\_with\_spectators\_;\_it\_is\_al  
ways\_an\_event\_at\_marseilles\_for\_a\_ship\_to\_come\_int  
o\_port\_,\_especially\_when\_this\_ship\_,\_like\_the\_phara  
on\_,\_has\_been\_built\_,\_rigged\_,\_and\_laden\_at\_the\_ol  
d\_phocee\_docks\_,\_and\_belongs\_to\_an\_owner\_of\_the\_  
city\_.

e	56
a	45
t	44
o	37
n	31
...	...

# BPE

- Find common adjacent sequences, build from bottom up.

on\_the\_24th\_of\_february\_,\_1815\_,\_the\_look\_-\_out\_at  
\_notre\_-\_dame\_de\_la\_garde\_signalled\_the\_three\_-\_m  
aster\_,\_the\_pharaon\_from\_smyrna\_,\_trieste\_,\_and\_na  
ples\_.

as\_usual\_,\_a\_pilot\_put\_off\_immediately\_,\_and\_round  
ing\_the\_château\_d'\_if\_,\_got\_on\_board\_the\_vessel\_be  
tween\_cape\_morgiou\_and\_rion\_island\_.immediately\_  
\_,\_and\_according\_to\_custom\_,\_the\_ramparts\_of\_fort\_s  
aint\_-\_jean\_were\_covered\_with\_spectators\_;\_it\_is\_al  
ways\_an\_event\_at\_marseilles\_for\_a\_ship\_to\_come\_int  
o\_port\_,\_especially\_when\_this\_ship\_,\_like\_the\_phara  
on\_,\_has\_been\_built\_,\_rigged\_,\_and\_laden\_at\_the\_ol  
d\_phocce\_docks\_,\_and\_belongs\_to\_an\_owner\_of\_the\_  
city\_.

e	56
a	45
t	44
o	37
n	31
...	...

th	14
he	11
an	10
nd	8
ar	7
...	...



# BPE

- Find common adjacent sequences, build from bottom up.

on\_**th**e\_24**th**\_of\_february\_,\_1815\_,\_**th**e\_look\_-\_out\_at\_  
notre\_-\_dame\_de\_la\_garde\_signalled\_**th**e\_**th**ree\_-\_mas  
ter\_,\_**th**e\_pharaon\_from\_smyrna\_,\_trieste\_,\_and\_napl  
es\_.

as\_usual\_,\_a\_pilot\_put\_off\_immediately\_,\_and\_round  
ing\_**th**e\_château\_d'if\_,\_got\_on\_board\_**th**e\_vessel\_bet  
ween\_cape\_morgiou\_and\_rion\_island\_.immediately\_,  
\_and\_according\_to\_custom\_,\_**th**e\_ramparts\_of\_fort\_sa  
int\_-\_jean\_were\_covered\_with**th**\_spectators\_;\_it\_is\_alw  
ays\_an\_event\_at\_marseilles\_for\_a\_ship\_to\_come\_into  
\_port\_,\_especially\_when\_**this**\_ship\_,\_like\_**th**e\_pharao  
n\_,\_has\_been\_built\_,\_rigged\_,\_and\_laden\_at\_**th**e\_old  
\_phocee\_docks\_,\_and\_belongs\_to\_an\_owner\_of\_**th**e\_ci  
ty\_.

e	56
a	45
t	44
o	37
n	31
...	...

th	14
he	11
an	10
nd	8
ar	7
...	...

# BPE

- Find common adjacent sequences, build from bottom up.

on\_the\_24th\_of\_february\_,\_1815\_,\_the\_look\_-\_out\_at\_  
notre\_-\_dame\_de\_la\_garde\_signalled\_the\_three\_-\_mas  
ter\_,\_the\_pharaon\_from\_smyrna\_,\_trieste\_,\_and\_napl  
es\_.

as\_usual\_,\_a\_pilot\_put\_off\_immediately\_,\_and\_round  
ing\_the\_château\_d'\_if\_,\_got\_on\_board\_the\_vessel\_bet  
ween\_cape\_morgiou\_and\_rion\_island\_.immediately\_,  
\_and\_according\_to\_custom\_,\_the\_ramparts\_of\_fort\_sa  
int\_-\_jean\_were\_covered\_with\_spectators\_;\_it\_is\_alw  
ays\_an\_event\_at\_marseilles\_for\_a\_ship\_to\_come\_into  
\_port\_,\_especially\_when\_this\_ship\_,\_like\_the\_pharao  
n\_,\_has\_been\_built\_,\_rigged\_,\_and\_laden\_at\_the\_old  
\_phocee\_docks\_,\_and\_belongs\_to\_an\_owner\_of\_the\_ci  
ty\_.

e	56
a	45
t	44
o	37
n	31
...	...

the	10
an	10
nd	8
ar	7
on	6
...	...

# BPE

- Find common adjacent sequences, build from bottom up.

on\_**the**\_24th\_of\_february\_,\_1815\_,\_**the**\_look\_-\_out\_at\_  
notre\_-\_dame\_de\_la\_garde\_signalled\_**the**\_three\_-\_mas  
ter\_,\_**the**\_pharaon\_from\_smyrna\_,\_trieste\_,\_and\_napl  
es\_.

as\_usual\_,\_a\_pilot\_put\_off\_immediately\_,\_and\_round  
ing\_**the**\_château\_d'\_if\_,\_got\_on\_board\_**the**\_vessel\_bet  
ween\_cape\_morgiou\_and\_rion\_island\_.immediately\_,  
\_and\_according\_to\_custom\_,\_**the**\_ramparts\_of\_fort\_sa  
int\_-\_jean\_were\_covered\_with\_spectators\_;\_it\_is\_alw  
ays\_an\_event\_at\_marseilles\_for\_a\_ship\_to\_come\_into  
\_port\_,\_especially\_when\_this\_ship\_,\_like\_**the**\_pharaon  
\_,\_has\_been\_built\_,\_rigged\_,\_and\_laden\_at\_**the**\_old\_p  
hocee\_docks\_,\_and\_belongs\_to\_an\_owner\_of\_**the**\_city  
\_.

e	56
a	45
t	44
o	37
n	31
...	...

the	10
an	10
nd	8
ar	7
on	6
...	...

# Some more specifics

- Iterations continue until one of 2 events happens:
  - We reach the pre-specified number of merges (often on the order of 16K, 32K, or 64K).
  - We no longer have any bigrams that occur  $> 1$

# Some more specifics

- Iterations continue until one of 2 events happens:
  - We reach the pre-specified number of merges (often on the order of 16K, 32K, or 64K).
  - We no longer have any bigrams that occur  $> 1$
- More merges approaches word-based model; fewer approaches character-based model.

# Some more specifics

- Iterations continue until one of 2 events happens:
  - We reach the pre-specified number of merges (often on the order of 16K, 32K, or 64K).
  - We no longer have any bigrams that occur  $> 1$
- More merges approaches word-based model; fewer approaches character-based model.
- Cross-word bigrams are not considered bigrams; likewise, word-internal breaks are indicated by a special symbol: @@ -> This is necessary to re-construct words afterwards.

# SentencePiece and WordPiece

- Similar methods to BPE:
- WordPiece has a different selection process - instead of selecting the character bigram with the highest frequency, it chooses the one that maximises the likelihood of the output sequence.
- SentencePiece does everything in Unicode, and works tokenization into the algorithm.

# IRRMGP

- Very similar to BPE, SentencePiece, etc.
- Instead of choosing most common bigram, it instead maximizes  $|x| * freq(x)$
- Jointly trying to find the longest sequence that is also frequent.



# IRRMGP

- Very similar to BPE, SentencePiece, etc.
- Instead of choosing most common bigram, it instead maximizes  $|x| * freq(x)$
- Jointly trying to find the longest sequence that is also frequent.

$|'language'| * 50 > |'es'| * 150$   
 $8 * 50 > 2 * 150$

language@@ s    vs.    languag@@ es

# MaxLen

- MaxLen finds the longest sequence that has already been seen in the text, and separates it out.

‘One language is fun, but many languages are better’

# MaxLen

- MaxLen finds the longest sequence that has already been seen in the text, and separates it out.

‘One *lan*guage is fun, but *man*y language s are better’

# Sequitur

Builds a grammar word-by-word:

# Sequitur

Builds a grammar word-by-word:

Read in one word - reduce it as much as possible: banana  $\rightarrow$  baXX; na $\rightarrow$ X

# Sequitur

Builds a grammar word:

Read in one word - reduce it as much as possible: banana  $\rightarrow$  baXX; na $\rightarrow$ X

Read in next word, expand grammar.

# Morfessor

- Read in all words
- Find segmentation that maximizes the probability of the segments.

# Morfessor

- Read in all words
- Find segmentation that maximizes the probability of the segments.

“^languages\$” -> “^language + s\$”

?

“^languages\$” -> “^languag + es\$”



# Morfessor

- Read in all words
- Find segmentation that maximizes the probability of the segments.

“^languages\$” -> “^language + s\$”

?

“^languages\$” -> “^languag + es\$”

Likely “^language + s\$” -> removing “-es” would likely create some strange morphemes for other words.