

Models In a Spelling Bee: Language Models Implicitly Learn the Character Composition of Tokens

Anonymous ACL-IJCNLP submission

Abstract

Standard pretrained language models operate on sequences of subword tokens without direct access to the characters that compose each token’s string representation. We probe the embedding layer of pretrained language models and show that models learn the internal character composition of whole word and subword tokens to a surprising extent, without ever seeing the characters coupled with the tokens. Our results show that the embedding layer of RoBERTa holds enough information to accurately spell up to a third of the vocabulary and reach high average character n gram overlap on all token types. We further test whether enriching subword models with additional character information can improve language modeling, and observe that this method has a near-identical learning curve as training without spelling-based enrichment. Overall, our results suggest that language modeling objectives incentivize the model to implicitly learn some notion of spelling, and that explicitly teaching the model how to spell does not enhance its performance on such tasks.

1 Introduction

Contemporary subword tokenization algorithms such as BPE (Sennrich et al., 2016) partition a string into contiguous spans of characters. Each span represents a frequent character n gram, from individual characters (a), through prefixes (uni) and suffixes ($tion$), and even complete words ($cats$). The tokenizer then converts each such span into a discrete symbol (a token) with no internal structure, effectively discarding the token’s orthographic information. Therefore, a model operating over sequences of subword tokens should be oblivious to the spelling of each token. In this work, we show that despite having no direct access to the subwords’ internal character composition, pretrained language models *do* learn some notion of spelling.

To examine what pretrained language models learn about spelling, we present the *SpellingBee* probe. SpellingBee is a generative language model that predicts the character composition of a token given only its (uncontextualized) vector representation from the pretrained model’s embeddings matrix. SpellingBee is trained on part of the model’s vocabulary, and then tested by spelling unseen token types. If the probe can successfully reconstruct the correct character sequence from an unseen token’s embedding, then there must be significant orthographic information encoded in the vector.

We find that the embedding layers of several pretrained language models contain surprising amounts of character information. SpellingBee accurately spells 31.8% of the held-out vocabulary for RoBERTa-Large (Liu et al., 2019), 32.9% for GPT2-Medium (Radford et al., 2019), and 40.9% for the Arabic language model AraBERT-Large (Antoun et al., 2020). A softer metric that is sensitive to partially-correct spellings (chrF) (Popović, 2015) shows a similar trend, with 48.7 for RoBERTa-Large and 62.3 for AraBERT-Large. These results are much higher than the baseline of applying SpellingBee to randomly-initialized vectors, which fails to spell a single token.

Given that subword models learn some notion of character composition to fulfill language modeling objectives, could they perhaps benefit from knowing the exact spelling of each token a priori? To that end, we reverse SpellingBee’s role and use it to pretrain the embedding layer of a randomly-initialized model, thus imbuing each token representation with its orthographic information before training the whole model on the masked language modeling objective. We compare the pretraining process of the character-infused model to that of an identical model whose embedding layer is randomly initialized (and not pretrained), and find that both learning curves converge to virtually identi-

cal values within the first 1000 gradient updates, a fraction of the total optimization process. This experiment suggests that while language models may need to learn some notion of spelling to optimize their objectives, they can quickly acquire all the character-level information they need without directly observing the composition of each token.

2 Spelling Bee

To measure how much a model knows the character composition of its tokens, we introduce SpellingBee, a generative probe that tries to spell out a token character-by-character. Specifically, SpellingBee probes the original model’s *embedding matrix*, since spelling is a property of token *types*, invariant to context. For example, given the embedding of the token *cats*, SpellingBee will try to generate the sequence $[c, a, t, s]$. We do so by modeling SpellingBee as a character-based language model, where the first token is a vector representation of the vocabulary item.¹

Training We split the vocabulary to train and test sets,² and use teacher forcing to train SpellingBee. In the example of *cats*, SpellingBee will compute the following probabilities:

$$\begin{aligned} P(x_1 = c \mid x_0 = \text{cats}) \\ P(x_2 = a \mid x_0 = \text{cats}, x_1 = c) \\ P(x_3 = t \mid x_0 = \text{cats}, x_1 = c, x_2 = a) \\ \vdots \end{aligned}$$

All of SpellingBee’s parameters are randomly initialized. The only parameters that are pretrained are the token embeddings (e.g. the representation of *cats* or *a*), which are taken from the original pretrained language model we intend to probe, and treated as constants; i.e. kept frozen during SpellingBee’s training.

Inference & Evaluation Once SpellingBee is trained, we apply it to the test set using greedy decoding. For each vocabulary item w in the test set, SpellingBee is given only the corresponding embedding vector e_w , and is expected to generate the character sequence w_1, \dots, w_n that defines w . We measure success on the test set using two metrics:

¹Some vocabularies have symbols for indicating preceding whitespaces (.) or that the next token is part of the same word (##). SpellingBee learns to predict these symbols too.

²We test various train/test splits to ensure the robustness of our findings. See Section 3 for more detail.

exact match (EM), and character *n*gram overlap score using *chrF* (Popović, 2015). While EM is strict, chrF allows us to measure partial success. We also report edit distance using Levenshtein distance ratio in Appendix A.

SpellingBee for Pretraining Embeddings

While we mainly use SpellingBee as a probe, a variation of our method could potentially imbue the embedding layer with character information before training a language model. We could train a probe with randomly-initialized embeddings (instead of pretrained embeddings from another model) to predict the spelling of *all* vocabulary items, and use these trained probe embeddings to initialize any target model’s embedding layer (instead of random initialization). We experiment with this method in Section 5, but find that it does not have any significant impact on the convergence of language models.

3 Experiment Setup

We begin with a series of probing experiments, where we apply SpellingBee to the embedding layer of various pretrained models.³

Pretrained Models We probe four pretrained models: RoBERTa-Base and Large (Liu et al., 2019), GPT2-Medium (Radford et al., 2019), and AraBERT-Large (Antoun et al., 2020). This set introduces some diversity in vocabulary, objective, and scale: the first three models are trained on English corpora, while AraBERT is trained on text in Arabic; GPT2 is an autoregressive language model, while the rest are masked language models; RoBERTa-Base consists of 125M parameters (with 768 dimensions per embedding), while the other models have approximately 350M parameters (with 1024 dimensions per embedding).

Control Since SpellingBee is a *trained* probe, we wish to establish the probe’s baseline performance when provided with inputs with no orthographic information. As an empirical *control*, we train and test SpellingBee on randomly-initialized vectors, in addition to the main experiments where we utilize the pretrained embedding layers.

Training & Testing Data We split the vocabulary into training and testing data using the following protocol. First, we randomly sample 1000 token types as test. We then filter the remaining

³Hyperparameters are detailed in Appendix E.

	Filter	RoBERTa		GPT2	AraBERT
		Base	Large	Medium	Large
EM	None	27.3	31.8	32.9	40.9
	Similarity	15.7	18.2	17.9	21.9
	Lemma	15.7	17.7	16.5	19.5
	Control	0.0	0.0	0.0	0.0
chrF	None	44.7	48.7	51.6	62.3
	Similarity	32.7	35.1	36.4	46.0
	Lemma	32.6	34.8	35.2	43.9
	Control	7.0	7.0	7.0	7.0

Table 1: The percent of token types that can be spelled out exactly (EM) from their embeddings by SpellingBee, and the n gram overlap between SpellingBee’s reproductions and the token types’ true spellings (chrF). The first three rows reflect different methods for filtering the training data, and the fourth represents the control experiment, which uses randomly initialized embeddings. All SpellingBee instances in this table are trained on 32000 examples.

vocabulary to eliminate tokens that may be too similar to the test tokens, and randomly sample 32000 training examples. We experiment with three filters: *none*, which do not remove tokens beyond the test-set tokens; *similarity*, which removes the top 20 most similar tokens for every token in test, according to the cosine similarity induced by the embedding vectors; *lemma*, which removes any token type that shares a lemma with a test-set token (e.g. if *diving* is in the test set, then *diver* cannot be in the training set).⁴ The lemma filter always applies the similarity filter first, providing an even more adversarial approach for splitting the data. To control for variance, we create 10 such splits for each model and filter, and report the averaged evaluation metrics over all 10 test sets.

4 Results

Main Result Table 1 shows how well SpellingBee can spell a vocabulary token using only its frozen pretrained embedding. We observe that SpellingBee is able to accurately recover the spelling of up to 40.9% of the test set, while the control is unable to spell even a single word correctly. A similar trend can be seen when considering the finer character n gram metric (chrF). Manually analyzing the predictions of the control baselines (see Appendix D) indicate that it primarily generates combinations of frequent character sequences, which mildly contributes to the chrF score,

⁴We lemmatize using NLTK’s WordNet lemmatizer (Bird and Loper, 2004) for English and Farasa’s Stemmer (Darwish and Mubarak, 2016) for Arabic.

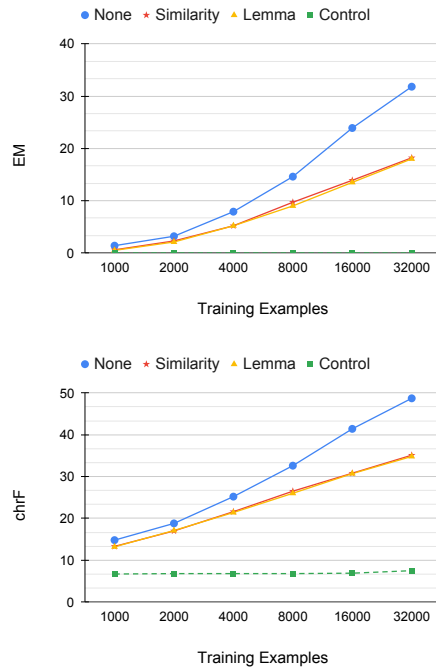


Figure 1: The amount of character information SpellingBee is able to extract from RoBERTa-Large, as measured by EM (top) and chrF (bottom), given different quantities of training examples.

but does not affect EM. These results are persistent across different models and filters, strongly indicating that the embedding layer of pretrained models contain significant amounts of information about each token’s character composition.

One may suggest that training SpellingBee over 32000 examples may leak information from the test set. For example, if *dog* was seen during training, then spelling out *dogs* might be easy. We thus consider the similarity and lemma filters, which remove such near-neighbors from the training set. While results are indeed lower (and probably do account for some level of information leakage), they are still considerably higher than the control, both in terms of EM and chrF. Results using the similarity and lemma filters are rather similar, suggesting that embedding-space similarity captures some information about each token’s lemma.

Finally, we find that the properties of pretrained models also seem to have a significant effect on the amount of spelling information SpellingBee can extract. Larger models tend to score higher in the probe, and the model trained on text in Arabic appears to have substantially higher EM and chrF scores than those trained on English corpora. One possibility is that Arabic’s rich morphology

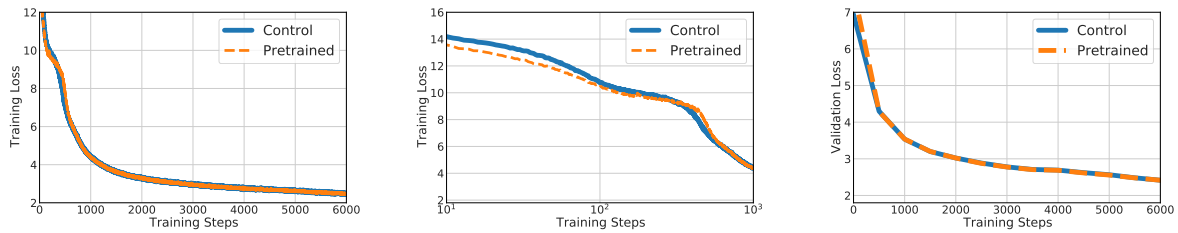


Figure 2: The overall training loss (left), first steps of training loss (center), and validation loss (right) of RoBERTa-Large, when training on the masked language modeling objective with embeddings pretrained by SpellingBee (*pretrained*) and randomly-initialized embeddings (*control*).

incentivizes the model to store more information about each token’s character composition, however it is also possible that AraBERT’s different vocabulary, which allocates shorter character sequences to each token type, might explain this difference (we discuss the link between sequence length and accuracy later in this section).

Overall, our probing experiments show that even though subword-based language models do not have direct access to spelling, they *can* and *do* learn a surprising amount of information about the character composition of each vocabulary token.

Probing with Less Training Data We further examine whether SpellingBee can extract information when trained on less examples. Figure 1 shows how well SpellingBee can spell RoBERTa-Large’s vocabulary when trained on varying amounts of data, across all filters. We find that more data makes for a better probe, but that even a few thousand examples are enough to train SpellingBee to extract significant character information from the embeddings, which *cannot* be extracted from randomized vectors (the control).⁵

5 Pretraining Language Models to Spell

Our probing experiments reveal that language models learn some partial notion of spelling, despite the lack of direct access to characters. Therefore, we hypothesize that learning to spell is beneficial for language models, and propose pretraining the embedding layer using a variant of the SpellingBee probe described in Section 2. Here, the goal is to imbue each embedding with enough information for SpellingBee to accurately generate its surface form, and then initialize the language model with the pretrained embeddings before it starts training on the language modeling objective.

⁵We provide additional analysis on spelling accuracy by subword frequency and length in Appendices B and C.

We apply this process to RoBERTa-Large, training the model’s embedding layer with SpellingBee using the same hyperparameter settings from Appendix E, with the key difference being that the embeddings are now tunable parameters (not frozen).⁶ We train RoBERTa-Large on English Wikipedia using the hyperparameter configuration of 24hBERT (Izsak et al., 2021), and cease training after 24 hours (~16000 steps). For comparison, we train exactly the same model with a randomly-initialized embedding layer.

Figure 2 shows the masked language modeling loss with and without pretrained embeddings. We see that the curves quickly converge into one. After only 1000 training steps, the difference between the validation losses never exceeds 0.01. This result indicates that the model does not utilize the character information injected into the tokens’ embeddings. Along with the results from Section 4, we conjecture that the model learns an implicit notion of spelling during pretraining, which is sufficient for masked language modeling, and does not benefit from explicitly adding orthographic information.

6 Conclusion

This work reveals that pretrained language models learn, to some extent, the character composition of subword tokens. We show that our SpellingBee probe can spell many vocabulary items using their uncontextualized embedding-layer representations alone. Trying to explicitly infuse character information into the model appears to have a minimal effect on the model’s ability to optimize its language modeling objective, suggesting that the model can independently learn all the character-level information it needs for the task.

⁶To verify that this process does indeed encode the tokens’ spellings into the embeddings, we apply a SpellingBee probe (using a different random initialization) to the learned embeddings, which yields 93.5% EM on held-out token types.

References

- 400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
- Wissam Antoun, Fady Baly, and Hazem Hajj. 2020. [AraBERT: Transformer-based model for Arabic language understanding](#). In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 9–15, Marseille, France. European Language Resource Association.
- Steven Bird and Edward Loper. 2004. [NLTK: The natural language toolkit](#). In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pages 214–217, Barcelona, Spain. Association for Computational Linguistics.
- Kareem Darwish and Hamdy Mubarak. 2016. [Farasa: A new fast and accurate Arabic word segmenter](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 1070–1074, Portorož, Slovenia. European Language Resources Association (ELRA).
- Peter Izsak, Moshe Berchansky, and Omer Levy. 2021. [How to train bert with an academic budget](#). *arXiv preprint arXiv:2104.07705*.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). *CoRR*, abs/1412.6980.
- Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *arXiv preprint arXiv:1907.11692*.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. [fairseq: A fast, extensible toolkit for sequence modeling](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Maja Popović. 2015. [chrF: character n-gram F-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. [Language models are unsupervised multitask learners](#). *OpenAI blog*, 1(8):9.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- 450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499

A Levenshtein Distance

Levenshtein distance (Levenshtein et al., 1966) is an edit distance metric that, given two strings, calculates the minimal number of changes needed to be done in order to make the two strings identical. Levenshtein distance *ratio* is the length-normalized version, which is computed by adding the sum of lengths of both strings to the edit distance and dividing by the same sum of lengths. We report the main experiment’s results using this ratio in Table 2.

Filter	RoBERTa		GPT2	AraBERT
	Base	Large	Medium	Large
None	69.7	72.7	74.4	83.6
Similarity	61.5	63.7	64.5	75.8
Lemma	61.4	63.3	63.7	74.8
Control	25.6	26.4	27.0	25.7

Table 2: Levenshtein distance ratio. The first three rows reflect different methods for filtering the training data, and the fourth represents the control experiment, which uses randomly initialized embeddings. All SpellingBee instances in this table are trained on 32000 examples.

B Spelling Accuracy by Frequency

We test whether pretrained models tend to store more spelling-related information in higher-frequency token types. We focus on RoBERTa-Large, and assign each token in the test set to its frequency quintile according to the number of times it appeared in the pretraining corpus – from the 10000 most frequent token types (top 20%) to those ranked 40000-50000 in the vocabulary (bottom 20%) – and measure the average performance of SpellingBee within each quintile. Figures 3 and 4 shows the results with and without the similarity filter. We observe that SpellingBee is indeed able to extract more information from higher-frequency token types, suggesting that the pretrained model has more information about their character composition.

C Spelling Accuracy by Length

We analyze the effect of token length on the probe’s ability to spell. A priori, it is reasonable to assume that it is easier for the probe to spell shorter tokens, since less information needs to be extracted from the embedding and there are less discrete decisions to be made while decoding. Indeed, Figure 5 shows that with the none filter most vocabulary

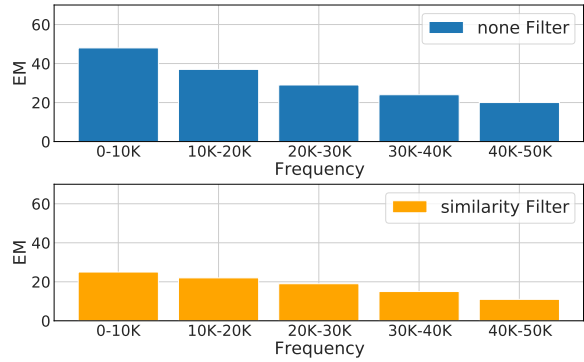


Figure 3: The EM scores of SpellingBee on RoBERTa-Large for each frequency quintile with the *none* filter (top) and the *similarity* filter (bottom).

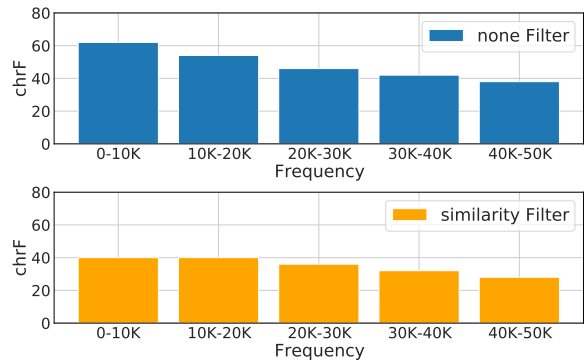


Figure 4: The chrF scores of SpellingBee on RoBERTa-Large for each frequency quintile with the *none* filter (top) and the *similarity* filter (bottom).

tokens with 2-4 characters can be accurately reproduced from their vector representations, while longer tokens are harder to replicate. This trend is particularly sharp when the similarity filter is applied, as the probe is hardly able to spell tokens with 6 or more characters accurately; having said that, the probe is able to generate many *partially correct* spellings, as measured by chrF (Figure 6). Perhaps a less intuitive result is the probe’s failure to spell single-character tokens. A closer look reveals that many of these examples are rare or non-alphanumeric characters (e.g. ζ and $\$$), which are probably very difficult for the probe to generate if it had not seen them during training. While these results show strong trends with respect to length, token length is also highly correlated with frequency, and it is not necessarily clear which of the two factors has a stronger impact on the amount and resolution of character-level information stored in the embedding layer of pretrained models.

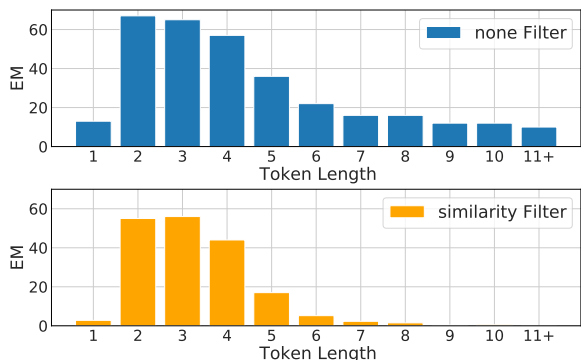


Figure 5: The EM scores of SpellingBee on RoBERTa-Large for each token length with the *none* filter (top) and the *similarity* filter (bottom). The rightmost column groups together tokens with length of 11 or above.

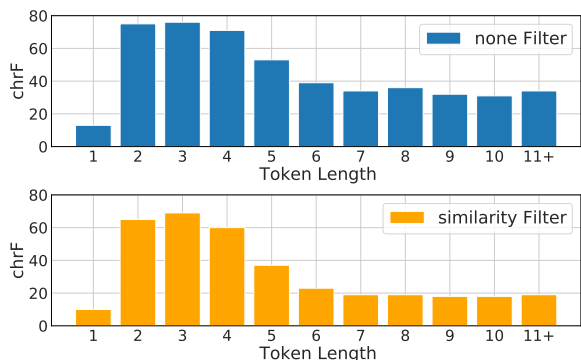


Figure 6: The chrF scores of SpellingBee on RoBERTa-Large for each token length with the *none* filter (top) and the *similarity* filter (bottom). The rightmost column groups together tokens with length of 11 or above.

D Manual Error Analysis

We manually analyze 100 random tokens that SpellingBee spelled incorrectly with the lemma filter to understand the nature of the spelling mistakes. Out of those 100 we display 20 mistakes in Table 3 alongside the spelling prediction of the control baseline. SpellingBee’s mistakes vary from single-character typos to completely different words. Having said that, the vast majority of mistakes have significant overlap with the correct spelling, such as shared prefixes and capitalization.

E Hyperparameters

We implement SpellingBee with a 6-layer encoder-decoder model, with 512 model dimensions. The model parameters are optimized with Adam (Kingma and Ba, 2015) for 1000 steps with up to 1024 tokens per batch, a learning rate of 5e-4, and a dropout rate of 0.1. These are the default hyperpa-

Token	SpellingBee	Control
_Issa	_Asey	_kinston
_Rhod	_Rob	_hoedn
Memory	_Mathinge	_entically
_metals	_metrys	_leaved
_Reed	_Redd	_fomparing
_break	_breach	_promoters
_summit	_mosump	_seasons
Catholic	Cravital	_tonversal
_cleanup	_lamed	_paclus
_Winner	_Womer	_purden
_LIM	_LUM	_Send
Copy	Cople	_providers
_voicing	_relicing	_walking
_Stab	_Stamb	_hovidors
_356	_353	_budiance
find	wive	_malding
_Psychic	_Syptanc	_joacter
_Looking	_Lowing	parging
CLOSE	DEFIC	_tuldence
_prolific	_promistic	_complexement

Table 3: Sampled SpellingBee errors with the lemma filter alongside the control baseline’s spelling for the same tokens. The underscore (.) represents a preceding whitespace.

rameters for training a transformer language model in Fairseq (Ott et al., 2019).