
Homework

Vlad-Cristian Matei

Departmentul Computer Science

National University of Science and Technology POLITEHNICA

Bucharest, Romania

vlad_cristian.matei@stud.acs.upb.ro

1 Dataset

For completing the task, I tried several preprocessing approaches for the dataset. First, I attempted to normalize the text, then formatted the text in the following structure: "[\${author}] \${text}". The two approaches were as follows: one involved taking the lines of text exactly as they appeared in the dataset, while the other involved merging consecutive lines from the same author, thereby constructing a direct dialogue with longer lines for each character.

```
0 [ king henry iv ] so shaken as we are, so wan ...
1 [ westmoreland ] my liege, this haste was hot ...
2 [ king henry iv ] it seems then that the tidin...
3 [ westmoreland ] this match'd with other did, ...
4 [ king henry iv ] here is a dear, a true indus...
```

Figure 1: Preprocessing

In the end, I proceeded with the second method, where I created longer, unified lines.

2 Tokenization

For completing the task, I used both proposed tokenizers, the character-level tokenizer and the imported one. Specifically, for the subword tokenizer, I used "gpt2" from the transformers.GPT2Tokenizer library.

```

tokenizer = SubTok()
print(tokenizer.pad_token, tokenizer.pad_token_id)
text = "Anne, has apples."
tokens = tokenizer.tokenize(text)
print("Tokens:", tokens)
ids = tokenizer.convert_tokens_to_ids(tokens)
print("Token IDs:", ids)
reconstructed_tokens = tokenizer.convert_ids_to_tokens(ids)
print("Reconstructed Tokens:", reconstructed_tokens)
reconstructed_text = tokenizer.detokenize(reconstructed_tokens)
print("Reconstructed Text:", reconstructed_text)
reconstructed_text_via_decode = tokenizer.decode(ids)
print("Reconstructed Text via Decode:", reconstructed_text_via_decode)

tokenizer_config.json: 0% | 0.00/26.0 [00:00<?, 7B/s]
vocab.json: 0% | 0.00/1.04M [00:00<?, 7B/s]
merges.txt: 0% | 0.00/456k [00:00<?, 7B/s]
tokenizer.json: 0% | 0.00/1.30M [00:00<?, 7B/s]
config.json: 0% | 0.00/665 [00:00<?, 7B/s]
<PAD> 50259

Tokens: ['<SOS>', 'Anne', ',', 'Ghas', 'Gapples', '.', '<EOS>']
Token IDs: [50257, 43227, 11, 468, 22514, 13, 50258]
Reconstructed Tokens: ['<SOS>', 'Anne', ',', 'Ghas', 'Gapples', '.', '<EOS>']
Reconstructed Text: Anne, has apples.
Reconstructed Text via Decode: Anne, has apples.

```

Figure 2: SubWord tokenization

3 Models

I trained multiple models to address the task. However, in the final form, I will present only four of them.

3.1 Small Model - Character Tokenizer

During the experiments, I settled on the following configuration: a batch size of 128, a maximum sequence length of 512, an embedding dimension of 384, 6 attention heads and layers, and a hidden dimension of 128, and Adam Optimizer with a learning rate of 1e-4.

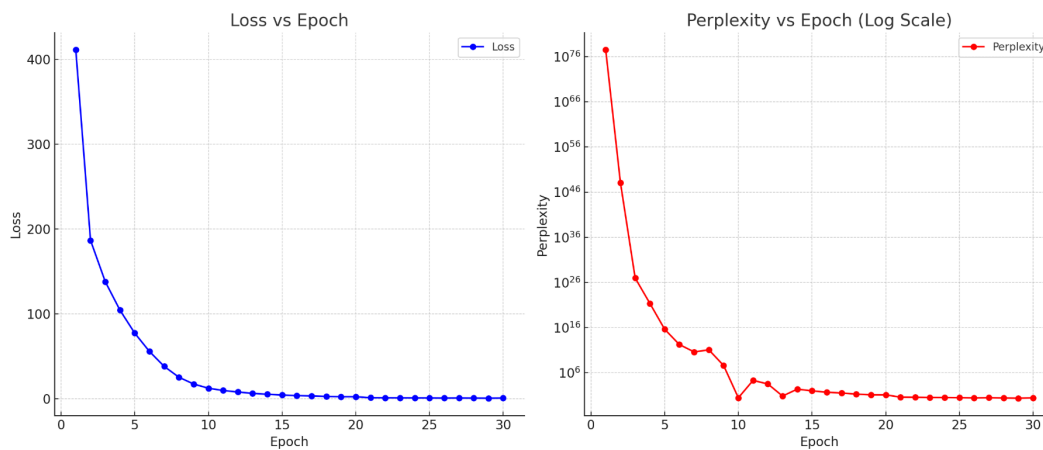


Figure 3: Small Model - Character Tokenizer

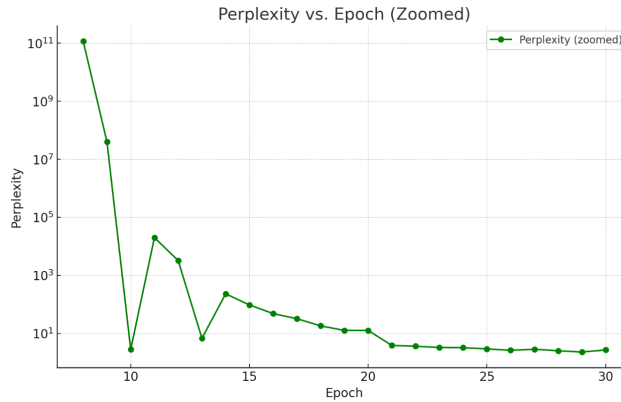


Figure 4: Small Model - Character Tokenizer

The evolution of the losses is shown in the figures below. As can be observed, there is a convergence, with the final loss on the test dataset being 1.3589. Additionally, the last perplexity value calculated in the final epoch is 2.69.

Qualitative Evaluation

```
Epoch 25/30, Loss: 1.0755, Perplexity: 2.9313
Predicted Text: [ rrr rrrr r e] i l l l r i eltser rshr ser ss sssr rr ir Joero ill lr irri erioarr xb r b rh bor b r r i l e bhe b
roabr ehre b e e breririrr brirebhlrr blrr r bil eberre b rh bier rbire bir b er b erb relrb ne e re blrr e rb erb r br berenrrr ee
erb rbrrre brrre b erbeerslrebr biel blrr rire bir be b erbes rber b b er bir eer b r r brb e s r b bee brrrrorber oerer br
rr bher rb rr borer rb ererbh b soeb l s r b soeberr e b er rersb er srre r lr be bher rbhlr ehr r b r rbr b r b rhberbheer?hb ?reb
er rb r l br b r r r tr
Epoch 26/30, Loss: 0.9611, Perplexity: 2.6146
Predicted Text: [ l l ald r ] aa aln aean aarte aar auu an aas hnue ang a n alt t aura e aoen o a a uagaerha gtoeh aae ga
nllon eane a e e aen anan an u aanl alun ea neae ne aerl wieu wue win w en w e w g el aaue ea al n eg w e wa dwu we e ue ee
e w uw e e w n neg w e weecarinewnd aiel wr l d u i e win ae a e a es aeng au a en wrl eanga ne d ana eus nd au aee aa ual aer oe an
au aieu a n ahuor a e e ain a rheae r a rheae n l a ena?l? a n ue n eu ae aheu ?aa l eiu a u an a n a uhae ainheu?ha ?hea
en aa ahul au a ru au
Epoch 27/30, Loss: 1.0481, Perplexity: 2.8521
Predicted Text: [ cir cuos yoo] s als a en ae hra ae handoor asud acdaaycdoanyns ouaoerooerer oar o uidoor oora rodn o wordw
ciliw ooce w orenwharwctwfa w efaalrin al r ndais eanenroadraa alauarauea alra a ar a araduaaraue gare warr gd w erwar wurwerenrue aes
erw uw s w e c sda w oruseer racewd wisa w ard uire wicawaowhorwer rwond awu w oc wiraaocdw cird wchouracd wu wosawatuatwarhaetet wr
au waeuatw t waeuaatwhetwacawaraeasarat warawatacahaowaeca a aawhec a ueachsu weowheuarwiraueia w u wacawacow uhaa winheu?hwa?awea
eca w hua wu waoul dwu
Epoch 28/30, Loss: 0.9260, Perplexity: 2.5245
Predicted Text: [ k e e ] h h e he e t h h a an a e e e a a h a e a e l e a e t a
a e a e a e a a a a a aht eae ea aie aie ai ae aea elae e eal aea a rae e r e e a
an e a a e a e e l e a a l l a l i e ai ae a enae naen a rae ainle a n a a e a rae a a laae eea a a a
e aa an ae aa ea aa n a eea a a eaeaa l a en rl a n r n a ce aah l eh c cn c n c ce che e r c reec en c
l c c atc r
Epoch 29/30, Loss: 0.8241, Perplexity: 2.2798
Predicted Text: [ ml o mo ] m ,ml r m emme m rmo mmeu mnometn muo ml, myl, monyn mu m e moel r m r m ,h,men m r r ,n e mhergm
lllml ehl m er rmol rmlnsa m an mhter m nt r mhn emerte agr moeu rmou mot m et m enm elrmo , r, mlt e m er m rk men, n , e
e,rm m e e nt teg m erme,,r t,mtg moe m rg uor, mot me m erm, rmlg mu m et mor e tgm terg atm eu tg mu mee mernelrme eer r m
ngaeu rm rl men r m er rmol m eemoer m eemerl m lelo l m lg nol sng m m e mh , eh ogm grml mol mg me mhn e r m room e
l m nl an m n rgau
```

Figure 5: Qualitative Evaluation - Small Model Character Tokenizer

Regarding qualitative evaluation, although the model seems to converge, I don't think it generates well. Even though it starts from the premise of being a model without pretraining, the results don't seem to illustrate anything very coherent. Different word forms are not being formed with each epoch; the model tends to develop a particular style in which it prioritizes certain letters more or less. However, there are instances where it seems to attempt to write different words.

It's not possible to talk about a BLEU score because appropriate words can't be identified, nor are the letters arranged in an intuitive, suitable order.

3.2 Big Model - Character Tokenizer

This time, I tried to scale up the model. Thus, the configuration looked as follows: batch size of 64, max sequence length of 512, embedding dimension of 512, attention heads and layers equal to 8, and hidden dimension of 1024. Additionally, I used the Adam optimizer with a learning rate of 1e-4. Everything was trained for 90 epochs, using a GPU P100, with the condition to fit within the 12 hours of resources provided by Kaggle.

Thus, the number of parameters increased from 786,432 to 25,165,824.

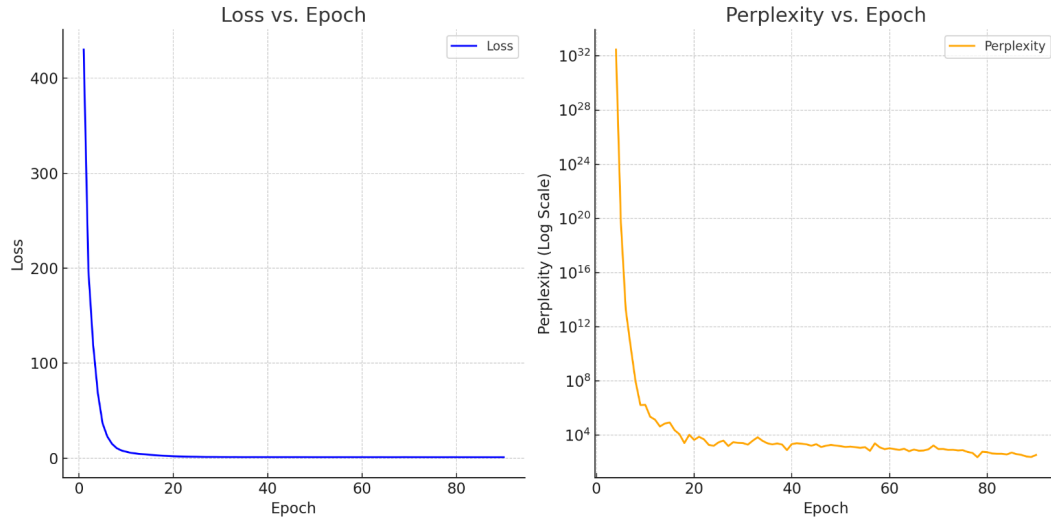


Figure 6: Big Model - Character Tokenizer

As can be seen in the figure above, the model saturates. Perhaps, if I had trained it for a longer period, I would have reached lower scores, as seen in Figure 7, but in a way, with this configuration, these were the results.

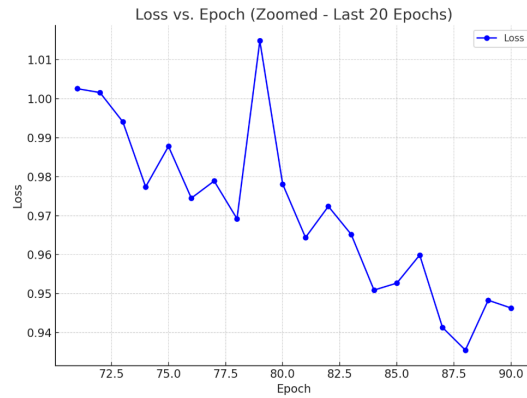


Figure 7: Big Model - Character Tokenizer - Loss

Qualitative Evaluation

```
Epoch 85/100, Loss: 0.9527, Perplexity: 509.6111
Predicted Text: f bi lee dn o o y rh yeshjeth m thmo w s wul ehsot w d mwa t w sash,woowha loel r he heoldher h e dsea lia dftaalt eit f e fttta wtuo
be lo wblila elewb: dliit ala beefd ip fiewp fiw p fircp f ec f e fdnl fano oa o f ceodhf e fa dfrofe or eeeo f rfoohe fihlode f e fo ol l fitefio efo dewi
ficfeef e f le fedderfwar ec fi aneacdf ce d ftf enledefnofeoe ne feb ee e fbeodaiooe ae teauo ee a e aiteaeboeae ebe aeboeae te aeeta eea tdeou at
ud ae f eue fa eelna dreed fterat fdu fe fi eul faloeefeta f u efu feeue dnuoe
Epoch 86/100, Loss: 0.9599, Perplexity: 391.9939
Predicted Text: [ ihifi-yl l d l ln winro wcan hal hrhol favhnoftn xnoo flo,fyloof nynth foofharwaroo a ha,oiolrhao ih egnoah oia ol oiaalltioia t foly fifo
luh t shhr nla frin hini h nrish rll hi firhif l hili h l h rh ,o fh fi n i h lin oh rh r hfuh rirvfil nish fhvno hloln h rhniror lihl ihit iho r ifiri
hil hooh rhirir l l nfoh l hir yl h loda hnh fr vo,hvft t hofe rhm k d r, hruf hi fr,h,rr,hef oia d dhir:hre hytir:d hre hode: oh: ryd d:ih roivf yv
tfo hooh f dhid : lfydah,fodhvihyupohr hodhit fr hyne h vyah d f ,hf hiof hf -
Epoch 87/100, Loss: 0.9413, Perplexity: 345.1668
Predicted Text: feds a l o w s montor n n san oin sau s surr doo lo on lldad smos y son e he i slts to tao tllien ant t t ton tnte
o t t omally mionsv ant m vle a ve me o owo wet w t w ur oolw u n v ulvnoow ve vom w w w oe n w w na umannoe w wen wa in wjowenlowalwoe ew
wenno w me oae n ome w t umale et v t v um o n oao w n unaoalwoe c v vo l a le viondic a vih v vien i c iomp pv ip c iouphl ip novipppi n p on
n il l pveivp e ov lo vimpion i hilvie io c ip novi v loi ipa sv i u
Epoch 88/100, Loss: 0.9355, Perplexity: 265.0152
Predicted Text: f me roore lrlow wy sir o wrothtoph? tophl fillalup offt lc llaft r l o pffo erpo voo h hh fh hush h e, ooh oho ,oeeoohcho h o oc crf f
qofahnr hifo,r hahoborne h rh hnef nhf, alm a o a e a fel o f s ae n s,oa e a jafae , of, es, a faeoe anons, a e a s, ocan,hn,hhsa,hoe fh , hhn
he h e h,c, hen , hfjh en hh e a n h n hvbhttr d ,htlt hr tre her kt l , hgrf htf h, d,heflr h t l hhd h geth t g h gethe d e h t d e h ld efl d tf
he h tf hh e thf h,f hd h d h f he hnt tfg h geth td hl fe,hfjh rf hfee
Epoch 89/100, Loss: 0.9483, Perplexity: 242.0783
Predicted Text: f ktn orve kee e all o eees eeb r ehb wuu wioaes esuo abo aeboohoses huoe e ehebee e esueh ee eeo e plisee hhe ghnoohsehnoshie o hto hno
u snouolhno geub gghneegs beogg h gheu ghua ,hnn heenehee hluad houaes a hd noughree ho ghuae a ouao esa heuhosheehnhng ehee hsaa hrenaangoohsdodhd lofh a
hohn oeeoeo oaro oenlo ofoeeoneh d eengemeto orh euf ro,hoes hotuoether he o hrsuuhueu h r heuos h e o hhr h reelon r h rhehe r doh ero d h oaro ouoo
ruo heah eu hh d ehua ohuao hr horohou hethhn eur horeeh erotot udohou suetohuoe
Epoch 90/100, Loss: 0.9463, Perplexity: 341.1633
```

Figure 8: Qualitative Evaluation - Big Model Character Tokenizer

I would say that, again, although word forms are not clearly shaped, the generations seem a bit more varied. A positive sign for generating diverse words that are appropriate in their context.

3.3 Small Model - Subword Tokenizer

For training using this tokenizer, I used the same text preprocessing strategy. The initial configuration for a 'small' model was: embedding dimension equal to 128, hidden dimension equal to 128, the number of attention heads and the number of layers equal to 2, max sequence length equal to 512, and batch size equal to 16. The learning rate remained at $1e-4$, and the optimizer was of the Adam type. The model has around 13.2 million parameters.

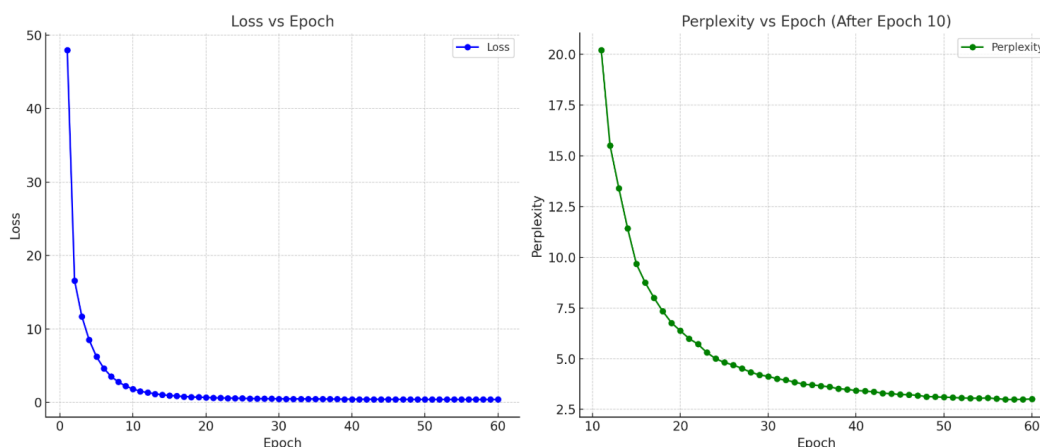


Figure 9: Small Model - Subword Tokenizer

As can be seen, the model converges again, and the perplexity decreases to values around 3. At this score, there is an expectation of some form of logic among the generated words. However, the BLUE score remains minimal, around 0.0001.

3.3.1 Qualitative Evaluation

```
Predicted Text: [ dmore ] i lorde, andward,,,ze, and i,, all king,,,orkight, and,wart, is, man,,ither,,,ius the incl, and,,, and i king,imer, and, king, the,sh
ire, the you, king, i,ouwer, and,, king,, the ish,,, and man, the,,, is, and the,o is, misuse, and,,, transformation, and, leshuomen,, garden,, not,, loyaltyire
s,,, thear

Epoch 57/60, Loss: 0.3985, Perplexity: 2.9918, BlueScore: 0.0001

Predicted Text: [ kingmoreland ] i lorde, and repro,,, seeking, and i,, the king,,,orkight, employ, thewart's, man, visitingretch,,, a,, and,, eve and i king,in
er, and, king, the,shire, specify pictures, king, i,ouwer, and, eve, king,, the ish,,, and man marks epid,,,work merchant and the,,o's unc misuse, and misuse of wh
ist transformation and, leshuomen, and much,, gone,, loyaltyires of,, epidar

Epoch 58/60, Loss: 0.3978, Perplexity: 2.9905, BlueScore: 0.0001

Predicted Text: [ firstmore ] i lorde, and,,, seeking, and i,, the king,,,orkight, the, thimwart's, man,,retch,,, the,, and,,, and i king,imer, and, king, th
e,shire, the pictures mistrust king, i,ouwer, and,, king,, the ikin,,, and man welcomes the,,, and, and the,arers's, misuse, and,, shameless absent, and, leshuom
en, and garden,, not and,,ires,,, thear

Epoch 59/60, Loss: 0.3955, Perplexity: 3.0032, BlueScore: 0.0001

Predicted Text: [ kingmoreland ] i lorde brow i opposition,,, seeking, i i,, the world,,,orkshright, and, athwart, is, man, whenceales,,, him,, and,, casting
and, world,imer, and, world, the,shire, the you, world, i turfoower, and casting, world,, the,sh,,, and man, epid funeral,,work, and the,,o is casting misuse, an
d,,, transformation, and, leshuomen,, full,, a,,,ires in,, thear
```

Figure 10: Qualitative Evaluation - Small Model - Subword Tokenizer

In terms of qualitative evaluation, it can be observed how the model generates the author's name very well, and also how it uses words from adjacent lexical fields or related domains close to the main domain. There are still issues, but one can see early signs of more qualitative generation.

3.4 Big Model - Subword Tokenizer

In this case, I tried to use a larger model, which was trained for 11 hours using a P100 GPU. It was configured with: 32 batch size, 384 embedding dimension, 4 attention heads and layers, 128 hidden dimension, and a maximum sequence length of 256. The learning rate was $3e-4$ with the Adam optimizer.

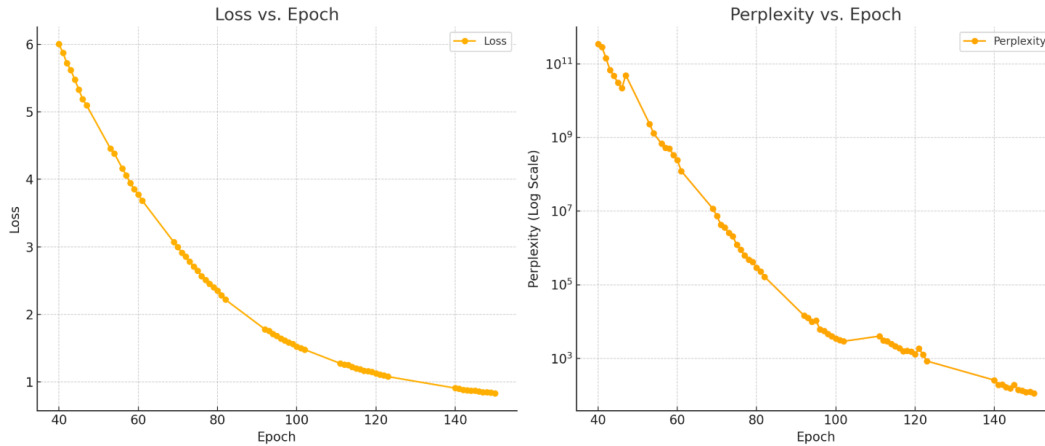


Figure 11: Big Model - Subword Tokenizer

Thus, the model size has increased to 20.1 million parameters. As can be seen, there is very good convergence. And if I had more resources, I would definitely have let the model train for more epochs. However, perhaps increasing the learning rate would have also helped.

3.4.1 Qualitative Evaluation

```
Epoch 146/150, Loss: 0.8578, Perplexity: 138.6421, BlueScore: 0.0001
Predicted Text: [ kingmoreland ] i lordge him and lodging as ael this, and, limitless clouds heels.'d
with asesternight oath and it athwart wallance, place gates roaked load,st you and, and quality deal c
ried and shade rab ruinimer, andents rab must clouds issire ] ire,, rabth spritessendower'd and a int
er rab,. clouds thesh regard mistrust, and place- wit sprites midnight,red, and my quality corpse, is
a'd, and'd, shameless,, and cried-shwomen,, i move, thrown, shores sitsold. was, heavenly times
Epoch 147/150, Loss: 0.8461, Perplexity: 132.9305, BlueScore: 0.0001
Predicted Text: [ tammoreland ] i lordge. and present was hot and the. and i limitless his vol,,, iest
ernight thousand and i athwarty lur, debt. theales load,,al,, and worst,hes and magical bits mortimer,
andents bits, his,shire, the, the bits y vent,endower, and never the bits,. the ish breeds seeks. and
retreat to the nap old ired govern and whose worst corpse a be never misuse. and misuse, shameless tra
nsformation, and death welshwomen,, i move,,, retold.' spoken, the,
Epoch 148/150, Loss: 0.8483, Perplexity: 122.1640, BlueScore: 0.0001
Predicted Text: [ protomoreland ] i necessge serve thorough bos for hot, the hears thorough i limits,
this question, in, iesternight thousand somew mine athwart came, sh, whenceitten load.. you news e and
worst, past and i irregular mortimer. and, irregular, this comesshire, the, the irregular irregular,,o
wer, and a lawful irregular with, this ish breeds taken, and barren, late weapon, letred, and whose wo
rst corpse, came a interim. and interim, shameless transformation, and your welshwomen, i be,,, more r
etold, faith, a,
Epoch 149/150, Loss: 0.8390, Perplexity: 123.9837, BlueScore: 0.0001
Predicted Text: [ benmoreland ] i lordge's and army! a, the, and i limits forbid the heels.'d, iestern
ight wounds i quiet athwart's,, mile? theales load, dick a,, and,, a and i fatal,imer, andents fatal,
the issire, the, thanks fatal mer, we,ower, and a the fatal,, the ishilla,, and goose, wit majesty's i
red govern and my, corpse a, a misuse, and misuse malice shameless transformation, and the welshwome
n,, i. i thou slander more retoldin false, the.
Epoch 150/150, Loss: 0.8323, Perplexity: 113.5560, BlueScore: 0.0001
```

Figure 12: Big Model - Subword Tokenizer

In terms of qualitative analysis, the model delivers very good results even with a perplexity score of 113. Yes, it's not a logically coherent text, but at least the author is generated very intuitively, and the rest of the words seem to form some sort of connection between them. By comparison, this model is not much larger than the previous one, but it seems to have more potential.