

Additional topics

# BPE - Byte Pair Encoding

## *Initial state:*

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

## *First step:*

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

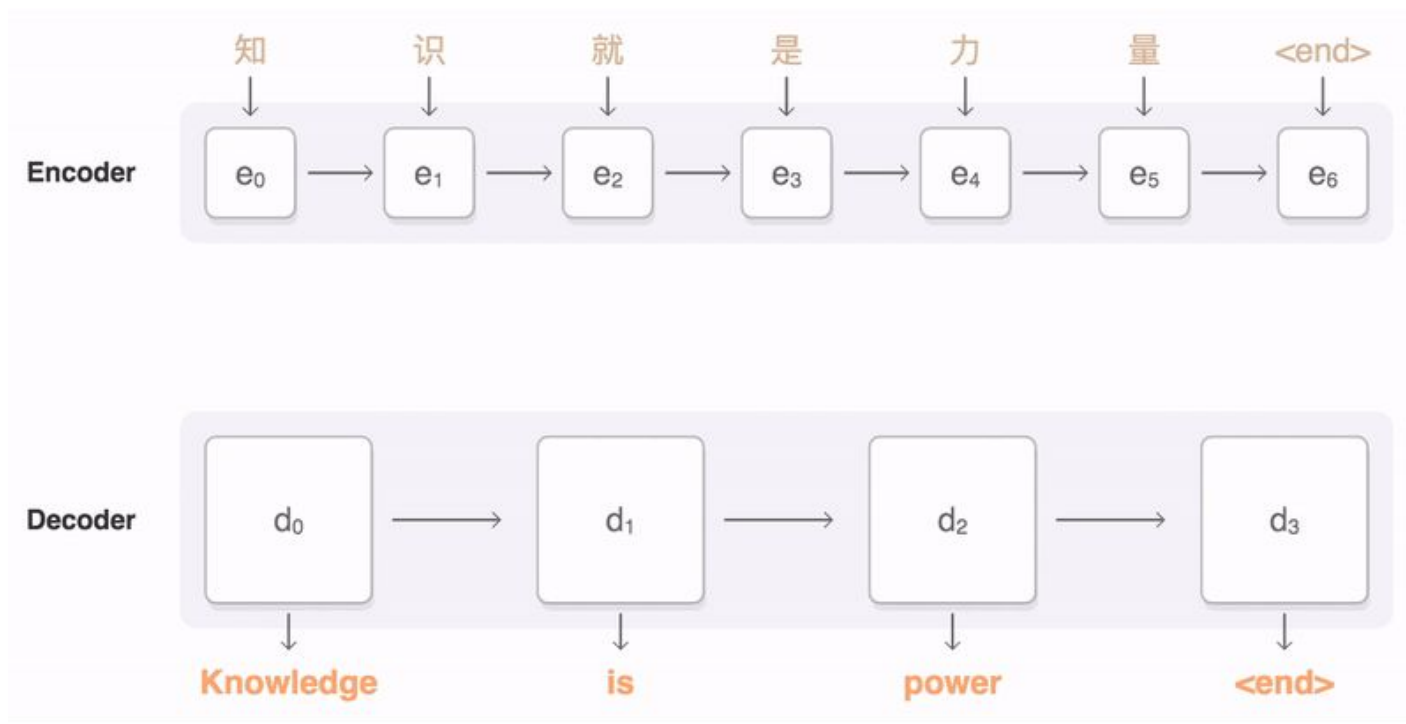
## *Second step:*

{'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w **es** t </w>': 6, 'w i d **es** t </w>': 3}

After each merge, there could be three scenarios, the number of tokens decreases by one, remains the same or increases by one. But in practice, as the number of merges increases, usually the number of tokens first increases then decreases.

[More info](#)

# Recap



# GS - Greedy search

```
# define a sequence of 10 words over a vocab of 5 words
```

```
data = [[0.1, 0.2, 0.3, 0.4, 0.5],  
        [0.5, 0.4, 0.3, 0.2, 0.1],  
        [0.1, 0.2, 0.3, 0.4, 0.5],  
        [0.5, 0.4, 0.3, 0.2, 0.1],  
        [0.1, 0.2, 0.3, 0.4, 0.5],  
        [0.5, 0.4, 0.3, 0.2, 0.1],  
        [0.1, 0.2, 0.3, 0.4, 0.5],  
        [0.5, 0.4, 0.3, 0.2, 0.1],  
        [0.1, 0.2, 0.3, 0.4, 0.5],  
        [0.5, 0.4, 0.3, 0.2, 0.1]]
```

```
Result = [4, 0, 4, 0, 4, 0, 4, 0, 4, 0]
```

[More info](#)

# BS - Beam search

[[4, 0, 4, 0, 4, 0, 4, 0, 4, 0], 0.025600863289563108]

[[4, 0, 4, 0, 4, 0, 4, 0, 4, 1], 0.03384250043584397]

[[4, 0, 4, 0, 4, 0, 4, 0, 3, 0], 0.03384250043584397]

