Analysing The Impact of Sequence Composition on Language Model Pre-Training

Anonymous ACL submission

Abstract

Most language model pre-training frameworks concatenate multiple documents into fixedlength sequences and use *causal masking* to compute the likelihood of each token given its context; this strategy is widely adopted due to its simplicity and efficiency. However, to this day, the influence of the pre-training sequence composition strategy on the generali-009 sation properties of the model remains underexplored. In this work, we find that applying causal masking can lead to the inclusion 011 of distracting information from previous doc-013 uments during pre-training, which negatively impacts the performance of the models on language modelling and downstream tasks. In intra-document causal masking, the likelihood of each token is only conditioned on the previous tokens in the same document, which eliminates potential distracting information from previous documents and significantly improves the performance. Furthermore, we find that 022 concatenating related documents can reduce some potential distractions during pre-training, and our proposed efficient retrieval-based se-025 quence construction method, BM25Chunk, can improve in-context learning (+11.6%), knowl-026 edge memorisation (+9.8%), and context utilisation (+7.2%) abilities of language models without sacrificing efficiency.

1 Introduction

034

042

Large Language Models (LLMs) are pre-trained on large amounts of documents by optimising a language modelling objective and show an intriguing ability to solve a variety of downstream NLP tasks (Brown et al., 2020; Biderman et al., 2023; Touvron et al., 2023; Jiang et al., 2023). Previous works emphasise the importance of pre-training data quality (e.g., Gunasekar et al., 2023; Lee et al., 2022; Tirumala et al., 2023; Soboleva et al., 2023) and diversity (e.g., Xie et al., 2023; Gao et al., 2021; Kaddour, 2023) to improve the generalisation properties of language models. However, the influence of the pre-training sequence composition strategy remains largely under-explored.

044

045

046

047

049

054

057

060

061

062

063

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

For most decoder-only language model pretraining pipelines (e.g., Shoeybi et al., 2019; Ott et al., 2019; Brown et al., 2020; Biderman et al., 2023; Geng, 2023; Liu et al., 2023b; Zhang et al., 2024), constructing a pre-training instance involves packing, which refers to the process of combining randomly sampled documents into a *chunk* that matches the size of the context window; and causal masking, which refers to predicting the next token conditioned on all previous tokens, including those from different documents in the chunk. An alternative to causal masking is *intra-document* causal masking, where the likelihood of each token is conditioned on the previous tokens from the same document; intra-document causal masking is not commonly used in existing open-source pretraining frameworks as it is argued to adversely impact pre-training efficiency (Brown et al., 2020; Pagliardini et al., 2023). However, to the best of our knowledge, there is no systematic analysis in the literature on how causal masking affects the generalisation properties of models despite its role in improving efficiency.

To analyse the impact of the packing and masking strategies on pre-training, we pre-train language models using intra-document causal masking (referred to as INTRADoc, Section 2.2) and compare them with models pre-trained via causal masking with several *packing* strategies by varying the relatedness of the documents in the pre-training chunks. Specifically, we analyse the results produced by a commonly used baseline method that randomly samples and packs documents (MIXChunk); a method that samples and packs documents from the same source based on their meta-information (UNIChunk); and our proposed efficient retrievalbased packing method, which retrieves and packs

related documents (BM25Chunk, Section 2.1).

087

090

093

095

101

102

103

104

105

106

107

109

110

111

112

113

114

115

116

117

118

119

121

122

123

124

125

126

127

128

129

130

131

Our experimental results indicate that using causal masking without considering the boundaries of documents can lead to the inclusion of distracting information from previous documents during pre-training (Section 3 and Section 5.1), negatively impacting the performance of the models in downstream tasks (Section 4). We observe that intradocument causal masking, which eliminates the potential distractions from irrelevant documents during pre-training, can significantly improve the performance of the model while increasing its runtime (+4% in our implementation, see Appendix A).

We also find that improving the relatedness of the documents in pre-training chunks can reduce some potential distractions from previous documents (e.g., UNIChunk avoids packing documents from different distributions, such as code and news text), which can improve the performance of causal masking models on a wide array of downstream tasks. Finally, we show that our proposed efficient retrieval-based packing method, BM25Chunk, can improve a model's language modelling (+6.8%), in-context learning (+11.6%), knowledge memorisation (+9.8%), and context utilisation (+7.2%) abilities using causal masking and thus without sacrificing pre-training efficiency.

> Our main contributions and findings can be summarised as follows:

- We systematically analyse and compare the models pre-trained using causal masking and intradocument causal masking; our experimental results reveal that using causal masking without considering the boundaries of documents can result in significant performance degradation (Section 3 and Section 4).
- We find that improving the relatedness of the documents in each pre-training chunk benefits causal masking models, and our proposed efficient retrieval-based packing method (BM25Chunk, Section 2.1) can improve the performance of language models significantly.
- We quantitatively analyse the attention distribution of the models during language modelling (Section 5.1), and investigate the burstiness property of pre-training chunks (Section 5.2); our findings indicate that models can be more robust to irrelevant contexts and obtain better performance when improving the relatedness of documents in pre-training chunks.

2 Packing and Masking Strategies for Pre-Training Sequence Composition

In this section, we formally introduce the pretraining data packing strategies, as well as causal masking and intra-document causal masking.

2.1 Packing Strategies

Let \mathcal{D}_i represent a corpus, such as Wikipedia, C4, or GitHub, and let $\mathcal{D} = \bigcup_s \mathcal{D}_s$ denote the dataset resulting from the union of such corpora. Furthermore, each corpus \mathcal{D}_s is defined as a set of documents $\mathcal{D}_s = \{d_1, \ldots, d_{|\mathcal{D}_s|}\}$, where each document d_i is defined as a sequence of tokens $d_i = (x_1, \ldots, x_{|d_i|})$.

A packing strategy involves first selecting a set of documents $\{d_i\}_{i=1}^n$ from \mathcal{D} , and then packing them into a chunk C with a fixed length |C| = L. Following Brown et al. (2020), we concatenate the documents $\{d_i\}_{i=1}^n$ by interleaving them with endof-sentence ([EOS]) tokens to construct a chunk. A packed sequence (or chunk) C is denoted as:

$$C = (d_1[\operatorname{EOS}]d_2[\operatorname{EOS}]\dots\operatorname{SPLIT}(d_n)), \quad (1)$$

where [EOS] is the end-of-sentence token, SPLIT() truncates the last document such that |C| = L, and the content of the chunk C will be removed from the dataset \mathcal{D} to avoid sampling the same documents multiple times.

In the following, we introduce three strategies to sample the documents $\{d_i\}_{i=1}^n$ from the dataset \mathcal{D} for composing each pre-training chunk, namely MIXChunk, UNIChunk, and BM25Chunk.

MIXChunk In MIXChunk (baseline), documents $d_i \in \mathcal{D}$ are sampled uniformly at random from the entire pre-training corpus \mathcal{D} :

$$d_i \sim \text{Uniform}(\mathcal{D}).$$

As a result, in MIXChunk, a chunk can contain documents from different source datasets, *e.g.*, Wikipedia and GitHub, as shown in Figure 1(a).

UNIChunk In UNIChunk, each chunk is composed of documents from a single source corpus D_s :

 $d_i \sim \text{Uniform}(\mathcal{D}_s), \text{ with } \mathcal{D}_s \subseteq \mathcal{D}.$

This helps to avoid packing documents from different distributions (such as code and news text) together. To construct a training batch, we sample sequences from each corpus \mathcal{D}_s in proportion to the number of tokens in \mathcal{D}_s .



(a) MIXChunk randomly samples documents from all corpora to construct pre-training sequences, which can pack documents from different sources. UNIChunk randomly samples documents from a single source to construct a sequence.



(b) The sequence construction process in BM25Chunk. The left part represents a document buffer that caches a set of documents randomly sampled from the corpus.

Figure 1: Packing strategies for pre-training chunks construction. (a) illustrates the compositions of MIXChunk and UNIChunk; (b) presents the sequence construction process of BM25Chunk.

BM25Chunk To improve the relevance of documents in pre-training chunks, we employ a BM25based retriever to construct pre-training chunks, referred to as BM25Chunk. Specifically, given a document $d_i \in \mathcal{D}_s$, we retrieve a sequence of documents $\{d_i\}_{i=1}^n$ by $d_{i+1} = \text{RETRIEVE}(d_i, \mathcal{D}_s)$; here, $\text{RETRIEVE}(d_i, \mathcal{D}_s)$ retrieves the most similar documents to d_i from \mathcal{D}_s based on BM25 scoring.

178

179

180

183

187

190

191

192

194

However, this retrieval process can be computationally inefficient due to the size of the pretraining corpus \mathcal{D}_s . To improve the efficiency of the retrieval step, we restrict the retrieval scope to a subset $\mathcal{B}_s \subseteq \mathcal{D}_s$ of the corpus \mathcal{D}_s , reducing the computational complexity of retrieval; the proposed approach is outlined in Figure 1(b). More formally, we introduce a document buffer $\mathcal{B}_s \subseteq \mathcal{D}_s$ that contains k documents uniformly sampled from \mathcal{D}_s , which serves as the retrieval source for constructing pre-training chunks:

$$d_1 \sim \text{Uniform}(\mathcal{B}_s), \quad d_{i+1} = \text{RETRIEVE}(d_i, \mathcal{B}_s).$$
 197

195

200

201

203

204

205

209

210

211

212

214

215

216

217

218

219

221

222

228

229

231

232

236

After retrieving a sequence of documents $\{d_i\}_{i=1}^n$ from the buffer \mathcal{B}_s for constructing a chunk, we refill the buffer by sampling new documents from documents from \mathcal{D}_s . The time complexity analysis and more details are presented in Appendix C.

2.2 Masking Strategies

Another core element of LLM pre-training is the *masking* strategy, which determines how next-token prediction distributions are conditioned on other tokens in the sequence.

Causal Masking In causal masking, each token in a sequence is predicted solely based on all preceding tokens in the sequence. More formally, given a chunk $C = (x_1, \ldots, x_{|C|})$ defined as in Equation (1), the likelihood of C is given by:

$$P(C) = \prod_{i=1}^{|C|} P(x_i \mid x_1, \dots, x_{i-1}),$$
213

where $P(x_i | x_1, ..., x_{i-1})$ denotes the probability of the token x_i given all preceding tokens $x_1, ..., x_{i-1}$ in the chunk. During pre-training, *causal masking* implies that, given a chunk C, the probability of each token in C will be conditioned on all preceding tokens, including those belonging to different documents. Causal masking is the *standard practice* when pre-training decoder-only LLMs (e.g., Shoeybi et al., 2019; Brown et al., 2020; Zhang et al., 2022; Biderman et al., 2023; Geng, 2023; Liu et al., 2023b; Zhang et al., 2024).

Intra-Document Causal Masking In intradocument causal masking, on the other hand, the probability of each token is conditioned on the previous tokens within the same document. More formally, given a chunk C defined as in Equation (1), the probability of each token d_{ij} belonging to document d_i is only conditioned on the preceding tokens within d_i :

$$P(C) = \prod_{i=1}^{n} \prod_{j=1}^{|d_i|} P\left(d_{ij} \mid d_{i1}, \dots, d_{i(j-1)}\right).$$
 233

We refer to models trained using intra-document causal masking as INTRADOC. The details of implementation are available in Appendix A.

L	Model	CommonCrawl	C4	Wikipedia	GitHub	StackExchange	Book	ArXiv	Avg.
2K	MIXChunk UNIChunk BM25Chunk INTRADoc	13.284 11.805 11.418 <u>11.631</u>	13.884 <u>13.650</u> 13.677 13.197	6.811 6.546 <u>6.237</u> 6.084	5.531 5.518 <u>4.585</u> 4.252	8.051 7.839 <u>7.623</u> 7.535	11.623 11.353 <u>11.253</u> 11.130	5.203 5.106 <u>5.059</u> 5.030	$\begin{array}{c} 9.172 \\ 8.831_{\downarrow 0.341} \\ \underline{8.550}_{\downarrow 0.622} \\ \textbf{8.410}_{\downarrow 0.883} \end{array}$
8K	MIXChunk UNIChunk BM25Chunk INTRADoc	9.645 9.478 <u>9.144</u> 8.994	14.424 14.190 <u>13.579</u> 13.173	7.010 6.897 <u>6.287</u> 6.073	7.496 7.006 <u>5.463</u> 5.010	8.634 8.456 <u>8.022</u> 7.894	11.337 11.117 <u>10.810</u> 10.701	4.911 4.812 <u>4.715</u> 4.705	$\begin{array}{c} 9.065 \\ 8.851_{\downarrow 0.214} \\ \underline{8.289}_{\downarrow 0.776} \\ \textbf{8.079}_{\downarrow 0.986} \end{array}$

Table 1: Evaluation of perplexity on SlimPajama's test set. The best score is highlighted in bold, and the second best is highlighted with an underline. L is the maximum length of the sequence for pre-training. Subscript \downarrow presents the PPL improvement over the *baseline* method MIXChunk.

3 Language Model Pre-Training

3.1 Settings

237

240

241

242

243

244

245

246

247

248

255

257

258

261

262

263

264

267

269

270

271

273

Pre-Training Corpora In this work, we use SlimPajama (Soboleva et al., 2023) as the pre-training corpus, which consists of seven sub-corpora, including CommonCrawl, C4, Wikipedia, GitHub, StackExchange, ArXiv, and Book. This allows us to investigate packing strategies in a mixed corpora setting. We sample documents with 150B tokens from SlimPajama as the pre-training corpus and ensure each subset maintains the same proportion of tokens as in the original dataset.

Pre-Training Models The model implementation is based on the LLaMA (Touvron et al., 2023) architecture with minor modifications to support intra-document causal masking. We pre-train 1.3B parameters models using context windows of 2,048 (referred to as 2K) and 8,192 (8K) tokens. We use the same set of documents with the difference in pre-training sequence composition to pretrain models, including causal masking models, i.e., MIXChunk, UNIChunk, and BM25Chunk, and intra-document causal masking models INTRADoc. More details are available in Appendix B.

Previous works (Brown et al., 2020; Pagliardini et al., 2023) argued that dynamic sequence-specific sparse masking reduces training efficiency. Compared to causal masking, we observe a 4.0% efficiency degradation on intra-document causal masking in our implementation, and the discussion on implementation is presented in Appendix A.

3.2 Results

For evaluating LLMs trained under different packing strategies, in this work, we compute the perplexity (PPL) of a held-out set of documents where each document is treated independently. The results are summarised in Table 1. We can see that BM25Chunk achieves the lowest PPL among the three causal masking models, yielding a lower average PPL compared to MIXChunk in the 2K (-0.62) and 8K (-0.78) settings. Furthermore, UNIChunk also yields a lower average PPL than the baseline MIXChunk (-0.34 and -0.21). These results indicate that increasing the relatedness of documents in a sequence can improve the language modelling ability of models. 274

275

276

277

278

279

280

281

282

283

284

287

290

291

293

294

295

299

300

301

302

303

304

305

306

307

308

309

310

When considering models trained via intradocument causal masking, we can see that IN-TRADoc achieves the lowest PPL compared to all models trained via causal masking. This indicates eliminating the potential distracting information from irrelevant documents during pretraining benefits the language modelling ability of models. Specifically, we observe that both BM25Chunk and INTRADoc obtain significantly lower PPLs on GitHub, where INTRADoc improves over UNIChunk in both the 2K (-1.3 PPL) and 8K (-2.0) models. For UNIChunk, though we avoided packing web text and code, its improvement over MIXChunk on GitHub is slight. This phenomenon could imply that *code pre-training* is more adversely affected by the distraction of unrelated context, and both intra-document causal masking and retrieval-based sequence construction strategy can alleviate this issue.

4 Experiments on Downstream Tasks

In the following, we evaluate the in-context learning, knowledge memorisation, and context utilisation abilities of the models.

4.1 In-Context Learning

Following Shi et al. (2023), we evaluate in-context learning abilities of the models using seven text classification datasets, namely SST2 (Socher et al., 2013), Amazon (Zhang et al., 2015), Yelp (Zhang

L	Model	SST2	Amazon	DBpedia	AGNews	Yelp	Hate	Offensive	Avg.
2K	MIXChunk UNIChunk BM25Chunk INTRADoc	$\begin{array}{c} 71.53_{\pm13.8} \\ \underline{77.61}_{\pm10.05} \\ 83.73_{\pm8.17} \\ 73.65_{\pm13.61} \end{array}$	$\begin{array}{c} 81.57_{\pm 15.7} \\ \underline{90.88}_{\pm 1.13} \\ \textbf{90.90}_{\pm 3.20} \\ 84.06_{\pm 12.68} \end{array}$	$\begin{array}{c} 40.87_{\pm 3.34} \\ 36.61_{\pm 2.15} \\ \textbf{50.16}_{\pm 2.61} \\ \underline{46.82}_{\pm 1.82} \end{array}$	$\begin{array}{c} \underline{74.98}_{\pm 0.99} \\ \overline{70.39}_{\pm 2.23} \\ 75.98_{\pm 2.73} \\ \overline{72.32}_{\pm 2.66} \end{array}$	$\begin{array}{c} 86.89_{\pm 4.81} \\ 91.16_{\pm 0.35} \\ \underline{91.67}_{\pm 3.68} \\ 91.91_{\pm 0.97} \end{array}$	$\begin{array}{c} 47.10_{\pm 7.51} \\ 46.20_{\pm 5.67} \\ \underline{48.58}_{\pm 5.26} \\ \textbf{55.72}_{\pm 3.47} \end{array}$	$\begin{array}{c} 41.82_{\pm 20.46} \\ 42.30_{\pm 14.92} \\ \underline{55.36}_{\pm 15.10} \\ \textbf{69.14}_{\pm 5.37} \end{array}$	63.54 65.02 70.91 <u>70.52</u>
8K	MIXChunk UNIChunk BM25Chunk INTRADoc	$\begin{array}{c} 76.01_{\pm 8.14} \\ \textbf{81.61}_{\pm 8.63} \\ \underline{80.87}_{\pm 6.16} \\ 72.38_{\pm 3.97} \end{array}$	$\begin{array}{c} 87.32 \pm 3.08 \\ 88.30 \pm 2.68 \\ \underline{91.39} \pm 1.30 \\ \textbf{93.25} \pm 0.91 \end{array}$	$\begin{array}{c} 45.94_{\pm 3.70} \\ 52.84_{\pm 2.36} \\ \underline{56.57}_{\pm 2.33} \\ \textbf{61.85}_{\pm 6.89} \end{array}$	$\begin{array}{c} 68.21_{\pm 6.21} \\ 63.16_{\pm 9.25} \\ \textbf{74.79}_{\pm 2.89} \\ \underline{72.49}_{\pm 4.72} \end{array}$	$\begin{array}{c} 79.06 {\scriptstyle \pm 9.99} \\ 83.45 {\scriptstyle \pm 6.41} \\ \underline{85.19} {\scriptstyle \pm 6.93} \\ \textbf{92.83} {\scriptstyle \pm 1.38} \end{array}$	$\begin{array}{c} 42.85 {\scriptstyle \pm 1.19} \\ 45.50 {\scriptstyle \pm 3.00} \\ \textbf{49.12} {\scriptstyle \pm 5.17} \\ \underline{46.20} {\scriptstyle \pm 3.26} \end{array}$	$\begin{array}{c} 37.03 \pm 14.28 \\ 46.84 \pm 16.78 \\ \underline{48.33} \pm 15.88 \\ \textbf{59.59} \pm 9.88 \end{array}$	$\begin{array}{c c} 62.43 \\ 65.96 \\ \underline{69.47} \\ 71.23 \end{array}$

Table 2: In-context learning performance evaluated by text classification accuracy across seven datasets. Accuracy and deviation (subscript) are calculated using different sets of demonstrations sampled by 16 random seeds.

L

Model



 $14.47_{\pm 0.75}$ MixChunk $6.19_{\pm 0.24}$ 10.33 UNIChunk 6.70 ± 0.26 15.53 ± 0.74 11.122K $\underline{15.57}_{\pm 0.65}$ BM25Chunk 11.34 $\frac{7.10}{1.00}$ $\overline{7.17}_{\pm 0.33}^{\pm 0}$ INTRADoc $16.04_{\pm 0.35}$ 11.60 MixChunk $5.08_{\pm 0.14}$ $10.90_{\pm 1.34}$ 7.99 $10.59_{\pm 1.10}$ **UNIChunk** $5.25_{\pm 0.37}$ 7.928K 11.09 ± 0.67 BM25Chunk $5.37_{\pm 0.43}$ 8.23 15.09±0.79 **6.89**±0.08 10.99 INTRADoc

NQ

Avg.

337

338

339

340

341

342

343

344

345

346

347

349

350

351

352

353

354

355

357

358

359

360

361

362

363

TQA

Figure 2: Average in-context learning accuracy using different numbers of few-shot demonstrations – the left and right figures show the results of 2K and 8K models.

et al., 2015), DBpedia (Lehmann et al., 2015), AG-News (Zhang et al., 2015), and TweetEval hate/offensive tweet detection tasks (Barbieri et al., 2020).

311

312

313

316

317 318

319

321

324

328

332

333

335

336

In Table 2, we report the in-context learning accuracy values of the models in few-shots learning settings, using 20 and 48 demonstrations for 2K and 8K models, respectively. We truncate the input sequences to fit within their respective context windows. For models pre-trained using causal masking, we can see that UNIChunk produces more accurate results than MIXChunk, while BM25Chunk yields a higher average accuracy than MIXChunk for 2K (+11.6%) and 8K (+11.3%) models. These results indicate that *increasing relatedness of the documents in pre-training chunks can improve the in-context learning abilities of the models*.

In Figure 2, we present the average accuracy using different numbers of few-shot demonstrations. We observe that BM25Chunk has an on-par accuracy with INTRADoc on the 2K setting; however, INTRADoc obtains a significantly higher accuracy compared to BM25Chunk on the 8K setting. It may imply that using a longer context window size can result in increased distractions for causal masking pre-training; meanwhile, constrained by the performance of the retrieval method, BM25Chunk

Table 3: Exact Match scores on closed-book closed-book QA tasks.

decreases the accuracy on the 8K setting. For 8K models, MIXChunk and UNIChunk obtain similar results to their corresponding 2K models, and they do not improve the accuracy when increasing the number of demonstrations. It might be due to the similar levels of distraction in both 2K and 8K settings using random packing strategies.

4.2 Knowledge Memorisation

We use two open-domain question-answering (ODQA) datasets, namely NaturalQuestions (NQ, Kwiatkowski et al., 2019) and TriviaQA (TQA, Joshi et al., 2017), to evaluate the knowledge memorisation properties of the models. We use 12 demonstrations for the 2K models and 48 demonstrations for the 8K models. In Table 3, we show the mean Exact Match (EM) scores calculated based on 5 different sets of demonstrations.

For models trained with causal masking, we can see that *increasing the relatedness of documents in pre-training chunks can improve the knowledge memorisation ability of models*. Compared to the baseline MIXChunk, BM25Chunk obtains +9.8% and +3.0% EM improvements on 2K and 8K models, respectively. We also note that intra-document causal masking significantly improves the knowledge memorisation ability, especially for 8K models, where INTRADoc improves EM by +12.3%

L	Model	RACE-h	RACE-m	SQuAD	HotpotQA	NQ-open	TQA-open	Avg.
2K	MIXChunk UNIChunk BM25Chunk INTRADoc	$\begin{array}{c} 32.34_{\pm 0.43}\\ \underline{34.01}_{\pm 0.52}\\ 33.17_{\pm 0.36}\\ 34.49_{\pm 0.56}\end{array}$	$\begin{array}{c} 42.77_{\pm 0.69} \\ 43.52_{\pm 0.44} \\ \underline{44.92}_{\pm 0.46} \\ 44.96_{\pm 0.59} \end{array}$	$\begin{array}{c} 36.70_{\pm1.79} \\ 37.33_{\pm2.31} \\ \underline{37.91}_{\pm1.84} \\ 39.91_{\pm1.48} \end{array}$	$\begin{array}{c} 7.32_{\pm 1.31} \\ 7.12_{\pm 1.35} \\ \textbf{10.30}_{\pm 0.42} \\ \underline{8.29}_{\pm 1.27} \end{array}$	$\begin{array}{c} 20.00_{\pm 0.46} \\ 21.16_{\pm 0.96} \\ \textbf{22.10}_{\pm 0.91} \\ \underline{21.66}_{\pm 0.85} \end{array}$	$\begin{array}{c} 42.72_{\pm 1.37} \\ 42.32_{\pm 1.10} \\ \textbf{46.24}_{\pm 0.63} \\ \underline{45.67}_{\pm 1.02} \end{array}$	30.31 30.91 <u>32.42</u> 32.49
8K	MIXChunk UNIChunk BM25Chunk INTRADoc	$\begin{array}{c} 31.66_{\pm 0.47} \\ 31.68_{\pm 0.94} \\ \underline{32.63}_{\pm 0.68} \\ 33.17_{\pm 0.37} \end{array}$	$\begin{array}{c} 41.57_{\pm 0.44} \\ 41.64_{\pm 0.55} \\ \underline{44.14}_{\pm 0.48} \\ \textbf{45.56}_{\pm 0.38} \end{array}$	$\begin{array}{c} 32.79_{\pm 1.56} \\ 34.94_{\pm 1.84} \\ \underline{39.45}_{\pm 1.05} \\ \textbf{41.32}_{\pm 2.28} \end{array}$	$\begin{array}{c} 10.53_{\pm 0.70} \\ 10.57_{\pm 1.13} \\ \textbf{14.46}_{\pm 0.93} \\ \underline{12.60}_{\pm 1.49} \end{array}$	$\begin{array}{c} 20.53_{\pm 0.58} \\ 21.76_{\pm 0.80} \\ \underline{22.17}_{\pm 1.02} \\ \textbf{22.25}_{\pm 0.13} \end{array}$	$\begin{array}{c} 40.53 {\scriptstyle \pm 1.03} \\ 39.60 {\scriptstyle \pm 1.77} \\ \underline{43.40} {\scriptstyle \pm 0.38} \\ \textbf{44.19} {\scriptstyle \pm 0.60} \end{array}$	29.60 30.03 34.54 <u>33.18</u>

Table 4: Evaluation results of machine reading comprehension and retrieval-augmented generation tasks.

and +37.5% over MIXChunk for 2K and 8K models, respectively. These results support our hypothesis that reducing the distractions deriving from concatenating multiple, potentially unrelated documents in pre-training chunks can improve the knowledge memorisation ability of the models.

364

365

366

367

370

373

375

376

386

390

391

394

399

400

401

402

4.3 Reading Comprehension and Retrieval-Augmented Generation

We evaluate the pre-trained models on a set of reading comprehension tasks, namely RACE (Lai et al., 2017), SQuAD (Rajpurkar et al., 2016), HotpotQA (Yang et al., 2018), and the following retrieval-augmented generation (RAG) tasks: NQ, TQA, and Multi-Document Question-Answering (MDQA, Liu et al., 2023a). For NQ and TQA, we use the top two passages retrieved by the dense retriever (Karpukhin et al., 2020; Izacard and Grave, 2021), denoted as NQ-open and TQA-open. Our results for RACE, SQuAD, and RAG tasks are summarised in Table 4, while the results on MDQA are available in Figure 3.

We can see that BM25Chunk produces more accurate results than MIXChunk and UNIChunk in all tasks and obtains the best average accuracy, showing that *increasing the relatedness of documents in pre-training chunks can improve the context utilisation ability*. Specifically, BM25Chunk obtains a significantly better accuracy on multi-hop QA task HotpotQA, showing it can better utilise multiple relevant information from the context.

INTRADoc obtains the best average accuracy in the 2K models and obtains the best scores in 5 of 6 tasks in the 8K models. It indicates that eliminating potential distractions from unrelated documents and *learning each document independently can improve context utilisation ability*. This finding is different from the ideas in previous works, which suggested that pre-training with multiple documents in one context (Shi et al., 2023) and adding distraction



Figure 3: Accuracy on Multi-Document Question-Answering (MDQA). The x-axis represents the position of the document that contains the answer. The y-axis presents the accuracy for a position.

in context during pre-training (Tworkowski et al., 2023) benefit context utilisation ability.

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

In MDQA, for each question, there are 30 documents provided in the context, where only one of them contains the answer to the question - MDQA is used to evaluate the ability of models to filter out irrelevant information and identify the relevant parts of a long context. This task has been used to analyse the *lost-in-the-middle* phenomenon in LLMs where they struggle to retrieve information stored in the middle of long contexts (Liu et al., 2023a). In the following, we analyse how the accuracy of models varies with the position of relevant information in the context. In these experiments, we focus on 8K models due to their ability to handle long contexts. The zero-shot results on MDQA are outlined in Figure 3. We observe that both BM25Chunk and INTRADoc tend to produce more accurate predictions than MIXChunk and UNIChunk when the relevant passage is located at the beginning or middle of the context. These results show that BM25Chunk and INTRADoc can better filter irrelevant context and locate relevant information; these results are further corroborated



Figure 4: Distracted attention proportions of models. The x-axis presents the token position of the second document; the y-axis presents the distraction proportion calculated by Equation (2). Figures (a) and (b) show the distraction proportion of the first and last layers. Figures (c) and (d) are the average distraction proportion over layers. In Figure (d), we separate documents by a newline token ("\n") and present the distraction proportion of INTRADoc. The results are averaged from 4096 examples. More analysis is presented in Appendix E.

by our experiments in Section 5.1 where we analyse the attention distribution of the models during the language modelling process.

5 Discussion and Analysis

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

5.1 Can Models Ignore Irrelevant Contexts Before the End-of-Sequence Token?

In the following, we analyse whether models can filter irrelevant context during language modelling by examining the attention score distributions over the context. Specifically, we concatenate two randomly sampled documents from the SlimPajama validation set, separate them by an end-of-sequence token [EOS], and check to which extent the attention distributions of the model focus on the irrelevant document in the sequence. More formally, we define the *distraction proportion* of the token in position p in the current document at layer l as:

$$\text{DISTRPROP}(l,p) = \sum_{i=1}^{|C_d|} a_{p,i}^l$$
(2)

where $|C_d|$ denotes the number of tokens in the irrelevant document, $a_{p,i}^l$ is the average multi-head attention scores to the *i*-th token in the irrelevant document C_d at layer *l*, and $\sum_{i=1}^{|C_d|+p} a_{p,i}^l = 1$. In our experiments, we set $|C_d| = 256$, and the results are outlined in Figure 4.

We can see that the latter positions have lower distraction proportions but remain 45%-52% average distraction proportion until the 256th token of the second document, as shown in Figure 4(c). We find that models trained via BM25Chunk (green line) tend to have lower distraction proportions than other causal masking models, showing that they can better recognise relevant information in the context, matching the results in Figure 3. The above analysis also demonstrates that during the pre-training, causal masking models can be distracted by unrelated documents in context, and the models can be more robust to irrelevant contexts when reducing distractions in pre-training sequences. 459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

Furthermore, in Figure 4(d), we compare IN-TRADoc and causal masking models using "\n" as the separator instead of [EOS], because [EOS] can only appear at the end of sequences during pre-training using intra-document causal masking. The results indicate that INTRADoc has the lowest distraction proportion compared to causal masking models; meanwhile, BM25Chunk consistently has a lower distraction proportion than MIXChunk and UNIChunk using "\n" as the separator. These results match the finding in Section 4.3, where IN-TRADoc and BM25Chunk can better recognise relevant information in the context.

5.2 Burstiness Property of Sequences

Chan et al. (2022); Han et al. (2023) found a positive correlation between the model's in-context learning ability and *burstiness* property of the training sequences. Here, burstiness refers to the phenomenon where certain types of tokens occur in clusters or bursts rather than being uniformly distributed across all documents. Burstiness is an inherent property of text; for example, a specific medical term might be frequently used in medical articles and rarely appear in general texts. Higher burstiness results in a lower Zipf's coefficient of token frequency *within a sequence* (Han et al., 2023).

Following Han et al. (2023), we use Zipf's coefficient to measure the burstiness property of pretraining sequences. Formally, let r denote the rank

L	Method	Zipf's Coeffeicient (α)	In-Context Learning (Acc.)	Knowledge Memorisation (EM)
2K	MIXChunk UNIChunk BM25Chunk	2.122 2.119 2.107	$63.54 \\ 65.02 \\ 70.91$	$10.33 \\ 11.12 \\ 11.34$
8K	MIXChunk UNIChunk BM25Chunk	$1.976 \\ 1.951 \\ 1.925$	$62.43 \\ 65.96 \\ 69.47$	7.99 7.92 8.23
2K 8K	INTRADoc INTRADoc	$2.119 \\ 1.952$	70.52 71.23	$11.60 \\ 10.99$

Table 5: Zipf's coefficients of token frequency in different data. In-context learning and knowledge memorisation abilities are evaluated in Section 4.

of a token in a sequence, and f is a frequency function that maps the rank r to the frequency of that token in the sequence. Then, according to Zipf's law, we have that $f(r; \alpha) \propto \frac{1}{r^{\alpha}}$, where $\alpha \in \mathbb{R}^+$ is the Zipf's coefficient; a lower α presents an increased burstiness property within the sequence.

In Table 5, we show the Zipf's coefficients α on different pre-training sequences. Our results show that, for causal masking approaches that use the same chunk size, a lower Zipf's coefficient, which denotes increased burstiness property, often results in more accurate results. However, INTRADoc can obtain significantly better results than UNIChunk with the same Zipf's coefficient. The above results indicate that, for causal masking approaches, *the correlation between higher burstiness and better performance could derive from reduced distractions in pre-training chunks*. We report additional evidence for the burstiness property in Appendix D.

Note that duplication in pre-training sequences can also result in increased burstiness property, which may negatively impact the performance of language models. We analyse the distinct n-gram phrases of pre-training sequences in Appendix D and will investigate the impact of duplication using different pre-training corpora in future work.

6 Related Works

Instance-Level Pre-training Data Composition GPT-3 (Brown et al., 2020) was pre-trained by packed documents with causal masking, with the idea that not adopting any dynamic masking can improve pre-training efficiency. Current opensource pre-training frameworks, such as MegatronLM (Shoeybi et al., 2019), FAIRSEQ (Ott et al., 2019), EasyLM (Geng, 2023), LLM360 (Liu et al., 2023b), also follow this strategy for pre-training. Levine et al. (2022) pairs similar sentences within the same sequence and Gu et al. (2023) packs documents that contain similar intrinsic tasks for continual pre-training, improving the in-context learning ability of models. Recently, Shi et al. (2023) emphasises that packing relevant documents can enhance language models' in-context learning and context utilisation; however, our findings indicate that packing documents can adversely affect performance, and learning each document independently using intra-document causal masking can reduce the distraction and improve the performance.

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

Distribution Properties of Pre-Training Data Chan et al. (2022) shows several data distribution properties can drive in-context learning ability, e.g., large numbers of long-tail classes, dynamic meanings of inputs, and Zipf's distribution of class frequency. Han et al. (2023) used a gradient-guided method to select small-scale data for continual pretraining, showing data exhibiting burstiness properties can enhance in-context learning performance.

Pre-training Data Quality Gunasekar et al. (2023) selected high-quality data to pre-train a small-size coding model, achieving comparable performance with larger models. Shin et al. (2022); Gao et al. (2021) emphasised the importance of pre-training data diversity. Lee et al. (2022); Tirumala et al. (2023); Soboleva et al. (2023); Abbas et al. (2023) showed the importance of data deduplication on models' generalisation. In our work, we use diverse and less duplicated pre-training dataset SlimPajama (Soboleva et al., 2023), highlighting the importance of pre-training sequence composition for language models' performance.

7 Conclusion

In this work, we investigate the impact of pretraining sequence compositions by pre-training models from scratch. First, we find causal masking can result in unrelated documents distracting language modelling pre-training and hurting the performance on downstream tasks; we show that intra-document causal masking can significantly improve the performance while decreasing the pretraining efficiency. Second, we find improving the relevance of documents in pre-training chunks for causal masking pre-training can reduce some potential distractions in chunks; our proposed efficient retrieval-based packing method BM25Chunk can improve the performance of language models significantly without reducing pre-training efficiency.

522

524

526

530

494

Limitations

580

599

603

608

610

611

612

613

614

615

616

617

618

619

621

622

625

626

627

629

- 581Efficiency of Intra-Document Causal Masking582We show that intra-document causal masking is583an effective method to improve the performance584while decreasing the pre-training efficiency. We use585FlashAttention2 (Dao, 2023) to implement intra-586document causal masking masking without sacrific-587ing too much efficiency (discussed in Appendix A).588Still, we do not propose a method to solve this589efficiency issue completely.
- 590 **Objective of Sequences Construction.** We dis-591 cuss sequence construction methods, showing the 592 importance of sequence compositions on the per-593 formance of models, but these methods lack an 594 objective during sequence construction. Since spe-595 cific data distribution properties may be related to 596 models' performance, we will explore using indica-597 tors of distributional properties to guide sequence 598 construction in future works.

Scaling The Size of Language Models. Limited by the computation resources, we cannot conduct experiments on larger models with more pretraining steps, and different results might be drawn when increasing the models at a specific scale. However, this work could be directly valuable for investigating pre-training relatively small models that aim at facilitating the use of language models under resource-constrained conditions.

References

- Amro Abbas, Kushal Tirumala, Daniel Simig, Surya Ganguli, and Ari S. Morcos. 2023. Semdedup: Dataefficient learning at web-scale through semantic deduplication. *CoRR*, abs/2303.09540.
- Francesco Barbieri, José Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. Tweeteval: Unified benchmark and comparative evaluation for tweet classification. In *Findings of the Association* for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020, volume EMNLP 2020 of *Findings of ACL*, pages 1644–1650. Association for Computational Linguistics.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. Pythia: A suite for analyzing large language models across training and scaling. *CoRR*, abs/2304.01373.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind

Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual. 630

631

632

633

634

635

636

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

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

- Stephanie Chan, Adam Santoro, Andrew K. Lampinen, Jane Wang, Aaditya Singh, Pierre H. Richemond, James L. McClelland, and Felix Hill. 2022. Data distributional properties drive emergent in-context learning in transformers. In *NeurIPS*.
- Tri Dao. 2023. Flashattention-2: Faster attention with better parallelism and work partitioning. *CoRR*, abs/2307.08691.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2021. The pile: An 800gb dataset of diverse text for language modeling. *CoRR*, abs/2101.00027.
- Xinyang Geng. 2023. Easylm: A simple and scalable training framework for large language models.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023. Pre-training to learn in context. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 4849– 4870. Association for Computational Linguistics.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. 2023. Textbooks are all you need. *CoRR*, abs/2306.11644.
- Xiaochuang Han, Daniel Simig, Todor Mihaylov, Yulia Tsvetkov, Asli Celikyilmaz, and Tianlu Wang. 2023. Understanding in-context learning via supportive pretraining data. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 12660–12673. Association for Computational Linguistics.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals,

799

and Laurent Sifre. 2022. Training compute-optimal large language models. CoRR, abs/2203.15556. Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. Trans. Mach. Learn. Res., 2022. Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021, pages 874-880. Association for Computational Linguistics. Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. arXiv preprint arXiv:2310.06825. Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. IEEE Transactions on Big Data, 7(3):535–547. Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke

689

701

705

710

711

712

713

714

715

716

717

718

719

720

721

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

- Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Jean Kaddour. 2023. The minipile challenge for dataefficient language models. *CoRR*, abs/2304.08442.
- 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.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard H. Hovy. 2017. RACE: large-scale reading comprehension dataset from examinations. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017, pages 785–794. Association for Computational Linguistics.

- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. Deduplicating training data makes language models better. In *Proceedings* of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 8424–8445. Association for Computational Linguistics.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N. Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick van Kleef, Sören Auer, and Christian Bizer. 2015. Dbpedia -A large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web*, 6(2):167–195.
- Yoav Levine, Noam Wies, Daniel Jannai, Dan Navon, Yedid Hoshen, and Amnon Shashua. 2022. The inductive bias of in-context learning: Rethinking pretraining example design. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023a. Lost in the middle: How language models use long contexts.
- Zhengzhong Liu, Aurick Qiao, Willie Neiswanger, Hongyi Wang, Bowen Tan, Tianhua Tao, Junbo Li, Yuqi Wang, Suqi Sun, Omkar Pangarkar, et al. 2023b. Llm360: Towards fully transparent open-source llms. *arXiv preprint arXiv:2312.06550*.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT* 2019: Demonstrations.
- Matteo Pagliardini, Daniele Paliotta, Martin Jaggi, and François Fleuret. 2023. Faster causal attention over large sequences through sparse flash attention. *CoRR*, abs/2306.01160.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100, 000+ questions for machine comprehension of text. In *Proceedings* of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016, pages 2383–2392. The Association for Computational Linguistics.
- Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Victoria Lin, Noah A Smith, Luke Zettlemoyer, Scott Yih, and Mike Lewis. 2023. Incontext pretraining: Language modeling beyond document boundaries. arXiv preprint arXiv:2310.10638.
- Seongjin Shin, Sang-Woo Lee, Hwijeen Ahn, Sungdong Kim, HyoungSeok Kim, Boseop Kim, Kyunghyun Cho, Gichang Lee, Woo-Myoung Park, Jung-Woo Ha, and Nako Sung. 2022. On the effect of pretraining corpora on in-context learning by a largescale language model. In *Proceedings of the 2022*

- Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022, pages 5168–5186. Association for Computational Linguistics.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. Megatron-lm: Training multi-billion parameter language models using model parallelism. *CoRR*, abs/1909.08053.
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. 2023.
 SlimPajama: A 627B token cleaned and deduplicated version of RedPajama.

810

811

812

813

814

815

816

817

818

823

825

830

835

840 841

844

845

847

850

851 852

855

- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIG-DAT, a Special Interest Group of the ACL, pages 1631–1642. ACL.
- Kushal Tirumala, Daniel Simig, Armen Aghajanyan, and Ari S. Morcos. 2023. D4: improving LLM pretraining via document de-duplication and diversification. *CoRR*, abs/2308.12284.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Szymon Tworkowski, Konrad Staniszewski, Mikolaj Pacek, Yuhuai Wu, Henryk Michalewski, and Piotr Milos. 2023. Focused transformer: Contrastive training for context scaling. *CoRR*, abs/2307.03170.
- Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V. Le, Tengyu Ma, and Adams Wei Yu. 2023. Doremi: Optimizing data mixtures speeds up language model pretraining. *CoRR*, abs/2305.10429.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018, pages 2369–2380. Association for Computational Linguistics.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. Tinyllama: An open-source small language model. arXiv preprint arXiv:2401.02385.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. OPT: open pre-trained transformer language models. *CoRR*, abs/2205.01068. 856

857

858

859

860

861

862

864

865

866

867

868

869

Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, pages 649–657. 870 871

872

873

874

876

882

896

900

901

902

904

A Implementation of Intra-Document Masking

We use FlashAttention2 (Dao, 2023) to implement intra-document causal masking. The pseudo-code is presented as follows:

Pseudo-code for intra-document causal masking

-
<pre># qkv_states: query, key and value # max_seqlen: max length of documents # cu_seqlens: boundaries of documents</pre>
<pre>qkv_states = qkv_project(hidden_states)</pre>
<pre>qkv_states = qkv_states.view(batch_size, seq_len, 3,</pre>
<pre>qkv_states = rotary_embed(qkv_states)</pre>
<pre>qkv_states = qkv_states.view(batch_size * seq_len, 3</pre>
<pre>attn = flash_attn_var_len_qkvpacked_func(qkv_states,</pre>
attn = attn.view(batch_size, seq_len, num_heads * head_dim) attn = output_project(attn)

In this implementation of intra-document causal masking, we first apply the rotary position embedding to the hidden states, ensuring INTRADoc uses the same position information that is used in causal masking for each document.

We observe a 4% pre-training speed decrease in our implementation compared to causal masking pre-training, testing on 128 80G A100 GPUs. Another choice to implement intra-document causal masking is using a binary attention bias matrix for masking tokens that belong to other documents. Compared to causal masking using FlashAttention2, we observe that it reduces efficiency by 32% in xFormers, which uses Triton to implement FlashAttention with the support for arbitrary masking matrix; besides, it reduces efficiency by 52% using the standard PyTorch implementation.

B Pre-Training Details

Hyperparameters In our experiments, we use the 1.3B model, which has 24 layers, a hidden size of 2048, and 16 attention heads. We use a batch size of 4 million tokens for both the models with 2K and 8K context window sizes and pre-train models using 150B tokens with 38400 steps, which costs 40 hours to pre-training a causal masking model using 128 80G A100 GPUs. We use Adam optimiser with $\beta_1 = 0.90$, $\beta_2 = 0.95$, a weight decay of 0.1, and a cosine learning rate scheduler. The peak learning rate is 3×10^{-4} , decreasing to 3×10^{-5} at the end.

Subset	# documents	Token proportion
CommonCrawl	42960927	52.2%
C4	76520211	26.7%
GitHub	5233374	5.2%
Books	47848	4.2%
ArXiv	383058	4.6%
Wikipedia	7044397	3.8%
StackExchange	7265708	3.3%

Table	6:	Pre-training	corpus.
14010	· ·		e or p ao.

Pre-Training Corpus We sample documents with 150B tokens sampled from SlimPajama for pre-training. All models are pre-trained using the same set of documents. In Table 6, we present the number of documents and the token proportions for each subset.

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

C Analysis of BM25Chunk

C.1 Time Complexity Analysis

In BM25, the similarity score between a query and a document is based on sparse representations, where each query and document is represented by the terms it contains; such sparse representations are stored in *inverted indices*, which map terms to the documents that contain them, along with necessary statistics such as the term frequency and the document frequency. The time complexity of computing similarities between a query and documents in BM25 using an inverted index is $\mathcal{O}(Q \times K)$, where Q denotes the number of tokens in the query, and K represents the number of total documents.

We restrict BM25Chunk's retrieval process within a document buffer rather than entire largescale corpora to improve efficiency. The buffer caches k documents, which enables similarity calculations between a term and documents to be at most k times. Since each query is a document, it could contain a large number of tokens; we remove the stop words and randomly sample q tokens to reduce the length. Therefore, the time complexity of sequence construction in BM25Chunk is reduced to $O(q \times k)$. In Figure 5, we test the sequence construction speed using different q and k.

C.2 Implementation Details

We group documents in batches of 5000K and build indexes within each group. The BM25 indexes of pre-training corpora with 150B tokens require 244GB storage memory. For both 2K and 8K lengths BM25Chunk, the document buffer size k is 3072, and the maximum length of query q is



Figure 5: Pre-training sequence construction speeds using different buffer sizes k and maximum query lengths q. Test on 16 CPU cores.

500. The data construction speed is 50.0K tokens per second using 16 CPU cores, and speeds using different settings are presented in Figure 5.

C.3 Ablation Studies

Effectiveness of Document Buffer BM25Chunk conducts retrieval within a document buffer, which may result in retrieving less relevant documents, so we conduct experiments on different document buffer sizes to investigate its effectiveness. We conduct ablation experiments using 0.3B models with a context window of 2048, trained with 13B tokens, the compute-optimal number of tokens according to Hoffmann et al. (2022). We present the PPL improvement over UNIChunk on the validation set of SlimPajama in Table 7. The results show that retrieving from different sizes of document buffers can improve PPL, indicating the effectiveness of retrieving from a small-scale document set. BM25Chunk with a buffer size of 4096 achieves the best result while increasing the size to 8192 does not improve the PPL.

Effectiveness of Retrieval BM25Chunk con-965 ducts multi-hop retrieval to retrieve a sequence of 966 documents, which could potentially help models 967 learn long-distance relationships across documents, 968 and this benefit has been revealed by its high accuracy on HotpotQA, a multi-hop QA task, as shown 970 in Section 4.3. An alternative choice is retrieving 971 multiple documents at once to fill a pre-training 972 chunk, and we present such one-hop retrieval in Ta-974 ble 7. The result indicates that BM25Chunk with multi-hop retrieval can improve the PPL more effec-975 tively. Besides, we experiment with random sam-976 pling documents from the buffers without retrieval; the result shows the effectiveness of retrieval. 978

Model (0.3B)	Document Buffer Size	Valid. PPL
MIXChunk	-	15.474
INTRADoc	-	$12.443_{\downarrow 3.031}$
	2048	$13.657_{\downarrow 1.817}$
BM25Chunk	4096	$12.528_{\downarrow 2.946}$
	8192	$12.684_{\downarrow 2.790}$
BM25Chunk		
w/o multi-hop retrieval	4096	$13.497_{\downarrow 1.977}$
w/o retrieval	4096	$14.241_{\downarrow 1.233}$
CONTRIEVERChunk	-	$13.720_{\downarrow 1.654}$

Table 7: PPL on the validation set of SlimPajama. Subscript_{\downarrow} is the PPL improvement over MIXChunk. The label "w/o multi-hop retrieval" means retrieving multiple documents at once to construct the sequence; "w/o retrieval" represents random sampling from document buffers, which is equivalent to UNIChunk.



Figure 6: Chunk frequency. The *x*-axis indicates the frequency rank of tokens; the *y*-axis presents the number of chunks containing a specific token.

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

Dense Retrieval Method An alternative retrieval method to BM25 is dense retrieval. We use Contreiver (Izacard et al., 2022) as the dense retriever and compare it with BM25. Following Shi et al. (2023), we embed pre-training documents to dense vectors using Contriever and use FAISS (Johnson et al., 2019) to accelerate the retrieval process instead of using the document buffer. Then, we construct pre-training chunks using the same process introduced in BM25Chunk. We present the result produced by the dense retrieval method in the last line of Table 7. We observe that the improvement of the dense retrieval method is less than BM25.

D Analysis of Data Distribution Properties

Chunk FrequencyIn addition to Zipf's coefficient, we analyse the burstiness property through994995

963

964

944



Figure 7: Average distraction proportions over layers. We compare results using different corpora (Wikipedia and GitHub), distraction length ($|C_d| = 256$ and 512), and the separator [EOS] and \n). The first row, (a) (b) (c) and (d), use [EOS] as the separator; the second row, (e) (f) (g) and (h) use \n . The first and the third columns, (a) (c) (e) and (g), have an irrelevant context length $|C_d|$ of 256; and the others are 512. The first two columns, (a) (b) (e) and (f), present the results of Wikipedia, and the last two columns, (c) (d) (g) and (h), present the results of GitHub. We present the baseline $y = |C_d|/(|C_d| + x)$ whose attention scores are uniformly distributed over all preceding tokens.

the chunk frequency of tokens. Specifically, chunk frequency refers to the number of chunks where a specific token appears. Given a corpus, if a specific token appears in fewer chunks, it indicates more concentrated occurrences in chunks containing the token, demonstrating a higher burstiness property. In Figure 6, we can see that low-frequency tokens appear in fewer chunks in BM25Chunk compared to MIXChunk and INTRADoc, indicating these lowfrequency tokens are gathered through the retrievalbased construction method.

996

997

1001

1002

1004

1005

1006

1010

1011

1013

1014

1016

1017

1021

Distinct N-gram The burstiness property can correlate to the duplication in a sequence, which may negatively affect models, e.g., models may tend to copy phrases from context. We use SlimPajama, a highly deduplicated dataset, as the pre-training corpus, which alleviates the duplication issue. We use the percentage of distinct n-grams within a sequence to analyse the duplication issue, as shown in Table 10. The results show that BM25Chunk only has a slightly lower percentage of distinct n-grams compared to MIXChunk and UNIChunk.

E Analysis of Distraction Proportions in Different Settings

In Figure 7, we report the average distraction proportion (defined in Equation (2)) over layers us-

Method	Δ PPL $\%$	Δ DistProp $\%$
MIXChunk UNIChunk BM25Chunk INTRADoc	$14.6\% \\ 15.3\% \\ 13.5\% \\ -0.7\%$	3.4% 4.6% 4.6% -0.6%

Table 8: The PPL and DISTPROP changes after replacing the separator [EOS] by \n . A positive value means PPL or DISTPROP increases (performance drops).

ing different settings. Specifically, we analyse distraction proportions in different settings by varying the *1*) modalities of corpus: text and code using Wikipedia and GitHub; *2*) the separator token: [EOS] and line break token $\langle n; 3 \rangle$ the length of distraction context, $|C_d| = 256$ and 512.

Comparing different separators [EOS] and \n , (a) (e), (b) (f), (c) (g), and (d) (h), we observe that causal masking models can obtain lower distraction proportions using [EOS], indicating causal masking models can benefit from [EOS] to ignore irrelevant context during pre-training. We present the impact of changing the separator from [EOS] to \n on PPL and distraction proportion in Table 8. The results show that PPL and DISTPROP increase after the replacement for causal masking models, while INTRADoc obtains better results using \n as the separator since it does not train on sequences

1039

1022

1024

1025

1026

L	Model	CommonCrawl	C4	Wikipedia	GitHub	StackExchange	Book	ArXiv	Avg.
	MIXChunk	0.5429	0.4950	0.6238	0.7665	0.5974	0.5001	0.6406	0.5952
эv	UNIChunk	0.5468	0.4984	0.6298	0.7709	0.6011	0.5033	0.6436	0.5991
2 K	BM25Chunk	0.5496	0.5021	0.6394	0.7782	0.6041	0.5050	0.6452	0.6034
	INTRADoc	0.5507	0.5048	0.6426	0.7793	0.6050	0.5062	0.6458	0.6049
	MIXChunk	0.5402	0.4867	0.6219	0.7443	0.5820	0.5042	0.6531	0.5903
8K	UNIChunk	0.5429	0.4888	0.6235	0.7483	0.5859	0.5065	0.6564	0.5932
	BM25Chunk	0.5489	0.4952	0.6391	0.7621	0.5919	0.5108	0.6599	0.6011
	INTRADoc	0.5506	0.4988	0.6443	0.7643	0.5936	0.5119	0.6597	0.6033

Table 9: Evaluation of next token accuracy on SlimPajama's test set.

L	Method	Distinct 2-gram %	Distinct 3-gram %	Distinct 4-gram %
2K	MIXChunk UNIChunk BM25Chunk INTRADoc	$\begin{array}{c} 71.84 {\scriptstyle \pm 14.68} \\ 71.84 {\scriptstyle \pm 15.07} \\ 71.49 {\scriptstyle \pm 15.21} \\ 80.35 {\scriptstyle \pm 15.26} \end{array}$	$\begin{array}{c} 84.06 {\scriptstyle \pm 14.47} \\ 84.17 {\scriptstyle \pm 14.74} \\ 84.00 {\scriptstyle \pm 14.91} \\ 89.01 {\scriptstyle \pm 13.07} \end{array}$	$\begin{array}{c} 89.02 {\scriptstyle \pm 13.16} \\ 89.16 {\scriptstyle \pm 13.26} \\ 89.07 {\scriptstyle \pm 13.41} \\ 92.61 {\scriptstyle \pm 11.34} \end{array}$
8K	MIXChunk UNIChunk BM25Chunk INTRADoc	$\begin{array}{c} 64.81_{\pm 12.84} \\ 64.57_{\pm 14.09} \\ 63.49_{\pm 14.63} \\ 79.88_{\pm 14.86} \end{array}$	$\begin{array}{c} 80.61_{\pm 13.69} \\ 80.61_{\pm 14.92} \\ 80.06_{\pm 15.64} \\ 88.90_{\pm 12.63} \end{array}$	$\begin{array}{c} 86.76_{\pm 12.76} \\ 86.88_{\pm 13.64} \\ 86.56_{\pm 14.31} \\ 92.61_{\pm 10.96} \end{array}$

Table 10: The percentages of the distinct n-grams in different pre-training sequences.

where documents are separated by [EOS] using intra-document causal masking.

Comparing Wikipedia (a) (b) (e) (f) and GitHub (c) (d) (g) (h), MIXChunk is more distracted by the irrelevant context in code generation.

Comparing different length distraction contexts, (a) (b), (c) (d), (e) (f) and (g) (h), models are more distracted when $|C_d|$ increases, while much better than the baseline of uniform distribution $y = |C_d|/(|C_d| + x)$.

Comparing INTRADoc (red line) and causal masking models, we observe that intra-document causal masking results in significantly lower distraction proportions in all cases. This phenomenon demonstrates that using causal masking without considering the boundaries of documents negatively impacts language modelling performance, and the models can be more robust to irrelevant contexts when increasing the relatedness of documents in pre-training chunks.

F Next Token Accuracy of Pre-Trained Language Models

In addition to PPL, we report the next token accuracy of pre-trained language models in Table 9.