

WordPiece

Espen Bang og Iver Omholt

Hva er WordPiece?

- WordPiece er en algoritme for å tokenisere tekst ved å dele ord opp i mindre biter (subwords)
- Opprinnelig utviklet hos Google (for machine translations, speech recognition)
- Brukes i liten grad i klassiske IR systemer. Mer vanlig i moderne IR systemer med NLP.
- Bidirectional Encoder Representations from Transformers (BERT)
- Likheter med Byte Pair Encoding (BPE)

Hva er formålet med WordPiece?

Addressere noen av begrensningene ved standard ordbasert tokenisering. Spesielt for NLP modeller.

- Whitespace, Rule-Based, N-gram tokenisering.
- Regelbasert vs. sannsynlighet

Ønsker å finne en balanse mellom å representere tekst med hele ord, og å dele ord opp i mindre deler når det trengs.

Hva oppnås med WordPiece?

Out-of-Vocabulary words

- Ved å dele opp ord i mindre biter kan man enklere behandle nye ord.
- NLP modeller må forstå ny tekst, feilstavelser, sjeldne eller sammensatte ord.

Redusere størrelse på vokabular

- Mer effektiv trening av maskinlæringsmodeller.

Generalisering

- Ved å gjenkjenne sub-words kan man generalisere bedre på tvers av ordformer

Fleksibilitet på tvers av språk

- Tysk, norsk, finsk

Eksempel

- happiness, joyfulness → happi, joyful, ness
- skrivebord, skrivebok, lesebord → skrive, bord, bok, lese

Algorithme

Corpus:

«all squares are rectangles, but not all rectangles are squares»

Algoritme

Corpus:

- all: 2
- rectangles: 2
- are: 2
- squares: 2
- but: 1
- not: 1

Algoritme

Corpus:

- all: 2
- rectangles: 2
- are: 2
- squares: 2
- but: 1
- not: 1

Splits:

a, ##l, ##l
r, ##e, ##c, ##t, ##a, ##n, ##g, ##l, ##e, ##s
a, ##r, ##e
s, ##q, ##u, ##a, ##r, ##e, ##s
b, ##u, ##t
n, ##o, ##t

Algorithme

Corpus:

- all: 2
- rectangles: 2
- are: 2
- squares: 2
- but: 1
- not: 1

Splits:

a, ##l, ##l
r, ##e, ##c, ##t, ##a, ##n, ##g, ##l, ##e, ##s
a, ##r, ##e
s, ##q, ##u, ##a, ##r, ##e, ##s
b, ##u, ##t
n, ##o, ##t

Vocabulary:

a, ##l, s, ##q, ##u, ##a, ##r, ##e, ##s, r, ##c, ##t, ##n, ##g, b, n, ##o

$$\text{Score} = \text{freq}(\text{AB}) / (\text{freq}(\text{A}) * \text{freq}(\text{B}))$$

Algoritme

Corpus:

- all: 2
- rectangles: 2
- are: 2
- squares: 2
- but: 1
- not: 1

Splits:

a, ##l, ##l
r, ##e, ##c, ##t, ##a, ##n, ##g, ##l, ##e, ##s
a, ##r, ##e
s, ##q, ##u, ##a, ##r, ##e, ##s
b, ##u, ##t
n, ##o, ##t

Vocabulary:

a, ##l, r, ##e, ##c, ##t, ##a, ##n, ##g, ##s, s, ##q, ##u, ##r, b, n, ##o

$$\text{Score} = \text{freq}(a, \text{\#\#l}) / (\text{freq}(a) * \text{freq}(\text{\#\#l}))$$

Algoritme

Corpus:

- all: 2
- rectangles: 2
- are: 2
- squares: 2
- but: 1
- not: 1

Splits:

a, \text{\#\#l}, \text{\#\#l}

r, \text{\#\#e}, \text{\#\#c}, \text{\#\#t}, \text{\#\#a}, \text{\#\#n}, \text{\#\#g}, \text{\#\#l}, \text{\#\#e}, \text{\#\#s}

a, \text{\#\#r}, \text{\#\#e}

s, \text{\#\#q}, \text{\#\#u}, \text{\#\#a}, \text{\#\#r}, \text{\#\#e}, \text{\#\#s}

b, \text{\#\#u}, \text{\#\#t}

n, \text{\#\#o}, \text{\#\#t}

Vocabulary:

a, \text{\#\#l}, r, \text{\#\#e}, \text{\#\#c}, \text{\#\#t}, \text{\#\#a}, \text{\#\#n}, \text{\#\#g}, \text{\#\#s}, s, \text{\#\#q}, \text{\#\#u}, \text{\#\#r}, b, n, \text{\#\#o}

$$\text{Score} = 2/(4*6) = 2/24 = 0.0833$$

Algoritme

Corpus:

- all: 2
- rectangles: 2
- are: 2
- squares: 2
- but: 1
- not: 1

Splits:

a, ##l, ##l
r, ##e, ##c, ##t, ##a, ##n, ##g, ##l, ##e, ##s
a, ##r, ##e
s, ##q, ##u, ##a, ##r, ##e, ##s
b, ##u, ##t
n, ##o, ##t

Vocabulary:

a, ##l, r, ##e, ##c, ##t, ##a, ##n, ##g, ##s, s, ##q, ##u, ##r, b, n, ##o

Algorithme

Corpus:

- all: 2
- rectangles: 2
- are: 2
- squares: 2
- but: 1
- not: 1

Splits:

al, ##l
r, ##e, ##c, ##t, ##a, ##n, ##g, ##l, ##e, ##s
a, ##r, ##e
s, ##q, ##u, ##a, ##r, ##e, ##s
b, ##u, ##t
n, ##o, ##t

Vocabulary:

a, ##l, r, ##e, ##c, ##t, ##a, ##n, ##g, ##s, s, ##q, ##u, ##r, b, n, ##o,
al

Algorithme

Vocabulary:

..., all, rectangle, are, square, but, not, ##es, #s, ...

Tokenisering

- notes = not ##es

Vocabulary:

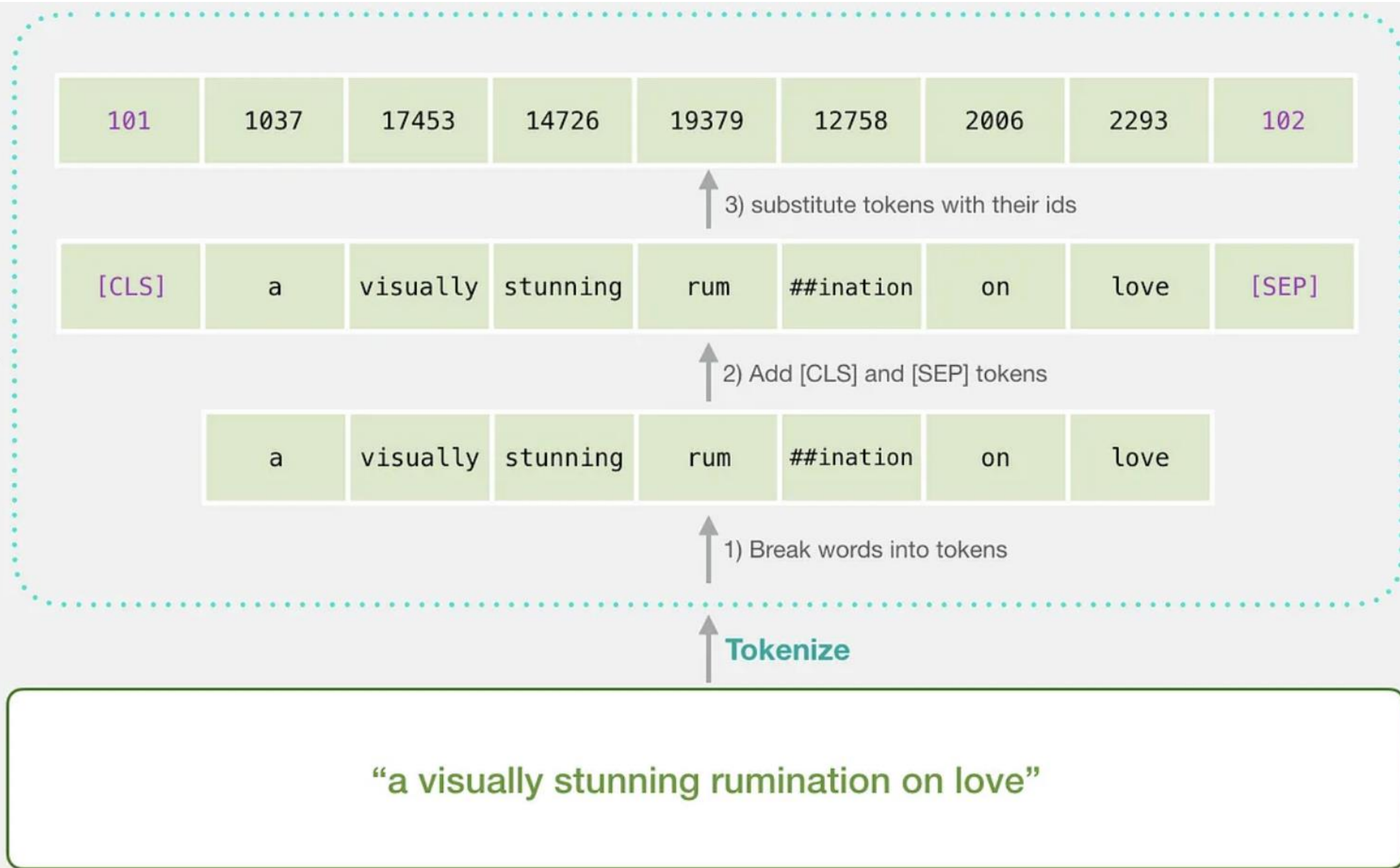
..., all, rectangle, are, square, but, not, #es, #s, ...

Tokenisering

- apples = ?? = [UNK]

Vocabulary:

..., all, rectangle, are, square, but, not, #es, #s, ...



Tokenization

DistilBertTokenizer

3) substitute tokens with their ids

2) Add [CLS] and [SEP] tokens

1) Break words into tokens

Tokenize

"a visually stunning rumination on love"