

Detokenization

Ki Hyun Kim

nlp.with.deep.learning@gmail.com

Tokenization

1. 영어 원문

- There's currently over a thousand TED Talks on the TED website.

2. tokenization을 수행하고, 기존 띄어쓰기와 구분을 위해 _ (U+2581) 삽입

- _There 's _currently _over _a _thousand _TED _Talks _on _the _TED _website .

3. subword segmentation을 수행, 공백 구분 위한 _ 삽입

- __There _'s __currently __over __a __thous and __TED __T al ks __on __the __TED __we b site _.

Detokenization

1. whitespace를 제거

- __There_'s__currently__over__a__thousand__TED__Talks__on__the__
__TED__website__.

2. __을 white space로 치환

- There_'s currently over a thousand TED Talks on the TED website__.

3. _를 제거

- There's currently over a thousand TED Talks on the TED website.



```
sed "s/ //g" | sed "s/__/ /g" | sed "s/_//g"
```