Unsupervised Morphological Tree Tokenizer

Anonymous ACL submission

Abstract

001 As a cornerstone in language modeling, tokenization involves segmenting text inputs into 003 pre-defined atomic units. Conventional statisti-004 cal tokenizers often disrupt constituent bound-005 aries within words, thereby corrupting semantic information. To address this drawback, we introduce morphological structure guidance to to-007 kenization and propose a deep model to induce character-level structures of words. Specifically, the deep model jointly encodes internal 011 structures and representations of words with a mechanism named MorphOverriding to ensure 012 the indecomposability of morphemes. By training the model with self-supervised objectives, our method is capable of inducing characterlevel structures that align with morphological rules without annotated training data. Based on the induced structures, our algorithm tokenizes words through vocabulary matching in a top-down manner. Empirical results indicate that the proposed method effectively retains 022 complete morphemes and outperforms widely adopted methods such as BPE and WordPiece 024 on both morphological segmentation tasks and language modeling tasks. The code will be released later.

1 Introduction

027

034

040

Tokenization, the initial step of language modeling, segments natural language into manageable units. While this process is crucial for representing natural language, research on new tokenization methods has remained limited, particularly in contrast to the rapid advancements in language model architectures and learning approaches. Currently, the de-facto tokenizers are BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012), which have been widely adopted by stateof-the-art language models such as GPT (Radford et al., 2019) and BERT (Devlin et al., 2019). However, numerous studies have challenged these methods (Bostrom and Durrett, 2020; Church, 2020;



Figure 1: BPE (top) tokenizes a word through a bottom-up greedy merging approach given pre-learned merge operations, while ours (bottom) tokenizes a word via a top-down vocabulary matching while traversing a global parse tree.

Hofmann et al., 2021; Minixhofer et al., 2023), arguing that they cannot adequately capture linguistic information. They often disrupt constituent boundaries within words, leading to unnatural and fragmented token representations. Figure 1(top) demonstrates an example where BPE fails to identify the appropriate boundaries in a word.

According to linguistic theories, both words and sentences are believed to have internal structures (Selkirk, 1982; Marvin, 2002; Cotterell and Schütze, 2015). While sentence-level grammar induction methods based on deep neural networks are highly effective, whether these methods can be applied equally well to words remains underexplored. In this work, we systematically evaluate neural grammar induction methods at the word level, propose a hypothesis explaining their suboptimal performance, and empirically validate this hypothesis. Building on these insights, we introduce the first effective unsupervised neural grammar induction model at the word level and present a more morphologically aligned tokenizer that leverages our model as shown in Figure 1(bottom).

Our approach draws inspiration from syntactic composition models (Maillard et al., 2017), where a sentence is encoded as a weighted sum over all composed root representations of its underlying binary parse trees via dynamic programming. Instead of composing a sentence from words, we apply

068

069

042

043

163

164

165

166

167

168

169

170

121

071composition models on characters in a word to in-
duce its morphological parse tree. To train the com-
position model, we propose two self-supervised
objectives akin to next token prediction and span
prediction that effectively leverage both contextual
information at the sentence level and semantic in-
formation at the subword level. Thus the model
can learn to assign higher probabilities to morpho-
logical constituents of a word and induce the un-
derlying morphological parse tree.

However, character sequences present a unique 081 challenge to composition models because morphemes, the smallest meaning-bearing units in a language (Jurafsky and Martin, 2009), are indecomposable. While we can represent a constituent by composing its sub-constituents in most cases, we cannot represent a subword by composing its components if the subword is a morpheme. For exam-088 ple, the meaning of windsurf can be decomposed to wind+surf, but wind is a morpheme whose meaning is not a function of its components. To 091 address the challenge, we propose a mechanism named MorphOverriding. During the bottom-up composition process in our model, upon identifying a subword that matches an entry in a heuristically constructed morpheme vocabulary, we compute the subword representation from both its components and the corresponding morpheme embedding, i.e., the model may learn to mix or override the com-099 position with the morpheme embedding. Our ex-100 periments show that such a mechanism is critical 101 in morphological structure induction. 102

Building upon the resolution of morphological structure induction, we introduce a novel tokeniza-104 tion algorithm named TreeTok, which includes both 105 vocabulary construction and word segmentation. 106 During vocabulary construction, TreeTok first uti-107 lizes a tree-based BPE variant to build an initial 108 vocabulary and then applies a tree-based Unigram 109 variant to prune the initial vocabulary to a speci-110 fied size. Because TreeTok operates in a top-down 111 manner, it does not need to retain all intermedi-112 ate tokens produced by merge operations in the 113 vocabulary as BPE does. By this means, we can 114 build a more compact vocabulary by pruning less 115 important subwords. During word segmentation, 116 we employ a lightweight parser with compact pa-117 rameters distilled from the composition model to 118 parse a word into a character-level binary tree and 119 then apply top-down vocabulary matching to en-120

hance the tokenizer's alignment to morphological structure, as illustrated in Figure 1.

In our experiments, we train TreeTok and baselines on the Wikitext-103 corpus (McClosky et al., 2006) and assess their performance on morphological segmentation tasks and language modeling tasks. Evaluation results indicate that TreeTok consistently outperforms BPE and WordPiece across all the tasks.

In conclusion, our contributions are three-fold:

- We conduct empirical study on character-level neural parsing, identifying its limitations and proposing a novel explanation—lack of MorphOverriding—to account for its suboptimal performance.
- Building on the MorphOverriding hypothesis, we introduce the first effective unsupervised neural model for character-level structure induction, addressing a critical gap in the field.
- We show that our character-level structure induction method can be integrated into mainstream tokenizers to significantly enhance their performance on morphological tasks.

2 Related Work

Subword Tokenizers. Subword tokenization, with typical methods such as BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012), has become customary in most NLP fields. BPE builds its vocabulary by repeatedly merging the most frequent subword unit pairs, whereas WordPiece selects pairs using the highest mutual information. During tokenization, BPE applies learned merge operations in the same order to new text initialized with characters while WordPiece iteratively finds the longest match in the vocabulary. Unigram (Kudo, 2018), another popular tokenizer, builds its vocabulary in the opposite direction: it starts with a large set of potential subwords and prunes them based on delta entropy in a unigram language model.

Our tokenizer aims to build upon the advantages of these effective statistical tokenizers and augment them with unsupervised induced tree structures.

Unsupervised Morphological Segmentation. In the line of work on unsupervised morphological segmentation, the most well-known model is Morfessor (Creutz and Lagus, 2002), along with its multiple variants (Creutz and Lagus, 2005; Grönroos et al., 2014, 2020). In Morfessor, an online search algorithm is utilized to apply a hierarchical word

splitting strategy with a Minimum Description 171 Length (MDL) (Rissanen, 1989) cost function. 172 However, its lack of explicit control over vocab-173 ulary size makes it unsuitable for use as a tokenizer. 174 In addition, although some studies (Ataman and Federico, 2018; Hou et al., 2023) find morpholog-176 ically motivated segmentation can improve data-177 driven tokenizers, most other studies (Machácek 178 et al., 2018; Domingo et al., 2019; Sälevä and Lignos, 2021) find no reliable improvement of such 180 methods over BPE. According to Gallé (2019), the 181 effectiveness of BPE lies in its superior compres-182 sion capability. A more detailed discussion can be 183 found in Mielke et al. (2021). 184

185

189

190

191

194

195

197

211

216

Some other studies try to model morphological structures using Bayesian PCFGs (Johnson et al., 2007) or a non-parametric Bayesian generalization of PCFGs (Johnson et al., 2006). However, they are pure statistical models and do not utilize modern neural methodologies.

Our method differs from previous unsupervised morphological methods in our characterbased structures, thereby possessing the superior compression capability of BPE. Meanwhile, our method leverages modern neural methodologies to better utilize contextual and intra-word semantic information.

198 **Composition Model.** In this work, we utilize a composition model to induce morphological struc-199 tures. Composition models jointly learn representations and structures of a symbol sequence by transforming text encoding into a combinatorial optimization problem. Maillard et al. (2017) proposes a CKY-like (Cocke, 1969; Kasami, 1966; Younger, 1967) encoder, in which each constituent is represented as a weighted average of the set of 206 composed representations computed from different 207 splits of the constituent. Drozdov et al. (2019) proposes a deep inside-outside encoder (Baker, 1979; Lari and Young, 1990), enabling the encoder to 210 learn underlying structures via an auto-encoding objective. Recently, a series of studies (Hu et al., 212 2024a,b) have been conducted to reduce the deep 213 inside-outside encoder complexity from cubic to linear, on which our work is based. 215

Methodology 3

Given a word $\mathbf{x} = \{x_1, x_2, ..., x_n\}$ where x_i is 217 the *i*-th character, we aim to parse it into a binary 218 tree and then tokenize it via top-down vocabulary 219 matching. The parser is a deep composition model



Figure 2: (a) The composition representation of asking $(i_{1,6})$ is a weighted sum over all subword pairs such as ask+ing $(\overline{i}_{1,6}^3)$ and as + king $(\overline{i}_{1,6}^2)$. (b) The composition function. Take ask $(\mathbf{i}_{1,3})$ as an example. $\mathbf{s}_{1,3}$ is $\mathbf{E}_{\mathbb{V}[\mathsf{ask}]}$ if $\mathsf{ask} \in \mathbb{V}$. Thus the representation of ask depends not only on its components but also on $\mathbf{E}_{\mathbb{V}[\mathsf{ask}]}$. However, if asking $\notin \mathbb{V}$, then $\mathbf{s}_{1,6}$ is \mathbb{E}_{empty} and the representation of asking $(\mathbf{i}_{1,6})$ only depends on the composition representation of its components.

capable of jointly modeling the internal structures and representations of words and is trained by optimizing self-supervised objectives. In the following sections, we sequentially introduce the composition model, training objectives, and the tree-based tokenization algorithm.

3.1 Composition Model for Word

For a given word x, we denote $\mathbf{i}_{i,j}$ as the representation of subword $\mathbf{x}_{i:j} = \{x_i, ..., x_j\}$. The inside pass (Drozdov et al., 2019) of a composition model computes a composition vector $\overline{\mathbf{i}}_{i,j}^k$ and a compatibility score $\bar{a}_{i,j}^k$ for each pair of sub-constituents (i, k) and (k + 1, j). The compatibility score indicates how likely these two sub-constituents are to be merged. The constituent representation $\mathbf{i}_{i,j}$ is computed as a weighted average over composition vectors of all possible pairs of sub-constituents as follows:

$$\bar{a}_{i,j}^{k}, \bar{\mathbf{i}}_{i,j}^{k} = f_{\alpha}(\mathbf{i}_{i,k}, \mathbf{i}_{k+1,j}), \\ \hat{w}_{i,j}^{k} = \frac{\exp(\bar{a}_{i,j}^{k})}{\sum_{k'=i}^{j-1} \exp(\bar{a}_{i,j}^{k'})}, \mathbf{i}_{i,j} = \sum_{k=i}^{j-1} \hat{w}_{i,j}^{k} \bar{\mathbf{i}}_{i,j}^{k}.$$
⁽¹⁾

The inside pass starts with characters by initializing $\mathbf{i}_{i,i}$ with character embeddings and recursively computes constituent representations bottom up following Equation 1. Representation $\mathbf{i}_{1,n}$ of the whole word \mathbf{x} is regarded as the word embedding EMB(x). f_{α} is the composition function implemented with a multi-layered Transformer. An example of the bottom-up composition process is depicted in Figure 2(a). In this work, we employ a pruned version of deep inside encoder (Hu et al., 2024b) as our backbone, which is easy to scale up, thanks to the logarithmic parallel time complexity and the linear space complexity.

The limitation of this approach is that the representation of any subword is always composed of 221

222

223

224

226

227

229

230

231

232

233

234

235

245 246 247

248

249

250

251

252

253

its component pairs, which is incompatible with the linguistic constraint that morphemes are the 256 smallest meaning-bearing units and should not be 257 decomposed further. Hence, we introduce MorphOverriding to enable a subword representation to disentangle from its component pairs when the subword is a morpheme. Specifically, we construct 261 a morpheme vocabulary \mathbb{V} heuristically using a statistical method (BPE in this work), in which each entry is associated with a learnable vector in a 264 morpheme embedding table **E**. When $\mathbf{x}_{i:i}$ hits the 265 vocabulary \mathbb{V} , we insert its morpheme embedding $\mathbf{s}_{i,j}$ into the computation of $\mathbf{i}_{i,j}$, making it possi-267 ble to mix or override the composition vector with 268 the morpheme embedding. Thus, the composition 269 vector and the compatibility scores can then be reformulated as:

 $l_2 = l_2$

273

274

275

276

278

279

284

285

287

292

293

296

$$\begin{split} \bar{a}_{i,j}^{n}, \mathbf{i}_{i,j}^{n} &= f_{\alpha}(\mathbf{i}_{i,k}, \mathbf{i}_{k+1,j}, \mathbf{s}_{i,j}) \,, \\ \mathbf{s}_{i,j} &= \begin{cases} \mathbf{E}_{\mathbb{V}[\mathbf{x}_{i:j}]} & \text{if } \mathbf{x}_{i:j} \in \mathbb{V} \\ \mathbf{E}_{empty} & \text{if } \mathbf{x}_{i:j} \notin \mathbb{V} \end{cases} \,, \end{split}$$

Figure 2(b) illustrates the composition function equipped with MorphOverriding. Our experiments demonstrate that this mechanism is crucial for character-level structure induction.

Tree induction. For a given span (i, j), the best split-point is k with the highest compatibility score $\bar{a}_{i,j}^k$. Thus, to derive a parse tree, we can recursively select the best split-points top-down starting from the root span (1, n). As the pruned inside-outside encoder produces a lightweight parser (Hu et al., 2022) with a compact parameter set as a byproduct, we use it for efficient inference during tokenization.

3.2 Training Objectives

The overall loss for training the composition model is the summation of an auto-encoding loss \mathcal{L}_{ae} and an auto-regression loss \mathcal{L}_{ar} . The auto-encoding loss is based on predicting each character or morpheme from the rest of a word, leveraging intraword information. The auto-regression loss is based on next token prediction that leverages contextual information to disambiguate the underlying structures of a word. Under these objectives, the composition model learns to assign proper scores to each split point of a subword.

297Auto-encoding Loss.Auto-encoding is a com-298mon practice of training a composition model. For299our character-level composition model, we try to300predict each character x_i based on its neighbor-301ing context representations $\mathbf{i}_{1,i-1}$ and $\mathbf{i}_{i+1,n}$ (Hu

et al., 2021). However, the auto-encoding objective turns out to be empirically ineffective when training our model probably because unlike word-level auto-encoding that requires selecting from tens of thousands of words in a vocabulary, here we only need to select from tens of characters, which is much less challenging. To enhance learning efficacy, we propose predicting both individual characters and morphemes in the vocabulary \mathbb{V} . For instance, given the word windsurf, we mask out wind and let the model uncover the masked morpheme based on the visible part _surf. Analogous to the inside pass, the outside pass computes each outside representation $o_{i,j}$ in a top-down manner based on context information outside span (i, j), whose details are described in Appendix A.3. we use $o_{i,j}$ to predict each subword $x_{i:j}$ that belongs to \mathbb{V} :

$$\mathcal{L}_{ae} = -\frac{1}{\mathcal{N}} \sum_{\mathbf{x}_{i:j} \in \mathbb{V}} \log \frac{\exp(\mathbf{o}_{i,j}^T \mathbf{E}_{\mathbb{V}[\mathbf{x}_{i:j}]})}{\sum_{k=1}^{|\mathbb{V}|} \exp(\mathbf{o}_{i,j}^T \mathbf{E}_k)},$$
 320

where \mathcal{N} is the total number of subwords belonging to the vocabulary 1

Auto-regression Loss. Given a sentence $S = {x^1, ..., x^m}$, whose word embedding is computed by the composition model, we feed the composed word embeddings into a causal language model and let it pick the correct next word from candidates built via in-batch sampling for each step. Let h_t denote the *t*-th hidden states of the causal language model and W denote a deduplicated vocabulary built on all input words in the same batch, we have the auto-regression loss defined as:

$$\mathcal{L}_{ar} = -\frac{1}{m} \sum_{t=1}^{m-1} \log \frac{\exp(\mathbf{h}_t \operatorname{EMB}(\mathbf{x}^{t+1}))}{\sum_{\mathbf{x} \in \mathcal{W}} \exp(\mathbf{h}_t \operatorname{EMB}(\mathbf{x}))} \,.$$

3.3 Tokenization

The proposed tree-based tokenization algorithm, TreeTok, includes segmentation and vocabulary construction procedures. As the latter depends on the former, we first discuss the segmentation procedure, followed by the vocabulary construction.

Segmentation Procedure. Given a constructed vocabulary, whose details are described later, we

333

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

323

324

325

326

327

328

329

330

331

332

334

- 335
- 336
- 337 338
- 339

341

¹Note that multiple subwords may be mutually exclusive, such as "asking" with "ask" and "king". Intuitively, it should be more reasonable to predict a constituent than other span candidates from the context. Therefore, we assign a constituency weight to each subword in the objective, as detailed in Appendix A.4.

parse each word into a morphological tree and seg-342 ment it via a top-down matching approach, as il-343 lustrated in Figure 1(bottom). Specifically, during 344 the top-down traversal of a parse tree, we retain a subword and backtrack if the subword matches an entry in the vocabulary. Note that unsupervised 347 structural learning is often imperfect, causing erroneous tokenization. For instance, an incorrect parse tree [[[book]e]d] may yield tokens book e d where e d should be merged. To address this issue, we propose a post-processing step to deal with mergeable pairs of segmented tokens. Specifically, we define the empirical probability of token t as $\frac{\text{COUNT}(t)}{T}$, where COUNT(t) is the frequency of t in the entire corpus and $T = \sum_{t \in \mathbb{V}} \text{COUNT}(t)$. 356 Therefore, the probability of a certain merge is the production of the probabilities of all tokens. We find the optimal merge by searching for the one with maximum probability among all potential merges via dynamic programming. Detailed 361 pseudo-code can be found in Appendix A.1.

Vocabulary Construction. One drawback of BPE and WordPiece is that they have to keep all intermediate "junk" tokens produced during the iterations of merge operations, which results in limited vocabulary space occupied by these meaningless tokens. For instance, if the corpus contains many occurrences of low and lower, the meaningless token 10 will be added to the vocabulary before low and will not be removed later. However, with the top-down matching framework, we don't need bottom-up merge operations to restore tokens, allowing us to prune unnecessary tokens and create 374 a more compact vocabulary. To build a compact vocabulary, we propose a vocabulary construction algorithm in which we employ a tree-based BPE-377 like algorithm to build a heuristic vocabulary and a tree-based Unigram algorithm to prune unnecessary subword units. Specifically, we initialize the token vocabulary with the character vocabulary and repeat the following steps to build a heuristic vocabulary given character-level tree structures of words:

363

364

371

373

375

381

391

- 1. Count adjacent token pairs that share the same parent in the tree structure, e.g., given [[b[o o]]k], only the pair (o, o) is counted.
- 2. Merge adjacent symbol pairs whose counts exceed a given threshold, e.g., $[[b[o o]]k] \rightarrow$ [[b oo]k].
- 3. Repeat 1-2 until there are no new symbol pairs.

In the pruning procedure, we start from the heuristic symbol vocabulary and prune it as follows:

392

393

394

395

396

397

398

399

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

- 1. Tokenize the corpus via the top-down matching according to the current vocabulary. The total entropy of the whole corpus is defined as $\mathcal{H}_{\mathbb{V}} = -\sum_{t \in \mathbb{V}} \frac{\operatorname{count}(t)}{T} \log \frac{\operatorname{count}(t)}{T}$ where $T = \sum_{t \in \mathbb{V}} \operatorname{count}(t)$.
- 2. For each token s, calculate the entropy gain after removing that word from the vocabulary denoted as $\Delta \mathcal{H}_s = \mathcal{H}_{\mathbb{V}/\{s\}} - \mathcal{H}_{\mathbb{V}}$. Intuitively, the higher $\Delta \mathcal{H}_s$ is, the more important s is.
- 3. Sort delta entropy of tokens and remove the lowest k% from \mathbb{V} . Repeat step 1-2 until $|\mathbb{V}|$ reaches the target vocabulary size.

In practice, we design a tree-based Viterbi algorithm (Viterbi, 1967) to implement the pruning procedure efficiently. The pseudo-code is presented in Appendix A.2.

Experiments 4

We evaluate the performance of TreeTok against the de-facto tokenizers such as BPE, WordPiece, and Unigram as primary baselines.

Training setups. For a fair comparison, we train all tokenizers from scratch on the lowercase version of the WikiText-103 corpus (Merity et al., 2017) without any word boundary marker and set the same vocabulary size of 30,000. For BPE, WordPiece, and Unigram, we use the implementation and default training paradigm provided by the HuggingFace library². Regarding the composition model, we train it with a context window of up to 512 characters. We use GPT2 implemented from HuggingFace³ as our causal language model when computing the auto-regression loss. We present detailed configurations of our model and training setup in Appendix A.5.

Evaluation datasets. We compare our tokenizer with other tokenizers for morphological alignment using two datasets with gold-standard morphological segmentation. One is from the Morpho Challenge 2010 Workshop (Kurimo et al., 2010) (Morpho), which contains 1,000 word forms with their segmentations corresponding to the surface forms of morpheme labels. The dataset contains instances of all kinds of morphological transformations, including inflection, derivation, and compounding. The other dataset is from Minixhofer et al. (2023)

²https://github.com/huggingface/tokenizers ³https://github.com/huggingface/transformers

	Morpho (Acc.) ↑	Compound (Acc.) ↑	$ \mathbb{V} $
BPE	19.50	62.98	30,000
WordPiece	26.20	62.19	30,000
Unigram	27.10	53.10	30,000
TreeTok	37.9	68.07	30,000

Table 1: Results on two morphological segmentation datasets. This table can be seen as a comparison between tree-enhanced BPE(TreeTok) and vanilla BPE/WordPiece/Unigram.

(Compound), which contains 759 compound words specifically designed to test the models' capabilities in decompounding. We also use these morphological segmentation datasets to evaluate the induced morphological parse trees.

In addition, we evaluate the tokenizers using statistical metrics that have been shown to strongly correlate with the performance on downstream tasks. These metrics are calculated on the validation set of WikiText-103.

4.1 Tokenization Quality

439

440

441

442

443

444

445

446

447

448

449

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

Metrics. We measure the performance of morphological segmentation via accuracy, i.e., the ratio of examples that are correctly segmented. We also consider a few statistical metrics that can directly assess the quality of tokenization, including Rényi Efficiency (Zouhar et al., 2023), average sentence-level perplexity, and average number of tokens per sentence. Rényi Efficiency is introduced by Zouhar et al. (2023) as a principled intrinsic measure of tokenization quality and is claimed to yield a Pearson correlation of 0.78 with BLEU (Papineni et al., 2002) on machine translation. Sentence-level perplexity is defined as $-\log p(\mathbf{s}) = -\sum_{i=1}^{n} \log p(s_i|s_{<i})$, where $\mathbf{s} = \{s_1, s_2, ..., s_n\}$ is a sentence with s_i being the *i*-th token. Since different tokenizers generate distinct segmentations leading to different numbers of tokens of the same word, sentence-level perplexity provides fairer evaluation compared with the default token-level perplexity $-\frac{1}{n}\log p(\mathbf{s})$.

According to Table 1, TreeTok significantly surpasses BPE, WordPiece, and Unigram on the two morphological segmentation datasets. The results demonstrate the efficacy of TreeTok in aligning with morphology.

Rényi efficiency & Perplexity. Table 2 reports the evaluation results in terms of Rényi efficiency and perplexity (PPL). TreeTok outperforms BPE and WordPiece on both Rényi and PPL. The improvements illustrate the benefits of TreeTok's

	Rényi↑	PPL↓	BLEU↑	avg. #tokens
BPE	44.66	107.76	26.55	26.58
WordPiece	44.54	110.97	-	26.60
Unigram	45.07	106.91	-	31.68
TreeTok	44.82	107.26	26.68	25.99

Table 2: Results for different tokenization models on Wiki-Text103 with 30,000 vocabulary size.

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

structural constraints and more compact vocabulary. The tree structure constraints enable the segmentation of words into more morphology-aligned tokens, while the compact vocabulary allows for the inclusion of meaningful morphemes by removing intermediate tokens in the pruning process during vocabulary construction, under a top-down matching framework. Unigram performs slightly better than TreeTok, but produces 22% more tokens on average. A possible explanation for the better performance of Unigram is that Unigram tends to produce inflectional suffixes such as "-ing" and "-ly", while other methods tend to retain entire words. This difference makes it easier for Unigram to share the same stems and affixes between different word forms, thus achieving better parameter sharing. However, under the Transformer architecture, an additional 22% number of tokens means extra inference steps and nearly 1.4 times the cost of self-attention. Such additional costs only bring marginal improvements as can be seen in the table.

We also note that TreeTok achieves the shortest average token length among all the tokenizers, which is desirable as Gallé (2019) shows that given a fixed vocabulary size budget, the fewer tokens a tokenizer needs to cover the test set, the better the translation.

Machine Translation. We conduct experiments on machine translation as a complementary. We use the fairseq framework ⁴ to train a Transformer on WMT14 English to German from scratch and measure the performance by calculating the BLEU score on the official test split.

We compare the model's performance when the tokenizer is BPE and TreeTok, respectively. We use the same model training setups. For the two tokenizers, the vocabulary size and basic characters are exactly the same. The results from Table 2 show that TreeTok is slightly better than BPE. Based on the results, TreeTok can improve alignment with morphology on top of BPE, while the new seg-

⁴https://github.com/facebookresearch/fairseq/ blob/main/examples/translation/README.md#wmt14english-to-german-convolutional

	Morpho	Compound	Word Seg.
	EN.	EN.	ZH.
Fast R2D2	67.69	48.96	
Neural PCFG	39.87	58.33	74.26
TreeTok	90.10	86.20	_
w/o context	70.00	63.02	—
w/o MorphOverriding	75.99	46.35	99.24
w/o span weights	89.42	78.39	_
w/o span loss	86.79	73.70	—

Table 3: Performance evaluation of our model, baseline models, and ablation studies on morphological segmentation, measured by morpheme recall rate. EN:English, ZH: Chinese.

mentation does not compromise downstream task performance.

4.2 Tree Structure Quality

521

522

523

524

525

529

531

532

533

537

539

540

541

542

543

544

546

549

Since tree structures play an important role in both vocabulary construction and segmentation, we evaluate the quality of trees induced by various composition models.

Metric. We use recall of morphemes (van den Bosch and Daelemans, 1999) in a tree to assess the quality of the tree structures against gold-standard segmentations, which is defined as the percentage of morphemes in the gold segmentation that can be found in the spans of the evaluated tree. We discard spans that are trivial for a tree (characterlevel and word-level spans) and report word-level recall (averaged over word samples).

Baselines. For baseline composition models, we include Fast-R2D2 (Hu et al., 2022), which is a variant of DIORA (Drozdov et al., 2019), and an efficient variant of neural PCFG (Yang et al., 2022).

We also include four variants of our composition model for an ablation study. In w/o context, we remove the auto-regression loss from our architecture so that each representation only contains information from individual words. In w/o MorphOverriding, we degenerate $s_{i,j}$ to the default empty embedding regardless of whether span $x_{i:j}$ hits the external vocabulary or not. In w/o span loss, for our auto-encoding loss, we only count loss from predicting characters instead of every subword span that hits the external vocabulary.

552Results and Discussions.As shown in Table 3,553our model outperforms all the other composition554models.Compared with Fast-R2D2, our main555differences lie in the training objectives and the556MorphOverriding mechanism.557validates the effectiveness of these improvements.



Figure 3: The effect of changing the vocabulary size learned by BPE. The initial results on both tasks show that the performance curve is a concave function where the maximum resides in the middle.

Our ablation experiments further analyze the contribution of these improvements to performance enhancement. Specifically, we have the following findings from each ablation.

Removing the auto-regression loss to prevent the model from getting feedback from contextual information significantly impacts the performance on both tasks, especially Morpho. We believe that contextual information can help the model capture the regularities of tenses and learn how to build composition representations for compound words. For example, consider how the context can help determine whether we should build the representation of asking as ask+ing or as+king. While either is a valid combination of morphemes, the former is more likely to be learned by our model since the context around asking often indicates the continuous tense or the gerund form, thus matching better with ing.

Removing MorphOverriding from the model results in a significant decrease of around 50% in performance on the decompounding task. The results consolidate our insight about conventional composition models violating the indecomposability of morphemes. Creating a morpheme's representation using its components' representation might make representations of disparate morphemes (e.g., wind and win) entangled together.

Removing the span loss also causes a performance drop on the two morphology tasks. This aligns well with the insight behind our design of morpheme-level loss, which augments the character-level loss by enhancing the learning of intra-word representations for most morphemes that are at an intermediate granularity.

In addition, we train both Neural-PCFG and our composition model on Chinese wiki (Xu and Lap-

594

558

559

original word	bed	commonly	windsurfing	tricycles	uniquenesses
BPE	bed	commonly	wind/sur/fing	tric/y/cles	uniqu/eness/es
Unigram	b/e/d	common/ly	wind/surf/ing	t/r/i/cycle/s	unique/ness/e/s
WordPiece	bed	commonly	winds/ur/fing	tric/y/cles	unique/ness/es
TreeTok	bed	commonly	wind/surf/ing	tri/cycles	unique/ness/es

Table 4: Example tokenizations.

ata, 2019) and evaluate the recall of word bound-595 596 ary against Penn Chinese Treebank (XUE et al., 2005). Our model can achieve a word boundary recall of 99.24% without MorphOverriding. Chinese is a language system that can be considered simplistic in terms of its internal structure of words where most of the time, each Chinese character (referred to as a hàn zì) represents one morpheme, and there are always explicit boundaries between 603 morphemes without orthographic changes during word formation from characters. These features make compositionality applicable in most cases, 606 thus alleviating the difficulty of modelling intraword structures. Hence, comparing with the poor 609 performance of w/o MorphOverriding on the English dataset (Compound), we can conclude that the difficulties of modelling the internal word structure vary greatly across languages, and MorphOverrid-612 ing is effective and necessary for languages with more challenging morphology structures.

597

598

604

610

611

615

616

617

618

619

621

625

627

631

632

634

635

Influence of Heuristic Vocabulary Size Additionally, we conduct experiments to investigate how the size of our heuristic morpheme vocabulary influences the performance of structure induction.

Figure 3 shows that the optimal size of an external vocabulary should be neither too large nor too small. According to our hypothesis that the compositional representation of subcomponents of a morpheme should be overridden by a high-level representation, ideally, the external vocabulary should contain all morphemes and only morphemes, because our model will trigger the soft morpheme overriding mechanism for every span that hits the external vocabulary. If BPE is used and the vocabulary is too small, many morphemes (especially longer standalone words) are excluded. Conversely, if it is too large, BPE merges across morphemes, creating spans larger than the smallest meaningbearing units.

Case Studies 4.3

To further examine the difference between tokenizers, we list their tokenizations in Table 4 and tree



Figure 4: Example tree structures induced by our composition model.

structures induced by our composition model in Figure 4.

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

Tokens produced by Unigram often include many characters. BPE and WordPiece often violate morpheme boundaries and tokenize words into some intermediate "junk" tokens introduced during the bottom-up vocabulary construction, such as fing, cles, and eness in Table 4.

TreeTok aligns significantly better with morphology. By merging the best of BPE and Unigram pruning, our vocabulary construction algorithm eliminates "junk" tokens. Meanwhile, top-down matching under linguistic constraints prevents excessive word fragmentation and morpheme boundary breaks.

In Figure 4, our model's high-level tree structures are generally accurate, although some lowlevel structures appear random, since MorphOverriding prioritizes the most reasonable high-level segmentations based on context, making lowerlevel details less important.

5 Conclusion

8

Our work introduces the first effective unsupervised neural model for character-level structure induction. We discovered that recognizing the indecomposability of morphemes is key, and to address this, we developed a composition model with a MorphOverriding mechanism alongside two selfsupervised objectives. TreeTok induces tree structures that closely match human-labeled morphology and consistently outperforms baselines like BPE and WordPiece across various tasks, offering new insights into unsupervised morphological segmentation.

6 Limitations

Our main limitation is that we need additional train-672 ing and inference overheads. Considering that the 673 composition model only needs to be trained once 674 and the overall time consumption is acceptable⁵, 675 we believe it is not a fatal flaw. Regarding inference cost, because a lightweight parser is produced as 677 a byproduct, it can be afforded by even CPU environments. According to Table 5, Treetok's average processing time per token is longer than other tokenizers. However, if we allow Treetok to tokenize in batches on a GPU in advance, this gap can be easily compensated. Furthermore, we can maintain a cache of high-frequency words to avoid repeated tokenization. In wikitext-103, the hit rate for a cache that stores the top 100000 frequent words is 98.11%, which means only 2% tokens need to be parsed on the fly. e.g. for 1000 tokens, it only needs 0.16s to parse. Furthermore, these results are based on single-core computation, and there is still room for multi-core acceleration.

Tokenizer	Avg Time/Token (s)
BPE	2.49e-05
WordPiece	2.33e-05
Unigram	2.54e-05
TreeTok(single processing)	1.98e-3

Table 5: Average processing time per token for different tokenizers with the same vocabulary size. The result is collected when TreeTok tokenizes on a CPU and only tokenize one sample at a time.

References

691

701

702

703

704 705

706

707

708

- Duygu Ataman and Marcello Federico. 2018. An evaluation of two vocabulary reduction methods for neural machine translation. In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas, AMTA 2018, Boston, MA, USA, March* 17-21, 2018 - Volume 1: Research Papers, pages 97–110. Association for Machine Translation in the Americas.
- James K. Baker. 1979. Trainable grammars for speech recognition. *Journal of the Acoustical Society of America*, 65.
- Kaj Bostrom and Greg Durrett. 2020. Byte pair encoding is suboptimal for language model pretraining. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4617–4624, Online. Association for Computational Linguistics.

Kenneth Ward Church. 2020. Emerging trends: Subwords, seriously? *Natural Language Engineering*, 26(3):375–382. 709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

- John Cocke. 1969. Programming Languages and Their Compilers: Preliminary Notes. New York University, USA.
- Ryan Cotterell and Hinrich Schütze. 2015. Morphological word-embeddings. In NAACL HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015, pages 1287–1292. The Association for Computational Linguistics.
- Mathias Creutz and Krista Lagus. 2002. Unsupervised discovery of morphemes. In *Proceedings of the ACL-*02 Workshop on Morphological and Phonological Learning, pages 21–30. Association for Computational Linguistics.
- Mathias Johan Philip Creutz and Krista Hannele Lagus. 2005. Inducing the morphological lexicon of a natural language from unannotated text. In *Proc. International and Interdisciplinary Conference on Adaptive Knowledge Representation and Reasoning* (*AKRR'05*), pages 106–113. International and Interdisciplinary Conference on Adaptive Knowledge Representation and Reasoning (AKRR'05) ; Conference date: 01-01-1800.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Miguel Domingo, Mercedes García-Martínez, Alexandre Helle, Francisco Casacuberta, and Manuel Herranz. 2019. How much does tokenization affect neural machine translation? In *Computational Linguistics and Intelligent Text Processing - 20th International Conference, CICLing 2019, La Rochelle, France, April 7-13, 2019, Revised Selected Papers, Part I*, volume 13451 of *Lecture Notes in Computer Science*, pages 545–554. Springer.
- Andrew Drozdov, Patrick Verga, Mohit Yadav, Mohit Iyyer, and Andrew McCallum. 2019. Unsupervised latent tree induction with deep inside-outside recursive auto-encoders. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 1129–1141, Minneapolis, Minnesota. Association for Computational Linguistics.
- Matthias Gallé. 2019. Investigating the effectiveness of BPE: The power of shorter sequences. In *Proceedings of the 2019 Conference on Empirical Methods*

⁵less than 1 day for $8 \times A100$ for WikiText-103

766

- 774
- 775 776
- 777
- 778 779
- 7

784

- 786
- 787 788
- 790 791
- 792

7

- 79
- 79
- 7
- 8

8

805

8

- 8
- 810
- 811 812

815 816

817 818

819

820 821 in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 1375–1381, Hong Kong, China. Association for Computational Linguistics.

- Stig-Arne Grönroos, Sami Virpioja, and Mikko Kurimo. 2020. Morfessor EM+Prune: Improved subword segmentation with expectation maximization and pruning. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3944–3953, Marseille, France. European Language Resources Association.
- Stig-Arne Grönroos, Sami Virpioja, Peter Smit, and Mikko Kurimo. 2014. Morfessor FlatCat: An HMMbased method for unsupervised and semi-supervised learning of morphology. In *Proceedings of COLING* 2014, the 25th International Conference on Computational Linguistics: Technical Papers, pages 1177– 1185, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. 2021. Superbizarre is not superb: Derivational morphology improves BERT's interpretation of complex words. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3594–3608, Online. Association for Computational Linguistics.
- Jue Hou, Anisia Katinskaia, Anh-Duc Vu, and Roman Yangarber. 2023. Effects of sub-word segmentation on performance of transformer language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023, pages 7413– 7425. Association for Computational Linguistics.
- Xiang Hu, Pengyu Ji, Qingyang Zhu, Wei Wu, and Kewei Tu. 2024a. Generative pretrained structured transformers: Unsupervised syntactic language models at scale.
- Xiang Hu, XinYu KONG, and Kewei Tu. 2023. A multigrained self-interpretable symbolic-neural model for single/multi-labeled text classification. In *The Eleventh International Conference on Learning Representations*.
- Xiang Hu, Haitao Mi, Liang Li, and Gerard de Melo. 2022. Fast-R2D2: A pretrained recursive neural network based on pruned CKY for grammar induction and text representation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2809–2821, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Xiang Hu, Haitao Mi, Zujie Wen, Yafang Wang, Yi Su, Jing Zheng, and Gerard de Melo. 2021. R2D2: Recursive transformer based on differentiable tree for

interpretable hierarchical language modeling. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4897–4908, Online. Association for Computational Linguistics. 822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

- Xiang Hu, Qingyang Zhu, Kewei Tu, and Wei Wu. 2024b. Augmenting transformers with recursively composed multi-grained representations. In *The Twelfth International Conference on Learning Representations*.
- Mark Johnson, Thomas Griffiths, and Sharon Goldwater. 2006. Adaptor grammars: A framework for specifying compositional nonparametric bayesian models. In *Advances in Neural Information Processing Systems*, volume 19. MIT Press.
- Mark Johnson, Thomas Griffiths, and Sharon Goldwater. 2007. Bayesian inference for PCFGs via Markov chain Monte Carlo. In Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference, pages 139–146, Rochester, New York. Association for Computational Linguistics.
- Daniel Jurafsky and James H. Martin. 2009. *Speech and Language Processing (2Nd Edition)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- Tadao Kasami. 1966. An efficient recognition and syntax-analysis algorithm for context-free languages. *Coordinated Science Laboratory Report no. R-257.*
- Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics.
- Mikko Kurimo, Sami Virpioja, Ville Turunen, and Krista Lagus. 2010. Morpho challenge 2005-2010: Evaluations and results. In Proceedings of the 11th Meeting of the ACL Special Interest Group on Computational Morphology and Phonology, pages 87– 95, Uppsala, Sweden. Association for Computational Linguistics.
- K. Lari and S.J. Young. 1990. The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer Speech & Language*, 4(1):35– 56.
- Dominik Machácek, Jonás Vidra, and Ondrej Bojar. 2018. Morphological and language-agnostic word segmentation for NMT. In *Text, Speech, and Dialogue - 21st International Conference, TSD 2018, Brno, Czech Republic, September 11-14, 2018, Proceedings,* volume 11107 of *Lecture Notes in Computer Science,* pages 277–284. Springer.

2017.

abs/1705.09189.

Technology.

Linguistics.

sentations.

enization in nlp.

Linguistics.

Linguistics.

Science.

Jean Maillard, Stephen Clark, and Dani Yogatama.

Tatjana Marvin. 2002. Topics in the Stress and Syntax

David McClosky, Eugene Charniak, and Mark Johnson.

2006. Effective self-training for parsing. In Proceed-

ings of the Human Language Technology Conference

of the NAACL, Main Conference, pages 152-159,

New York City, USA. Association for Computational

Stephen Merity, Caiming Xiong, James Bradbury, and

Sabrina J. Mielke, Zaid Alyafeai, Elizabeth Salesky,

Colin Raffel, Manan Dey, Matthias Gallé, Arun Raja,

Chenglei Si, Wilson Y. Lee, Benoît Sagot, and Sam-

son Tan. 2021. Between words and characters: A brief history of open-vocabulary modeling and tok-

Benjamin Minixhofer, Jonas Pfeiffer, and Ivan Vulić.

2023. CompoundPiece: Evaluating and improving

decompounding performance of language models.

In Proceedings of the 2023 Conference on Empiri-

cal Methods in Natural Language Processing, pages

343–359, Singapore. Association for Computational

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-

Jing Zhu. 2002. Bleu: a method for automatic evalu-

ation of machine translation. In Proceedings of the

40th Annual Meeting of the Association for Compu-

tational Linguistics, pages 311-318, Philadelphia,

Pennsylvania, USA. Association for Computational

Alec Radford, Jeffrey Wu, Rewon Child, David Luan,

Jorma Rissanen. 1989. Stochastic complexity in statisti-

Jonne Sälevä and Constantine Lignos. 2021. The effec-

tiveness of morphology-aware segmentation in low-

resource neural machine translation. In Proceedings

of the 16th Conference of the European Chapter of

the Association for Computational Linguistics: Student Research Workshop, EACL 2021, Online, April

19-23, 2021, pages 164-174. Association for Com-

Mike Schuster and Kaisuke Nakajima. 2012. Japanese and korean voice search. 2012 IEEE International

cal inquiry. In World Scientific Series in Computer

models are unsupervised multitask learners.

Dario Amodei, and Ilya Sutskever. 2019. Language

Richard Socher. 2017. Pointer sentinel mixture mod-

els. In International Conference on Learning Repre-

of Words. Ph.D. thesis, Massachusetts Institute of

and syntax with unsupervised tree-lstms. CoRR,

Jointly learning sentence embeddings

- 886
- 892

- 897
- 900 901
- 902 903 904

905 906

- 907 908 909
- 910 911 912

913 914

915

916 917 918

919

- 920 921
- 922 923
- 924

926

927

930 931

Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 5149-5152.

putational Linguistics.

Elisabeth Selkirk. 1982. The Syntax of Words. Linguistic inquiry monographs. MIT Press.

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

- 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.
- Antal van den Bosch and Walter Daelemans. 1999. Memory-based morphological analysis. In Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics, pages 285-292, College Park, Maryland, USA. Association for Computational Linguistics.
- Andrew J. Viterbi. 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. IEEE Trans. Inf. Theory, 13(2):260-269.
- Yumo Xu and Mirella Lapata. 2019. Weakly supervised domain detection. Transactions of the Association for Computational Linguistics, 7:581-596.
- NAIWEN XUE, FEI XIA, FU-DONG CHIOU, and MARTA PALMER. 2005. The penn chinese treebank: Phrase structure annotation of a large corpus. Natural Language Engineering, 11(2):207–238.
- Songlin Yang, Wei Liu, and Kewei Tu. 2022. Dynamic programming in rank space: Scaling structured inference with low-rank HMMs and PCFGs. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4797-4809, Seattle, United States. Association for Computational Linguistics.
- Daniel H Younger. 1967. Recognition and parsing of context-free languages in time n3. Information and control, 10(2):189-208.
- Vilém Zouhar, Clara Meister, Juan Gastaldi, Li Du, Mrinmaya Sachan, and Ryan Cotterell. 2023. Tokenization and the noiseless channel. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5184-5207, Toronto, Canada. Association for Computational Linguistics.

A Appendix

975

976

977 978

A.1 Pseudo-codes of tokenization

Algorithm 1 Tokenize

```
1: Input: string x, parse tree root r, vocabulary \mathbb{V}
 2:
     procedure TOKENIZE(\mathbf{x}, r, \mathbb{V})
 3:
          t \leftarrow []
                                          ▷ tokenized subword units list
 4:
          stack \leftarrow [r]
 5:
          while |stack| > 0 do
 6:
               c \leftarrow \text{POP}(stack)
 7:
                i, j \leftarrow c.i, c.j
               \bar{\mathbf{x}} \leftarrow \mathbf{x}_{i:j}
 8:
               if \bar{\mathbf{x}} \in \mathbb{V} then
 9:
10:
                     APPEND(t, \bar{\mathbf{x}})
                else if i < j then
11:
                                                     ▷ Non-terminal nodes
12:
                     PUSH(stack, c.right)
13:
                     PUSH(stack, c.left)
14:
          t \leftarrow \text{PostMerge}(t, \mathbb{V})
                                          > post processing if over-split
15:
16:
          return t
```

Algorithm 2 Post-Merge Algorithm

1:	Input: tokens t , vocab2entropy \mathbb{V}
2:	procedure POSTMERGE(\mathbf{t}, \mathbb{V})
3:	$n \leftarrow \text{length of } \mathbf{t}$
4:	if $n \leq 1$ then
5:	$\mathbf{t}_{ ext{MERGE}} \leftarrow \mathbf{t}$
6:	else
7:	$\mathcal{H}[n][n]$ init with ∞ \triangleright Best entropy
8:	$s[n][n]$ init with [] \triangleright Best segments
9:	for $i \leftarrow 0$ to $n - 1$ do \triangleright Base case
10:	$\mathcal{H}_{i,i} \leftarrow \mathbb{V}[x_i]$
11:	$s_{i,i} \leftarrow [x_i]$
12:	for $h \leftarrow 1$ to $n - 1$ do \triangleright Iterate tree height
13:	for $i \leftarrow 0$ to $n - h - 1$ do
14:	$j \leftarrow i + h$
15:	$k_{ ext{best}} \leftarrow -1$
16:	$m \leftarrow \text{concatenate } t_i \dots t_j$
17:	$\mathcal{H}_{ ext{best}} \leftarrow \operatorname{Get}(\mathbb{V}, m, \infty)$
18:	for $k \leftarrow i$ to $j - 1$ do
19:	if $\mathcal{H}_{i,k} + \mathcal{H}_{k+1,j} \leq \mathcal{H}_{\text{BEST}}$ then
20:	$k_{ ext{best}} \leftarrow k$
21:	$\mathcal{H}_{ ext{best}} \leftarrow \mathcal{H}_{i,k} + \mathcal{H}_{k+1,j}$
22:	if $k_{\text{BEST}} \neq -1$ then
23:	$s_{i,j} \leftarrow s_{i,k_{\text{BEST}}} + s_{k_{\text{BEST}}+1,j}$
24:	else
25:	$s_{i,j} \leftarrow [m] \qquad \qquad \triangleright \text{Merge}$
26:	$\mathcal{H}_{i,j} \leftarrow \mathcal{H}_{ extsf{best}}$
27:	$\mathbf{t}_{ ext{MERGE}} \leftarrow s_{0,n-1}$
28:	return t _{MERGE}

A.2 Pseudo-codes of vocab construction

Please refer to Algorithm 3 for details.

Algorithm 3 Vocabulary Construction

```
1: Input: tree-freq pair list T, vocab size k, pruning rate \alpha
 2: procedure VOCABULARY CONSTRUCTION(T, k, \alpha)
          procedure E-STEP(T, \mathbb{V})
 3:
               \mathbb{V}' \leftarrow \text{DICT}()
 4:
                                          ▷ E-step: Update vocab freq
 5:
               for \{root, freq\} \in T do
                    \_, seg \leftarrow TREEVITERBI(root, V, null)
 6:
 7:
                    for token \in seg do
                         \mathbb{V}'[token] \leftarrow \mathbb{V}'[token] + freq
 8:
               return V
 9:
10:
11:
          procedure M-STEP(T, \mathbb{V})
12:
               l \leftarrow \text{DICT}()
                                           ▷ M-step: Update delta loss
13:
               for \{root, freq\} \in T do
14:
                                                 ▷ word-level delta-loss
                    l_{word} \leftarrow \text{DICT}()
15:
                     \_, seg \leftarrow \text{TREEVITERBI}(root, \mathbb{V}, l_{word})
16:
                    for token \in seg do
17:
                         loss \leftarrow l_{word}[token]
18:
                         l[token] \leftarrow l[token] + loss * freq
19:
               return V
20:
21:
          \mathbb{V} \leftarrow \text{INITVOCAB}(T)
                                                ▷ Init with a BIG vocab
          while |\mathbb{V}| > k do
22:
23:
                \mathbb{V} \leftarrow \text{E-STEP}(T, \mathbb{V})
                                                 ▷ Estimate token count
               \mathbb{L} \leftarrow \text{M-STEP}(T, \mathbb{V})
24:
                                               ▷ Maximize delta losses
               Remove \min(|\mathbb{V}| - k, \lfloor \alpha |\mathbb{V}| \rfloor) of the
25:
26:
               tokens t with lowest L_t from \mathbb{V}
27:
          return \mathbb V
```

Algorithm 4 TreeViterbi

1: Input: parse tree root r, vocabulary \mathbb{V} , delta loss dict l2: **procedure** TREEVITERBI(r, V, l)3: $w \gets r.token$ 4: if r.i = r.j then 5: $s \leftarrow \text{GET}(\mathbb{V}, w, \infty) \mathrel{
ightarrow} \text{Infinity entropy if } w \notin \mathbb{V}$ return s, [w]6: 7: else $s_L, w_L \leftarrow \text{TreeViterbi}(r.left, \mathbb{V}, l)$ 8: $s_R, w_R \leftarrow \text{TREEVITERBI}(r.right, \mathbb{V}, l)$ 9: 10: $s \leftarrow \operatorname{Get}(\mathbb{V}, w, \infty)$ 11: if *l* then ⊳ Enter in M-step 12: $l[w] \leftarrow l[w] + MAX(s_L + s_R - s, 0)$ ▷ Record delta loss: Entropy increase 13: 14: if $s_L + s_R > s$ then 15: return s, [w]16: else 17: **return** $s_L + s_R, w_L + w_R$

Algorithm 5 Vocabulary Initialization

1:	Input: tree-freq pair list T , threshold k
2:	procedure INITVOCAB (T, k)
3:	$\mathbb{V} \leftarrow All \text{ character freq}$
4:	$n \leftarrow \mathbb{V} $
5:	while True do
6:	$\mathbb{V}^{'} \leftarrow CountBigrams(T,\mathbb{V})$
7:	Prune all the entries in $\mathbb{V}^{'}$ with freq less than k
8:	$\mathbb{V}.MERGE(\mathbb{V}') \qquad \triangleright \text{ Add new items in } \mathbb{V}' \text{ to } \mathbb{V}$
9:	if $ \mathbb{V} = n$ then
10:	break
11:	$n = \mathbb{V} $
12:	return V

Algorithm 6 Count Bigrams

1:	Input: tree-freq pair list T, vocabulary \mathbb{V}
2:	procedure COUNTBIGRAMS (T, \mathbb{V})
3:	$\mathbb{V}' \leftarrow \text{DICT}()$ \triangleright Store new merges
4:	procedure $RECURCOUNT(r, f)$
5:	if r.left & r.right then
6:	$hit_L \leftarrow \text{RecurCount}(r.left, f)$
7:	$hit_R \leftarrow \text{RecurCount}(r.right, f)$
8:	if hit_L and hit_R then
9:	if $r.token \in \mathbb{V}$ then
10:	return True
11:	else ,
12:	$\mathbb{V}[r.token] \leftarrow f \triangleright$ Merge: new entry
13:	return False
14:	else
15:	return False
16:	else
17:	return True
18:	for $\{root, freq\} \in T$ do
19:	Recurcount(root, freq)
20:	return $\mathbb{V}^{'}$

The outside computation is akin to the inside pass

but in a top-down manner. we denote the outside

representation and score of a given span as $\bar{\mathbf{o}}_{i,j}^k$ and

 $\bar{b}_{i,j}^k$ respectively, whose parent span is (i, k) or (k,

 $\check{w}_{i,j}^k = \frac{\exp(\bar{b}_{i,j}^k)}{\sum_{k'>i,k'< i} \exp(\bar{b}_{i,j}^{k'})}, \mathbf{o}_{i,j} = \sum_{k>j,k< i} \check{w}_{i,j}^k \bar{\mathbf{o}}_{i,j}^k.$

An intuitive idea is that the larger the probability of a span's existence, the greater its weight. A span exists if its parent span exists and the span is an im-

mediate child of its parent span. Therefore, we can

recursively estimate the existence probability of

each span top-down (Hu et al., 2023) and formalize

 $p_{i,j} = \sum_{k < i} p_{k,j} \hat{w}_{k,j}^{i} + \sum_{k > j} p_{i,k} \hat{w}_{i,k}^{j}, \, p_{1,n} = 1,$

 $\mathcal{L}_{ae} = -\frac{1}{\sum p_{i,j}} \sum_{\mathbf{x}_{i,i} \in \mathbb{V}} p_{i,j} \log \frac{\exp(\mathbf{o}_{i,j}^T \mathbf{E}_{\mathbb{V}[\mathbf{x}_{i:j}]})}{\sum_{k=1}^{|\mathbb{V}|} \exp(\mathbf{o}_{i,j}^T \mathbf{E}_k)}.$

Our composition function uses 4 layers of Trans-

use 128-dimensional embeddings with 4 attention

For span representations, we

The neural outside pass

$$\begin{split} \bar{\mathbf{o}}_{i,j}^k &= \begin{cases} f_\beta(\mathbf{o}_{i,k},\mathbf{i}_{j+1,k}) & \text{ if } k>j\\ f_\beta(\mathbf{o}_{k,j},\mathbf{i}_{k,i-1}) & \text{ if } k<i \end{cases},\\ \bar{b}_{i,j}^k &= \begin{cases} \phi_\beta(\mathbf{o}_{i,k},\mathbf{i}_{j+1,k}) & \text{ if } k>j\\ \phi_\beta(\mathbf{o}_{k,j},\mathbf{i}_{k,i-1}) & \text{ if } k<j \end{cases}, \end{split}$$

the auto-encoding loss as follows:

A.5 Experimental Setup and Hyperparameters

former layers.

j) for k > j or k < i.

A.4 Span weights

979 980

A.3

981 982

984

985

98

30

989

990 991

99

99

994

99

997

998

heads, 512-dimensional hidden layer representations, and a vocabulary size of 7835. This vocabu-1001 lary is built from concatenating 1903 most frequent 1002 characters in the training set of wikitext-103 and 1003 a 10,000-entry BPE dictionary, excluding all char-1004 acters. To guide the composition function, our 1005 lightweight parser is a 4-layer Transformer model 1006 that uses 64-dimensional embeddings with 4 attention heads and 128-dimensional hidden layer 1008 representations. For the causal language model, we use a 3-layer GPT2 equipped with 128-dimensional 1010 embeddings and 4 attention heads and follow the 1011 original configuration for the rest of the hyperpa-1012 rameters. 1013

1014

1015

1016

1017

1018

1019

Our composition models are trained on 8 PPUs with a learning rate of 1e-2 for the light-weight parser and 5e-4 for the rest. The batch size is 8×128 , and for each sample, we limit the context window to 512 characters (whitespace included). The total number of training steps is ten times the number of sentences in Wikitext-103.