# TemporalWiki: A Lifelong Benchmark for Training and Evaluating Ever-Evolving Language Models

**Anonymous ACL submission**

## Abstract

Language Models (LMs) become outdated as the world changes; they often fail to perform tasks requiring recent factual information which was absent or different during training, a phenomenon called *temporal misalignment*. This is especially a challenging problem because the research community still lacks a coherent dataset for assessing the adaptability of LMs to frequently-updated knowledge corpus such as Wikipedia. To this end, we introduce TEMPORALWIKI, a lifelong benchmark for ever-evolving LMs that utilizes the difference between the consecutive snapshots of Wikipedia and Wikidata for training and evaluation, respectively. The benchmark hence allows one to periodically track an LM's ability to retain previous knowledge and acquire new or updated knowledge at each point in time. We also find that training an LM on the *diff* data with an adapter achieves similar or better perplexity than on the entire snapshot in our benchmark with 12 times less computational cost, which verifies that factual knowledge in LMs can be safely updated with minimal training data via continual learning. The dataset and the code will be available at www.omitted.link.

## 1 Introduction

Large Language Models (LMs) pretrained on a vast amount of text corpus have shown to be highly effective when finetuned or prompted to perform various downstream tasks (Raffel et al., 2019; Brown et al., 2020; Sanh et al., 2021; Wei et al., 2021). However, most of the datasets used to evaluate these LMs are static benchmarks; the train and test data are both from similar points in time. On the other hand, in the real world, factual knowledge is frequently changed, added, or deprecated. For example, suppose a language model is asked what the most dominant coronavirus variant is (Figure 1). The answer would have been the *Delta variant* in the fall of 2021 but has changed to the *Omicron variant* near the end of 2021. If LMs remain unchanged and are not periodically trained to cope with the changing world, they will be outdated very quickly. This means downstream tasks that directly depend on or are finetuned from the LM will suffer from *temporal misalignment* (Luu et al., 2021; Lazaridou et al., 2021), which refers to the misalignment in time between the train and test data.

Temporal misalignment becomes a critical problem, especially when using language models for knowledge-intensive tasks such as closed-book question answering (Roberts et al., 2020; Petroni et al., 2021; Jang et al., 2021) since they rely solely on the knowledge stored in their parameters. Furthermore, LMs augmented with retrieval mechanism (Guu et al., 2020; Lewis et al., 2020; Borgeaud et al., 2021) often suffer from *hallucination* even if they successfully retrieve up-to-date information (Zhang and Choi, 2021; Chen et al., 2021), meaning temporal misalignment still needs to be addressed in such semi-parametric LMs.

Recently, Lazaridou et al. (2021); Jang et al. (2021) have explored updating the internal knowledge of LMs through continual pretraining on new and updated data as a solution for mitigating temporal misalignment. However, these datasets are still *static* in nature: as the world changes, they will eventually get outdated as well. In order to comprehensively measure the capability of ever-evolving LMs on addressing temporal misalignment, automated periodic evaluation of the LMs is crucial.

In this paper, we introduce TEMPORALWIKI, a *lifelong* benchmark for training and evaluating ever-evolving LMs in a periodic and automated manner, shown in Figure 1. The corpora used for updating LMs are constructed by comparing articles from consecutive Wikipedia snapshots and retrieving only *changed* information, which we name as TWIKI-DIFFSETS. The evaluation datasets are constructed in a similar manner by comparing Wikidata snapshots that correspond to the Wikipedia
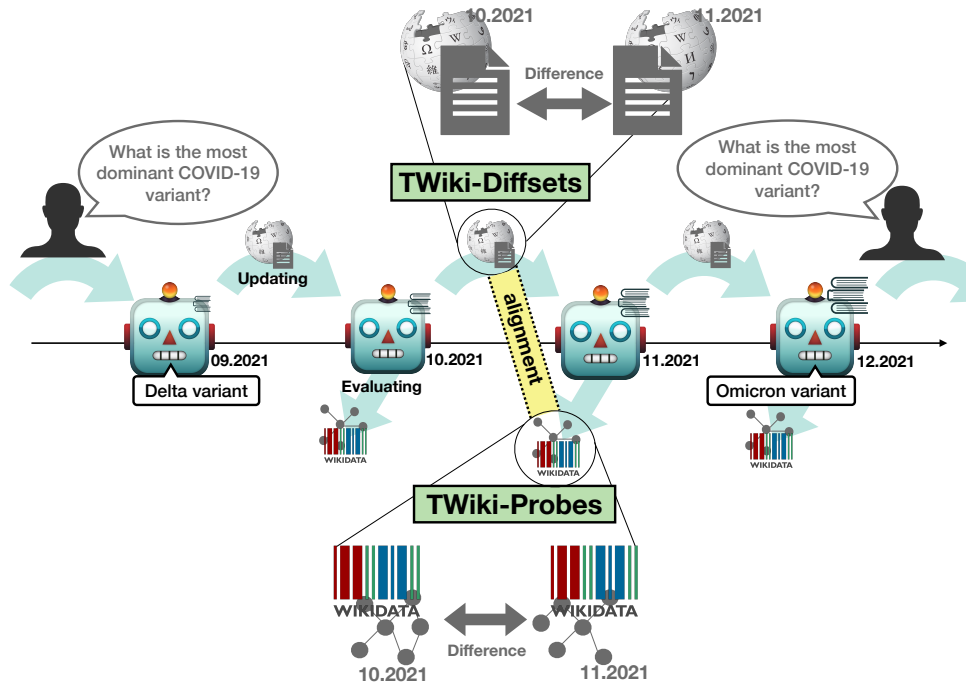
Figure 1: An overview of using TEMPORALWIKI, consisting of TWIKI-DIFFSETS and TWIKI-PROBES to train and evaluate ever-evolving LMs, respectively. Differences between Wikipedia snapshots at different points in time are used for temporal language modeling and categorized factual instances in the corresponding Wikidata snapshots are used for temporal evaluation.

snapshots in time and categorizing each factual instance into UNCHANGED, UPDATED, or NEW. Since Wikidata updates may not exactly align with Wikipedia updates, we only retain factual instances that can be grounded to articles in Wikipedia, ensuring the quality of the data and name the resulting evaluation dataset as TWIKI-PROBES. The whole benchmark creation process is done without any human annotation, thus allowing it to be automated and *lifelong* as new Wikipedia and Wikidata snapshots are released by Wikimedia[1] on a monthly basis.

Through TEMPORALWIKI, we aim to tackle the following research questions: How can we train ever-evolving LMs efficiently and automate the evaluation of each update? How does updating LMs only on new and updated data from Wikipedia compare to updating LMs on entire Wikipedia snapshots, especially in scenarios with multiple updates? How problematic is catastrophic forgetting (McCloskey and Cohen, 1989) when LMs are updated only on new and updated data, and how can we effectively mitigate catastrophic forgetting? Our main contributions are summarized as follows:

- We introduce TEMPORALWIKI, a *lifelong* benchmark for ever-evolving LMs. Unlike previous *static* benchmarks, TEMPORALWIKI

is responsive to the *dynamic* changes in the world and can be utilized to automatically train and evaluate ever-evolving LMs on each Wikipedia and Wikidata snapshot update.

- We find that continually training LMs only on the updated and the new portion of Wikipedia, which we call *temporal language modeling*, is much more computationally efficient than updating LMs on entire Wikipedia snapshots as well as being more effective in terms of stability-plasticity trade-off. It is still a challenging task especially when multiple updates are required due to catastrophic forgetting.

- As a competitive baseline for temporal language modeling, we implement an adapter-based continual learning approach that mitigates forgetting while bolstering the learning of new knowledge, thus providing an overall enhancement in terms of both stability and plasticity. We hope that TEMPORALWIKI will foster future research on continual learning methods for the temporal aspect of ever-evolving LMs.

## 2 Background

Recent works have introduced the need to tackle the issue of temporal misalignment, which refers

---

[1] https://commons.wikimedia.org/

2

to machine learning models showing poor performance due to misalignment in time between the train and test data. Temporal misalignment can be caused either by (1) the dynamic nature of language (Röttger and Pierrehumbert, 2021; Hombaiah et al., 2021; Rosin et al., 2021) or (2) the update of factual information (Chen et al., 2021; Dhingra et al., 2021; Jang et al., 2021).

Luu et al. (2021) have emphasized the effect of temporal misalignment on 8 different NLP downstream tasks, asserting that misalignment between the train and test sets of the downstream tasks causes severe performance degradation which can be mitigated by fine-tuning on the corpus from the target period. Agarwal and Nenkova (2021) have argued this to be less of a concern when utilizing representations from pretrained LMs and show that self-labeling on the downstream task is more effective than continued pretraining on more recent data for temporal adaptation. Note that these works have focused on misalignment caused by the dynamic nature of language on tasks that are not knowledge-intensive, such as text classification.

Others have tackled the problem of temporal misalignment caused by the update of factual knowledge. Lazaridou et al. (2021) have shown that LMs deteriorate significantly in performance when there is a misalignment in time between the pretraining data and the downstream task and argued ever-evolving LMs are necessary. Dhingra et al. (2021) have proposed explicitly including time information during pretraining as a potential solution. Jang et al. (2021); Jin et al. (2021) have implemented continual learning methods to mitigate catastrophic forgetting that occurs during continued pretraining on new data.

Despite the recent surge of community interest in the need for ever-evolving LMs, the community still lacks widely-available resources to train and evaluate such LMs. Previous works have introduced benchmarks comprised of data sources from Twitter feeds (Osborne et al., 2014; Yogatama et al., 2014), recent news articles (Jang et al., 2021), and arXiv papers (Lazaridou et al., 2021) where the temporal adaptability of LMs and the effectiveness of different methodologies of updating LMs can be evaluated. However, these data sources are domain-specific and inherently *static*. On the other hand, Wikipedia and Wikidata are great sources of tracking the dynamic change of world knowledge in diverse domains. 120K volunteer editors make 120

updates to the English Wikipedia per minute and add hundreds of new article entries every day (Logan IV et al., 2021)[2]. TEMPORALWIKI leverages the dynamic nature of Wikipedia and Wikidata to provide a *lifelong* benchmark for developing and maintaining ever-evolving LMs.

## 3 TemporalWiki

In this section, we delve into the process of creating TEMPORALWIKI, which is comprised of training corpora (TWIKI-DIFFSETS) and evaluation datasets (TWIKI-PROBES) constructed from comparing the consecutive snapshots of Wikipedia and Wikidata, respectively. In Section 3.1, we first describe the process of constructing the training corpora from Wikipedia snapshots. Then in Section 3.2, we describe the process of generating the evaluation datasets from Wikidata snapshots and their alignment with Wikipedia. In Section 3.3, we describe the quality control applied to the evaluation datasets. Lastly, in Section 3.4, we briefly discuss the current limitations of TEMPORALWIKI.

---

**Algorithm 1** Generating TWIKI-DIFFSETS

---

**Require:** Wikipedia snapshots $WP_{prev}$ and $WP_{recent}$ where $WP_{recent}$ is more recent.
  $D$ := An empty array to store new and updated data.
  *$article$ in $WP$ has attributes $id$ and $text$
  **for all** $article\ a_r \in WP_{recent}$ **do**
    **if** $a_r.id = a_p.id$ for some $article\ a_p \in WP_{prev}$ **then**
      $D$.append(GETDIFF($a_p, a_r$))
    **else**
      $D$.append($a_r$)
    **end if**
  **end for**

  **function** GETDIFF($a_p, a_r$)
  $Diff$ := An empty string to append difference between $text$ in two $article$s.
  **for all** paragraph $p_r \in a_r.text$ **do**
    **if** $p_r$ have *no* matching sentences with any paragraph $p_p \in a_p.text$ **then**
      $Diff \leftarrow Diff + p_r$
    **else if** $p_r$ have *some* matching and *some* different sentences with any paragraph $p_p \in a_p.text$ **then**
      $Diff \leftarrow Diff + sentences$ that differ between $p_r$ and $p_p$.
    **end if**
  **end for**
  **return** $Diff$

---

### 3.1 Generating Corpora for Temporal Language Modeling from Wikipedia

In terms of computational resources, it is highly inefficient to train an LM on the entire Wikipedia snapshot every time the LM requires updates since

---

[2]https://en.wikipedia.org/wiki/Wikipedia:Statistics

most part of Wikipedia has not changed since the previous snapshot. Moreover, it is not certain whether updating the LM on the entire Wikipedia snapshot is the best approach for updating the factual knowledge stored in the LM. Therefore, we compare the differences between consecutive Wikipedia snapshots in order to use only updated and new text for training. We call these subsets TWIKI-DIFFSETS. Algorithm 1 shows the procedure for generating them.

As shown in Algorithm 1, a single TWIKI-DIFFSET is generated by getting the differences (similarly to `git diff`) between two consecutive Wikipedia snapshots. If an article with a new unique id is included in the recent snapshot, we append the entire article to TWIKI-DIFFSET. For an article having an existing id in the previous snapshot, we compare the two articles by paragraphs and add new or updated sentences to TWIKI-DIFFSETS. Detailed Statistics are shown in Section 4.

### 3.2 Generating Evaluation Datasets from Wikidata

In our work, the main objective for continually pre-training LMs is to add and update the *factual* knowledge stored in the implicit parameters of LMs. The success of an LM update can be evaluated by quantifying the stability-plasticity dilemma (Mermillod et al., 2013): the dilemma of artificial and biological neural systems having to sacrifice either *stability*, ability to retain learned knowledge, or *plasticity*, ability to obtain new knowledge. In order to evaluate whether each update is successful, we need evaluation datasets that can quantify the amount of *new* and *updated* knowledge successfully gained (plasticity) and the amount of knowledge that remains *unchanged* as intended after the LM update (stability). Therefore, we categorize factual instances from Wikidata snapshots that are temporally aligned with Wikipedia snapshots and call the resulting datasets TWIKI-PROBES.

Wikidata snapshots are structured knowledge graphs that store factual information in the form of (`Subject`, `Relation`, `Object`) such as (`Barack Obama`, `born-in`, `Hawaii`). These factual instances can be used to probe the LM for factual knowledge (Petroni et al., 2019). Through Algorithm 2, we distinguish each factual instance into either one of the three categories: UN-CHANGED, UPDATED, or NEW.

---

**Algorithm 2** Generating TWIKI-PROBES

**Require:** Wikidata snapshots $WD_{prev}$ and $WD_{recent}$ where $WD_{recent}$ is more recent.
$Un, Up, N$ := Arrays that store UNCHANGED, UPDATED, and NEW factual instances, respectively.
  **for all** fact $(s_r, r_r, o_r) \in WD_{recent}$ **do**
    $\mathbb{P} \leftarrow \{(s, r, o) \mid s = s_r \text{ where } (s, r, o) \in WD_{prev}\}$
    **if** $\mathbb{P} = \emptyset$ **then**
      $N$.append$(s_r, r_r, o_r)$
    **else if** $r_r \notin \mathbb{P}$ **then**
      $N$.append$(s_r, r_r, o_r)$
    **else if** $r = r_r$ and $o = o_r$ for some$(s, r, o) \in \mathbb{P}$ **then**
      $Un$.append$(s_r, r_r, o_r)$
    **else**
      $Up$.append$(s_r, r_r, o_r)$
    **end if**
  **end for**

---

As shown in Algorithm 2, given two consecutive Wikidata snapshots, a single TWIKI-PROBE is constructed which is used to evaluate an LM updated with TWIKI-DIFFSET. We categorize instances with new `Relation` into NEW, instances with the same `Relation` but a new `Object` into UPDATED, and the others into UNCHANGED.

### 3.3 Quality Control for Evaluation Data

We further apply several quality control steps to the categorized factual instances from Section 3.2 (Algorithm 2) to best represent the actual change of knowledge from the LM update.

**Aligning with TWIKI-DIFFSETS** We ensure correct alignment of UPDATED and NEW factual instances with articles in TWIKI-DIFFSETS and UNCHANGED factual instances with articles from the entire Wikipedia since Wikidata updates do not necessarily entail Wikipedia updates and vice versa. In order to do this, we crawl information from Wikipedia article pages to find the mapping to the corresponding Wikidata entity id and store the information as a dictionary.

Then, for each factual instance from UPDATED and NEW, we check if the `Subject` id can be mapped to an article id from TWIKI-DIFFSETS using the dictionary of id mappings. For each instance from UNCHANGED, we check if the `Subject` id can be mapped to an article id from Wikipedia. For a given factual instance, if `Subject` id is successfully mapped to an article id, we finally check if the `Object` exists in the text of the article. Figure 2 shows an example of a successful alignment. Finally, we remove duplicate instances and instances containing `Object` which has $> 5\%$ overlap on the same evaluation subset.
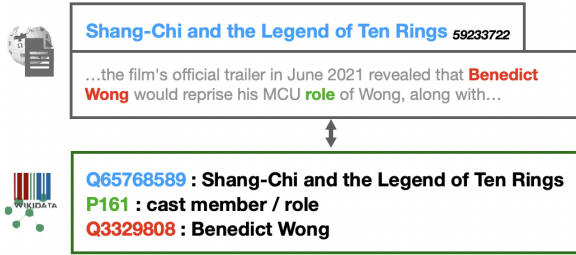
4

Figure 2: An example of a successful alignment between an NEW factual instance from TWIKI-PROBES and an article from TWIKI-DIFFSETS. The alignment is considered successful because for the given factual instance, the Subject matches the title of the Wikipedia page and the Object exists in the article.

|  | # of Articles | # of Tokens |
|---|---|---|
| WIKIPEDIA-08 | 6.3M | 4.6B |
| TWIKI-DIFFSET-0809 | 306.4K | 347.29M |
| WIKIPEDIA-09 | 6.3M | 4.6B |
| TWIKI-DIFFSET-0910 | 299.2K | 347.96M |
| WIKIPEDIA-10 | 6.3M | 4.7B |
| TWIKI-DIFFSET-1011 | 301.1K | 346.45M |
| WIKIPEDIA-11 | 6.3M | 4.6B |
| TWIKI-DIFFSET-1112 | 328.9K | 376.09M |
| WIKIPEDIA-12 | 6.3M | 4.7B |

Table 1: Statistics of TWIKI-DIFFSETS. The two digits indicate the month of the year 2021 that the Wikipedia snapshot was obtained from. The four digits for WIKI-DIFFSET indicate the months of the two snapshots being compared. For instance, TWIKI-DIFFSET-0809 indicates the difference between August (08) and September (09).

### 3.4 Limitations of TEMPORALWIKI

One aspect that is not covered in this work is *knowledge deletion*. While maintaining Wikipedia and Wikidata, volunteer editors not only update or add new information but also *delete* information that is incorrect or misinformed. As removing the misinformation and bias stored in LMs is an important issue and necessary for truly ever-evolving LMs, future work should address this aspect.

### 4 Dataset Statistics

TEMPORALWIKI is constructed from 08.2021 to 12.2021, and its statistics are discussed below.

**Training Corpora Statistics** Statistics of Wikipedia snapshots and TWIKI-DIFFSETS constructed from comparing the snapshots are shown in Table 1. An interesting aspect of TWIKI-DIFFSETS is that the amount of information being updated and added (i.e. number of tokens in each subset) is similar for each month.

**Evaluation Dataset Statistics** Statistics of TWIKI-PROBES divided into categories before and after quality control at each time step are shown in

| Month | Initial Categorization | | | → | Alignment | | |
|---|---|---|---|---|---|---|---|
| | Un | Up | N | | Un | Up | N |
| 0809 | 514,017 | 807,320 | 401,952 | | 10,133 | 785 | 1,544 |
| 0910 | 544,708 | 747,022 | 449,784 | | 10,625 | 802 | 1,819 |
| 1011 | 460,228 | 1,037,248 | 535,530 | | 10,544 | 580 | 1,162 |
| 1112 | 463,623 | 906,002 | 747,707 | | 10,580 | 850 | 2,622 |

Table 2: Detailed Statistics of TWIKI-PROBES during construction phase. **Un**, **Up**, and **N** represents UNCHANGED, UPDATED, and NEW factual instances, respectively. For each subset, alignment with Wikipedia corpus is processed to ensure quality of the dataset. For **Un**, we randomly sample 0.1% of the factual instances after Algorithm 2 because majority of factual instances were initially categorized as **Un**.

Table 2. A single Wikidata snapshot is comprised of 93 million entities where there are around 30 facts for each entity which amounts to roughly *2.8 billion factual instances*. Since most instances from Algorithm 2 are categorized into UNCHANGED, we randomly sample 0.1% of its original size and then apply alignment for quality control, reducing the number of unchanged instances to around 10k.

For further analysis, we break down the entity types of Subject and Object and observe a similar proportion of each entity category for each month of TWIKI-PROBES (Appendix A). We also find that the distribution of Relation is skewed in the decreasing order of NEW, UPDATED, and UNCHANGED (Appendix B).

## 5 Experiments with TEMPORALWIKI

In this section, we train and evaluate ever-evolving LMs with TEMPORALWIKI, which consists of TWIKI-DIFFSETS and TWIKI-PROBES. Section 5.1 describes the experimental settings. Section 5.2 describes the baseline methodologies for updating LMs. Section 5.3 shows evaluation results on the training corpora. Section 5.4 presents evaluation results on TWIKI-PROBES.

### 5.1 Settings

We continue pretraining GPT-2-Large (Radford et al., 2019) (774M parameters), which serves as our baseline language model (LM). We compare the baseline performances between updating GPT-2 with TWIKI-DIFFSETS and updating it with entire Wikipedia snapshots and evaluate each update using TWIKI-PROBES. We also implement an adapter-based continual learning method for mitigating *catastrophic forgetting* that occurs when updating GPT-2 with only TWIKI-DIFFSETS. See the detailed configuration in Appendix C.

## 5.2 Baseline Models

Here we describe four baseline methods used for training and evaluation, namely INITIAL, FULL, DIFF and DIFF-CL, as shown in Table 3 and 4.

**Initial** As the initial model checkpoint, we first bring pretrained GPT-2 from Radford et al. (2019), continue pretraining it on the 08.2021 Wikipedia snapshot for four epochs in total (around 546K global steps), and denote it as INITIAL.

**Full** We start from INITIAL and continue pretraining it on the entire Wikipedia snapshot of each month in a sequential manner. For example, after training on the 09.2021 Wikipedia snapshot from INITIAL, we continue training it on the 10.2021 Wikipedia snapshot and move on to the next snapshot. We denote the resulting model as FULL. We iterate through the training data only once, which corresponds to an average of 4.6 billion token updates (140K global steps) for each month.

**Diff** We start from INITIAL and continue pretraining it on TWIKI-DIFFSETS in a sequential manner. We denote the resulting model as DIFF. Similarly to FULL, we iterate through the training data only once, which is an average of 347 million token updates (12K global steps) for each month.

**Diff-CL** Since catastrophic forgetting may occur when updating LMs with TWIKI-DIFFSETS, we also experiment with applying a competitive adapter-based continual learning method, K-Adapters (Wang et al., 2021), which is a method of freezing the original parameters and adding additional adapters (an increase of 103M parameters) to the LM. We denote the resulting LM as DIFF-CL.[3]

## 5.3 Intrinsic Evaluation

We first perform intrinsic evaluation by measuring the perplexity of the baseline models on their training corpora. For each month, we measure the model's perplexity on TWIKI-DIFFSETS and NON-TWIKI-DIFFSETS, where the latter refers to the subset of the month's entire Wikipedia snapshot that does not include the data from TWIKI-DIFFSETS. We sample 10,000 input instances from each subset with a fixed length of 512 and measure the perplexity on proper noun tokens determined by

---

[3]We add the additional parameters once for the updates from 08.2021. Exploring the optimal interval to add parameters for ever-evolving LMs is left for future work.
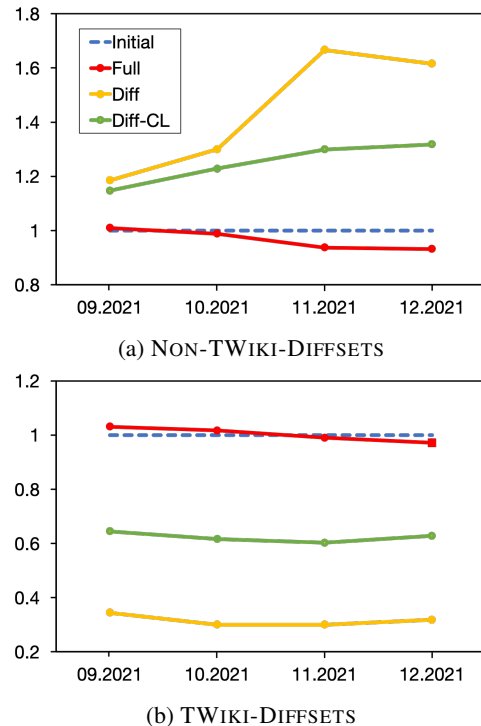


(a) NON-TWIKI-DIFFSETS

(b) TWIKI-DIFFSETS

Figure 3: Relative pronoun perplexity of FULL, DIFF, and DIFF-CL compared to INITIAL on TWIKI-DIFFSETS and NON-TWIKI-DIFFSETS for each month. Lower ratio indicates better performance.

a Part-of-Speech (POS) tagger (Honnibal and Montani, 2017) as in (Lazaridou et al., 2021), which can be considered as a proxy for tokens containing factual knowledge. Therefore, the result on NON-TWIKI-DIFFSETS is meant to indicate the performance on unchanged knowledge, while the result on TWIKI-DIFFSEETS corresponds to updated and new knowledge. Figure 3 shows the relative perplexity of FULL, DIFF, and DIFF-CL compared to INITIAL (i.e., dividing each model by INITIAL, and thus the lower the better).

Results on NON-TWIKI-DIFFSETS show that the relative perplexity of DIFF increases rapidly while that of FULL remains constant as time goes, which implies that forgetting occurs when the LM is trained with TWIKI-DIFFSETS. The relative perplexity of DIFF-CL increases more slowly than DIFF, which means that applying continual learning mitigates catastrophic forgetting.

On the other hand, the results on TWIKI-DIFFSETS show the opposite trend: the relative perplexity of DIFF is much lower than FULL. One thing to note is that the perplexity of FULL is very similar to that of INITIAL on TWIKI-DIFFSETS, which suggests that updating LMs on entire Wikipedia snapshots hinders the effective learning of *changed* data compared to DIFF, despite

6

| | Time | TWiki-Probes-0809 | | | TWiki-Probes-0910 | | | TWiki-Probes-1011 | | | TWiki-Probes-1112 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | **Un** | **Up** | **N** | **Un** | **Up** | **N** | **Un** | **Up** | **N** | **Un** | **Up** | **N** |
| INITIAL | 0 hours | 370.42 | 312.32 | 367.68 | 333.39 | 355.88 | 429.43 | 336.14 | 349.18 | 410.94 | 342.48 | 373.27 | 386.56 |
| FULL | ~24 hours | 364.94 | 309.05 | 366.69 | 357.87 | 369.68 | 458.45 | **306.25** | 329.29 | 365.66 | 348.00 | 354.95 | 357.17 |
| DIFF | ~2 hours | 395.33 | **246.49** | 301.05 | 384.12 | **268.58** | 353.96 | 418.19 | **272.18** | 337.45 | 384.08 | **297.75** | 315.19 |
| DIFF-CL | ~2 hours | **326.54** | 255.02 | **297.64** | **305.26** | 285.92 | **344.42** | 311.44 | 287.01 | **321.77** | **318.66** | 314.37 | **320.90** |

Table 3: Zero-shot perplexity of LMs measured on TWIKI-PROBES where each of **Un, Up, N** represents UNCHANGED, UPDATED, and NEW factual instances, respectively. **Time** represents the average training time of a single update under the setting described in Section 5.1. The descriptions of each baseline models are explained in Section 5.2.
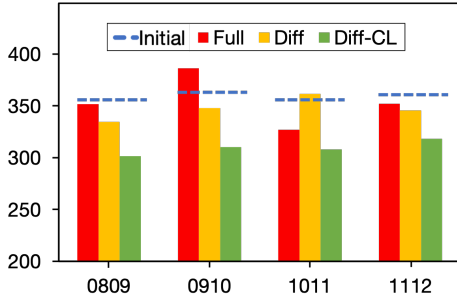


Figure 4: Weighted overall perplexity of TWIKI-PROBES. We weigh and sum the perplexity with equal importance placed on stability and plasticity. A lower score indicates better performance.

both having seen the same instances of TWIKI-DIFFSETS during training for the same number of iterations. DIFF-CL shows higher overall perplexity than DIFF on TWIKI-DIFFSETS due to less number of trainable parameters compared to DIFF.

## 5.4 Extrinsic Evaluation on TWIKI-PROBES

Performing only intrinsic evaluation on the training corpora is not sufficient because the intrinsic evaluation itself only tests the capability of the LMs for memorization (McCoy et al., 2021). Through extrinsic evaluation with TWIKI-PROBES (Section 3.2), we specifically focus on evaluating the factual knowledge from each month.

**Zero-shot**  We use TWIKI-PROBES to measure the zero-shot average perplexity of the updated LMs on each fact instance, shown in Table 3. DIFF and DIFF-CL show better overall performance on UPDATED and NEW factual instances than INITIAL in all months, bolstering the results from intrinsic evaluation. For UNCHANGED, however, DIFF suffers from *catastrophic forgetting*, showing consistent performance degradation as the number of updates increases. In contrast, DIFF-CL shows surprising results on UNCHANGED, outperforming even FULL for three out of the four months. This means that DIFF-CL has effectively mitigated much of the catastrophic forgetting during temporal language modeling. Moreover, it surpasses DIFF

on NEW factual instances, showing that the continual learning method does not hinder the LM from effectively learning new knowledge.

Placing equal importance on *stability* (UNCHANGED) and *plasticity* (UPDATED and NEW), we show the weighted sum of the perplexity on UNCHANGED, UPDATED and NEW data from Table 3 with weights of 0.5, 0.25, 0.25, respectively, in Figure 4. The figure shows that DIFF-CL is the most effective method of updating the LMs over all the time periods while being around 12 times more computationally efficient than FULL. DIFF also outperforms FULL in all months but 1011, showing that temporal language modeling is an effective approach for stability-plasticity trade-off.

We note that, as also shown in previous works (Lazaridou et al., 2021), results in Table 3 present an overall high perplexity (>200) because the sentences in TWIKI-PROBES are not natural sentences; they are factual phrases *synthetically* generated from a naive concatenation of Subject, Relation, and Object. We aim to address this issue via light-tuning, as discussed below.

**Light-tuning**  To alleviate the distributional shift that causes high zero-shot perplexity, we *light-tune* the LMs on 500 instances randomly sampled from WikiData that do not overlap with instances from TWIKI-PROBES (details in Appendix D). Unlike fine-tuning, *light-tuning* lets the LM only learn the input and output distribution of the task, avoiding the problem of test-train overlap pointed out by Lewis et al. (2021). Table 4 shows the results of light-tuning, which demonstrate a similar trend as the zero-shot performance. We also report light-tuning results with the F1 score metric in Appendix E. Although light-tuning avoids the problem of test-train overlap, results are largely affected by the sampled instances for tuning, so a zero-shot evaluation setting is preferred for reliability.

**Effect of Temporal Misalignment**  We quantify the effect of temporal misalignment on each

| | TWiki-Probes-0809 | | | TWiki-Probes-0910 | | | TWiki-Probes-1011 | | | TWiki-Probes-1112 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **Un** | **Up** | **N** | **Un** | **Up** | **N** | **Un** | **Up** | **N** | **Un** | **Up** | **N** |
| INITIAL | **95.93** | 99.68 | 101.05 | **91.34** | 114.19 | 108.12 | **91.23** | 115.78 | 122.3 | **92.96** | 121.42 | 116.72 |
| FULL | 103.81 | 108.41 | 108.56 | 97.3 | 122.27 | 115.58 | 95.08 | 119.21 | 125.2 | 96.79 | 119.66 | 116.72 |
| DIFF | 105.95 | **84.92** | **86.3** | 110.89 | **95.59** | **99.51** | 119.86 | **104.11** | 121.67 | 116.54 | 116.69 | 114.11 |
| DIFF-CL | 98.46 | 94.58 | 99.48 | 102.41 | 113.89 | 105.59 | 99.18 | 111.29 | **115.73** | 98.42 | **113.94** | **112.94** |

Table 4: Light-tuning perplexity of LMs measured on TWIKI-PROBES where each of **Un, Up, N** represents UNCHANGED, UPDATED, and NEW factual instances, respectively.
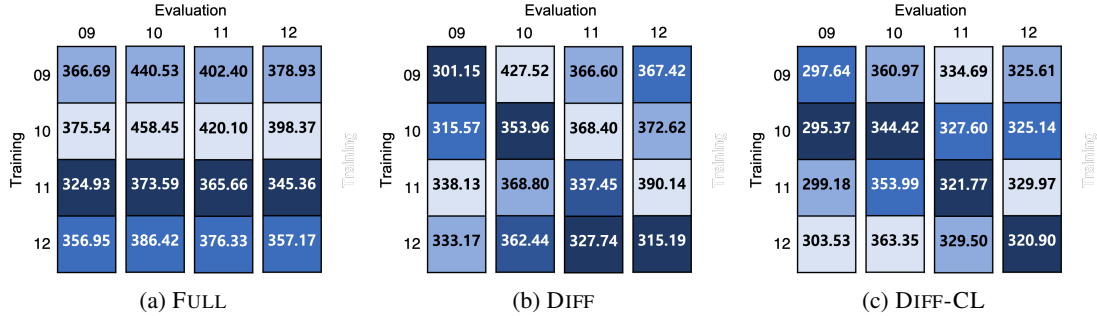


Figure 5: The zero-shot perplexity of the LMs updated and evaluated on various time intervals of NEW of TWIKI-PROBES, showing the effect of temporal misalignment. The better the results, the darker the performance is colored. The color is compared within the same method and also the same evaluation set.

method by training the LMs and evaluating their zero-shot perplexity on NEW instances of TWIKI-PROBES with various time intervals of training and evaluation. As shown in Figure 5, FULL method is mostly influenced by the number of training updates and not much by whether there is temporal alignment. Since FULL is continuously pretrained on the entire Wikipedia corpus in each month, it would have likely seen the data containing NEW factual instances multiple times, leading to lower perplexity as training steps increases.[4] For DIFF and DIFF-CL, there is a general trend of strong performance when there is temporal alignment (diagonal entries), outperforming FULL with much fewer global training steps. It is important to note that DIFF-CL shows robustness against temporal misalignment, i.e. the perplexity does not increase much even when the training and evaluation months do not match, compared to DIFF which suffers from a more severe perplexity spike.

## 6 Conclusion

In this paper, we provide some answers to the four proposed questions in Section 1. (1) *How can we train ever-evolving LMs efficiently and automate the evaluation of each update?* We introduce TEMPORALWIKI, a lifelong benchmark that can be used for training and evaluating ever-evolving language models (LMs) in an automated manner. It consists of TWIKI-DIFFSETS as the training corpora for temporal language modeling and TWIKI-PROBES as the evaluation datasets for measuring the stability-plasticity trade-off at each LM update. (2) *How does updating LMs only on new and updated data from Wikipedia compare to updating LMs on entire Wikipedia snapshots, especially in scenarios with multiple updates?* Through experiments on TEMPORALWIKI, we show that updating LMs on TWIKI-DIFFSETS leads to better acquisition of *new* and *updated* knowledge than updating on entire Wikipedia snapshots with much less computational cost. (3) *How serious is catastrophic forgetting when LMs are updated only on new and updated data?* Temporal language modeling is still a challenging problem, as we observe more forgetting of previous knowledge not contained in TWIKI-DIFFSETS as the number of LM updates increases. However, results still show an overall enhancement in terms of stability and plasticity compared to updating with entire Wikipedia snapshots, showing that temporal language modeling can be also an effective alternative. (4) *How can we mitigate catastrophic forgetting?* We find that an adapter-based continual learning method can mitigate forgetting without hindering the learning of new knowledge, thus achieving the best overall performance.

---

[4]Although directly training INITIAL on the whole Wikipedia corpus of a specific month can be an alternative, we exclude it here because it would only learn the knowledge of the specific month and thus inappropriate for ever-evolving LMs.

# References

Oshin Agarwal and Ani Nenkova. 2021. Temporal effects on pre-trained models for language processing tasks. *arXiv preprint arXiv:2111.12790*.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2021. Improving language models by retrieving from trillions of tokens. *arXiv preprint arXiv:2112.04426*.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *NeurIPS*.

Wenhu Chen, Xinyi Wang, and William Yang Wang. 2021. A dataset for answering time-sensitive questions. In *NeurIPS*.

Bhuwan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. 2021. Time-aware language models as temporal knowledge bases. *arXiv preprint arXiv:2106.15110*.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. Realm: Retrieval-augmented language model pre-training. In *ICML*.

Spurthi Amba Hombaiah, Tao Chen, Mingyang Zhang, Michael Bendersky, and Marc Najork. 2021. Dynamic language models for continuously evolving content. In *KDD*.

Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.

Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Jungkyu Choi, and Minjoon Seo. 2021. Towards continual knowledge learning of language models. *arXiv preprint arXiv:2110.03215*.

Xisen Jin, Dejiao Zhang, Henghui Zhu, Wei Xiao, Shang-Wen Li, Xiaokai Wei, Andrew Arnold, and Xiang Ren. 2021. Lifelong pretraining: Continually adapting language models to emerging corpora. *arXiv preprint arXiv:2110.08534*.

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. 2021. Mind the gap: Assessing temporal generalization in neural language models. In *NeurIPS*.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *NeurIPS*.

Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel. 2021. Question and answer test-train overlap in open-domain question answering datasets. In *EACL*.

Robert L Logan IV, Alexandre Passos, Sameer Singh, and Ming-Wei Chang. 2021. Fruit: Faithfully reflecting updated information in text. *arXiv preprint arXiv:2112.08634*.

Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Karishma Mandyam, and Noah A Smith. 2021. Time waits for no one! analysis and challenges of temporal misalignment. *arXiv preprint arXiv:2111.07408*.

Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. *Psychology of learning and motivation*.

R Thomas McCoy, Paul Smolensky, Tal Linzen, Jianfeng Gao, and Asli Celikyilmaz. 2021. How much do language models copy from their training data? evaluating linguistic novelty in text generation using raven. *arXiv preprint arXiv:2111.09509*.

Martial Mermillod, Aurélia Bugaiska, and Patrick Bonin. 2013. The stability-plasticity dilemma: Investigating the continuum from catastrophic forgetting to age-limited learning effects. *Frontiers in Psychology*.

Miles Osborne, Ashwin Lall, and Benjamin Van Durme. 2014. Exponential reservoir sampling for streaming language models. In *ACL*.

Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, et al. 2021. Kilt: a benchmark for knowledge intensive language tasks. In *NAACL*.

Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? In *EMNLP*.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*.

Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In *EMNLP*.

Guy D Rosin, Ido Guy, and Kira Radinsky. 2021. Time masking for temporal language models. *arXiv preprint arXiv:2110.06366*.

Paul Röttger and Janet B Pierrehumbert. 2021. Temporal adaptation of bert and performance on downstream document classification: Insights from social media. In *Findings of EMNLP*.

Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*.

Leslie N Smith. 2018. A disciplined approach to neural network hyper-parameters: Part 1–learning rate, batch size, momentum, and weight decay. In *CVPR*.

Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Cuihong Cao, Daxin Jiang, Ming Zhou, et al. 2021. K-adapter: Infusing knowledge into pre-trained models with adapters. In *Findings of ACL*.

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.

Dani Yogatama, Chong Wang, Bryan R Routledge, Noah A Smith, and Eric P Xing. 2014. Dynamic language models for streaming text. *TACL*.

Michael JQ Zhang and Eunsol Choi. 2021. Situatedqa: Incorporating extra-linguistic contexts into qa. In *EMNLP*.

## A    Details of Entity Types of `Subject` and `Relation`

Figure 6 shows the ratio of different entity types of `Subject` and `Relation` of UNCHANGED, UPDATED, and NEW.
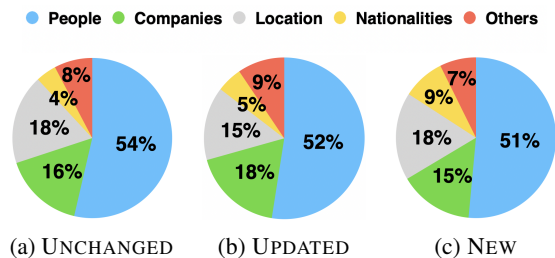


Figure 6: Entity types of `Subject` and `Object` in TWIKI-PROBES.

## B    Details of `Relation` Distribution

The distribution of `Relation` for UNCHANGED, UPDATED, and NEW factual instances in TWIKI-PROBES is shown in Figure 7.

## C    Continual Pretraining and Light Tuning Configuration

For continual pretraining of LMs, we use 8 V100 GPUs with a global batch size of 64 and a fixed input sequence length of 512 for each update. We use the max learning rate of 1e-4 and one cycle learning rate scheduling policy (Smith, 2018). For light-tuning, the training is done for only 1 epoch with a learning rate of 1e-5 and batch size of 32. Input and output sequence lengths are equal to 25. For DIFF-CL, we unfreeze the whole parameters for tuning, following Jang et al. (2021).

## D    Light-Tuning Data

We sample 500 instances from WikiData for each time step that do not overlap with instances from TWIKI-PROBES for each factual instance category. During sampling, we keep the distribution of each `Relation` proportional to the original distribution. Table 5 shows the size and distribution of `Relation` of light-tuning datasets.

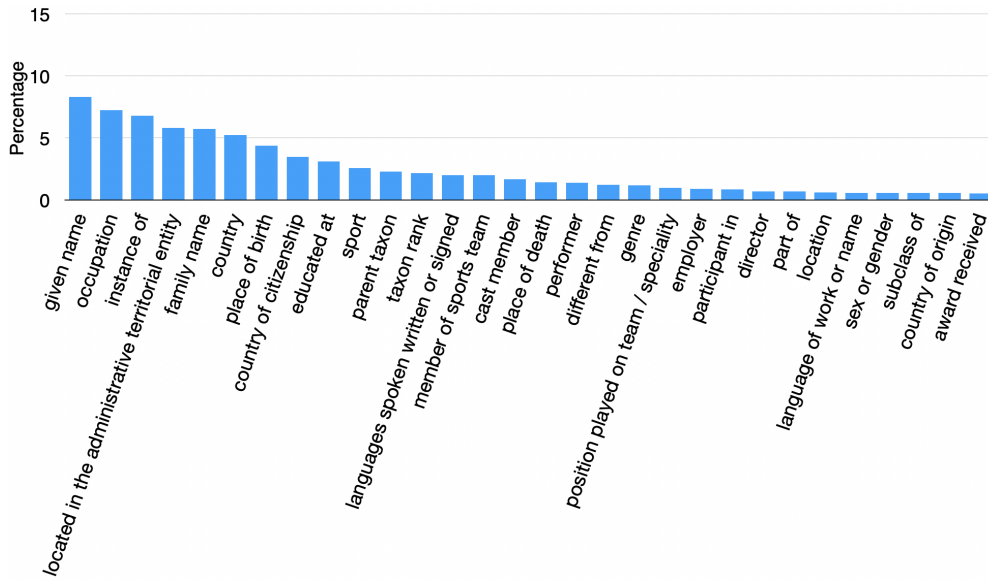|  | Size | # of Relation | Maximum Repetition of Relation | # of Subject |
|---|---|---|---|---|
| UNCHANGED | 500 | 115 | 38 | 500 |
| UPDATED | 500 | 114 | 32 | 496 |
| NEW | 500 | 145 | 50 | 497 |

Table 5: Statistics of the data used for Light-Tuning

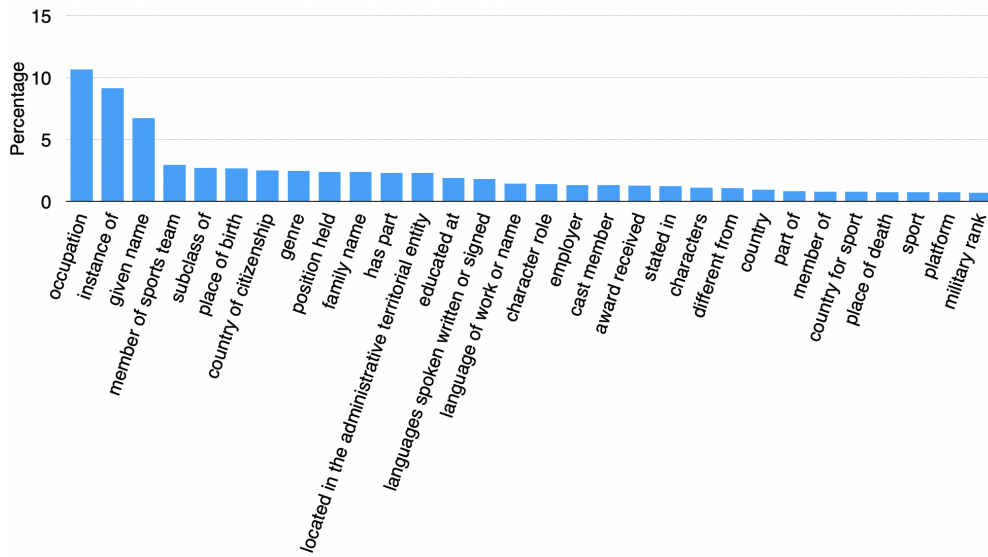## E    F1 score results of light-tuning with TWIKI-PROBES

Many knowledge-intensive tasks such as closed-book question answering (Roberts et al., 2020; Petroni et al., 2021; Jang et al., 2021) or slot filling (Petroni et al., 2021) use accuracy, EM, or F1 score to evaluate the task. We also show the F1 score on TWIKI-PROBES in Table 6. Overall trend is consistent with the perplexity metric: INITIAL shows good performance for UNCHANGED while DIFF and DIFF-CL shows better results for UPDATED and NEW. There are some cases that INITIAL performs best for UPDATED. This is due to the small evaluation set size (<1,000) and low absolute F1 score of UPDATED.

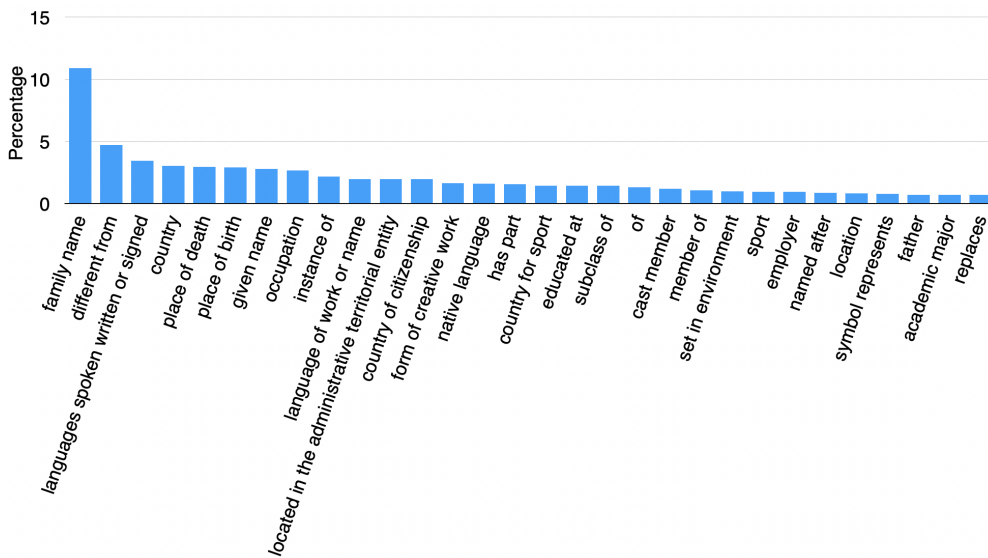| | TWiki-Probes-0809 | | | TWiki-Probes-0910 | | | TWiki-Probes-1011 | | | TWiki-Probes-1112 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **Un** | **Up** | **N** | **Un** | **Up** | **N** | **Un** | **Up** | **N** | **Un** | **Up** | **N** |
| INITIAL | **13.50** | **4.99** | 13.32 | 12.95 | **4.11** | 17.57 | **13.12** | 3.93 | 12.12 | 13.04 | 3.98 | 13.58 |
| FULL | 12.97 | 4.91 | 13.08 | 12.66 | 3.66 | 16.04 | 12.96 | 3.70 | 11.11 | **13.38** | 3.60 | 11.98 |
| DIFF | 13.32 | 4.86 | **14.66** | **13.14** | 3.93 | **18.39** | 11.38 | **5.39** | 12.42 | 13.13 | **4.11** | **13.94** |
| DIFF-CL | 11.38 | 4.86 | 11.60 | 10.87 | 3.93 | 16.47 | 11.09 | 3.81 | **12.50** | 11.42 | 3.60 | 12.07 |

Table 6: F1 score result of LMs on TWIKI-PROBES after light-tuning.

(a) UNCHANGED

(b) UPDATED

(c) NEW

Figure 7: TWiki-Probes distribution of the top 30 most frequent `Relation`.