## **Generation with Dynamic Vocabulary**

### **Anonymous ACL submission**

#### Abstract

Vocabulary is a crucial component of language models. Traditional language models generate text by selecting tokens from a fixed vocabulary. In this paper, we introduce a novel dynamic setting for the vocabulary. Under this setting, vocabulary can include arbitrary text spans on demand. These text spans act as basic bricks, akin to tokens in referred the fixed vocabulary. Our proposed model can be deployed in a way of plug-and-play. Extensive experimental results demonstrate that our approach yields superior generation quality. For instance, compared to the standard language model, the MAUVE metric increases from 20.47 % to 25.69%. We also demonstrate that dynamic vocabulary can be effectively applied to different domains in a training-free manner, and it also helps to generate reliable citations in question answering tasks (substantially enhancing citation results without compromising answer accuracy).

### 1 Introduction

011

012

014

037

041

Vocabulary, which defines basic bricks (tokens) for composing new sentences, bridging different languages (Stahlberg, 2020; Lample and Conneau, 2019; Liu et al., 2020; Koehn and Knowles, 2017), and alleviating harmful generations (Kirk et al., 2022; Weidinger et al., 2021; Banko et al., 2020), is essential for language models. In modern development, vocabularies are often obtained by training tokenizers with a pre-defined vocabulary size on the specific corpus (Sennrich et al., 2015; Radford et al., 2019), and once built, they are kept unchanged in the following model construction and deployment(Dagan et al., 2024).

Although it is enough for basic language modeling, this *static* setting makes vocabulary be quietly ignored (Dagan et al., 2024) in advanced generation tasks (Gao et al., 2023; Rozière et al., 2024; Fried et al., 2023). For example, it can not be augmented with new phrases for better adapting an unseen domain (Koehn and Knowles, 2017; Jin et al., 2020; Dynamic vocabulary constructed based on input text





Figure 1: Generation process of our proposed dynamic vocabulary. The model's vocabulary dynamically changes based on the input text, with phrases serving as basic blocks that are directly input and output.

Chen et al., 2022) or verbatim reference text spans for better inline evidence generation (Menick et al., 2022; Gao et al., 2023). To bring vocabulary back to the stage, it is natural to ask whether prior constraints posted by tokenization corpus and fixed vocabulary size can be relaxed.

Here, we explore vocabulary in a new *dynamic* setting. Instead of being a fixed token table, dynamic vocabulary is required to be able to include *arbitrary text spans* on demand. This setup brings new challenges to the language model. On the input side, using a single embedding layer is no longer feasible as the full table can not be enumerated. On the output side, the model needs a stronger next-token predictor as the model allows multiple oracles (tokenized to different granularity) for a single string.

In this work, we build a dynamic vocabulary by building a dynamic phrase encoder. Akin to the embedding layer, the encoder maps arbitrary text spans (called *phrases* in the following sections) to the input space of language models. It can be trained with existing language models in the same

self-supervised manner, despite that multiple to-065 kens (in the original static vocabulary) can be input or output at a single step. Though the paradigm is almost unchanged, supporting dynamic tokens needs non-trivial modification on data curation. Specifically, we find that, to prevent the learned model from either biased towards full static token 071 outputs or towards full new phrase outputs, it is crucial to make the two properly interleaved in training samples. We also show that without informative negative samples, the token encoder is hard to learn. We thus develop two retrieval-based and generation-based methods for accelerating the learning of the dynamic phrase encoder.

> The obtained dynamic vocabulary can be deployed in a way of plug-and-play: the underlying architecture of language models are kept, and those new on-demand phrases can be used as ordinary tokens during the generation. To evaluate the dynamic vocabulary, we investigate three exemplar applications, including basic language modeling, domain adaptation, and generating citations for question answering. Results show that the new flexibility of vocabulary both improve basic generation performances (e.g., stronger fluency and diversity scores on WikiText-103 (Merity et al., 2016) with lower latency) and provide a new tool to handle advanced language modeling tasks (e.g., generating more accurate citations with QA scores also increased).

### 2 The Approach

084

086

097

099

102

103

104

105

106

107

108

109

110

111

112

113

114

## 2.1 Problem Definition

Given a language model LM, denote V as its vocabulary, and  $x = x_1, x_2, ..., x_n$  as a tokenized sentence according to  $V(x_i \text{ is a token in } V)$ . A dynamic vocabulary  $V' = V \cup P$  augments V with arbitrary phrases (text spans) P. The same sentence x now can be tokenized to a different sequence  $x'_1, x'_2, ..., x'_m$ , where  $x'_i \in V'$ . The usage of dynamic vocabulary V' is identical to the vanilla static vocabulary V: the language model LM can accept any token in V' as input and choose output tokens from V'.

Supporting arbitrary phrase set P and integrating V' with language models are two cruxes to implement dynamic vocabularies. For the first one, it is possible to support new phrases by fine-tuning the language model with V', but it requires updating the model when P changes which can hardly be used in real applications. We will also see that,

for the second crux, simply replacing V with V' fails to learn the language model due to the decoding ambiguity introduced by P. We elaborate our solutions in the following sections.

#### 2.2 Dynamic Phrase Encoder

Instead of fine-tuning the language model for every possible P to support arbitrary phrase sets, we build a parametric encoder for those dynamic phrases. Once the encoder is learned, it can be deployed with the model.

Specifically, the dynamic phrase encoder is built with a casual Transformer. To get the representation of a phrase  $p \in P$ , it first tokenizes  $p = w_1, w_2, ..., w_s$  according to the static vocabulary V, and after going through several casual Transformer layers followed by an MLP, the hidden vector of the last token  $\mathbf{h}_s$  is the vector representation of p.

The above setting is different from existing works in three ways (Lan et al., 2023; Teehan et al., 2024). First, it is common to use Transformer encoder (full attention) to build the phrase encoder, while we apply Transformer decoders (casual masking). The choice is mainly guided by efficient negative sampling (see Section 2.4 for further details).

Second, the dynamic phrase encoder adopts the same tokenizer of LM (which is used to build the static vocabulary V). Sharing tokenizers means the language model doesn't need to load additional vocabularies and tokenizers during inference. <sup>1</sup>

Third, following the same idea of consistent treatment, we also use a non-contextualized representation of phrases, which makes the new phrases more like the original tokens in V. Contextualized representations can also be used (Joshi et al., 2020; Lan et al., 2023), but it means that, besides the phrases themselves, the contexts of them should also be included in the dynamic vocabulary.

To summarize, the considerations above aim to make the dynamic phrase encoder aligns with the embedding layer as much as possible: both of them map tokens (phrases) into the input space of the language model, one by lookup operations, and another by running the phrase encoder.

## 2.3 Inference with Dynamic Vocabulary

In testing time, the new dynamic vocabulary can be used as the ordinary vocabulary. We take an

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

115

116

117

118

<sup>&</sup>lt;sup>1</sup>As a comparison, the phrase encoder in COG (Lan et al., 2023) is BERT, and one should load both the BERT vocabulary and GPT-2 vocabulary when testing.

161auto-regressive language model LM as an exam-162ple. For a set of new phrases P,  $^2$  we run the163learned dynamic phrase encoder to get representa-164tions of its phrases, denoted by a matrix P. The165language model's input and output embedding ma-166trices  $W_{emb,in}, W_{emb,out}$  are expanded with these167embeddings,

 $\mathbf{W}_{\mathrm{emb,in}}' = [\mathbf{W}_{\mathrm{emb,in}}, \mathbf{P}],$ 

 $\mathbf{W}'_{\text{emb out}} = [\mathbf{W}_{\text{emb.out}}, \mathbf{P}].$ 

At each auto-regressive decoding step, the lan-

guage model LM outputs a hidden vector  $\mathbf{h}_{< i}$  rep-

resenting current prefix  $x'_{< i}$ , the probability of next

 $Z = \sum_{k' \in V} \exp(\mathbf{h}_{\langle i} \cdot \mathbf{e}_{\text{out}}^{k'}) + \sum_{k' \in P} \exp(\mathbf{h}_{\langle i} \cdot \mathbf{e}_{\text{out}}^{k'}),$ 

where  $\mathbf{e}_{\mathrm{out}}^k$  is the k-th column of  $\mathbf{W}_{\mathrm{emb,out}}'$ . When

the *i*-th token is selected, no matter whether it is

a token in V or a phrase in P, its embedding is

looked up from  $\mathbf{W}'_{\mathrm{emb,in}}$  as the input of the next

 $\mathbb{P}(x'_i = k | x'_{< i}) = Z^{-1} \exp(\mathbf{h}_{< i} \cdot \mathbf{e}^k_{\text{out}})$ 

168

169

172

173

token is

decoding step.<sup>3</sup>

174

176 177

178

179

180

184

185

188

189

190

191

192

193

194

195

196

197

# 2.4 Training with Dynamic Vocabulary

**Building Samples** To train the dynamic phrase encoder, we follow the same self-supervision regime as the training of language models. The key difference here is that, besides tokens in V, we need to organize phrases (text spans) in a training sample for learning the phrase encoder. In particular, 1) the diversity of training time in-domain phrases would influence the generalization of the learned phrase encoder, and 2) the distribution of phrases in samples would influence the how the language model switches between tokens and phrases. For building phrases, we test the following two methods.

• *"real" phrases.* We can use classical chunking algorithms to recognize phrases in a sentence. The result phrases can recognized as a single

grammatical unit or as a common word collocation. Here, we follow Lan et al. (2023) to use an unsupervised chunker forward maximum matching (FMM). Basically, FMM recognizes phrases that frequently appear in a support corpus and as long as possible. The algorithm (and other external chunkers) may need additional time costs to compile samples (e.g., in our experiments, FMM needs  $\approx 15$  hours to build its phrase table).

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

231

232

233

234

235

236

237

238

239

240

241

242

• *Ngrams*. Another candidate set of phrases is ngrams, which is much simpler to build than involving external chunkers. Though a ngram may not carry a meaning, it could be a stronger learning target for the phrase encoder: the connections between ngrams and its contexts are more complex than "real" phrases (as they usually follow the simple patterns which are used to extract them). We study two settings, ngrams of words and ngrams of tokens (denoted by N-words and N-ids respectively). Taking N-words as an example, a word tokenizer <sup>4</sup> first recognizes words in a sentence, then randomly sequences of 2-5 consecutive words are grouped into phrases.

Next, given a sentence and a set of candidate phrases, we need to determine the distribution of phrases. One may build samples with full ngrams phrases, but they could be both hard to learn (the learning ignores the prior knowledge of original vocabulary V in the model), and hard to apply (the setting is rare in applications). In our practice, to accelerate learning and prevent unnecessary data bias, it is crucial to make phrases and tokens properly interleaved in training samples. Therefore, we control the interval between two phrases to be at least five tokens.

**Negative Phrases** After building training samples, we can directly optimize the log-probability defined in Equation 1, which requires the correct next token in  $V' = V \cup P$  has the largest logit than other tokens in V and P (negative tokens). However, the number of phrases in the training set would be large, and it is prohibitive to include all of them in the loss function. <sup>5</sup> A common workaround is to include only in-batch and prebatch phrases in P (Gao et al., 2021). Unfortu-

<sup>&</sup>lt;sup>2</sup>The phrase set P can change at each decoding step. Here, for simplicity, we assume it is kept unchanged during testing, and we can run the dynamic phrase encoder only once.

<sup>&</sup>lt;sup>3</sup>When decoding a phrase, another option adopted by (Joshi et al., 2020; Lan et al., 2023) is to unfold tokens in the phrase and input them individually. Despite the inconsistency between input and output vocabulary (our experiments indicate a negative influence on performances), this setting may also slow the decoding speed (or generate shorter texts given a fixed length budget) even if it can predict a phrase.

<sup>&</sup>lt;sup>4</sup>N-words uses the word tokenizer in the NLTK toolkit, and N-ids uses GPT-2's tokenizer.

<sup>&</sup>lt;sup>5</sup>It is worth noting that all training time phrases are dropped after learning the encoder. For ngram phrases (N-words and N-ids), phrases are built on the fly in the batching process, and there is no global training time P.



Figure 2: The overall architecture of our proposed dynamic vocabulary. During training, there are four sources of negative phrases: pre-batch, corpus-retrieval, self-retrieval, and generation. Phrases are embedded by the dynamic phrase encoder with an additional linear layer. The hidden layer of the last token serves as the phrase embedding. In the model input layer, phrases are treated as a basic brick without splitting into tokens.

nately, it doesn't help learning the phrase encoder. 243 Specifically, we find that the model struggles to 244 correctly transit from a phrase token to an ordinary token and vice versa. More concretely, when pre-246 dicting a phrase  $p = w_1, w_2, ..., w_s$ , the dynamic 247 phrase encoder has trouble on distinguish p from 248 1) phrases which are prefixes of that phrase (e.g.,  $w_1w_2$  and  $w_1w_2w_3$ ) and 2) phrases which have p as their prefix (e.g.,  $pw_{s+1}$  and  $pw_{s+1}w_{s+2}$ ). There-251 fore, we also manually add the above phrases to Pin each batch (we call them informative negative phrases).

255

257

259

261

262

263

For the first type, we can simply enumerate all prefixes of *p*. For the second type, we develop *retrieval-based* and *generation-based* methods for getting successor tokens of *p*,

- retrieval-based continuation finds appearances of p in a support corpus and takes p and its successor tokens there as negative phrases (corpus-retrieval).
   <sup>6</sup> One simplification is only considering p's successor tokens in the current sample (self-retrieval).
- generation-based continuation, instead of searching corpus, tries to get synthetic negative phrases by employing a language model. <sup>7</sup> The model is

prompted with p and the following generations are included in P (generation).

268

269

271

272

273

274

275

277

278

279

281

284

285

286

287

289

290

292

293

Finally, regarding getting embeddings of these informative negative phrases, recall that we adopt an causal Transformer as the phrase encoder and use the hidden state of the final token to represent p, the embeddings of negative phrases could be efficiently obtained by feeding the longest phrase to the encoder.

**Loss Functions** The first part of the training loss is defined by Equation 1 (with negative samples added to P), which we denote by  $L_p$ . We also add a special setting of  $L_p$  in the loss (denoted by  $L_t$ ), in which  $P = \emptyset$  (i.e., the vanilla language modeling). It helps to maintain generation ability with the static vocabulary V.

We can further align the above two settings by requiring their next token distributions at each token position are close (measured by KL divergence). Concretely, given a sentence x, recall that (Section 2.1) the oracle of training  $L_p$  is  $x'_1, x'_2, ..., x'_m$ , the oracle of training  $L_e$  is  $x_1, x_2, ..., x_n$ . Assume a function  $\sigma$  which aligns  $x'_i$  to a token position in  $L_e$ 's oracle: if  $x'_i$  is a token in V, it is mapped to the same token position, otherwise,  $x'_i$  is mapped to its last token's position.

$$L_{kl} = \frac{1}{m} \sum_{i=0}^{m} \mathrm{KL}(\mathbb{P}(x'_i | x'_{< i}) || \mathbb{P}(x_{\sigma(x'_i)} | x_{< \sigma(x'_i)})).$$
294

<sup>&</sup>lt;sup>6</sup>Due to the time complexity of matching phrases, we only adopt corpus-retrieval when phrases are obtained by FMM, and keep the efficiency of Ngram phrases.

<sup>&</sup>lt;sup>7</sup>Here we use GPT-2, stronger models can also be applied.

386

344

345

The final loss function is  $L = L_p + L_t + L_{kl}$ .

### **3** Experiments

### 3.1 Setups

296

301

303

304

305

311

314

315

317

319

321

326

328

**Configurations** As a default configuration, we use GPT-2 (Radford et al., 2019) as the initial model of both the language model and dynamic phrase encoder. In testing time, we construct phrases for each test sample. By default, we follow COG (Lan et al., 2023) to use DPR (Karpukhin et al., 2020) to retrieve the top-k (k = 32) relevant documents from the training set We consider all n-grams of top-k documents as candidate phrases. Please refer to Appendix B for more details.

**Baselines** We compare the proposed method with the following state-of-the-art models as baselines:

**Transformer** (Vaswani et al., 2023) is the standard token-level language model. We fine-tune the pre-trained GPT2 in our experiments.

**KNN-LMs** (Khandelwal et al., 2020) extends a pre-trained neural language model by linearly interpolating it with a k-nearest neighbors(KNN) model.

**RETRO** (Borgeaud et al., 2022) is a retrievalenhanced transformer that combines a frozen Bert retriever, a differentiable encoder, and a chunked cross-attention mechanism.

**COG** (Lan et al., 2023) decomposes text generation into a series of copy-and-paste operations. It first retrieves semantically relevant documents and then considers all n-grams within them as candidate phrases.

**MWT** (Gee et al., 2023) propose to expand vocabulary with top-k frequent n-grams in support corpus. Rather than expanding vocabulary dynamically, it still focuses on building a static vocabulary.

**Metrics** Following previous works (Lan et al., 330 2023; Cao et al., 2024), we utilize four automatic 331 evaluation metrics to measure the quality of the 332 generated texts: (i) MAUVE (Pillutla et al., 2021) 333 measures the distribution similarity between the reference text and generated text; (ii) Rep-n (Welleck 335 et al., 2019) reflects the repetition at different ngram levels in the generated text; (iii) Diversity 337 (Welleck et al., 2019) evaluates the variety of gen-339 erated content; and (iv) Perplexity measure the difficulty in predicting the next word in a sequence. 340 In addition, we also compare the average time cost 341 of different methods to decode a continuation con-342 sisting of 128 tokens given a prefix of 32 tokens, 343

referred to as **latency**. The details for these metrics can be found in Appendix C

In the following experiments, we verify the effectiveness of our methods in three exemplar applications: basic language, domain adaptation, and generating citations for question answering.

### 3.2 Basic Language Modeling

We employ the WikiText-103 (Merity et al., 2016) dataset for training and conduct a language modeling task using its test set to evaluate our model. For each test sample, we provide the first 32 tokens as a context prefix, and both the baselines and our proposed model are tasked with generating the subsequent 128 tokens. For a fair comparison, we specify that our method generates text of the same length as other baselines (128 tokens after tokenization by GPT2Tokenizer (Radford et al., 2019)).

The results are listed in Table 1. We find that,

- Notably, our model outperforms standard Transformer with 5.22% improvements in MAUVE, indicating the high quality of the generated text. While MWT gets 24.74% MAUVE, it is a static vocabulary. Meanwhile, our model achieve 47.44% diversity, which significantly exceeds the baselines.
- Additionally, we delve into the discussion of generation latency. Given that phrase representations are pre-computed, encoding times are excluded. As seen, our model exhibits a clear advantage in faster text generation, nearly 0.99s to generate 128 ids. This is because a single phrase encompasses several tokens, translating to fewer decoding steps for identical-length output.
- However, the perplexity of our model is marginally higher than that of the Transformer. This discrepancy could potentially stem from the fact that during testing, the input prefixes are strictly composed of tokens from a fixed vocabulary, whereas the model is not subjected to such constraints during training, which results in an inconsistency between the training and testing data distributions, potentially leading to the observed difference in perplexity scores.

we also evaluate the generation results under 387nucleus sampling (p = 0.95), with detailed metrics 388provided in the appendix A. 389

Model	MAUVE $\uparrow$	Rep-2↓	Rep-3↓	Rep-4↓	Diversity $\uparrow$	Latency(s)↓	$\mathbf{PPL}\downarrow$
Transformer	20.47	41.96	36.82	33.74	24.30	1.10	3.60
RETRO	19.59	43.78	38.58	35.35	22.33	4.43	3.96
KMM-LM*	19.92	43.79	38.76	35.69	22.13	10.36	3.48
COG	21.61	34.77	30.67	28.35	32.41	1.04	7.89
MWT	24.74	33.78	26.72	22.76	37.48	1.13	5.58
Ours	25.69	27.77	20.80	17.08	47.44	0.99	8.03

Table 1: The automatic evaluation on the test set of WikiText-103. \* indicates that we directly utilize the results from the COG paper for KNN-LM due to limited GPU memory. Additionally, our method retrieves only 32 text documents for phrase segments during evaluation, whereas COG retrieves 1024. MWT Gee et al. (2023) apply MWT to encoder-only model but we implement MWT with GPT-2

Comparison	Better	No Prefer	Worse
Transformer	0.61	0.05	0.34
MWT	0.58	0.02	0.40
COG	0.58	0.08	0.34

Table 2: GPT evaluation on WikiText-103.

390

392

395

400

401

402

403

404

405

406

407

408

409

**GPT4 Evaluation.** To ensure the reliability of the aforementioned automated evaluation and the quality of the generated text, we employ GPT-4 (gpt-4-0125-preview) (Achiam et al., 2023) for assessment. Although human evaluation is considered the gold standard for assessing human preferences, it is slow and costly. Zheng et al. (2023) have demonstrated that strong LLMs, such as GPT-4, can match most human preferences well. Specifically, we randomly sample 100 cases and evaluate the results of the Baselines and our model. GPT-4 is asked to evaluate the generated texts by considering fluency, coherence, informativeness, and grammar. Detailed implementations and prompt can be found in Appendix D.

As depicted in Table 2, our proposed method significantly outperforms the Transformer with Better cases of 61 and only 34 cases of slight inferiority, which suggests that our model possesses enhanced generation capabilities.

Sequence Compression Sequence compression 410 reflects the length of text that a model can accom-411 modate within the same window size. Following 412 Dagan et al. (2024), we measure the two compres-413 sion metrics: Normalized Sequence Length (NSL) 414 and the average number of Bytes per Token. For-415 mally, we define NSL as the token count of an 416 417 encoded sequence from a tokenizer T. Given that our proposed model does not incorporate a genuine 418 tokenizer, we record the tokens and phrases output 419 of each sample and regard it as the tokenization 420 result of our proposed method. We use tokenizers 421

of both GPT and MWT on the textual outputs of our model.

Model	$\mathbf{NLS}\downarrow$	UTF-8 Bytes ↑
Transformer	127.72	4.28
MWT	114.84	4.77
Ours	101.38	5.54

Table 3: Compression on WikiText-103. Since COG, KNN-LM, and RETRO do not modify the model's tokenizer or input vocabulary, the compression results are the same with the Transformer.

As shown in the table 3, our proposed model holds the highest information content per token, averaging 101.38 tokens or phrases per sequence and 5.54 UTF-8 bytes per token, and necessitates fewer tokens or phrases to generate the same text. This is a natural consequence of the fact that the dynamically added phrases inherently consist of at least two tokens.

**Negative Samples** As mentioned above, we have designed four different negative sampling strategies to investigate the influence of negative examples on the generation results. The results shown in Table 4 indicate that the choice of negative sampling method significantly impacts the fluency and quality of the generated text.

- Specifically, compared with the other negative sampling methods, the in-batch and pre-batch negative sampling methods result in a markedly higher PPL, approximately 10 and 3 points higher in the FMM setting, irrespective of the phrase segment method used. This observation suggests that increasing the phrases that are prefixes of gold phrase or vice versa is crucial for generating fluent text.
- Furthermore, for any given negative sampling method, the PPL in the FMM setting is con-

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

Negative Samples	<b>MAUVE</b> ↑	Diversity ↑	PPL
FMM			
in-batch	21.95	57.92	16.48
pre-batch	22.28	48.91	9.02
generation	22.87	42.19	6.34
corpus-retrieval	21.98	41.32	6.40
self-retrieval	21.65	41.67	6.39
self-retrieval + generation	21.25	42.40	6.62
N-words			
in-batch	25.42	64.82	18.04
pre-batch	23.98	61.80	14.60
generation	24.99	49.03	8.51
self-retrieval	24.83	48.46	8.13
self-retrieval + generation	25.69	47.44	8.03
N-ids			
in-batch	25.69	67.53	22.25
pre-batch	23.66	61.16	14.83
generation	23.91	46.40	8.07
self-retrieval	23.64	48.38	8.36
self-retrieval + generation	24.85	47.08	8.21

Table 4: The automatic evaluation on different negative samples and training samples. During testing, each phrase is constrained to 2-8 tokens. Here, the pre-batch method contains prefixes of gold phrases as well and the number of preceding batches is set to 1.

sistently lower than that of the N-words and Nids methods. This phenomenon occurs because phrases obtained through FMM possess a relatively well-defined meaning.

 Interestingly, the average MAUVE values for the N-words and N-ids are approximately 1% higher than that of FMM. The result indicates that sample building has a substantial influence on the text quality Cao et al. (2024).

#### 3.3 Domain Adaptation

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

The plug-and-play property of our proposed model motivates us to explore the model's performance on a specific domain in a training-free manner. Specifically, we investigate the performance of the model trained on the WikiText-103 dataset but tested on another one. The LawMT (Koehn and Knowles, 2017) dataset is an English-German translation dataset in the legal domain. Consistent with prior work (He et al., 2021a; Alon et al., 2022; Lan et al., 2023), we utilize the English portion of this dataset.

470 Results. As shown in table 5, our model achieves
471 26.35 % MAUVE score and 82.99 % Diversity
472 score and even outperforms the transformer fur473 ther fine-tuned on LawMT datasets with 3.29 % on
474 MAUVE, indicating that our model can generate
475 high-quality text in the specific domain without
476 further fine-tuning.

Model	MAUVE ↑	Diversity $\uparrow$	Latency(s)↓	PPL
Transformer w/o FT	22.97	72.12	1.03	3.21
Transformer w/ FT	23.06	80.21	1.02	3.54
RETRO	19.07	72.68	5.72	3.78
KMM-LM*	23.32	19.85	-	-
COG	19.46	81.93	1.39	6.74
MWT	24.55	77.45	1.10	5.38
Ours	26.35	82.99	1.09	7.61

Table 5: The automatic evaluation on Law-MT. In this experiment, we retrieve 512 documents for each sample. To guarantee a fair comparison, we also evaluate the performance of the Transformer model both with and without further fine-tuning on LawMT.

In addition, we evaluate the sequence compression ratio and conduct a GPT Evaluation on LawMT. The details can be found in Appendix D, E. 477

478

479

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

#### 3.4 Generation with Citations

Given that we can customize our dynamically expanding vocabulary as needed, and that the phrases always have a corresponding document, a straightforward downstream task would be citation generation. ASQA (Stelmakh et al., 2022) is a long-form QA dataset and we utilize the dataset processed by Gao et al. (2023), where potential candidate documents have already been provided for each query. We extract phrases from the provided documents. We first label each document with a unique ID marker starting from 1 and add the corresponding marker to each phrase, such as "dynamic vocabulary[1]". Therefore, during generation, we can easily know which document the generated phrase belongs to, thereby solving the citation task.

**Results** We evaluate the generated results from two perspectives: citation quality and QA answer accuracy. The detailed definitions of the metrics can be found in Gao et al. (2023). We provide the model with the top-k articles and leverage incontext learning to instruct it to cite accordingly.

The experimental outcomes demonstrate a significant boost in the citation capability of our model with citation recall and precision surpassing TinyL-lama baseline by 9.14% and 27.76%, respectively. However, phrase collections have a significant impact on the citation results. The phenomenon occurs potentially due to the property of the n-grams approach, which yields an extensive collection of phrases. Consequently, there is a higher likelihood of encountering suitable phrases that align with the context compared with the parsing method.

Furthermore, our model exhibits a superior QA

Model(shot-1)	Citation_rec	Citation_prec	QA-EM	QA-F1	Rouge-L
TinyLlama	0.62	1.54	6.00	8.78	25.43
ours					
w/ n-grams	9.76	29.30	8.88	11.83	30.06
w/ parsing	2.94	9.17	9.87	13.06	30.16
w/o phrases	0.00	0.00	8.81	11.81	29.60

Table 6: The automatic evaluation on ASQA. In this experiment, we opt for TinyLlama as the language model to imbue the model with in-context learning capabilities. All baseline models are configured in a one-shot setting, with the number of candidate documents set to 3. Parsing denotes that we use Stanza parser (Qi et al., 2020) to extract phrases from candidate documents, which ensures that the phrases possess a relatively complete and well-defined meaning.

performance compared to Tinyllama, with an Exact Match score of 9.87% and an F1 of 13.06%. Due to our model's further training on the WikiText-103 dataset, which is derived from a subset of Wikipedia articles and the property of ASQA tasks that necessitate Wikipedia-based information, our model's QA performance is expected to be superior when not utilizing phrases (i.e., ours w/o phrases) as compared to directly employing TinyLLama.

### 4 Related Work

515

516

517

518

519

520

521

522

524

525

526

527

532

533

534

535

537

538

539

541

542

543

544

546

547

**Tokenizer** Tokenizer is an essential component of language models (Dagan et al., 2024; Mielke et al., 2021), responsible for transforming raw text into a sequence of tokens. Byte-Pair Encoding (BPE) is commonly used to build tokenizer (Radford et al., 2019; Liu et al., 2019; Lewis et al., 2019; He et al., 2021b) and, there exist other tokenization algorithms, such as Unigram (Kudo, 2018) and WordPiece tokenization used in BERT (Devlin et al., 2019). However, these tokenizations are limited to subwords or whole words. Kumar and Thawani (2022) and Gee et al. (2023) generalize the BPE algorithm to multi-words and multi-tokens separately. Whereas these approaches necessitate training the tokenizer and remain static.

COG (Lan et al., 2023), which employs a "dynamic vocabulary", retrieves related documents based on the input text and expanded vocabulary with phrases extracted from these documents. However, COG only employs dynamic vocabulary in the output module and splits phrases into tokens in the input. In this paper, we treated phrases as atomic units same as tokens, and dynamically expanded vocabulary both in input and output layers.

549Sequence CompressionLanguage models are550constrained by the limited length of input se-551quences they can process. Increasing this length552results in a prohibitive computational overhead. A

series of techniques have been proposed to compress sentences into one or a few tokens or latent representations (Qin and Van Durme, 2023; Chevalier et al., 2023; Bulatov et al., 2022; Mu et al., 2024). MWT (Gee et al., 2023) enhances compression by retraining the tokenizer, incorporating the most frequent n-grams of a support corpus into the vocabulary. In contrast to the static vocabulary of MWT, our method dynamically adapts the model's vocabulary to the input text, resulting in a more flexible and efficient adaptation.

## 5 Conclusion

In this paper, we propose a novel approach for dynamically adjusting the model's vocabulary based on input text. It is a plug-and-play approach that can be performed simultaneously with pre-training tasks. We investigated standard language modeling, domain adaptation, and citation generation, and discussed the impact of different training samples and negative phrase construction methods on the quality of generated text. Our experimental results show that our proposed model can rapidly generate high-quality, high-compression text compared to baselines.

### 6 Limitations

In this paper, we propose a method to dynamically expand the vocabulary based on the input text. While our approach can improve generation speed and increase the effective length of the generated text, our model does not modify the underlying tokenizer. As a result, it cannot reduce the token numbers for known input information like prompts or questions. The dynamic vocabulary is, therefore, limited to the subsequent content generated by the model.

Furthermore, to obtain embedding representations for phrases, a dynamic phrase encoder is nec-

588

589

553

554

555

556

557

558

559

697

698

643

644

590

596

- 599

- 604

- 610
- 611 612 613
- 614 615 616
- 617 618
- 619
- 623

625

- 627
- 631 632

637

- 641 642

essary. This encoder has a more intricate structure compared to the model's linear embedding layer and requires additional memory allocation during implementation.

Lastly, our method relies on external techniques, such as a retriever, to obtain relevant documents and extract phrases from them during testing. This adds complexity to the preparation process.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.
- Uri Alon, Frank F. Xu, Junxian He, Sudipta Sengupta, Dan Roth, and Graham Neubig. 2022. Neuro-symbolic language modeling with automatonaugmented retrieval. Preprint, arXiv:2201.12431.
- Michele Banko, Brendon MacKeen, and Laurie Ray. 2020. A unified taxonomy of harmful content. In Proceedings of the Fourth Workshop on Online Abuse and Harms, pages 125-137, Online. Association for Computational Linguistics.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. Improving language models by retrieving from trillions of tokens. Preprint, arXiv:2112.04426.

- Aydar Bulatov, Yuri Kuratov, and Mikhail Burtsev. 2022. Recurrent memory transformer. In Advances in Neural Information Processing Systems.
- Bowen Cao, Deng Cai, Leyang Cui, Xuxin Cheng, Wei Bi, Yuexian Zou, and Shuming Shi. 2024. Retrieval is accurate generation. Preprint, arXiv:2402.17532.
- Zhivu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2022. Finqa: A dataset of numerical reasoning over financial data. Preprint, arXiv:2109.00122.
- Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. 2023. Adapting language models to compress contexts. Preprint, arXiv:2305.14788.
- Gautier Dagan, Gabriel Synnaeve, and Baptiste Rozière. 2024. Getting the most out of your tokenizer for pre-training and domain adaptation. Preprint, arXiv:2402.01035.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. Preprint, arXiv:1810.04805.
- Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen tau Yih, Luke Zettlemoyer, and Mike Lewis. 2023. Incoder: A generative model for code infilling and synthesis. Preprint, arXiv:2204.05999.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Dangi Chen. 2023. Enabling large language models to generate text with citations. Preprint, arXiv:2305.14627.
- Leonidas Gee, Leonardo Rigutini, Marco Ernandes, and Andrea Zugarini. 2023. Multi-word tokenization for sequence compression. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track. Association for Computational Linguistics.
- Leonidas Gee, Andrea Zugarini, Leonardo Rigutini, and Paolo Torroni. 2022. Fast vocabulary transfer for language model compression. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Industry Track, pages 409– 416, Abu Dhabi, UAE. Association for Computational Linguistics.
- Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. 2021a. Efficient nearest neighbor language models. Preprint, arXiv:2109.04212.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021b. Deberta: Decodingenhanced bert with disentangled attention. Preprint, arXiv:2006.03654.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. Preprint, arXiv:2009.13081.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. Span-BERT: Improving pre-training by representing and predicting spans. Transactions of the Association for Computational Linguistics, 8:64–77.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for opendomain question answering. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6769-6781, Online. Association for Computational Linguistics.

- 700 701 702 703 704 705 706 707 708 709 710 711 712
- 712 713 714 715 716 717 718 719 720 721
- 722 723 724 725 726 727 728 729 730 731
- 729 730 731 732 733 734 735 736 737
- 7
- 741
- 742
- 743 744
- 745 746
- 747 748 749
- 750 751
- 751
- 752 753

- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through memorization: Nearest neighbor language models. *Preprint*, arXiv:1911.00172.
- Hannah Kirk, Abeba Birhane, Bertie Vidgen, and Leon Derczynski. 2022. Handling and presenting harmful text in NLP research. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 497–510, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. *Preprint*, arXiv:1706.03872.
  - Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. *Preprint*, arXiv:1804.10959.
  - Dipesh Kumar and Avijit Thawani. 2022. Bpe beyond word boundary: How not to use multi word expressions in neural machine translation. In *Proceedings* of the Third Workshop on Insights from Negative Results in NLP, pages 172–179.
  - Guillaume Lample and Alexis Conneau. 2019. Crosslingual language model pretraining. *Preprint*, arXiv:1901.07291.
  - Tian Lan, Deng Cai, Yan Wang, Heyan Huang, and Xian-Ling Mao. 2023. Copy is all you need. *Preprint*, arXiv:2307.06962.
  - Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *Preprint*, arXiv:1910.13461.
  - Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pretraining for neural machine translation. *Preprint*, arXiv:2001.08210.
  - 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. *Preprint*, arXiv:1907.11692.
  - Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. *Preprint*, arXiv:1711.05101.
- Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nat McAleese. 2022. Teaching language models to support answers with verified quotes. *Preprint*, arXiv:2203.11147.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *Preprint*, arXiv:1609.07843.

Sabrina J. Mielke, Zaid Alyafeai, Elizabeth Salesky, Colin Raffel, Manan Dey, Matthias Gallé, Arun Raja, Chenglei Si, Wilson Y. Lee, Benoît Sagot, and Samson Tan. 2021. Between words and characters: A brief history of open-vocabulary modeling and tokenization in nlp. *Preprint*, arXiv:2112.10508. 754

755

758

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

790

791

792

793

794

795

796

798

799

800

801

802

803

804

805

- Jesse Mu, Xiang Lisa Li, and Noah Goodman. 2024. Learning to compress prompts with gist tokens. *Preprint*, arXiv:2304.08467.
- Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. Mauve: Measuring the gap between neural text and human text using divergence frontiers. *Preprint*, arXiv:2102.01454.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations.*
- Guanghui Qin and Benjamin Van Durme. 2023. Nugget: neural agglomerative embeddings of text. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code Ilama: Open foundation models for code. *Preprint*, arXiv:2308.12950.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.
- Felix Stahlberg. 2020. Neural machine translation: A review. *Journal of Artificial Intelligence Research*, 69:343–418.
- Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. ASQA: Factoid questions meet long-form answers. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 8273–8288, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ryan Teehan, Brenden Lake, and Mengye Ren. 2024. College: Concept embedding generation for large language models. *arXiv preprint arXiv:2403.15362*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *Preprint*, arXiv:1706.03762.

807 808

809

810

811 812

813

814

815

816

817

818 819

820

821

822

823 824

825

826

827

828

829

831

- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023. Large language models are not fair evaluators. *Preprint*, arXiv:2305.17926.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. 2021. Ethical and social risks of harm from language models. *Preprint*, arXiv:2112.04359.
  - Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2019. Neural text generation with unlikelihood training. *Preprint*, arXiv:1908.04319.
  - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.

# 834

- 836
- 837

838

839

841

842

845

852

853

857

858

865

871

873

875

876

877

## A Full Results

We show the full results of our experiments in Tables 7, 8, 9, 10.

## **B** More IMPLEMENTATION DETAILS

The training of our proposed model was carried out on two NVIDIA RTX 3090 GPUs, each with 24GB of memory, over a total of 400,000 training steps. During the training process, we implemented a gradient accumulation step of 2, with a batch size of 4. We also used a linear learning rate schedule with a warmup, alongside the AdamW optimizer (Loshchilov and Hutter, 2019), maintaining the default beta values. The initial learning rate was set at 5e-5. Additionally, we applied gradient clipping with a clipping value of 1.0 to ensure training stability. When conducting nucleus sampling, we set the *p* to 0.95.

The experiments of MWT in paper (Gee et al., 2023) were conducted on encoder-only models such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019). In our implementation, we modified the foundation model to GPT2 (Radford et al., 2019), a decoder-only model, and added the top 10000 most frequent 2-grams to the original GPT2 Tokenizer. The embeddings for newly added words were initialized using Fast Vocabulary Transfer(FVT) (Gee et al., 2022). MWT was trained for a total of 150000 steps on the WikiText103 dataset.

# C More Details of Automatic Evaluation

In this section, we provide a detailed introduction to the automatic evaluation metrics.

• MAUVE Pillutla et al. (2021) measures how closely the token distribution in the generated text matches that in human-written text across the entire test set. We follow prior work and leverage the GPT2-large model to generate the scores. In our implementation, the scaling factor is set as 2.0.

• **Rep-n** Welleck et al. (2019) measures the repetition at different n-gram levels in the generated text. It is defined as  $100 \times (1.0 - \frac{|uniquen-gram(x)|}{|totaln-gram(x)|})$ . Higher Rep-n represents the severe degeneration problem in generations.

Diversity Welleck et al. (2019) evaluates the variety of generated content, which is formu-

lated as  $\prod_{n=2}^{4} (1 - \frac{Rep - n}{100})$ . More informative generations get higher Diversity scores.

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

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

925

926

• **Perplexity** is a measure of the uncertainty or difficulty in predicting the next word in a sequence. A lower perplexity score indicates that the model is more certain about its predictions.

# **D** GPT4 Evaluation

Although human evaluation is considered the gold standard for assessing human preferences, it is slow and costly. Zheng et al. (2023) have demonstrated that strong LLMs, such as GPT-4, can match most human preferences well , achieving over 80%agreement, which is the same level of agreement between humans. Therefore, LLM-as-a-judge is an interpretable approach to approximating human preferences. We random sample 100 cases and evaluate the results of the Baselines and our model. GPT-4 is asked to evaluate the generated texts by considering fluency, coherence, informativeness, and grammar. Owing to GPT4's sensitivity to the order of the two candidate sentences (Wang et al., 2023), we adhere to the approach employed in Wang et al. (2023) and determine the final result by calculating the average of the outcomes from interchanging the order of the candidate sentences.

Figure 3 shows the detailed prompt used for GPT-4. Despite the template emphasizing that the order should not affect the results (red text), large language models still exhibit a significant positional bias. Therefore, for each triplet (prefix, <generation\_1>, <generation\_2>), we include another corresponding triplet (prefix, <generation\_2>, <generation\_1>). This is done to mitigate the impact of the order of the two generations on GPT evaluation.

Here are the full results of our evaluation using GPT-4 shown in Table 11. It can be seen that our model is capable of producing generations that are comparable or even superior to the baselines.

# E Sequence Compression On LawMT

Analogous to the section 3.2, we calculate the compression ratio of LawMT. The conclusion aligns with those from section 3.2, indicating that our model could yield the highest information density per token. And for an equal number of tokens, our model encompasses a longer effective text length.

Model	Decoding	MAUVE $\uparrow$	Rep-2↓	Rep-3↓	Rep-4	<b>Diversity ↑</b>	Latency(s)↓	PPL
Transformor	greedy	20.47	41.96	36.82	33.74	24.30	1.10	3.60
11 ansior mer	nucleus	25.05	5.40	1.44	0.51	92.76	1.15	31.01
ретро	greedy	19.59	43.78	38.58	35.35	22.33	4.43	3.96
KEIKO	nucleus	20.77	5.83	1.91	0.83	91.61	5.43	39.74
KWW I W*	greedy	19.92	43.79	38.76	35.69	22.13	10.36	3.48
IX1V11V1-1/1V1	nucleus	22.50	3.33	0.69	0.21	95.8	10.42	78.01
000	greedy	21.61	34.77	30.67	28.35	32.41	1.04	7.89
COG	nucleus	25.96	5.43	1.53	0.67	92.50	1.06	36.66
	greedy	24.74	33.78	26.72	22.76	37.48	1.13	5.58
Gritinivi	nucleus	25.66	4.18	0.90	0.29	94.68	1.17	55.02
Ours	greedy	25.69	27.77	20.80	17.08	47.44	0.99	8.03
Ours	nucleus	24.34	4.59	1.03	0.28	94.16	1.00	51.38

Table 7: The automatic evaluation on the test set of WikiText-103. \* denotes that the results are obtained from COG (Lan et al., 2023) paper. For each sample, the first 32 tokens are provided and models are tasked with generating the subsequent 128 tokens. We can observe that our proposed model achieves the best scores in most metrics.

Negative Samples	Decoding	MAUVE $\uparrow$	Rep-2↓	Rep-3↓	Rep-4	Diversity $\uparrow$	Latency(s)↓	PPL
FMM								
nra hatah	greedy	22.28	26.90	20.07	16.29	48.91	0.95	9.02
pre-baten	nucleus	20.59	4.62	1.07	0.35	94.03	0.88	56.28
generation	greedy	22.87	31.17	23.82	19.55	42.19	1.20	6.34
generation	nucleus	20.33	4.35	1.01	0.31	94.39	1.06	49.51
	greedy	21.98	31.47	24.39	20.26	41.32	1.12	6.40
corpus-retrievar	nucleus	20.52	4.36	1.00	0.32	94.38	1.08	51.60
colf notmicrypl	greedy	21.65	31.33	24.15	20.00	41.67	1.15	6.39
sen-reurevar	nucleus	20.63	4.37	1.00	0.35	94.34	1.04	49.93
colf retrievel + concretion	greedy	21.25	30.89	23.73	19.57	42.40	1.16	$\bar{6}.\bar{6}2$
sen-reuleval + generation	nucleus	20.34	4.24	0.96	0.29	94.57	1.04	52.27
N-words								
nno hotoh	greedy	23.98	19.58	13.63	11.02	61.80	1.16	14.60
pre-batch	nucleus	23.60	5.71	1.82	0.92	91.73	1.11	47.17
conception	greedy	24.99	26.72	19.95	16.41	49.03	0.94	8.51
generation	nucleus	24.85	4.64	1.07	0.31	94.04	0.94	50.65
colf notmicryol	greedy	24.83	27.21	20.23	16.54	48.46	0.96	8.13
sen-reurevar	nucleus	24.51	4.57	1.05	0.33	94.12	0.94	51.85
colf retrievel + concretion	greedy	25.69	27.77	20.80	17.08	47.44	0.99	8.03
sen-reurevar + generation	nucleus	24.34	4.59	1.03	0.28	94.16	1.00	51.38
N-ids								
	greedy	23.66	19.81	13.96	11.36	61.16	1.12	14.83
pre-batch	nucleus	22.84	5.17	1.52	0.67	92.77	0.92	54.52
	greedy	23.91	28.12	21.45	17.82	46.40	0.99	8.07
generation	nucleus	24.50	4.41	0.97	0.29	94.38	0.96	53.98
a alf matrices 1	greedy	23.64	27.29	20.33	16.49	48.38	1.02	8.36
sen-retrieval	nucleus	23.85	4.43	0.94	0.27	94.41	0.88	55.76
colf rational L concretion	greedy	24.85	27.85	21.04	17.36	47.08	1.01	8.21
sen-reureval + generation	nucleus	23.91	4.41	0.96	0.28	94.40	0.98	53.03

Table 8: The automatic evaluation on different negative samples with greedy and nucleus sampling (top-p: 0.95) decoding algorithms on the WikiText103 dataset. The constructions of training samples and negative phrases have a significant influence on the generated text.

Model	Decoding	MAUVE ↑	Rep-2↓	Rep-3↓	Rep-4	<b>Diversity ↑</b>	Latency(s)↓	PPL
Transformor w/o FT	greedy	22.97	13.36	9.69	7.84	72.12	1.03	3.21
ITalisioniller w/orri	nucleus	24.15	4.05	1.62	0.80	93.64	1.05	31.48
Transformer w/ FT	greedy	23.06	9.74	6.45	5.00	80.21	1.02	3.54
ITalisiofiller w/ F I	nucleus	25.12	4.36	1.73	0.87	93.17	1.08	14.94
DETDO	greedy	19.07	13.19	9.34	7.66	72.68	5.72	3.78
KEIKU	nucleus	21.26	3.30	1.18	0.55	95.03	5.54	57.40
KMM I M*	greedy	23.32	-	-	-	19.85	-	-
KIVIIVI-LAVI	nucleus	24.75	-	-	-	94.60	-	-
202	greedy	19.46	9.29	5.68	4.24	81.93	1.39	6.74
00	nucleus	24.45	4.57	1.58	0.72	93.25	0.89	32.01
	greedy	24.55	11.59	7.34	5.46	77.45	1.10	5.38
	nucleus	22.68	3.15	1.01	0.39	95.49	1.16	68.55
Ours	greedy	26.35	9.26	5.21	3.52	82.99	1.09	7.61
Guis	nucleus	24.80	3.63	1.17	0.48	94.78	0.93	60.70

Table 9: The automatic evaluation on LawMT. We directly retrieve 512 documents for each sample in this experiment. Our proposed model even outperforms the Transformer further fine-tuned on the LawMT corpus.

Negative Samples	Decoding	MAUVE ↑	Rep-2↓	Rep-3↓	Rep-4	<b>Diversity ↑</b>	Latency(s)↓	PPL
FMM								
	greedy	25.00	8.71	4.76	3.16	84.20	0.98	8.26
pre-batch	nucleus	23.19	3.71	1.19	0.50	94.66	0.83	60.34
ganaration	greedy	22.87	11.00	6.76	4.85	78.96	1.26	6.17
generation	nucleus	22.50	3.50	1.13	0.48	94.95	1.07	65.26
Patriavel complex	greedy	23.00	10.45	6.36	4.53	80.06	1.21	6.11
Retrieval-samples	nucleus	23.24	3.43	1.01	0.46	95.07	1.02	68.26
colf notmoryal	greedy	23.41	10.98	6.80	4.92	78.89	1.20	6.11
sen-reurevar	nucleus	23.22	3.48	1.05	0.43	95.10	0.98	67.14
colf retrievel + concretion	greedy	24.15	10.50	6.31	4.49	80.08	1.22	6.24
sen-reuleval + generation	nucleus	22.55	3.40	1.16	0.53	94.98	1.04	69.40
N-words								
pro hotoh	greedy	26.15	6.53	3.11	1.92	88.82	0.61	14.40
pre-baten	nucleus	25.15	4.07	1.41	0.61	94.00	0.53	45.79
generation	greedy	26.35	9.26	5.21	3.52	82.99	1.09	7.61
generation	nucleus	24.66	3.53	1.16	0.48	94.89	0.92	62.58
colf ratrioval	greedy	23.65	8.92	4.88	3.29	83.87	1.04	8.05
sen-reurevar	nucleus	24.71	3.54	1.09	0.42	95.00	0.81	62.51
colf retrievel + concretion	greedy	26.35	9.26	5.21	3.52	82.99	1.09	7.61
sen-reurevar + generation	nucleus	24.80	3.63	1.17	0.48	94.78	0.93	60.70
N-ids								
nra hatah	greedy	25.08	6.70	3.14	1.87	88.68	0.62	14.49
pre-baten	nucleus	23.93	4.25	1.46	0.65	93.74	0.43	47.94
generation	greedy	22.55	9.24	5.21	3.55	82.98	1.04	8.03
generation	nucleus	23.14	3.59	1.14	0.49	94.85	0.85	61.89
colf ratrioval	greedy	24.63	9.46	5.43	3.71	82.44	1.05	7.86
sen-iculeval	nucleus	24.19	3.58	1.11	0.44	94.94	0.78	63.87
colf ratriaval L gaparation	greedy	23.18	9.31	5.25	3.59	82.85	1.07	7.57
sen-reurevar + generation	nucleus	24.63	3.57	1.10	0.46	94.93	0.87	60.32

Table 10: The automatic evaluation on different negative samples with greedy decoding and nucleus sampling(top-p: 0.95) on the LawMT dataset.

You are a helpful and precise assistant for checking the quality of the text. [Prefix] {prefix} [The Start of Assistant 1's Generation] {Generation\_1} [The End of Assistant 1's Generation] [The Start of Assistant 2's Generation] {Generation\_2} [The End of Assistant 2's Generation] [System]

We would like to request your feedback on the performance of two AI assistants in response to the user prefix displayed above.Please rate the fluency, coherence, informativeness, and grammar. Each assistant receives an overall score on a scale of 1 to 10, where a higher score indicates better overall performance.

Please first provide a comprehensive explanation of your evaluation, avoiding any potential bias and ensuring that the order in which the responses were presented does not affect your judgment. Then, output two lines indicating the scores for Assistant 1 and 2, respectively.

Output with the following format: Evaluation evidence: <your evluation explanation here> Score of the Assistant 1: <score> Score of the Assistant 2: <score>

Figure 3: The GPT evaluation template with three slot {prefix}, {Generation\_1} and {Generation\_2}.

Comparison (VS)	Better	No Prefer	Worse
WikiText103			
Transformer	0.61	0.05	0.34
MWT	0.58	0.02	0.40
COG	0.58	0.08	0.34
LawMT			
Transformer	0.46	0.02	0.52
MWT	0.67	0.07	0.26
COG	0.50	0.05	0.45

Table 11: GPT evaluation on WikiText-103. Due to the sensitivity of GPT-4 to the order of two candidates, we got the final result by calculating the average scores by changing the order of the two candidates.

Model	NLS	UTF-8 Bytes
WikiText103		
Transformer	127.72	4.28
MWT	114.84	4.77
Ours	101.38	5.54
LawMT		
Transformer	128.79	5.22
MWT	124.94	5.39
Ours	105.38	6.53

Table 12: Compression on WikiText-103 and LawMT. Our model compresses text in a larger margin than MWT in the specific domain.