

data drive tokenizer: tokenizer dengan proses training

Pasangan Byte

ganti XX jadi A

taro di lookup

xd jadiin B

karakter yang paling sering muncul di merging

Paling sering ketemu itu ""o"

jangan lupa tambahin ke vocab.

lo yang tadi di merge ke merge Rule

ekarang paling frequent alo

$2 + 1 = 3$, karena ada 2 sebagai okurensi hakim dan 1 di baki.

dari fase training cara tokenisasi nya agak berbeda.

yang pertama gak pake pager, cuman kedua, dst.

