003 004

010 011

012

013

014

015

016

017

018

019

021

025

026 027 028

029

043

044

TIC-LM: A MULTI-YEAR BENCHMARK FOR CONTINUAL PRETRAINING OF LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) are trained on data crawled over many years from the web. We investigate how quickly LLMs become outdated as the world evolves with time and how to best update them with newer data. Specifically, we simulate a world where the latest dump of Common Crawl (CC), the most prominent public source of pre-training data, is used every month to *continually* train an LLM. We design various dynamic evaluations from the CC data, Wikipedia, StackExchange, and code documentations to measure continual learning metrics such as forgetting and forward transfer. Notably, our TIC-CC training data is more than 100× larger compared with prior continual learning benchmarks for language modeling. We discover that recent DataComp-LM (Li et al., 2024a) models trained on data before 2023 have already become outdated, incurring up to 45% larger noun-perplexity on 2024 Wikipedia articles compared to pre-2023 articles. Further, we use our setup to evaluate the effectiveness of several large-scale continual learning methods and find that replaying older data is most effective for combating forgetting: for previously seen CC dumps, it can reduce the regret on held-out loss by 60% compared to other optimizer and loss-based interventions. However, some domains evolve more quickly than others, favoring different trade-offs between mixing old and new data.

1 INTRODUCTION

Large language models (LLMs) rely on massive amounts of data, a major portion of which comes 031 from large-scale web-crawls that have been running over the past 10–20 years. Common Crawl (CC), the most well-known source of such data, has been active since 2007 and continues to release monthly 033 dumps of data. While typically many (or all) previous dumps are combined together to train LLMs 034 from scratch (Penedo et al., 2023; Li et al., 2024a), the vast costs and inherent knowledge cutoffs of LLMs raise natural questions about how they can be most effectively updated as future dumps are released. In this work, we introduce a benchmark for Time-Continual Learning of Language Models 037 (TiC-LM) and investigate how to continually train LLMs over many months and years. Taking inspiration from the recent TiC-CLIP (Garg et al., 2024) work, our goal is to find efficient alternatives to training LLMs from scratch by reusing and updating prior pre-trained models. Overall, we seek to answer the following: 040

- How quickly does a pre-trained LLM become outdated? We observe that Gemini, Gemini-2, and DCLM models are outdated on 2024 data by 34%, 28%, and 45%, respectively (Fig. 2).
 - Can continual pretraining reach the performance of training from scratch when fixing the number of tokens? We demonstrate a variety of methods that shrink the gap, but this proves to be an open and challenging problem to be studied in future using our benchmark.
- Do forgetting and forward transfer vary across domains such as Wikipedia, News, etc? Yes.
 For example, although data replay methods generally avoid forgetting, they hurt performance on domains that change rapidly (Fig. 5).
- Our contributions in TiC-LM are: (1) introducing a large-scale continual pretraining benchmark
 for language modeling, and (2) evaluating continual learning strategies. TiC-LM centers around
 TiC-CommonCrawl (TIC-CC), a massive time-stratified set of training and evaluation data created
 using 114 CC dumps spanning 2013–2024 including evaluation subsets TIC-CC-WIKI and TIC-CC-NEWS. TIC-CC contains 2.9T tokens, more than 100× larger than prior continual continual
 learning benchmark for LLMs. TiC-LM also contains domain-specific evaluations sourced from



Figure 1: **TiC-LM experiment setup.** We simulate a setup where each Common Crawl dump D_0, \dots, D_T is revealed one-at-a-time. An LLM f_0 is first pre-trained for B tokens on the initial month D_0 and then continually updated for a fixed budget of B/T tokens in each following month (which may optionally include replaying any older data). The goal is for each monthly model f_1, \dots, f_T to perform well on both standard static downstream tasks as well as dynamic evaluations that evolve over time, requiring the balancing of learning (gray/red) with preventing forgetting (blue).

Table 1: Comparison with continual learning benchmarks for LLMs. CLS: classification, SUM:
Summarization KB: Knowledge Base, QA: Question-Answering, LM: Language Modeling. Acc.:
Accuracy, Ppl.: Perplexity, Tok.: Tokens, Art.: Articles.

Benchmark	Domain	Task	Metric	CL Train	Time-CL	Years	Timesteps	# CL Train	# Eval Samples
Gururangan et al. (2020)	Science, News, Reviews	CLS	micro/macro-F1	1	X-Task CL	_	_	0.3M	140k
Luu et al. (2022)	Tweet,Science,News,Reviews	CLS/SUM	F1/Rouge-L	1	1	2013-2022	4–7	695k	695k
Chrono. Tweet(2022)	Science, Tweet	CLS	micro/macro-F1	1	1	2014-2020	4	25M	4M
TempEL (2022)	Wikipedia	KB	EL Acc.	1	1	2013-2022	10	_	92k
TemporalWiki (2022a)	Wikipedia	KB	Noun Ppl.	×	1	2021	4	23B Tok.	36k
StreamingQA (2022)	News	QA	Acc.	1	1	2007-2020	12	99k Art.	46k
EvolvingQA (2024)	Wikipedia	QA	EM/F1	1	1	2007-2020	6	_	46k
TIQ (2024)	Wikipedia	QA	Precision/Rank	1	1	1801-2025	_	6k QA	4k
TAQA (2024)	Wikipedia	QA	F1	✓	1	2000-2023	_	9k QA	11k
TIC-CC (All/Wiki/News)	Generic Web	LM	Ppl.	1	1	2013-2024	114	2.9T Tok.	2.7M
TIC-WIKI	Wikipedia	KB	Noun Ppl.	×	1	2014-2024	62	_	10M
TIC-STACKE	Code, Math, English,	KB / QA	Ppl.	×	1	2008-2024	187	_	3.65M
TIC-CODEDOCS	Code	LM	Ppl.	×	1	2017-2024	16	_	6.5k

outside Common Crawl including TiC-Wikipedia (TIC-WIKI), TiC-StackExchange (TIC-STACKE), and TIC-CODEDOCS spanning 2008–2024. Using our benchmark, we evaluate several continual learning baselines and find that cyclic learning rate schedules and data replay can be effective for balancing learning on new data and preventing forgetting. We also find that different domains evolve at different rates, favoring different methods (e.g., benefiting from more or less replay).

2 RELATED WORK

066

067

068

069

071

085

087

880

090

091 092

Learning new capabilities from multiple, sequentially observed, distributions has long been an active area of ML research (Wang et al., 2024). More recently, several works have studied this setting 094 for LLMs (Wu et al., 2024), targeting improvements on: (1) general capabilities (via updating 095 on improved datasets) (Parmar et al., 2024; Ibrahim et al., 2024; Gupta et al., 2023); (2) specific 096 domains (Jin et al., 2022; Gururangan et al., 2020; Chen et al., 2024); (3) newer data as the world evolves (Jin et al., 2022; Jang et al., 2022b;a; Lazaridou et al., 2021; Nylund et al., 2024; Loureiro 098 et al., 2022; Qin et al., 2022; Liška et al., 2022). Works in this third category have demonstrated that in many domains, the performance of LLMs decay as training and evaluation sets grow farther 100 apart in time, motivating the need for methods to *efficiently* and *non-distruptively* adapt to temporal 101 distribution shifts. Our work scales up these previous efforts to more closely match current LLM 102 training practices. While older works typically focus on continual training runs involving individual 103 sources (e.g., news, Wikipedia, and social media) and <10 timesteps, we consider training on a 104 generic web-crawl (i.e., Common Crawl) spanning 114 different months. In turn, the generality 105 of our training data allows us to go beyond single-domain evaluations. We provide an extended discussion of related works in Appx. E. Table 1 summarizes our proposed datasets compared with the 106 most related time-continual benchmarks. With 2.9T tokens, TIC-CC is the largest and most diverse 107 continual learning benchmark for language modeling.



122

123

125

126

127

128

129

130

131 132 133

134 135

Figure 2: (Left) Performance of a model trained on DCLM-Baseline, which contained data up to
2022. Notably, the loss increases significantly on TIC-CC-WIKI and TIC-CC-NEWS subsets after
2022 data cutoff. (Right) Performance of the same DCLM-Baseline model as well as two versions of
Gemma (GemmaTeam et al., 2024) on our TIC-WIKI dynamic evaluation. Performance of the DCLM
model is about 45% worse on the latest evaluation data compared to the preceeding data. For the
Gemma series, the older Gemma-7b and newer Gemma-2-9b are 34% and 28% worse, respectively.



Figure 3: We plot the total number of tokens per month in TIC-CC (left) as well as the proportion of those tokens coming from our TIC-CC-WIKI and TIC-CC-NEWS subsets (right).

3 TIC-COMMONCRAWL (TIC-CC): MORE THAN A DECADE OF WEB DATA

136 We create a large *time-stratified* dataset of 2.9T tokens based upon Common Crawl, a free and open corpus of web-crawled data that has been online since 2007. CC releases new snapshots of the web 137 roughly every month. Each dump creates a representative snapshot of the web by sampling a limited 138 number of pages from each domain. Sampled URLs change from month to month independent 139 of whether they appeared in the previous dumps. We collect all dumps between May-2013 and 140 July-2024, resulting in 114 corresponding splits that we refer to by the month of their release date. 141 For each split, we then apply a pre-processing pipeline based on that of DataComp-LM (Li et al., 142 2024a). Notably, we do not perform any operations on a particular month that depend on future 143 months to retain causality and temporal order. 144

Data processing. We build upon the existing pipeline from DataComp-LM (Li et al., 2024a), 145 starting with DCLM-Pool (Li et al., 2024a), which contains all CC dumps between May-2013 146 and Dec-2022 and parsed to extract plain text from webpages via the open-source resiliparse 147 library (Bevendorff et al., 2018; 2021). We split this data by month and reuse the same download 148 and processing scripts to extend DCLM-Pool until July-2024. Next, we follow DCLM-Baseline's 149 pipeline by applying heuristic filters from RefinedWeb (Penedo et al., 2023) and a fuzzy-deduplication 150 step which we modify to run only within each month rather than non-causal global deduplication. 151 Alternatively, similar to TiC-CLIP, one could deduplicate data globally but keep the earliest occurrence 152 of each document. We avoid this deduplication for two reasons: (1) fuzzy deduplication across 153 months may not always be helpful, potentially removing near-duplicates such as Wikipedia pages where a few key facts have changed but most of the text is the same, (2) it allows for exploring the 154 benefits/pitfalls of such data-centric interventions as part of method design. Also, we do not use the 155 final classifier-based filter in DCLM-Baseline, as this classifier was trained on data from all months. 156

Finally, we leverage the fact that DCLM-Pool was randomly partitioned into ten equally-sized chunks
to construct held-out sets for loss-based evaluations (Sec. 4.1). In Fig. 3, we show the number of
tokens we have for each month of the dataset. In total, the dataset spans 29T tokens, with individual
months ranging between 100B to 500B tokens. We use smaller subset of 220B tokens from a single
global shard with 2.9T for our training while future work can expand to the full 2.9T/29T tokens. For
more details about the data pipeline see Appx. A.

¹⁶² 4 EVALUATIONS

163 164

In this section, we will discuss various time-continual evaluations that are designed both with
 and independent of CC data. As our focus is on continual pretraining, we focus on evaluations
 without instruction-tuning. We introduce three sets of novel evaluations: TIC-CC, TIC-CC-WIKI,
 TIC-CC-NEWS, TIC-WIKI, TIC-STACKE, and TIC-CODEDOCS.

168 Static downstream evaluations. We focus on pre-trained base models without any instruction 169 fine-tuning and evaluate our models on a variety of suitable downstream zero-shot and few-shot 170 tasks. Specifically, we use the CORE evaluations from the DCLM benchmark (Li et al., 2024a) 171 which includes 22 zero-shot and few-shot in-context learning tasks. These evaluations, which include 172 benchmarks such as ARC-Easy (Clark et al., 2018) and Hellaswag (Zellers et al., 2019), assess general 173 capabilities of base models via a variety of world knowledge and natural language understanding 174 tasks. While these evaluations are not designed to be time-dependent, we use them to assess (1) whether continually trained models match the general capabilities of models trained on all dumps; (2) 175 if they benefit from different months of Common Crawl. 176

Perplexity metrics. We employ three distinct perplexity metrics for different evaluations:

184 185

186 187 188

189

190 191 $ppl_{token} = \exp\left(\frac{\sum_{d \in \mathcal{D}} \sum_{t \in T_d} -\log P(t|c_{< t})}{\sum_{d \in \mathcal{D}} |T_d|}\right),$ (1) where \mathcal{D} is a set of documents, T_d is the set of tokens in document d, and $c_{< t}$ is the context prior to

$$ppl_{answer} = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \exp\left(-\log P(a_q|c_q)\right) , \qquad (2)$$

where Q is a set of question-answer pairs, a_q is the gold answer for question q, and c_q is the context.

$$\operatorname{ppl}_{\operatorname{noun}} = \exp\left(\frac{\sum_{d \in \mathcal{D}} \sum_{n \in N_d} -\log P(n|c_{< n})}{\sum_{d \in \mathcal{D}} |N_d|}\right),$$
(3)

where \mathcal{D} is the set of documents in a snapshot, N_d is the set of proper noun tokens (tagged as NNP or NNPS by a POS tagger) in document d, and $c_{< n}$ is the context prior to noun n.

192 4.1 TIC-COMMONCRAWL (TIC-CC) EVALUATIONS

193 194 CC data is a consistent, albeit partial, snapshot of the web over years that does not require special 195 processing of the history per website. We compute token-perplexity (ppl_{token}) on three monthly 196 subsets of our CC data which were held out from training:

- TIC-CC: Held-out documents coming from the full distribution for each month of TIC-CC.
- TIC-CC-WIKI: Pages in TIC-CC from English Wikipedia (i.e., whose URLs contain either the domain en.wikipedia.org or simple.wikipedia.org).
- TIC-CC-NEWS: Pages in TIC-CC from a set of news sites based on WMT competitions (Barrault et al., 2020).
- 201 202 203

197

199

200

4.2 TIC-WIKIPEDIA (TIC-WIKI)

204 Our TIC-CC-WIKI evaluation in Sec. 4.1 is based on sampled Wikipedia pages existing in each 205 CC dump which is a representative set of Wikipedia but not all of it. We create TIC-WIKI, a more 206 comprehensive evaluation from full dumps of Wikipedia while utilizing the knowledge graph from 207 Wikidata. TIC-WIKI allows us to construct question/answer and factual evaluations as well as split the performance over changed/unchanged knowledge. We build upon TemporalWiki (Jang et al., 208 2022a), which generates evaluations from four consecutive monthly snapshots of English Wikipedia 209 and Wikidata. Our TIC-WIKI evaluation captures a broader spectrum of knowledge evolution, we 210 extend the evaluation timespan to a full decade (2014-2024) and improve upon the matching of 211 Wikipedia/Wikidata (see Appx. B.1 for more details). 212

- To evaluate performance on TIC-WIKI diffsets, we adopt the approach of Lazaridou et al. (2021) and Jang et al. (2022a), calculating the average perplexity of proper nouns (ppl_{noun}) identified by a Part-of-Speech tagger (Honnibal & Montani, 2017). This method serves as a proxy for assessing
- factual knowledge changes, as proper nouns often contain key factual information.

2164.3TIC-STACKEXCHANGE (TIC-STACKE)217

We design another question/answering evaluation based on the historical data from StackExchange.StackExchange has 182 communities that share knowledge by posting questions and answers. We measure answer-perplexity (ppl_{answer}) on high-quality answers from selected sites by collecting answers that have been accepted by the question author (using the accepted answer timestamp to bin examples by month). The resulting evaluation contains examples from 2008–2024. We provide details of TIC-STACKE data processing in Appx. B.2.

- 224 225
- 4.4 TIC-CODE DOCUMENTATIONS (TIC-CODEDOCS)

Our TIC-CODEDOCS evaluation is based on code documentations from popular open-source Python libraries: NumPy (Harris et al., 2020) and PyTorch (Ansel et al., 2024). For NumPy, we use documentations from 16 major releases, ranging from version 1.13.0 (June 2017) to version 2.1.0 (August 2024). For PyTorch, we use documentations of major releases ranging from version 1.8.0 (March 2021) to version 2.4.0 (July 2024).

231 We build the documentation directly from each library's git repository. The process involves the 232 following steps: (1) Identify the commit tagged for the major release, (2) revert to that specific 233 commit, (3) install necessary dependencies, (4) build the project from source, (5) generate HTML 234 documentation from the source. We then convert all HTML pages to raw text by extracting the 235 main body of the documentation pages. This approach ensures that template-related elements of the 236 pages such as the index and footer are not included in the final text, focusing solely on the relevant documentation content. We evaluate the model's code understanding using perplexity (ppl_{token}), 237 calculated across entire snapshots of code docs. 238

239 240

241

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

5 CONTINUAL LEARNING BASELINE METHODS

The goal for TiC-LM methods is to match the performance of the *Oracle* which trains on all data (114 months) starting from random initialization for the full token budget. We consider methods from three categories: optimization-based, data replay, and regularization. Aside from average metrics across all timesteps, methods should balance forgetting and forward transfer metrics (defined in Sec. 6).

Optimization-based methods. In non-continual settings, LLMs are often trained with a cosinedecayed learning rate schedule which requires knowledge of total training steps ahead of time. In a continual setup, however, the number of total tokens grows over time and we care about the performance after each month. We benchmark the following optimization approaches in our work:

- *Cyclic Cosine* is the simplest alternative which applies cosine decay *within* each training month using the same maximum learning rate and warmup for each round. This was found to be most effective in TiC-CLIP (Garg et al., 2024).
- *Cyclic Cosine* + *AR (autoregressive)* is similar to cyclic cosine decay except the maximum learning rate in each cycle decays across months, regressed from a single-cycle cosine decay and shown to offer improvements by Roth et al. (2024).
- *Rsqrt* (*reciprocal*- $\sqrt{}$) are *infinite* schedules that decay the learning rate slowly in a global training run and branch off of this trajectory with linear cooldowns (Zhai et al., 2022). To keep training steps fixed compared to other methods, we follow Roth et al. (2024) and implement a version that maintains only a single trajectory by re-warming up from the previous cooldown.
- *Schedule-Free* is an optimizer proposed by Defazio et al. (2024) which aims to circumvent the need for defining a learning rate schedule by using iterate averaging and has achieved promising results in i.i.d. non-continual settings.

Data replay methods. Data replay is a classical continual learning strategy to prevent forgetting, whereby in each training round, the model is fed a mixture of data from both older and the current timesteps (Lopez-Paz & Ranzato, 2017; Rebuffi et al., 2017; Chaudhry et al., 2018). Defining a replay method therefore boils down to *how* the mixture ratios are specified. We consider the following replay strategies based on the best-performing strategies in TiC-CLIP (Garg et al., 2024):

• For the current timestep t, we allocate a ratio $0 \le \alpha_t \le 1$ of the monthly token budget B_t to data from the current month, seeing $\alpha_t B_t$ tokens from that month.

270 • For previous months, we redistribute the remaining $(1 - \alpha_t)B_t$ tokens equally, i.e., each month 271 contributing $\frac{1-\alpha_t}{t-1}B_t$ tokens to this round's training set. 272

In particular, when $\alpha_t = 1/t$, we see an equal number of tokens from all observed months. We also 273 consider setting $\alpha_t = 1/2$ which always allocates half the token budget to the current month. The 274 general downside of replay-based methods is the cost of retaining old data. This can be particularly 275 challenging if old data expires and needs to be removed. Methods with larger values of α_t are less 276 affected by such limitations. 277

Regularization-based methods. These methods alter the training objective instead of the data, 278 generally by adding a regularization term which encourages newer model updates to stay close to the model weights learned after the previous month. Following TiC-CLIP, we try two notable methods: LwF (Li & Hoiem, 2018) and EWC (Kirkpatrick et al., 2017).

- LwF adds an additional loss term based on KL divergence which penalizes differences in model outputs between the previous checkpoint and the current model.
- *EWC* attempts to slow down updates to particular model parameters which are highly influential for performing well on older months as measured by the (approximate) Fisher information matrix.

Because both LwF and EWC involve extra loss terms and model copies, it is important to note that they induce larger GPU memory footprints and run-times compared to optimizer and replay-based methods. That being said, we do not try to adjust the token counts to account for this given that our re-implementations may not be optimally efficient.

290 291 292

279

281

283 284

285 286

287

288

289

EXPERIMENTS 6

293 **Training details.** For all runs, we train 3B parameter language models using OpenLM. Unless otherwise indicated, each method observes a fixed total number of 220B tokens, equivalent to 4x 295 the Chinchilla optimal amount. ¹ We further assume that current practitioners are (a) likely to have 296 access to more than enough data to train initial models; (b) unlikely to wait to obtain non-trivial 297 performance. Hence, we front-load the total token budgets such that *half* is allocated to training on 298 the first month, May-2013. Then, the remaining 110B tokens are split equally among the other 113 299 continual timesteps. For more realistic hyperparameter selection in our continual setup (Cha & Cho, 2024), we only use the first 10 of these timesteps for tuning (see Appx. C for more details). 300

301 **Evaluation metrics.** Each continual run actually produces a $T_t \times T_e$ matrix of evaluations E where 302 T_t, T_e are the total number of training/evaluation timesteps, $E_{i,j}$ is the performance of the model 303 trained after training on data up to month i and evaluated on the month j. To control for inherent difficulty differences across evaluation months, we measure the regret $R_{i,j} = E_{i,j} - E_j^*$ where E_j^* is the performance of the Oracle trained on all months (May-2013–July-2024) on month j. We subtract 304 305 E_j^* instead of $E_{j,j}$ since if $E_{j,j}$ is bad it may lead to misleadingly good forward/backward metrics. 306

Following Garg et al. (2024), we consider the following summary metrics (assuming $T_t = T_e = T$).

- In-distribution (ID) performance: averages along the matrix diagonal, i.e., ∑_{i=1}^T = R_{i,i}/T.
 Backward transfer: averages the lower triangular of R, i.e., ∑_{i=1}^T ∑_{j<i} R_{i,j}/T(T-1)/2</sub>.
- Forward transfer: averages the upper triangular of R analogously to backward transfer.

313 For some downstream evaluations, the train/evaluation periods do not exactly align $(T_t \neq T_e)$, making the definition of ID more nuanced. For such evaluations, we define a_i as the index of the nearest 314 evaluation timestep that comes before the training timestep i. We then count R_{i,a_i} towards the ID 315 average only if no other training timestep is closer to a_i (i.e., $a_i \neq a_{i-1}$). Otherwise, we count $R_{i,j}$ 316 towards backward and forward transfer when $j < a_i$ and $j \ge a_i$ respectively. 317

318 319

307

308

309 310 311

312

6.1 TIC-CC HELD-OUT EVALUATIONS

Tab. 2 and Fig. 4 show results on our TIC-CC hold-out sets. Overall, we observe the various methods incur different trade-offs between forgetting and plasticity. We summarize the main findings below:

321 322 323

320

¹Here, following Li et al. (2024a), the token counts are given by $20 \times$ number of parameters \times Chinchilla multiplier with a $1 \times$ multiplier being a near-optimal compute allocation found by Hoffmann et al. (2022).



Figure 4: Cyclic Cosine demonstrates delayed forgetting while replay avoids forgetting in exchange for lower plasticity. We plot the difference in log-perplexity (log ppl_{token}) between continual checkpoints and the Oracle. While we train on all 114 months, we evaluate on a subset of the months which are roughly annually spaced. Overall, we observe that training sequentially 346 with the Cyclic Cosine method (top) leads to strong ID performance (along the diagonal) but also significant forgetting on TIC-CC and TIC-CC-NEWS. Meanwhile, adding replay can reduce this forgetting by sacrificing ID performance, especially for later model checkpoints.

Table 2: Loss-based evaluations for various methods at the 3B-4x scale. We report log-perplexity values relative to the Oracle. While various optimizer (top) and regularization-based (bottom) methods trade-off backward transfer with in-distribution performance, replay (middle) is required to obtain the least amount of forgetting. Bold values are within one standard deviation of the best in each column, with standard deviations estimated from three runs of Cyclic Cosine.

M-4h - J	TIC-CC↓			TIC	-CC-WI	KI↓	TIC-CC-NEWS ↓		
Method	Backward	ID	Forward	Backward	ID	Forward	Backward	ID	Forward
Cyclic Cosine (std)	0.072 (0.000)	0.027 (0.000)	0.161 (0.000)	0.038 (0.000)	0.032 (0.000)	0.074 (0.000)	0.058 (0.000)	0.015 (0.000)	0.109 (0.000)
Cyclic Cosine + AR	0.058	0.040	0.166	0.032	0.031	0.074	0.041	0.017	0.110
Cyclic Rsqrt	0.065	0.030	0.162	0.033	0.030	0.073	0.049	0.015	0.108
Schedule-Free	0.065	0.036	0.164	0.036	0.033	0.076	0.049	0.017	0.110
Replay ($\alpha_t = 1/t$)	0.023	0.074	0.178	0.020	0.036	0.078	0.005	0.035	0.117
Replay ($\alpha_t = 1/2$)	0.024	0.042	0.167	0.024	0.031	0.074	0.013	0.019	0.111
Replay ($\alpha_t = 1/t$) + AR	0.026	0.083	0.181	0.019	0.037	0.079	0.004	0.039	0.119
Replay ($\alpha_t = 1/2$) + AR	0.025	0.055	0.171	0.022	0.032	0.076	0.009	0.022	0.112
LwF	0.072	0.027	0.161	0.038	0.032	0.074	0.058	0.015	0.109
EWC	0.061	0.032	0.162	0.031	0.029	0.071	0.046	0.014	0.108

366

367

368

343

344

345

347

348 349

350

351

352

353

> Cyclic Cosine offers the best plasticity but also the most forgetting. We find the optimal maximum learning rate for ID performance to be $30 \times$ smaller than what was used for the initialization (see Tab. 5 in Appx. C). Further, based upon the gaps between Forward and ID metrics, TIC-CC-WIKI appears to change more slowly than TIC-CC-NEWS, while both evolve less rapidly than TIC-CC.

369 Alternative LR schedules and EWC can improve Backward at the cost of ID but replay is 370 required to further reduce forgetting. As shown in the heatmaps in Fig. 12, all non-replay methods 371 still result in significant forgetting at later checkpoints, while Replay ($\alpha_t = 1/t$) reaches a Backward 372 metric of 0.023 on TIC-CC, 60% smaller than the best non-replay approach. Between the two 373 replays, $\alpha_t = 1/2$ offers slightly less Backward improvement but much better ID, likely because 374 $\alpha_t = 1/t$ decreases the ratio of new data over time, scaling poorly to the >100 timesteps in our setup. 375 This differs from TiC-CLIP's findings, where different replays behave more similarly but $10 \times$ fewer rounds are used. However, even while $\alpha_t = 1/2$ can achieve good balance between Backward and ID 376 for the first few years, plasticity becomes an increasing issue across the larger timescales in TIC-CC, 377 likely due to later months still being underrepresented as they have been replayed fewer times.

Table 3: Selected downstream evaluations at 3B-4x scale. For all dynamic evaluations, we report
perplexity values relative to the Oracle with log-scaling. Meanwhile, CORE is an average of the
accuracies of 22 downstream zero/few-shot tasks used in DataComp-LM (Li et al., 2024a), evaluated
only on the final model checkpoint (with score relative to Oracle in parentheses). Bold values are
within one standard deviation (estimated with 3 runs of Cyclic Cosine) of the best in each column.

M-411	T1C-W1K1-Diff↓			TIC-STA	CKOVER	FLOW ↓	TIC-STACKE-CAT7↓		
Method	Backward	ID	Forward	Backward	ID	Forward	Backward	ID	Forward
Cyclic Cosine (std)	0.033 (0.000)	0.052 (0.000)	0.085 (0.000)	0.041 (0.002)	0.078 (0.002)	0.156 (0.003)	0.045 (0.001)	0.050 (0.001)	0.071 (0.000)
Cyclic Cosine + AR	0.033	0.054	0.087	0.032	0.077	0.159	0.035	0.044	0.068
Cyclic Rsqrt	0.031	0.051	0.084	0.034	0.076	0.158	0.039	0.046	0.069
Schedule-Free	0.035	0.055	0.087	0.038	0.079	0.160	0.045	0.050	0.072
Replay ($\alpha_t = 1/t$)	0.038	0.063	0.091	0.075	0.121	0.191	0.036	0.052	0.072
Replay ($\alpha_t = 1/2$)	0.032	0.055	0.086	0.055	0.094	0.170	0.038	0.049	0.070
Replay ($\alpha_t = 1/t$) + AR	0.039	0.063	0.092	0.066	0.119	0.193	0.031	0.050	0.072
Replay ($\alpha_t = 1/2$) + AR	0.033	0.057	0.088	0.047	0.096	0.176	0.032	0.046	0.071
LwF	0.033	0.053	0.085	0.037	0.075	0.155	0.044	0.048	0.070
EWC	0.030	0.051	0.083	0.033	0.077	0.162	0.035	0.043	0.067

Method	TIC-CODEDOCS-NUMPy ↓ Backward ID Forward			TIC-CODI Backward	Static Evals. ↑ Core (DCLM)		
Cyclic Cosine (std)	0.073 (0.004)	0.096 (0.003)	0.072 (0.002)	0.057 (0.002)	0.025 (0.001)	0.217 (0.002)	48.5 (-2.1)
Cyclic Cosine + AR	0.054	0.074	0.062	0.084	0.052	0.228	48.5 (-2.1)
Cyclic Rsqrt	0.066	0.092	0.071	0.062	0.029	0.220	49.0 (-1.6)
Schedule-Free	0.069	0.100	0.080	0.084	0.051	0.236	48.8 (-1.8)
Replay ($\alpha_t = 1/t$)	0.054	0.046	0.057	0.175	0.138	0.275	48.9 (-1.7)
Replay ($\alpha_t = 1/2$)	0.058	0.066	0.060	0.099	0.069	0.237	49.0 (-1.6)
Replay ($\alpha_t = 1/t$) + AR	0.040	0.045	0.057	0.197	0.169	0.277	49.0 (-1.6)
Replay ($\alpha_t = 1/2$) + AR	0.034	0.050	0.052	0.129	0.098	0.246	49.2 (-1.4)
LwF	0.076	0.104	0.073	0.058	0.028	0.214	48.5 (-2.1)
EWC	0.055	0.081	0.067	0.070	0.040	0.222	48.9 (-1.7)

6.2 DOWNSTREAM EVALUATIONS

Tab. 3 presents results for several of our downstream evaluations, while Fig. 5 and Appx. D show
 corresponding evaluation matrices. Here, we observe broadly that methods exhibit trade-offs *across the different evaluations*, showcasing the inherent challenges of performing well on a variety of
 domains while training on a sequence of generic web-data dumps (in which the coverages of said
 domains may also vary over time). We summarize the key findings below:

EWC is the best method for adapting to new knowledge in TIC-WIKI-Diff. This differs from the results of Garg et al. (2024) and Jin et al. (2022) which found EWC to have little positive impact in their settings. It also differs from TIC-CC-Wiki, where EWC could not match replay in terms of Backward performance. One explanation for this is that by isolating new/changed segments of Wikipedia, TIC-WIKI-Diff places strong pressure on seeing newer data. Indeed, we see in Tab. 8 in Appx. D that measuring performance on the *unchanged* segments returns replay's superiority on Backward, though the improvements are much smaller. To further shed light on this discrepancy, we also see from Fig. 5 (left) that unlike for TIC-CC-WIKI, peak performance on each TIC-WIKI month is often years after that month is seen, even when no replay is used in Cyclic Cosine. This suggests that the knowledge TIC-WIKI captures each month is also successfully learned from CC dumps that were crawled quite a bit later, which could be due to: (1) delayed alignment between CC's crawls of Wikipedia and TIC-WIKI's more comprehensive coverage of Wikipedia edits; (2) measuring TIC-CC loss on all tokens versus focusing on specific segments and proper nouns in TIC-WIKI to capture only factual knowledge rather than irrelevant nuances (e.g., page formatting).

Replay of earlier data tends to benefit more on domains expected to evolve slowly. In Tab. 3, we show both the performance of different methods on TIC-STACKOVERFLOW as well as an average over a subset of seven other large StackExchange sites excluding StackOverflow (TIC-STACKE-CAT7). Earlier CC dumps (i.e., before Feb-2016) are the most useful for TIC-STACKE-MATH leading to improvements from both replay and AR schedules as observed in Fig. 5. In contrast, for TIC-STACKOVERFLOW, there not only exists a larger distribution shift over time, but seeing less old data improves all three summary metrics. Similar to TIC-STACKE, we observe a sharp difference between two evaluations within TIC-CODEDOCS: TIC-CODEDOCS-NUMPY and TIC-CODEDOCS-PYTORCH. As shown in Tab. 3, continual runs involving replay perform better on all metrics for the



Figure 5: **Replay helps on TIC-STACKE-MATH but hurts on domains where new data is crucial** (TIC-STACKOVERFLOW). We show heatmaps the Cyclic Cosine (left) and Replay ($\alpha_t - 1/2$) + AR methods (right) evaluated on TIC-WIKI, TIC-STACKE-MATH, and TIC-STACKOVERFLOW. The purple dotted lines trace out when the training and evaluation timestamps are closest to one another.

former whereas the opposite is true for latter. This is likely due to NumPy being an older library first
released in 1995 compared to PyTorch in 2016. Based on the corresponding heatmaps in Appx. D,
it appears that models improve on NumPy in earlier years (i.e., 2013-2016) before forgetting this
knowledge when training on the following four to five years. This suggests the bulk of NumPy-related
content appeared in earlier years before decreasing, thereby *necessitating* replay for models to retain
this knowledge. In contrast for PyTorch, replay harms performance since it shifts more weight to
earlier CC dumps, three years of which were before PyTorch even first released.

471 Static evaluation of methods with an unbiased initialization matches the Oracle. Table 3 presents 472 evaluations on the CORE set of tasks from Li et al. (2024a). Initially, we observe that most continual 473 methods perform similarly and a sizable gap to the Oracle remains. Indeed, the initialization trained on 474 the May-2013 already achieves an average of 48.5, the same as the final checkpoint of Cyclic Cosine. Meanwhile, the interventions that mitigate forgetting can somewhat help, closing the remaining gap 475 to the Oracle by 33%. The two possible explanations for the remaining 67% are that the Oracle 476 benefits from: (1) having less restricted access to data throughout its training (i.e., the continual 477 phase of our runs is at fault); (2) being trained from scratch rather than starting from a model biased 478 towards the oldest data (i.e., the initial training on the May-2013 is at fault). To investigate, we run an 479 additional oracle variant which starts from the same May-2013 initialization but trains on an equal 480 mix of the remaining 113 months all at once. This model achieves a performance of only 48.9, below 481 our best continual runs which suggests that (2) is likely more at fault than (1).

482 483 484

6.3 EFFICIENCY OF CONTINUAL TRAINING

In this section, we perform a case study on the practical utility of current continual methods by measuring the potential compute savings offered by Replay ($\alpha_t = 1/2$). While the results in previous

Method	Tokens	TIC-CC	TIC-CC-Wiki	TIC-CC-News	TIC-WIKI-Diff	TIC-WIKI-Unchang
Replay	220B	0.024	0.019	0.017	0.013	0.017
Replay + AR	220B	0.027	0.017	0.014	0.014	0.016
Replay	330B	0.011	0.009	0.010	0.001	0.006
Replay + AR	330B	0.008	0.003	0.002	-0.001	-0.001
Replay	440B	0.004	0.004	0.007	-0.004	0.000
Demlerr + AD	440P	-0.002	-0.006	-0.005	-0.009	-0.010
Keplay + AK	440B	-0.002	-0.000	0.000	01003	01010
Method	Tokens	-0.002	ACKOVERFLOW	TIC-STACKE-Cat7	TiC-CD-PyTore	h TiC-CD-NumP
Method Replay	Tokens 220B	TIC-ST	ACKOVERFLOW	TIC-STACKE-Cat7	TiC-CD-PyTore	h TiC-CD-NumP
Method Replay Replay + AR	Tokens 220B 220B	TIC-STA	ACKOVERFLOW 0.034 0.027	ТIC-STACKE-Cat7 0.028 0.022	TiC-CD-PyTorc 0.052 0.082	h TiC-CD-NumP 0.047 0.023
Method Replay Replay + AR Replay	Tokens 220B 220B 330B	TIC-ST	ACKOVERFLOW 0.034 0.027 0.014	ТІС-STACKE-Cat7 0.028 0.022 0.020	TiC-CD-PyTorc 0.052 0.082 0.003	h TiC-CD-NumP 0.047 0.023 0.035
Method Replay Replay + AR Replay Replay + AR	Tokens 220B 220B 330B 330B		ACKOVERFLOW 0.034 0.027 0.014 0.017	ТІС-STACKE-Cat7 0.028 0.022 0.020 0.014	TiC-CD-PyTorc 0.052 0.082 0.003 0.040	h TiC-CD-NumP 0.047 0.023 0.035 0.024
Method Replay Replay + AR Replay + AR	Tokens 220B 220B 330B 330B 440B		ACKOVERFLOW 0.034 0.027 0.014 0.017 0.007	ТIC-STACKE-Cat7 0.028 0.022 0.020 0.014 0.150	TiC-CD-PyTorc 0.052 0.082 0.003 0.040 -0.025	h TiC-CD-NumF 0.047 0.023 0.035 0.024 0.018

Table 4: Replay-based approaches are competitive with re-training Oracles while 62% cheaper. We scale up two methods that use the ($\alpha_t = 1/2$) version of replay by increasing the monthly token budgets for the continual phase of training. The three scales considered correspond to 220B / 330B / 440B total tokens seen during initialization and continual training. We measure sub-optimality relative to a *series* of Oracle models trained roughly every two years (requiring 1.16T training tokens altogether) and report averages of all Backwards and ID elements of the resulting evaluation matrices.

sections indicate that all continual methods under-perform the Oracle, this Oracle is quite strong for two reasons: (1) for most entries $E_{i,j}$ in our evaluation matrices, it has seen considerably more tokens than model checkpoint i; (2) while it is compute matched with continual runs, it would not actually be obtainable until the last month; if one wanted to consistently re-train new models like the Oracle over this 11 year span, the costs would increase linearly with the number of desired updates.

Thus, to more practically measure the cost effectiveness of continual training, we now consider an 511 alternative baseline of a *series* of Oracle models with different data cutoffs. We then measure the 512 sub-optimality of any checkpoint relative to the *series* of Oracles by subtracting the performance 513 of the most recent Oracle (instead of always the Jul-2024 Oracle). Specifically we consider seven 514 cutoff dates roughly spaced two years apart (i.e., May-2013, Jan-2015, Jan-2017, Jan-2019, Jan-2021, 515 Jan-2023, Jul-2024). Each corresponding Oracle is then token matched with the continual checkpoints 516 corresponding to its cutoff date: e.g., the Jan-2019 Oracle is trained on data coming from all months 517 up to Jan-2019 and for the number of tokens seen by a 220B continual run's Jan-2019 checkpoint. In 518 total all seven oracles together require 1.16T tokens. Given that this now costs more than $5 \times$ our 519 current continual runs, we also consider increasing the compute budget for continual runs by $1.5 \times$ 520 and $2\times$ (by increasing the monthly budget for the continual phase while keeping the initialization the same). These longer runs still cost considerably less than the Oracle series (seeing 330B and 440B 521 tokens respectively), while also being able to update models every month instead of every two years. 522

523 Matching the Oracles with 62% less compute. In Tab. 4, we report the average of the elements that 524 would have appeared in the Backwards and ID elements of the corresponding matrix. Overall, we 525 observe that scaling up Replay ($\alpha_t = 1/2$) + AR to 440B tokens becomes competitive with the series 526 of Oracles. Despite requiring 62% less compute, it surpasses re-training Oracles on many evaluations 527 (i.e., all TIC-CC and TIC-WIKI subsets) while closing the gap significantly on most others.

528 529

530

499

500

501

502

504 505

7 CONCLUSION

531 We introduce a benchmark for continual LLM pretraining spanning more than a decade of times-532 tamped data. TiC-CommonCrawl (TIC-CC) consists of training and evaluation data spanning more than 100 months. We also introduce new TiC evaluations, TiC-Wikipedia (TIC-WIKI), TiC-534 StackExchange (TIC-STACKE), and TIC-CODEDOCS. Using these assets, we clearly observe models 535 need to be continually trained to stay up to date but that the ideal update frequency varies based upon the domain, motivating the need for methods that prevent forgetting. To this end, we compared 536 various baseline strategies for continual pretraining, finding that simple cyclic learning rate schedules 537 and data-replay shrink the gap to an Oracle that trains on all data. However, completely closing the 538 gap remains an open and challenging problem to be studied by future work on our benchmark.

540 REFERENCES

571

580

581

582

583

Oshin Agarwal and Ani Nenkova. Temporal effects on pre-trained models for language processing tasks. *Transactions of the Association for Computational Linguistics*, 10:904–921, 2022.

- 544 Jason Ansel, Edward Z. Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Aniali Chourdia, Will 546 Constable, Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael 547 Gschwind, Brian Hirsh, Sherlock Huang, Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, 548 Mario Lezcano, Yanbo Liang, Jason Liang, Yinghai Lu, C. K. Luk, Bert Maher, Yunjie Pan, 549 Christian Puhrsch, Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi, Helen Suk, Shunting 550 Zhang, Michael Suo, Phil Tillet, Xu Zhao, Eikan Wang, Keren Zhou, Richard Zou, Xiaodong 551 Wang, Ajit Mathews, William Wen, Gregory Chanan, Peng Wu, and Soumith Chintala. Pytorch 2: Faster machine learning through dynamic python bytecode transformation and graph compilation. 552 In Rajiv Gupta, Nael B. Abu-Ghazaleh, Madan Musuvathi, and Dan Tsafrir (eds.), Proceedings of 553 the 29th ACM International Conference on Architectural Support for Programming Languages and 554 Operating Systems, Volume 2, ASPLOS 2024, La Jolla, CA, USA, 27 April 2024- 1 May 2024, pp. 555 929-947. ACM, 2024. doi: 10.1145/3620665.3640366. URL https://doi.org/10.1145/ 556 3620665.3640366.
- Yogesh Balaji, Mehrdad Farajtabar, Dong Yin, Alex Mott, and Ang Li. The effectiveness of memory replay in large scale continual learning. *arXiv preprint arXiv:2010.02418*, 2020.
- Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, 561 Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, 562 Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, 563 Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. Findings of the 2020 conference on machine translation (WMT20). In Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, 565 Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco 566 Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, 567 Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri (eds.), 568 Proceedings of the Fifth Conference on Machine Translation, pp. 1–55, Online, November 2020. 569 Association for Computational Linguistics. URL https://aclanthology.org/2020. 570 wmt-1.1.
- Himanshu Beniwal, Mayank Singh, et al. Remember this event that year? assessing temporal
 information and reasoning in large language models. *arXiv preprint arXiv:2402.11997*, 2024.
- Janek Bevendorff, Benno Stein, Matthias Hagen, and Martin Potthast. Elastic ChatNoir: Search Engine for the ClueWeb and the Common Crawl. In Leif Azzopardi, Allan Hanbury, Gabriella Pasi, and Benjamin Piwowarski (eds.), *Advances in Information Retrieval. 40th European Conference on IR Research (ECIR 2018)*, Lecture Notes in Computer Science, Berlin Heidelberg New York, March 2018. Springer.
 - Janek Bevendorff, Martin Potthast, and Benno Stein. FastWARC: Optimizing Large-Scale Web Archive Analytics. In Andreas Wagner, Christian Guetl, Michael Granitzer, and Stefan Voigt (eds.), *3rd International Symposium on Open Search Technology (OSSYM 2021)*. International Open Search Symposium, October 2021.
- Sungmin Cha and Kyunghyun Cho. Hyperparameters in continual learning: A reality check. arXiv preprint arXiv:2403.09066, 2024. URL https://arxiv.org/abs/2403.09066.
- ⁵⁸⁷ Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- Jie Chen, Zhipeng Chen, Jiapeng Wang, Kun Zhou, Yutao Zhu, Jinhao Jiang, Yingqian Min,
 Wayne Xin Zhao, Zhicheng Dou, Jiaxin Mao, Yankai Lin, Ruihua Song, Jun Xu, Xu Chen,
 Rui Yan, Zhewei Wei, Di Hu, Wenbing Huang, and Ji-Rong Wen. Towards effective and efficient
 continual pre-training of large language models. arXiv preprint arXiv:2407.18743, 2024. URL
 https://arxiv.org/abs/2407.18743.

594 595 596	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>arXiv preprint</i> , 2018. https://arxiv.org/abs/1803.05457.
597 598 599 600	Aaron Defazio, Xingyu, Yang, Harsh Mehta, Konstantin Mishchenko, Ahmed Khaled, and Ashok Cutkosky. The road less scheduled. arXiv preprint arXiv:2405.15682, 2024. URL https: //arxiv.org/abs/2405.15682.
601 602 603 604 605	Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus Hutter, and Joel Veness. Language modeling is compression. In <i>The Twelfth International Conference on</i> <i>Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net, 2024. URL https://openreview.net/forum?id=jznbgiynus.
606 607 608	Bhuwan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. Time-aware language models as temporal knowledge bases. <i>Transactions of the Association for Computational Linguistics</i> , 10:257–273, 2022.
609 610 611 612 613 614	 Felix Drinkall, Eghbal Rahimikia, Janet B. Pierrehumbert, and Stefan Zohren. Time machine GPT. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (eds.), <i>Findings of the Association for</i> <i>Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024</i>, pp. 3281–3292. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-NAACL. 208. URL https://doi.org/10.18653/v1/2024.findings-naacl.208.
615 616 617	Mehrdad Farajtabar, Navid Azizan, Alex Mott, and Ang Li. Orthogonal gradient descent for continual learning. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 3762–3773. PMLR, 2020.
618 619 620	Bahare Fatemi, Mehran Kazemi, Anton Tsitsulin, Karishma Malkan, Jinyeong Yim, John Palowitch, Sungyong Seo, Jonathan Halcrow, and Bryan Perozzi. Test of time: A benchmark for evaluating llms on temporal reasoning. <i>arXiv preprint arXiv:2406.09170</i> , 2024.
621 622 623 624 625	Saurabh Garg, Mehrdad Farajtabar, Hadi Pouransari, Raviteja Vemulapalli, Sachin Mehta, Oncel Tuzel, Vaishaal Shankar, and Fartash Faghri. Tic-clip: Continual training of clip models. In <i>The Twelfth International Conference on Learning Representations (ICLR)</i> , 2024. URL https://openreview.net/forum?id=TLADT8Wrhn.
626 627 628 629 630	GemmaTeam, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy,
631 632 633 634	Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej
635 636 637 638 639	Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed,
640 641 642 643	Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024. URL https://arxiv.org/abs/2403.08295.
645 646	Kshitij Gupta, Benjamin Thérien, Adam Ibrahim, Mats Leon Richter, Quentin Gregory Anthony, Eugene Belilovsky, Irina Rish, and Timothée Lesort. Continual pre-training of large language

Eugene Belilovsky, Irina Rish, and Timothée Lesort. Continual pre-training of large language
 models: How to re-warm your model? In *Workshop on Efficient Systems for Foundation Models* @
 ICML2023, 2023. URL https://openreview.net/forum?id=pg7PUJe0T1.

- Wes Gurnee and Max Tegmark. Language models represent space and time. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.
 OpenReview.net, 2024. URL https://openreview.net/forum?id=jE8xbmvFin.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don't stop pretraining: Adapt language models to domains and tasks. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 8342–8360, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740. URL https://aclanthology.org/2020.acl-main.740.
- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2.
 URL https://doi.org/10.1038/s41586-020-2649-2.
- 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, and Laurent Sifre. Training compute-optimal large language models. *arXiv preprint arXiv:2302.00487*, 2022. URL https://arxiv.org/abs/2203.15556.
- Matthew Honnibal and Ines Montani. spacy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. *To appear*, 7(1):411–420, 2017.
- Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual learning scenarios: A categorization and case for strong baselines. *arXiv preprint arXiv:1810.12488*, 2018.
- Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L. Richter, Quentin Gregory Anthony, Eugene
 Belilovsky, Timothée Lesort, and Irina Rish. Simple and scalable strategies to continually pre-train
 large language models. *Trans. Mach. Learn. Res.*, 2024, 2024. URL https://openreview.
 net/forum?id=DimPeeCxKO.

683

684

685

- Joel Jang, Seonghyeon Ye, Changho Lee, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, and Minjoon Seo. Temporalwiki: A lifelong benchmark for training and evaluating ever-evolving language models. *arXiv preprint arXiv:2204.14211*, 2022a.
- Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Jungkyu Choi, and Minjoon Seo. Towards continual knowledge learning of language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022b. URL https://openreview.net/forum?id= vfsRB5MImo9.
- Zhen Jia, Philipp Christmann, and Gerhard Weikum. Tiq: A benchmark for temporal question an swering with implicit time constraints. In *Companion Proceedings of the ACM on Web Conference* 2024, pp. 1394–1399, 2024.
- Xisen Jin, Dejiao Zhang, Henghui Zhu, Wei Xiao, Shang-Wen Li, Xiaokai Wei, Andrew O. Arnold, and Xiang Ren. Lifelong pretraining: Continually adapting language models to emerging corpora. In Marine Carpuat, Marie-Catherine de Marneffe, and Iván Vladimir Meza Ruíz (eds.), *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*, pp. 4764–4780. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.NAACL-MAIN.351. URL https://doi.org/10.18653/v1/2022. naacl-main.351.

- Jungo Kasai, Keisuke Sakaguchi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A
 Smith, Yejin Choi, Kentaro Inui, et al. Realtime qa: what's the answer right now? *Advances in Neural Information Processing Systems*, 36, 2024.
- Yujin Kim, Jaehong Yoon, Seonghyeon Ye, Sangmin Bae, Namgyu Ho, Sung Ju Hwang, and Se-Young Yun. Carpe diem: On the evaluation of world knowledge in lifelong language models. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pp. 5401–5415. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.NAACL-LONG.302. URL https://doi.org/10.18653/v1/2024. naacl-long.302.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A
 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming
 catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114 (13):3521–3526, 2017.
- Angeliki Lazaridou, Adhi Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d'Autume, Tomas Kocisky, Sebastian Ruder, et al. Mind the gap: Assessing temporal generalization in neural language models. *Advances in Neural Information Processing Systems*, 34:29348–29363, 2021.
- Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Gadre, Hritik Bansal, Etash
 Guha, Sedrick Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation of training
 sets for language models. *arXiv preprint arXiv:2406.11794*, 2024a.
- Yucheng Li, Yunhao Guo, Frank Guerin, and Chenghua Lin. Evaluating large language models for generalization and robustness via data compression. *arXiv preprint arXiv:2402.00861*, 2024b.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(12):2935–2947, 2018. doi: 10.1109/TPAMI.2017.2773081.
- Adam Liška, Tomáš Kočiský, Elena Gribovskaya, Tayfun Terzi, Eren Sezener, Devang Agrawal, Cyprien de Masson d'Autume, Tim Scholtes, Manzil Zaheer, Susannah Young, Ellen Gilsenan-McMahon Sophia Austin, Phil Blunsom, and Angeliki Lazaridou. Streamingqa: A benchmark for adaptation to new knowledge over time in question answering models. *arXiv preprint arXiv:2205.11388*, 2022.
- Vincenzo Lomonaco, Lorenzo Pellegrini, Pau Rodriguez, Massimo Caccia, Qi She, Yu Chen, Quentin Jodelet, Ruiping Wang, Zheda Mai, David Vazquez, et al. Cvpr 2020 continual learning in computer vision competition: Approaches, results, current challenges and future directions. *Artificial Intelligence*, 303:103635, 2022.
- David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning.
 Advances in neural information processing systems, 30, 2017.
- Daniel Loureiro, Francesco Barbieri, Leonardo Neves, Luis Espinosa Anke, and Jose Camachocollados. TimeLMs: Diachronic language models from Twitter. In Valerio Basile, Zornitsa Kozareva, and Sanja Stajner (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pp. 251–260, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-demo.25. URL https://aclanthology.org/2022.acl-demo.25.
- Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *arXiv preprint arXiv:2308.08747*, 2023.
- Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Karishma Mandyam, and Noah A. Smith. Time
 waits for no one! analysis and challenges of temporal misalignment. In Marine Carpuat, MarieCatherine de Marneffe, and Iván Vladimir Meza Ruíz (eds.), *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*, pp. 5944–5958.

756 757 758	Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.NAACL-MAIN.435. URL https://doi.org/10.18653/v1/2022.naacl-main.435.
759 760 761	Sanket Vaibhav Mehta, Darshan Patil, Sarath Chandar, and Emma Strubell. An empirical investigation of the role of pre-training in lifelong learning. <i>Journal of Machine Learning Research</i> , 24(214): 1–50, 2023.
762 763 764	Seyed Iman Mirzadeh, Mehrdad Farajtabar, and Hassan Ghasemzadeh. Dropout as an implicit gating mechanism for continual learning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops</i> , pp. 232–233, 2020a.
765 766 767 768	Seyed Iman Mirzadeh, Mehrdad Farajtabar, Razvan Pascanu, and Hassan Ghasemzadeh. Understand- ing the role of training regimes in continual learning. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 33:7308–7320, 2020b.
769 770 771 772	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. Fast model editing at scale. In <i>The Tenth International Conference on Learning Representations, ICLR 2022,</i> <i>Virtual Event, April 25-29, 2022.</i> OpenReview.net, 2022a. URL https://openreview.net/ forum?id=ODcZxeWfOPt.
773 774 775	Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. Memory- based model editing at scale. In <i>International Conference on Machine Learning</i> , pp. 15817–15831. PMLR, 2022b.
776 777 778 779	Seyed Mahed Mousavi, Simone Alghisi, and Giuseppe Riccardi. Is your llm outdated? benchmarking llms & alignment algorithms for time-sensitive knowledge. <i>arXiv preprint arXiv:2404.08700</i> , 2024.
780 781 782 783 784 785	Kai Nylund, Suchin Gururangan, and Noah A. Smith. Time is encoded in the weights of finetuned language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pp. 2571–2587. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.141. URL https://doi.org/10.18653/v1/2024.acl-long.141.
786 787 788	Jupinder Parmar, Sanjev Satheesh, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Reuse, don't retrain: A recipe for continued pretraining of language models. <i>arXiv preprint</i> <i>arXiv:2407.07263</i> , 2024. URL https://arxiv.org/abs/2407.07263.
789 790 791 792 793	Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. <i>arXiv</i> preprint, 2023. https://arxiv.org/abs/2306.01116.
794 795 796	Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In <i>Computer Vision–ECCV 2020: 16th European Conference</i> , <i>Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16</i> , pp. 524–540. Springer, 2020.
797 798 799 800 801 802	Yujia Qin, Jiajie Zhang, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. ELLE: Efficient lifelong pre-training for emerging data. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pp. 2789–2810, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/ v1/2022.findings-acl.220. URL https://aclanthology.org/2022.findings-acl. 220.
803 804 805 806 807	Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabsa, Mike Lewis, and Amjad Almahairi. Progressive prompts: Continual learning for language models. In <i>The Eleventh International Con- ference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net, 2023. URL https://openreview.net/forum?id=UJTgQBc91
808 809	Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In <i>Proceedings of the IEEE conference on</i> <i>Computer Vision and Pattern Recognition</i> , pp. 2001–2010, 2017.

810 Guy D Rosin, Ido Guy, and Kira Radinsky. Time masking for temporal language models. In 811 Proceedings of the fifteenth ACM international conference on Web search and data mining, pp. 812 833-841, 2022. 813 Karsten Roth, Vishaal Udandarao, Sebastian Dziadzio, Ameya Prabhu, Mehdi Cherti, Oriol Vinyals, 814 Olivier Hénaff, Samuel Albanie, Matthias Bethge, and Zeynep Akata. A practitioner's guide to 815 continual multimodal pretraining. arXiv preprint arXiv:2408.14471, 2024. 816 817 Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray 818 Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. arXiv preprint 819 arXiv:1606.04671, 2016. 820 Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye 821 Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual 822 learning. In International conference on machine learning, pp. 4528–4537. PMLR, 2018. 823 824 Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. arXiv preprint arXiv:1904.07734, 2019. 825 826 Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry W. Wei, Jason Wei, Chris Tar, Yun-Hsuan 827 Sung, Denny Zhou, Quoc V. Le, and Thang Luong. Freshilms: Refreshing large language models 828 with search engine augmentation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), 829 Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and 830 virtual meeting, August 11-16, 2024, pp. 13697–13720. Association for Computational Linguistics, 831 2024. URL https://aclanthology.org/2024.findings-acl.813. 832 Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: 833 Theory, method and application, 2024. URL https://arxiv.org/abs/2302.00487. 834 835 Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. 836 Continual learning for large language models: A survey, 2024. URL https://arxiv.org/ abs/2402.01364. 837 838 Klim Zaporojets, Lucie-Aimée Kaffee, Johannes Deleu, Thomas Demeester, Chris Develder, and 839 Isabelle Augenstein. Tempel: Linking dynamically evolving and newly emerging entities. Advances 840 in Neural Information Processing Systems, 35:1850–1866, 2022. 841 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine 842 really finish your sentence? In Annual Meeting of the Association for Computational Linguistics 843 (ACL), 2019. https://aclanthology.org/P19-1472. 844 845 Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. 846 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 847 pp. 12104–12113, June 2022. 848 Bowen Zhao, Zander Brumbaugh, Yizhong Wang, Hannaneh Hajishirzi, and Noah A. Smith. Set 849 the clock: Temporal alignment of pretrained language models. In Lun-Wei Ku, Andre Mar-850 tins, and Vivek Srikumar (eds.), Findings of the Association for Computational Linguistics, ACL 851 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, pp. 15015-15040. Asso-852 ciation for Computational Linguistics, 2024. URL https://aclanthology.org/2024. 853 findings-acl.892. 854 855 Jonathan Zheng, Alan Ritter, and Wei Xu. NEO-BENCH: evaluating robustness of large language models with neologisms. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings 856 of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long 857 Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pp. 13885–13906. Association 858 for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.749. URL https: 859 //doi.org/10.18653/v1/2024.acl-long.749. 860 861 Da-Wei Zhou, Fu-Yun Wang, Han-Jia Ye, and De-Chuan Zhan. Pycil: A python toolbox for class-862 incremental learning, 2023. 863

A DATASET CONSTRUCTION

We build upon the existing pipeline and assets from DataComp-LM (Li et al., 2024a) to build our dataset, only altering steps that rely on global operations across months.

Initial pool and temporal splitting. We start with DCLM-Pool (Li et al., 2024a) which contains all CC dumps between May-2013 and December-2022. The only pre-processing that has been done on this pool is to parse the HTML (contained in WARC files of CC) into plaintext for each webpage via the open-source resiliparse library (Bevendorff et al., 2018; 2021)². In DCLM-Pool, documents are naturally grouped together into files based upon the CC dump, which is indicated by the file prefix ³. To split the data by month, we simply group files that share the same prefix. Since DCLM-Pool contains data up to December-2022, we also follow their exact download and extraction scripts to obtain more recent data until July-2024.

Data preprocessing and tokenization. Next, we follow DCLM-Baseline's filtering procedure which starts with their implementation of heuristic filters from RefinedWeb. We apply these filters independently on each page with no change. However, we have to modify deduplication that removes nearly identical pages given a similarity threshold. Instead of applying deduplication globally as in DCLM-Baseline, we apply the same deduplication method *only within* each month. Finally, we also skip the final classifier-based filtering in DCLM-Baseline, as their classifier was trained on data that comes from all months, including examples generated by recent LLMs such as GPT-4.

Data sampling and held-out sets. DCLM-Pool was partitioned randomly into 10 equally sized "global shards". Within our monthly splits, we also maintain the original global shard assignments. For our training scales, using just one of these global shards within each month is sufficient. Notably though, when we construct evaluation sets such as in (Sec. 4.1), we make sure to sample from a different global shard than the one used for training. This ensures the evaluation data is a sampled from the same distribution as the training data while also being *mostly* held out. Notably, since we do not deduplicate across globals shards or months, there could be overlap between training and eval sets across months. However, we observe from Fig. 6 that potential data leakages are unlikely significantly change relative losses values (compared to the Oracle). For each validation set, we cap the maximum number of tokens to 16.7M which corresponds to 8192 sequences for our context length of 2048. For some months of TIC-CC-WIKI and TIC-CC-NEWS, we end up with less than this amount, but the smallest are 5M and 12M respectively.



Figure 6: Findings from TIC-CC are robust to potential data leakages. We create a decontaminated version of our TIC-CC loss-based evaluation by deduplicating each month's evaluation set using a Bloom Filter pre-populated by the corresponding training set. Overall, across all the methods, checkpoints, and evaluation months we observe strong correlations between using the pre-decontamination (x-axis) and post-decontamination (y-axis) losses (relative to the Oracle).

²We use readability for parsing code documentations in our TIC-CODEDOCS (https://github.com/mozilla/readability).

³In DCLM-Pool, each file always starts with CC-MAIN-YYYYMM where YYYYMM indicates the dump month.

918 B DETAILS OF EVALUATIONS

920 B.1 TIC-WIKI

931

932

933

934

935

936

937

938

939

940

941

942

945

946

947

948

949

951

922We construct TIC-WIKI from Wikipedia and Wikidata which are sister projects from the non-profit923Wikimedia Foundation. Wikidata is a structured knowledge graph that stores the structured data of924Wikipedia and other sister projects. Data on Wikidata is represented in the form of statements in925the form of property-value about an item in the simplest form. For example, "Mount Everest is the926highest mountain in the world" is represented as Earth (Q2) (item) \rightarrow highest point (P610) (property)927 \rightarrow Mount Everest (Q513) (value) 4. The triplet (item, property, value) can also be referred to as928(subject, relation, object).

TemporalWiki dataset generation. TemporalWiki constructs evaluations from monthly snapshots
 of English Wikipedia and Wikidata through the following steps:

- 1. Generate TWiki-Diffsets by identifying changes and additions between consecutive snapshots of Wikipedia. For new articles, the entire article is added to the Diffset while for existing articles, only the changed or new paragraphs are added.
- 2. Construct TWiki-Probes by processing two consecutive snapshots of Wikidata. Statements are categorized into changed if the property/value has changed or categorized into unchanged otherwise.
- 3. Align TWiki-Diffsets with Wikidata by ensuring changed statements exist in TWiki-Diffsets and unchanged statements exist in Wikipedia.
- 4. Heuristic filtering by removing statements where the subject or object is a substring of the other or the object is more than 5 words. Moreover, a single subject is limited to maximum 1% and relation/object is limited to maximum 5% of the total statements.
- 943 TIC-WIKI extends TemporalWiki in various ways:
 - 1. Expanding the timespan from four months to a decade (2014-2024), thus capturing a broader spectrum of knowledge evolution.
 - 2. We improve the matching process of Wikipedia and Wikidata dumps, and enhance the robustness of data parser to format changes over time.
- 950 B.1.1 DATA PREPROCESSING

Wikidata and Wikipedia dumps. Wikimedia releases regular dumps^{5,6}, but only retains data for 952 the most recent 4 months. To access historical data, we utilized the Internet Archive⁷. The earliest 953 available dump dates back to November 2014. It is important to note that the archived dumps do not 954 cover every month, with several months missing from the record. In our study, we made use of all 955 available monthly dumps. The filenames of the dumps include the specific date of month that has 956 been collected on, which is typically the 1st or 20th of the month, though this can vary. We include 957 only one dump per month if multiple dumps are available. We check for the first date if not available 958 look for 20th and if neither we start from begining the monthh and check for the first available date in 959 that month.

Data cleanup. We utilize WikiExtractor ⁸ to clean up the Wikipedia data. This step extracts the main content and removes extraneous and non-essential characters.

963 Wikipedia diffsets. To construct consecutive diffs of Wikipedia, we developed a method comparing
 964 snapshots of articles from consecutive dumps. For comparing two snapshots of an article, we first
 965 remove extraneous whitespace and standardize formatting by preprocessing the text. This involves
 966 removing empty lines, stripping newline characters, and creating a normalized version of each line
 967 where punctuation is removed and text is converted to lowercase.

^{968 &}lt;sup>4</sup>https://www.wikidata.org/wiki/Help:About_data

^{969 &}lt;sup>5</sup>https://dumps.wikimedia.org/wikidatawiki/

^{970 &}lt;sup>6</sup>https://dumps.wikimedia.org/enwiki/

^{971 &}lt;sup>7</sup>https://archive.org

⁸https://github.com/attardi/wikiextractor

Afterward, we use a two-level comparison: first at the paragraph level, then at the sentence level for
 changed paragraphs. We utilize Python's difflib.SequenceMatcher to compare the normal ized versions of paragraphs and sentences. This hierarchical method, coupled with normalization,
 captures substantial edits while filtering out minor or stylistic changes.

We extract and store both changed and unchanged content separately. Changed content includes replaced paragraphs with modified sentences and newly inserted paragraphs. Unchanged content preserves paragraphs and sentences that remain identical between versions. New articles are treated as entirely changed content. This approach allows us to focus on meaningful content changes while
maintaining the context of unchanged information, providing a comprehensive view of how Wikipedia articles evolve over time. Algorithms 1 and 2 describe the process of constructing Wikipedia diffs and changed/unchanged content.

983 Wikidata diffsets. Next, we extract changed and unchanged Wikidata statements of the form 984 (subject, relation, object) from each consecutive dump. Identical triplets in both dumps are marked 985 as unchanged. Triplets in the new dump not present in the old are categorized as new, with the 986 exception that if a subject entity has more than 10 triplets, the algorithm randomly samples 10 to 987 represent it. When a triplet has the same subject and relation as one in the old dump but a different object and the old and new objects differ only in case (upper/lowercase), the triplet is classified as 988 unchanged; otherwise, it is categorized as new. Triplets from the old dump not found in the new 989 one are implicitly considered removed, but importantly, these are not included in the output sets of 990 changed or unchanged triplets. Throughout this process, the algorithm filters out triplets with overly 991 long object values (more than 5 words) and ensures no duplicates are added. This approach efficiently 992 tracks Wikidata evolution, capturing nuanced changes while managing the volume of data for new 993 entities. Algorithm 3 describes the process of triplet extraction. 994

Wikipedia historical dumps. It is possible to reconstruct each version of Wikipedia using the large
 history files Wikipeida provide ⁹. There are more than 200 historical dumps of English Wikipedia,
 each sized more than 2GB. Combined together, these files include all revisions and all pages of
 Wikipeida.

For Wikidata, Wikimedia does not provide historical diff files as Wikipedia except for the last three months ¹⁰. Wikidata file names are formatted similar to wikidatawiki-20190101-pages-articles.xml.bz2 and available at URLs similar to https://dumps.wikimedia.org/wikidatawiki/20240401/.

Each Wikidata dump is approximately 140GB whereas each Wikipeida dump is less than 22GB.
Therefore, it is possible to make a version of Wikipedia that keeps track of all changes which results in 200 files of 2GB. But as far as we know there are no such files for Wikidata.

- Using the dumps from archive.org has several advantages:
 - We are sure that we do not leak information from previous timesteps.
 - There exists a Wikidata dump close to each Wikipedia dump to be aligned.
 - We can use Wiki-Extractor for filtering and remove Wikipeida editorial discussions.

To illustrate the characteristics of our generated dataset, we present key statistics in the following figures. Figure 7 shows the number of Wikipedia pages with significant changes between consecutive database dumps over time. This graph provides insight into the volume and temporal distribution of our data generation process, highlighting periods of higher and lower content modification as well as distribution of our dumps.

1017

1019

1008

1009

1010

1011

1018 B.2 TIC-STACKE

1020 B.2.1 DATA PREPROCESSING

1021TIC-STACKE spans data from July 2008 through April 2024. The data was sourced from archive.
org using the April 2024 dump of StackExchangeEach category in the dump comes with two key1023

^{1024 &}lt;sup>9</sup>https://dumps.wikimedia.org/enwiki/latest/ file names containing

¹⁰²⁵ pages-meta-history.

¹⁰https://dumps.wikimedia.org/wikidatawiki/

$\frac{26}{Al}$	gorithm 1 Construct Wikipedia Consecutive Diffs
27 —	Input: oldSnapshot_newSnapshot
28 2:	Output: changedContent, unchangedContent
29 3:	oldArticles \leftarrow ReadArticles(oldSnapshot)
³⁰ 4:	newArticles \leftarrow ReadArticles(newSnapshot)
³¹ 5:	changedContent $\leftarrow \emptyset$, unchangedContent $\leftarrow \emptyset$
³² 6:	for each articleId in newArticles.keys do
33 7:	if articleId in oldArticles then
34 8:	$oldText \leftarrow NormalizeText(oldArticles[articleId].text)$
35 9 :	$newText \leftarrow NormalizeText(newArticles[articleId].text)$
36 10:	$changed \leftarrow ExtractChangedContent(oldText, newText)$
37 11:	unchanged \leftarrow ExtractUnchangedContent(oldText, newText)
38 12:	Add (articleId, changed) to changedContent
39 13:	Add (articleId, unchanged) to unchangedContent
40 14:	else
41 15:	Add (articleId, newArticles[articleId].text) to changedContent
42 16:	return changedContent, unchangedContent
43 .	anithm 2 Extract Changed Content
44 AI	
45 1:	Input: oldText, newText
46 2:	Output: changedContent
47 3:	oldParagraphs \leftarrow SplitIntoParagraphs(oldText)
4: 48 -	$here a Regraphs \leftarrow SplitIntoParagraphs(new Text)$
49 5	changedContent $\leftarrow \emptyset$
50 7	if LaDifferent (ald Dara, new Para) then
51 0.	II ISDITICICIII(OluPata, flewPata)
52 0·	$rac{1}{2}$
- 9. 53 10:	for each (oldSent, newSent) in Zip(oldSentences, newSentences) do
3 10. ⊿ 11.	if IsDifferent(oldSent_newSent) then
- 12·	Add newSent to changedContent
a 13:	return changedContent
7	
$\frac{1}{8}$ Al	gorithm 3 Wikidata Triplet Extraction and Categorization
i9 Re	quire: oldDump, newDump
o En	sure: unchanged, new
1	unchanged $\leftarrow \{\}$
2	$new \leftarrow \{\}$
3	newEntities $\leftarrow \{\}$
, i	for all triplet \in newDump do
5	if triplet \in oldDump then
	Add triplet to unchanged
,	else if hasSameSubjectPredicate(triplet, oldDump) then
(oldObject
8	if equalsIgnoreCase(triplet.object, oldObject) then
9	Add triplet to unchanged
U	else
1	Add triplet to new
2	
3	II ITIPIELSUDJECT \notin OLDUMP THEN
1	Add implet to newEntities[triplet.subject]
	else
5	Add inplet to new
	sample new Entity I riplets (new Entitles, new)
3	interLongObjects(unchanged, new)
)	return unchanged new
-	



Figure 7: Number of Wikipedia pages with significant Changes between consecutive archive.org dumps.

1099

1102

files: Post.xml and PostHistory.xml. Post.xml contains information on how answers and questions relate to each other and includes the latest text for each post entry. PostHistory.xml records the changes to each post, whether it is a question or an answer.

106 To construct our dataset, we first build the graph of question-answer relationships based on the 107 Post.xml. We then use PostHistory.xml to reconstruct exact snapshots of posts at specific 108 timestamps. This allowed us to capture the state of each post at the end of each month, ensuring our 109 data reflected the actual content available at those points in time.

1110 We construct binary classification tasks from StackExchange content. For each question, we extract 1111 two responses: the solution accepted by the original author and an alternative option. Our goal is to 1112 create clear distinctions in answer quality, so we implement rigorous selection criteria. Specifically, we requir the accepted solution to have received at least four times the number of upvotes as the 1113 alternative. For the alternative, we choose the response with the lowest upvote count that was posted 1114 before the accepted answer. This strict filtering, while effective in creating distinct quality differentials, 1115 significantly reduced our sample size across most categories. To maintain robust evaluation metrics 1116 while preserving data volume, we introduce an additional metric: the average perplexity of accepted 1117 answers, calculated without applying the strict upvote ratio filter. This approach allows us to include 1118 more samples in our analysis while still capturing meaningful performance trends 1119

We applied this process consistently across all categories of StackExchange, allowing for comprehensive evaluation. In total, we processed 174 out of 182 categories in stackexchange data, of which we focus on stackoverflow in this work as well as a group of seven categories: apple, codereview, electronics, english, gaming, math, and worldbuilding. Some categories had insufficient questions in a single month to provide statistically significant results. In such cases, we combined data from consecutive months, ensuring that each time frame contains at least 500 questions.

1126 The full set of sites includes:

3dprinting, academia, ai, android, anime, apple, arduino, astronomy, aviation, avp, beer, bicycles, bioacoustics, bioinformatics, biology, bitcoin, blender, boardgames, bricks, buddhism, cardano, chemistry, chess, christianity, civicrm, codegolf, codereview, coffee, cogsci, computer-graphics, conlang, cooking, craftcms, crafts, crypto, cs, cseducators, cstheory, datascience, dba, devops, diy, drones, drupal, dsp, earthscience, ebooks, economics, electronics, elementaryos, ell, emacs, engineering, english, eosio, esperanto, ethereum, expatriates, expressionengine, fitness, free-lancing, french, gamedev, gaming, gardening, genai, genealogy, german, gis, graphicdesign, ham, hardwarerecs, health, hermeneutics, hinduism, history, homebrew, hsm, interpersonal, iot, iota, islam,

1134 italian, japanese, joomla, judaism, korean, langdev, languagelearning, latin, law, lifehacks, linguistics, 1135 literature, magento, martialarts, materials, math, matheducators, mathematica, mechanics, meta, 1136 moderators, monero, money, movies, music, musicfans, mythology, networkengineering, opendata, 1137 opensource, or, outdoors, parenting, patents, pets, philosophy, photo, physics, pm, poker, politics, 1138 portuguese, proofassistants, puzzling, quant, quantumcomputing, raspberrypi, retrocomputing, reverseengineering, robotics, rpg, rus, russian, salesforce, scicomp, scifi, security, sharepoint, sitecore, 1139 skeptics, softwareengineering, softwarerecs, solana, sound, space, spanish, sports, sqa, stackoverflow, 1140 stats, stellar, substrate, sustainability, tex, tezos, tor, travel, tridion, ukrainian, unix, ux, vegetarianism, 1141 vi, webapps, webmasters, windowsphone, woodworking, wordpress, workplace, worldbuilding, and 1142 writers. 1143

1144

1145 B.2.2 ANALYSIS OF STACKEXCHANGE DATA

 This section presents an analysis of question-answer patterns across the top 20 categories of StackExchange, with a focus on StackOverflow, Mathematics, and English Language & Usage.

Overall category distribution. Figure 8 shows the distribution of questions across the top 20 StackExchange categories.



1182

(a) StackOverflow

(c) English Language & Usage

Figure 9: Number of questions per month in StackOverflow, Mathematics and English Language & Usage.

(b) Mathematics





Figure 10: Character Count Distribution in StackOverflow, Mathematics and English Language & Usage Questions.

Answer patterns. Figure 11 presents the distribution of answer counts per question for StackOverflow, Mathematics, and English Language & Usage.



Figure 11: Distribution of Answer Counts per Question in Mathematics and English Language & Usage.



1242 C HYPERPARAMTER TUNING

1243

In general we follow the configurations used in DataComp-LM (Li et al., 2024a) unless further specified. For our Oracle and initialization trained on May-2013, we exactly follow their hyperparameters given that these were also standard pre-training runs from scratch.

For our various continual methods, we do perform additional hyperparameter tuning using the first 10 TiC-CC training sets and held-out validation sets. Following Cha & Cho (2024), we limit the tuning to an early set of months given that it would be impossible for a practitioner to be able to tune based upon data they have not seen far in the future. We discuss the tuning and hyperparamter choices for specific methods in more detail below.

Cyclic Cosine. We mainly tuned the maximum learning rate in each cycle, trying values between 1253 1e-3 and 3e-5, as shown in Tab. 5. On our tuning set, the best setting across the board was 1e-4. 1254 When carrying out these tuning runs to completion on all 113 timesteps, we do observe an important 1255 difference in behavior. While 1e-4 continues to offer the best ID performance and strictly dominates 1256 all higher settings, lowering it further can be used to trade-off Backward and ID performance. The 1257 smallest fixed max learning rate, 3e-5 results in a similar yet overall worse performance profile to 1258 using an an AR meta-schedule. This makes sense given the AR schedule roughly can be considered 1259 to decrease the maximum learning rate at a 1/t rate; since our setup involves over 100 months, AR 1260 schedules set the maximum learning rate very close to the minimum of 3e-5 in most rounds. Overall, 1261 we find that learning rates do need to be lowered by at least 30x compared to the the May-2013 initialization (which used 3e-3). This is in contrast to Ibrahim et al. (2024); Gupta et al. (2023) which 1262 both suggest re-warming up to a similar learning rate as the initial pre-training or Parmar et al. (2024) 1263 who start from the minimum learning rate of the pre-trained model. We suspect this is due to the 1264 difference in setup (i.e., these works use only 2 or 3 training rounds of comparable sizes and face 1265 distribution shifts related to data quality and language rather than temporal evolution). 1266

MarilD	TIC-CC	(Tuning]	Months)	(TIC-CC All Months)			
Max LK	Backward	ID	Forward	Backward	ID	Forward	
1e-3	0.103	0.086	0.118	0.197	0.083	0.209	
3e-4	0.019	0.016	0.051	0.125	0.041	0.178	
1e-4	0.002	0.005	0.039	0.072	0.027	0.161	
5e-5	0.002	0.006	0.039	0.062	0.034	0.163	
3e-5	0.004	0.009	0.040	0.060	0.042	0.165	
AR Schedule	0.002	0.008	0.043	0.058	0.040	0.166	

1274 1275

1276

1286 1287

1290 1291

1293 1294 1295

Rsqrt. We tuned both the maximum learning rate within the same range as Cyclic Cosine as well as the cooldown length, choosing between 50 and 400. Our final run continued to use 1e-4 for the maximum learning rate and 400 for the cooldown, though there did not appear to be much difference when compared to smaller values such as 200 or 100 on the tuning months.

Schedule-Free. We continued to use warmup but given that Schedule-Free makes more drastic changes to optimization (i.e. using a different optimizer versus simply a different learning rate schedule), we re-tuned both the learning rate and weight decay. Interestingly, 1e-4 as the maximum learning rate continued to work best for us, though we found it helped slightly to drop the weight decay from 0.033 to 0.01.

		8		
Max LR	WD	TIC-CC Backward	Ionths) Forward	
1e-3	0.033	0.1025	0.0856	0.1178
5e-4	0.033	0.0448	0.0373	0.0713
3e-4	0.033	0.0206	0.0183	0.0532
5e-5	0.033	0.0053	0.0105	0.0406
1e-4	0.067	0.0049	0.0080	0.0406
1e-4	0.033	0.0044	0.0077	0.0404
1e-4	0.010	0.0042	0.0075	0.0403
1e-4	0.005	0.0042	0.0075	0.0403
1e-4	0.0001	0.0044	0.0077	0.0404

LwF. Following the original paper (Li & Hoiem, 2018), we used a temperature parameter of T = 2. 1297 We mainly tuned the regularization weight λ trying values between 0.1 and 1.0 and settling upon 0.3. 1298 However, overall we found using LwF either resulted in little difference (when using a small λ) or 1299 started to decrease all metrics (when using a larger λ).

EWC. We fixed the number of iterations used to estimate the Fisher matrix to 100 and similar to LwF, we focused on tuning the weight given to the EWC regularization term. Overall, we found that fairly high values were needed to overcome the small values in the approximate Fisher matrix (coming from small second order moment terms). We found that $\lambda = 10^7$ performed best when tuning between 10^1 and 10^9 , as shown in Tab. 7. The only other setting we tried that is not strictly dominated by this choice was $\lambda = 10^6$, which resulted in slightly better ID performance but significantly worse backward transfer.

λ	TIC-CC Backward	C (Tuning N ID	lonths) Forward
100	0.0025	0.0050	0.0394
10^{1}	0.0025	0.0050	0.0394
10^{4}	0.0025	0.0050	0.0394
10^{5}	0.0025	0.0049	0.0394
10^{6}	0.0021	0.0047	0.0391
10^{7}	0.0013	0.0050	0.0389
10^{8}	0.0107	0.0178	0.0462
10^{9}	0.0286	0.0400	0.0586

Table 7: **Tuning** λ **for EWC**

1350 D EXTENDED RESULTS



1352 D.1 TIC-COMMONCRAWL (TIC-CC) VALIDATION SETS

Figure 12: Evaluation matrix heatmaps for selected methods on our TIC-CC evaluations.



D.2 TIC-WIKIPEDIA (TIC-WIKI)



Table 8: Comparing TIC-WIKI-Diff versus TIC-WIKI-Unchanged

Mathad	T1C-W1K1-Diff↓			T1C-W1K1-Unchanged↓		
Method	Backward	ID	Forward	Backward	ID	Forward
Cyclic Cosine (std)	0.033 (0.000)	0.052 (0.000)	0.085 (0.000)	0.039 (0.000)	0.052 (0.000)	0.072 (0.000)
Cyclic Cosine + AR	0.033	0.054	0.087	0.035	0.051	0.074
Cyclic Rsqrt	0.031	0.051	0.084	0.035	0.050	0.070
Schedule-Free	0.035	0.055	0.087	0.040	0.055	0.074
Replay ($\alpha_t = 1/t$)	0.038	0.063	0.091	0.035	0.056	0.074
Replay ($\alpha_t = 1/2$)	0.032	0.055	0.086	0.034	0.053	0.072
Replay ($\alpha_t = 1/t$) + AR	0.039	0.063	0.092	0.034	0.055	0.077
Replay ($\alpha_t = 1/2$) + AR	0.033	0.057	0.088	0.033	0.052	0.074
LwF	0.033	0.053	0.085	0.039	0.053	0.072
EWC	0.030	0.051	0.083	0.034	0.050	0.069



D.3 TIC-STACKEXCHANGE (TIC-STACKE)

Figure 14: Evaluation matrix heatmaps for various methods on the Math site of TIC-STACKE.



Figure 15: Heatmaps for various methods on the StackOverflow site of TIC-STACKE.

Table 9: Average over an extended set of TIC-STACKE evaluations that we refer to as TIC-STACKE-CAT7. This includes the following sites: apple, codereview, electronics, english, gaming, math, and worldbuilding. Overall, we find that a combination of replay and AR meta-schedules does the most to reduce forgetting while EWC performs best on ID and Forward evaluations.

M-41 J	TIC-STACKE-CAT7↓			
Method	Backward	ID	Forward	
Cyclic Cosine (std)	0.045 (0.001)	0.050 (0.001)	0.071	
Cyclic Cosine + AR	0.035	0.044	0.068	
Replay ($\alpha_t = 1/t$)	0.036	0.052	0.072	
Replay ($\alpha_t = 1/2$)	0.038	0.049	0.070	
Replay ($\alpha_t = 1/t$) + AR	0.031	0.050	0.072	
Replay ($\alpha_t = 1/2$) + AR	0.032	0.046	0.071	
LwF	0.044	0.048	0.070	
EWC	0.035	0.043	0.067	





¹⁷⁸² E EXTENDED RELATED WORK

1783 1784

Temporal knowledge evaluations. Language models are expected to have an understanding of 1785 time to answer questions about specific time periods and generally reason about time. Various 1786 benchmarks have been proposed to evaluate temporal knowledge of LLMs. TemporalWiki (Jang 1787 et al., 2022a) evaluates the capability of models to update factual knowledge. TemporalWiki is 1788 constructed from the difference between four consecutive snapshots of Wikipedia and Wikidata. Our 1789 TIC-WIKI evaluation expands and improves on TemporalWiki in various ways (see Appx. B.1). 1790 StreamingQA (Liška et al., 2022) consists of human written and generated questions from 14 years 1791 of news articles. The evaluation is either open-book where a model receives a collection of news 1792 articles that contain the answer, or closed-book where the model is first fine-tuned on the training set containing the documents and then tested. We evaluate our TiC checkpoints on StreamingQA 1793 both in open/closed-book setups and find that there is high ambiguity in the questions that evaluates 1794 reasoning more than temporal knowledge understanding. TempEL (Zaporojets et al., 2022) evaluates 1795 entity linking performance across 10 yearly snapshots of Wikipedia. Entity linking is the task of 1796 mapping anchor mentions to target entities that describe them in a knowledge base. In comparison, 1797 our TIC-WIKI evaluates general language and knowledge understanding. TempLAMA (Dhingra et al., 2022) constructs an evaluation for factual queries from Wikidata. They focus on temporally 1799 sensitive knowledge with known start and end dates in a specific Wikidata snapshot. Notably, they propose TempoT5 to jointly model text and timestamp which allows a language model to answer 1801 temporal questions that change over time such "Who is the president". EvolvingQA (Kim et al., 2024) 1802 is also a benchmark for training and evaluating on Wikipedia over time where a LLM automatically 1803 generates question-answers from 6 months of articles in 2023. We avoid using any LLMs for generating our evaluations to prevent transfer of bias. TIQ (Jia et al., 2024) benchmark consists of 1805 10k questions-answers based on significant events for the years 1801–2025.

1806 **Temporal generalization.** Beyond understanding the past, LLMs need to be prepared for the future. 1807 Li et al. (2024b) observes performance deterioration of public LLMs on Wikipedia, news, code, and 1808 arXiv papers after their training data cutoff date. They particularly use compression rate achieved by 1809 treating an LLM as a general input compressor using arithmetic coding (Delétang et al., 2024). Our 1810 comprehensive evaluations on CommonCrawl, Wikipedia, news articles, StackExchange, and code evaluations verifies their results and more comprehensively shows that the rate of deterioration is 1811 domain-specific. DyKnow (Mousavi et al., 2024) evaluations also reaffirm that LLMs private and 1812 open-source LLMs have outdated knowledge by asking them questions constructed using Wikidata. 1813 They also observe LLMs output inconsistent answers in response to prompt variations and current 1814 knowledge editing methods do not reduce outdatedness. TAQA (Zhao et al., 2024) further demonstrate 1815 that pretrained LLMs mostly answer questions using knowledge from years before their pretraining 1816 cutoff. They construct question/answers from Wikipedia for years 2000-2023 and propose three 1817 methods to improve the temporal alignment of models. Similar observations have been made in 1818 RealTimeQA (Kasai et al., 2024) and TempUN (Beniwal et al., 2024). These works further solidify 1819 the need for continuously updating models with continual pretraining.

1820 Temporal understanding. General temporal understanding involves reasoning based on the relation 1821 between existing knowledge. Test of Time (Fatemi et al., 2024) benchmark evaluates temporal 1822 reasoning, logic, and arithmetics by constructing a synthetic dataset. Their goal is to reduce the chance of factual inconsistency in the evaluation using synthetic data. Our benchmark is designed 1824 to be fully realistic based on real data and timestamps to understand the challenges of large-scale 1825 continual pretraining in practice. Gurnee & Tegmark (2024) find that LLMs learn a representation of 1826 space and time with individual neurons that encode spatial and temporal coordinates. They construct datasets of named entities and find that linear probing LLMs performs well on predicting spatial 1827 and temporal coordinates. Nylund et al. (2024) proposed time vectors that specify a direction in the 1828 model's weight space that improve performance on text from a specific time period. 1829

Temporal domain-specific evaluations. We can further analyze the temporal understanding of a model based on the performance on specific domains with varying rates of change. Luu et al. (2022) studied temporal misalignment such as quantifying temporal degradation of domain-specific finetuning in four domains: social media, science, news, and food reviews. They observed significant temporal degradation in domains such as news, social media, and science but less in food reviews.
Gururangan et al. (2020) studied domain-adaptive pretraining and task-adaptive pretraining on unlabeled data for four domains in science, news, and reviews. They observe domain/task-adaptive 1836 pretraining improves performance on the new domain but do not evaluate forgetting on previous domains. Agarwal & Nenkova (2022) studies the temporal model deterioration on future evaluations. 1838 They find that the deterioration is task-dependent and domain-adaptive pretraining does not help 1839 hypothesizing that limited pretraining data is detrimental in continual pretraining. Jin et al. (2022) 1840 domain-incremental pretraining for four scientific domains as well as temporal pretraining on social media over 6 years. They focus on the impact on downstream performance after fine-tuning. They 1841 observe distillation-based approaches are the most effective in retaining dowstream performance 1842 for tasks related to earlier domains. Overall, the gap between different continual learning methods 1843 remained small that can be due to the small scale of pretraining. In comparison, our TIC-CC training 1844 is simulating large-scale pretraining. 1845

Domain/task-continual learning for LLMs. In domain/task continual learning, the model is 1846 presented with a sequence of tasks with predefined labels (Hsu et al., 2018; Van de Ven & Tolias, 1847 2019; Zhou et al., 2023). Each task comes with its training and test sets. In contrast with continual 1848 pretraining, the model needs to support a growing set of labels while compared with temporal 1849 continual learning, the order of tasks are often arbitrary (e.g., Split-CIFAR, Perm-MNIST). Prominent 1850 methods in this domain are regularization-based methods (Kirkpatrick et al., 2017; Mirzadeh et al., 1851 2020a;b; Farajtabar et al., 2020), replay-based methods that often perform superior (Lomonaco et al., 1852 2022; Balaji et al., 2020; Prabhu et al., 2020), and architecture-based methods that adapt the model over time (Schwarz et al., 2018; Rusu et al., 2016). Continual learning for language models has also been dominated by domain/task continual works. Jin et al. (2022) proposed benchmarks for 1855 continually training models on a sequence of research paper domains as well as chronologicallyordered tweet streams. Razdaibiedina et al. (2023) proposed learning a new soft prompt for each task 1857 and pass soft prompts for all seen tasks to the model which provides adaptability while preventing catastrophic forgetting. Luo et al. (2023) studied continual learning for instruction tuning and observed catastrophic forgetting, especially for larger models. Mehta et al. (2023) showed that 1859 generic pretraining implicitly reduces catastrophic forgetting during task incremental finetuning.

1861 Continual pretraining of LLMs. Recent work have studied continual pretraining of foundation 1862 models at large-scale. TiC-CLIP (Garg et al., 2024) proposed a benchmark of training and evaluation of image-text foundation models and demonstrated the deterioration of existing foundation models 1863 on new data. Lazaridou et al. (2021) studied time-stratified language pretraining on WMT, news, and 1864 arXiv up to 2019 and observed the models become outdated quickly on news data that holds even 1865 for models of various sizes. They study dynamic evaluation as a continual pretraining method that 1866 trains on a stream of chronologically ordered documents and observed that models can be updated. 1867 However, they did not explore the impact on forgetting and scalability of the method to more generic 1868 pretraining data over years. Jang et al. (2022b) proposed continual knowledge learning as a new 1869 problem and suggested that parameter expansion is necessary to retain and learn knowledge. They 1870 focus on one-step continual pretraining where models are pretrained on C4/Wikipedia data up to 1871 2020 and then trained once more on recent news articles. They find adapter methods perform better 1872 than regularization and replay methods. Adapter methods are not directly applicable in our multi-1873 year continual pretraining setup where we train in more than 100 steps on large-scale data. Gupta 1874 et al. (2023) proposed simple recipes for continual pretraining of LLMs such as utilizing cyclical learning rate schedules with warmup and ablated on hyperparameters such as warmup duration when 1875 continuing the pretraining on a fixed pair of pretraining datasets. 1876

1877 Time-aware training. Orthogonal to continual pretraining, one can modify the training or fine-tuning of a model to include explicit information about time. TempLAMA (Dhingra et al., 2022) proposed prepending a time prefix to each example during training which gives the model the flexibility to respond to time-sensitive questions. They train models on news articles where the time can be reliably extracted. Drinkall et al. (2024) proposed training a series of models with sequential data cutoffs dates to avoid data contamination with benchmark and private data. The observe no difference across time on static downstream evaluations when training models on news and Wikipedia

Factual editing and retrieval augmented generation (RAG). Another set of works aims to address the staleness of pretrained LLMs without further standard pretraining. One approach is to
surgically edit the facts a model "knows" by identifying and updating the relevant weights within a
model (Mitchell et al., 2022a). Another is to store edits in an explicit memory and learn to reason
over them (Mitchell et al., 2022b). Retrieval augmented generation (RAG) pairs an LLM with new
data sources to retrieve the most relevant document for a query. Generally, continual pretraining
and RAG are orthogonal approaches to generate up to date responses. RAG methods increase the

cost at inference time without changing the model while continual pretraining is the opposite. Fresh LLMs (Vu et al., 2024) proposes a QA benchmark and argues that fast-changing knowledge requires a retrieval-based solution compared with slow-changing knowledge. Continual pretraining can be
 crucial in reducing the cost of RAG by utilizing retrieval only on knowledge that changes faster than the rate of continual pretraining.

1896 F FUTURE WORK

Tokenizer. As the data changes over the years, new words appear in the language that would benefit from temporal adaptation of the tokenizer (Zheng et al., 2024). In this work, we fixed the tokenizer and did not change it across models. One important challenge that changing the tokenizer introduces is that the perplexity of models with different vocabularies will not be directly comparable. Future work would need to either focus on non-perplexity evaluations (Delétang et al., 2024) or normalize perplexity by a mapping between vocabularies of a checkpoint to the reference oracle model.

Joint training of text and timestamp. TIC-CC training data has monthly timestamp corresponding to the crawl time that could be passed as context to the LLM during training and evaluation. TempoT5 (Dhingra et al., 2022) and TempoBERT (Rosin et al., 2022) explored temporal language modeling for example by prefixing the input with "Year: " which helps resolve ambiguity in knowledge that has time-dependent answers such as "Who is the president".