# Time Waits for No One!
# Analysis and Challenges of Temporal Misalignment

**Kelvin Luu**[1]   **Daniel Khashabi**[2]   **Suchin Gururangan**[1]
**Karishma Mandyam**[1]   **Noah A. Smith**[1,2]

[1]University of Washington   [2]Allen Institute for AI

{kellu,sg01,krm28,nasmith}@cs.washington.edu,
danielk@allenai.org

## Abstract

When an NLP model is trained on text data from one time period and tested or deployed on data from another, the resulting *temporal misalignment* can degrade end-task performance. In this work, we establish a suite of eight diverse tasks across different domains (social media, science papers, news, and reviews) and periods of time (spanning five years or more) to quantify the effects of temporal misalignment. Our study is focused on the ubiquitous setting where a pretrained model is optionally adapted through continued domain-specific pretraining, followed by task-specific finetuning. We establish a suite of tasks across multiple domains to study temporal misalignment in modern NLP systems. We find stronger effects of temporal misalignment on task performance than have been previously reported. We also find that, while temporal adaptation through continued pretraining can help, these gains are small compared to task-specific finetuning on data from the target time period. Our findings motivate continued research to improve temporal robustness of NLP models.[1]

## 1 Introduction

Changes in the ways a language is used over time are widely attested (Labov, 2011; Altmann et al., 2009; Eisenstein et al., 2014); how these changes will affect NLP systems built from text corpora, and in particular their long-term performance, is not as well understood.

This paper focuses on *temporal misalignment*, i.e., where training and evaluation datasets are drawn from different periods of time. In today's pretraining-finetuning paradigm, this misalignment can affect a pretrained language model—a situation that has received recent attention (Jaidka et al., 2018; Lazaridou et al., 2021; Peters et al., 2018; Raffel et al., 2020; Röttger and Pierrehumbert, 2021)—or the finetuned task model, or both. We

suspect that the effects of temporal misalignment will vary depending on the genre or domain of the task's text, the nature of that task or application, and the specific time periods.

We focus primarily on measuring the extent of temporal misalignment on task performance. We consider eight tasks, each with datasets that span at least five years (§2.4), ranging from summarization to entity typing, a subproblem of entity recognition (Borthwick, 1999). Notably, these task datasets span four different domains: social media, scientific articles, news, and reviews. We introduce an easily interpretable metric that summarizes the rate at which task performance degrades as function of time.

Our research questions are:

($Q_1$) *how does temporal misalignment affect downstream tasks over time?*
($Q_2$) *how does sensitivity to temporal misalignment vary with text domain and task?*
($Q_3$) *how does temporal misalignment affect language models across domains and spans of time?*
($Q_4$) *how effective is temporal adaptation, or additional pretraining on a target year, in mitigating temporal misalignment?*

We find that temporal misalignment affects both language model generalization and task performance. We find considerable variation in degradation across text domains (§3.2) and tasks (§3.1). Over five years, classifiers' $F_1$ score can deteriorate as much as 40 points (political affiliation in Twitter) or as little as 1 point (Yelp review ratings). Two distinct tasks defined on the same domain can show different levels of degradation over time.

We explore domain adaptation of a language model, using temporally selected (unannotated) data, as a way to curtail temporal misalignment (Röttger and Pierrehumbert, 2021). We find that this does not offer much benefit, especially relative

---

[1]Data and code are available here.

to performance that can be achieved by finetuning on temporally suitable data (i.e., from the same time period as the test data). We conclude that temporal adaptation should not be seen as a substitute for finding temporally aligned labeled data.

The evidence and benchmarks we offer motivate careful attention to temporal misalignment in many applications of NLP models, and further research on solutions to this problem.

**Contributions.** To facilitate the study of temporal misalignment phenomenon on downstream applications, we compile a suite of eight diverse tasks across four important language domains. We define an interpretable metric that summarizes temporal misalignment of a model on a task with timestamped data. Our experiments reveal key factors in how temporal misalignment affects NLP model performance.

## 2 Methodology Overview

We begin by defining the scope of our study.

### 2.1 Learning Pipeline

We consider a process for building an NLP model that is in widespread use by the research community, illustrated in Fig. 1. First, a (neural network) language model (LM) is pretrained on a large text collection that is not necessarily selected for topical or temporal proximity to the text of the target application (our focus is on GPT-2; Brown et al., 2020). Second, the LM is optionally adapted by continued training on a collection strategically curated for closer proximity to the target (Beltagy et al., 2019); this stage is often referred to as domain-adaptive pretraining (DAPT; Gururangan et al., 2020). Finally, the model is finetuned to minimize a task-specific loss, using labeled data representative of what the model is expected to be exposed to in testing or deployment.
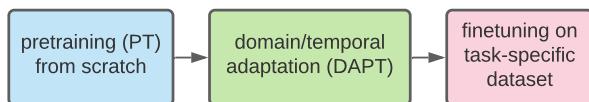


Figure 1: A typical modeling pipeline in NLP.

We study two ways in which temporal misalignment might affect the pipeline's performance as well as straightforward ways to mitigate them.

**Task Shift and Temporal Finetuning** The relationship between text inputs and target outputs may change over time. To the extent that this occurs, annotated datasets used to train NLP systems in the finetuning stage will become stale over time. Due to this temporal misalignment, performance will degrade after deployment, or any in evaluations that use test data temporally distant from the training data. We seek to quantify this degradation across a range of text domains and tasks.

**Language Shift and Temporal Domain Adaptation** Changes in language use can cause a pretrained LM, which commonly serves as the backbone for most modern NLP models, to become stale over time (Lazaridou et al., 2021), regardless of the end task. Lazaridou et al. (2021) explored *temporal adaptation*, continuing LM training on new text data. This is essentially the same procedure as DAPT, where the data is selected by time period. Their work focused on the LM alone, not downstream tasks; we consider both here.

Röttger and Pierrehumbert (2021), the closest to our work, studied temporal adaptation in conjunction to finetuning for a classification task over Reddit data. They conclude that temporal adaptation does not help any more than normal DAPT. We corroborate this work and extend it by studying a wider variety of tasks over a longer span of time periods and thus are better able to draw generalizations from our results.

We believe that the two kinds of shift—task shift and language shift—are difficult to disentangle, and we do not attempt to do so in this work. Instead, we aim to quantify the effect of temporal misalignment on a range of NLP tasks, as well as the benefits of these two strategies.

### 2.2 Evaluation Methodology

Our experiments are designed to measure the effect of temporal misalignment on task performance. To do so, for each task, we fix a test set within a given time period, $T_{test}$. We vary the time period of the training data, allowing us to interpret differences in performance as a kind of "regret" relative to the performance of a model trained on data temporally aligned with $T_{test}$.[2] We consider multiple different test periods for each task. We also seek to control the effect of training dataset size. We partition training data into time periods of roughly

---

[2]This setup avoids a confound of varying test set difficulty that we would encounter if we fixed the model and compared its performance across test datasets from *different* time periods.

the same size and always train on a single partition, keeping the training set size of each time period constant within each task. We expect that performance could be improved by accumulating training data across multiple time periods, but that would make it more difficult to achieve our research goal of quantifying the effect of temporal misalignment on performance.

## 2.3 Quantifying Temporal Degradation

Understanding temporal misalignment requires evaluating a model's performance across data with a range of different timestamps, which makes it difficult to compare various models in terms of their misalignment. We define a metric for temporal degradation (TD) which summarizes the expected speed of model degradation due to temporal misalignment on a task as a single value. In high-level terms, the TD score measures the average rate of performance deterioration (of perplexity, $F_1$, or Rouge-L) for each timestep of difference between that the train and evaluation sets. Higher TD scores imply greater levels of performance deterioration due to misalignment.

Let $S_{t' \to t}$ indicate the performance a model trained on timestep $t'$ data and evaluated on timestep $t$. We define $D(t' \to t)$ as:

$$D(t' \to t) = -(S_{t' \to t} - S_{t \to t}) \times \text{sign}(t' - t).$$

$D(t' \to t)$ is a modified difference in performance between two models.[3] Fig. 2 illustrates $D$ as a function of consecutive training time periods.

We find a line of best fit for $D(t' \to t)$ for all $t'$ using least-squares regression. The slope of this line is $\text{TD}(t)$, the TD score for evaluation time period $t$. The final TD score is the average of the $\text{TD}(t)$ across all evaluation time periods $t$. Further details can be found in Appendix A.

## 2.4 Domains, Tasks, and Datasets

We describe the eight tasks and four domains used for this study. Three (out of eight) of the tasks are newly defined in this work, and all tasks required nontrivial postprocessing. We provide examples and detailed statistics in Table 1.

---

[3]Without the modification, a task with degradation would have have positive performance gaps both $t' > t$ and $t' > t$; the function would not be monotone and the rate of change would be harder to approximate. The modification yields a simpler visual understanding of the deviations over time.
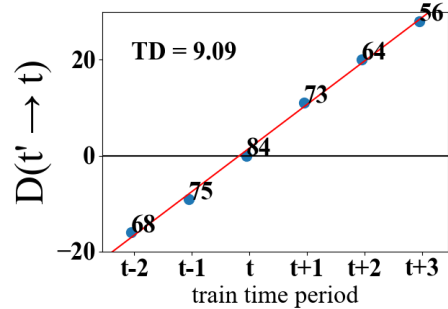


Figure 2: An example calculation of the TD score for a particular timestep $t$ (discussed in Section 2.3). The plotted markers represent $D(t' \to t)$ ($y$-axis) as a function of train time period $t'$ ($x$-axis). The annotated numbers on each blue dot are the raw evaluation scores $S_{t' \to t}$, not to be confused with the $y$ values. The red line is the line of best fit and its slope is the TD score for evaluation timestep $t$. In this example, we would expect to see, on average, 9.09 points of deterioration for each year of misalignment. The final TD score is averaged across all evaluation timesteps.

**Domain 1: Twitter**  Social media platforms like Twitter have been mined to study aspects of language change over time, such as the introduction or diffusion of new words (Eisenstein et al., 2014; Tamburrini et al., 2015; Wang and Goutte, 2017). We collect unlabeled data for domain adaptation by extracting a random selection of 12M tweets, spread semi-uniformly from 2015 till 2020.[4] We experiment with two tasks on Twitter data:

*Political affiliation classification (**POLIAFF**)*  We collect English tweets dated between 2015 and 2020 from U.S. politicians with a political affiliation (*Republican* or *Democrat*). We omit any politician who changed parties over this time period or identified as independent. We consider the downstream task of detecting political affiliations, i.e., given a text of a single tweet we predict the political alignment of its author at the time of the tweet. This task can be useful for studies that involve an understanding of ideologies conveyed in text (Lin et al., 2008; Iyyer et al., 2014).

*Named entity type classification (**TWIERC**)*  We use the Twitter NER dataset from Rijhwani and Preoţiuc-Pietro (2020). The dataset contains tweets dated from 2014 to 2019, each annotated with the mentions of named entities and their types (*Person*, *Organization*, or *Location*). We consider the task of typing a given mention, which is a subproblem of named entity recognition.

---

[4]Collected via the Twitter API.

| Domain | Task | Time Range | Size | Example |
|---|---|---|---|---|
| Twitter | political affiliation classification | 2015-2019 | 120k | **Input:** *History will note that Trump didn't merely fiddle while the planet burned but tried to throw the Arctic National W...* **Output:** *Democrat (vs Republican)* |
| | entity type classification | 2014-2019 | 8k | **Input:** *entity: Finola, tweet: Two 64-year olds enjoying their first birthday together in 40+ years. My twin sister, Finola, and I.* **Output:** *Person* |
| Science | mention type classification | 1980-2016 | 8k | **Input:** *mention: deep Long Short-Term Memory (LSTM) subnetwork, abstract: In this paper, we study the problem of online action detection from the streaming skeleton data .... by leveraging the merits of the deep Long Short-Term Memory (LSTM) subnetwork, the proposed model ...* **Output:** *Method* |
| | venue classification | 2009-2020 | 16k | **Input:** *Rank K Binary Matrix Factorization (BMF) approximates a binary matrix by the product of two binary matrices of lower rank, K...* **Output:** *AAAI (vs ICML)* |
| News | media frame classification | 2009-2016 | 20k | **Input:** *You think you have heard the worst horror a gun in the wrong hands can do, and then this.You think there could not have been anywhere more tragic for it to happen...* **Output:** *Gun Control (15 possible frames)* |
| | publisher classification | 2009-2016 | 67k | **Input:** *A Muslim woman said Sunday that her viral article explaining why she voted for Donald Trump has angered her liberal pals as well as other Muslims.* **Output:** *FoxNews (vs NYTimes or WaPost)* |
| | summarization | 2009-2016 | 330k | **Input:** *The Consumer Financial Protection Bureau is demanding PayPal return $15 million to consumers and pay a $10 million fine for ...* **Output:** *The CFPB alleges many customers unwittingly signed up for PayPal Credit* |
| Food Reviews | review rating classification | 2013-2019 | 126k | **Input:** *What a beautiful store and amazing experience! Not only the atmosphere, but the people...* **Output:** *4 (out of 5)* |

Table 1: The tasks from four domains studied in this paper, with examples. See Section 2.4 for more details.

**Domain 2: Scientific Articles** Scientific research produces vast amounts of text with great potential for language technologies (Wadden et al., 2020; Lo et al., 2020); it is expected to show a great deal of variation over time as ideas and terminology evolve. For adaptation to this domain, we collect unlabeled data from science documents available in Semantic Scholar's corpus,[5] which yields 650k documents, spread over a 30-year period (Ammar et al., 2018). For this domain, we study two tasks:

*Mention type classification (*SCIERC*)* We use the *SciERC* dataset from Luan et al. (2018) which contains entity-relation annotations for computer science paper abstracts for a relatively wide range of years (1980s to 2019). We subdivide the annotated data into time periods with roughly equal-sized numbers of papers (1980–1999, 2000–2004, 2005–2009, 2010–2016). The task is to map a mention of a scientific concept to a type (*Task*, *Method*, *Metric*, *Material*, *Other-Scientific-Term*, or *Generic*).

*AI venue classification (*AIC*)* We also examine temporal misalignment on the task of identifying whether a paper was published in AAAI or ICML. We group the data into roughly equal-sized time periods (2009–2011, 2012–2014, 2015–2017, and 2018–2020). This task is, loosely, a proxy for topic classification and author disambiguation applications (Subramanian et al., 2021).

**Domain 3: News Articles** News articles make up a significant part of the data commonly used

to train LMs (Dodge et al., 2021). News articles convey current events, suggesting strong temporal effects on topic. For adaptation, we use 9M articles from the Newsroom dataset (Grusky et al., 2018), ranging from 2009–2016.[6] We experiment with three tasks on news articles:

*Newsroom summarization (*NEWSUM*)* The Newsroom dataset provides a large quantity of high-quality summaries of news articles (Grusky et al., 2018). We group articles by years for this task (2009–2010, 2011–2012, 2013–2014, 2015–2016). Note that this task, unlike the other tasks considered here, is not a document classification task.

*Publisher classification (*PUBCLS*)* The Newsroom dataset also provides metadata, such as publication source. We take the documents published by the 3 most prolific publishers (Fox News, New York Times, and Washington Post) and train models to classify documents among them. We bin the years (2009–2010, 2011–2012, 2013–2014, 2015–2016). This task is a proxy for applications that seek to infer fact provenance (Zhang et al., 2020). We note that, unlike in our other tasks, we downsample to ensure that the labels are equally balanced.

*Media frames classification (*MFC*)* "Framing" often refers to the emphasis or deemphasis of different social or cultural issues in the media's presentation of the news (Entman, 1983). Card et al. (2015) provide a dataset of news articles annotated
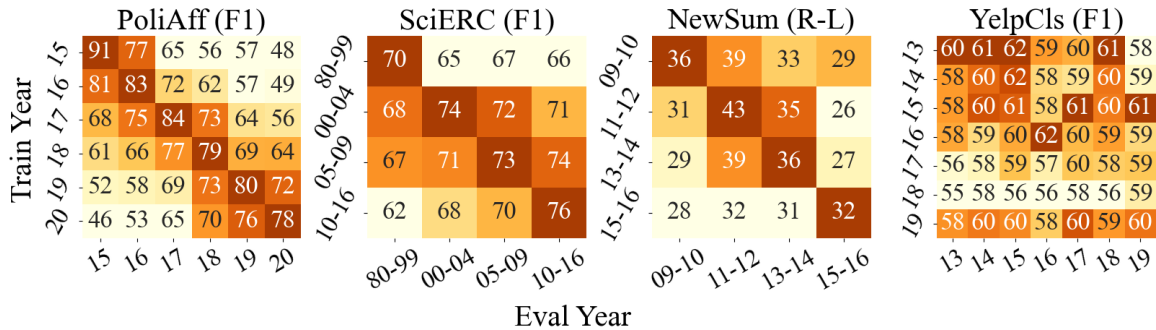
with framing dimensions. We predict the primary frame of a document, treating the problem as a 15-way classification task. We bin by timestamp (2009–2010, 2011–2012, 2013–2014, 2015–2016).

**Domain 4: Food Reviews**   Food and restaurant reviews have been widely studied in NLP research. We considered this domain as a possible contrast to those above, expecting less temporal change. Using data from the Yelp Open Dataset,[7] we consider one task:

*Review rating classification (*YELPCLS*)* This is a conventional sentiment analysis task, mapping the text of a review to the numerical rating given by its author (Pang et al., 2002; Dave et al., 2003). We partition the data by year (2013 to 2019) and ensure that each timestep has a roughly equal amount of reviews.

## 3   Empirical Results and Analysis

In this section, we summarize our experimental analysis, resulting from more than 500 experiments. In our experiments, we primarily explore the effect of temporal misalignment on GPT2 (Brown et al., 2020), a LM often used for generation.[8] We report the macro $F_1$ score for classification tasks and *Rouge-L* (Lin, 2004) for NEWSUM.

We first focus on quantifying temporal misalignment in end tasks. As a preliminary analysis, we investigate how the marginal distribution over labels changes over time. We then study how temporal misalignment affects performance of GPT2 models in downstream tasks with temporal finetuning ($Q_1$,$Q_2$). We find that the amount of performance degradation can vary by task; in some cases the degradation can be severe.

We then study how temporal misalignment affects LMs. As a first step, we analyze how vocabularies change over time in our datasets. We then experiment with ($Q_3$) how temporal misalignment affects upstream language modeling and ($Q_4$) how effective temporal adaptation, or additional pretraining on a target year, is in mitigating misalignment. We find that while LMs are affected by misalignment, temporal domain adaptation is not enough to mitigate temporal misalignment.

Details on temporal domain adaptation and finetuning, and an extended version of our results, can be found in Appendices B and D, respectively.

Label Distribution Change (KL-Div)



Figure 3: KL divergence between label distributions over time for a subset of tasks. See Appendix D for full results. For each cell, we compare the distribution of labels to that of the first time period; e.g., the 2017 PO-LIAFF cell contains the KL-divergence between the label distributions of POLIAFF in 2017 and 2015. While most tasks see little change over time, POLIAFF and MFC see a large shift.

### 3.1   Temporal Misalignment in Tasks

How much does misalignment affect task performance? We find that it depends on the task.

**Label Distribution Drift**   We first investigate how task datasets undergo changes in the marginal distribution over labels due to time. For each task and each test period, we calculate the KL divergence between the label distributions in that period and the first test period. Full results are reported in Fig. 3. In three cases, we detected notable label distribution drift: POLIAFF, AIC, and MFC.[9] In POLIAFF, Republican tweets outnumbered Democratic ones by over a 2:1 ratio in 2015, but the reverse held by 2020. This observation shows that, regardless of the properties of NLP models, the nature of many tasks changes over time, if only because the output distribution changes.

**Finetuning**   As described in §2.4, for each task, we create training and evaluation sets associated with different time periods. We finetune GPT2 on each of the task's training sets and evaluate each on two evaluation sets. Note that there is no domain adaptation here.

Fig. 4 shows our results on downstream tasks (with no domain adaptation). To get more reliable estimates, each number in this heatmap is an average of five independent experiments with different random seeds. A summary of the fine-tuning re-

---

[7]*https://www.yelp.com/dataset*
[8]In our preliminary results, we found that BERT, RoBERTa, and GPT2 models showed similar patterns.

[9]For other tasks, it is possible that the data collection/annotation procedures suppressed label distribution changes that would be visible in data "from the wild."

## PoliAff (F1) — Train Year (rows) × Eval Year (columns)

| Train \ Eval | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|
| 15 | 91 | 77 | 65 | 56 | 57 | 48 |
| 16 | 81 | 83 | 72 | 62 | 57 | 49 |
| 17 | 68 | 75 | 84 | 73 | 64 | 56 |
| 18 | 61 | 66 | 77 | 79 | 69 | 64 |
| 19 | 52 | 58 | 69 | 73 | 80 | 72 |
| 20 | 46 | 53 | 65 | 70 | 76 | 78 |

## SciERC (F1)

| Train \ Eval | 80-99 | 00-04 | 05-09 | 10-16 |
|---|---|---|---|---|
| 80-99 | 70 | 65 | 67 | 66 |
| 00-04 | 68 | 74 | 72 | 71 |
| 05-09 | 67 | 71 | 73 | 74 |
| 10-16 | 62 | 68 | 70 | 76 |

## NewSum (R-L)

| Train \ Eval | 09-10 | 11-12 | 13-14 | 15-16 |
|---|---|---|---|---|
| 09-10 | 36 | 39 | 33 | 29 |
| 11-12 | 31 | 43 | 35 | 26 |
| 13-14 | 29 | 39 | 36 | 27 |
| 15-16 | 28 | 32 | 31 | 32 |

## YelpCls (F1)

| Train \ Eval | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|
| 13 | 60 | 61 | 62 | 59 | 60 | 61 | 58 |
| 14 | 58 | 60 | 62 | 58 | 59 | 60 | 59 |
| 15 | 58 | 60 | 61 | 58 | 61 | 60 | 61 |
| 16 | 58 | 59 | 60 | 62 | 60 | 59 | 59 |
| 17 | 56 | 58 | 59 | 57 | 60 | 58 | 59 |
| 18 | 55 | 58 | 56 | 56 | 58 | 56 | 59 |
| 19 | 58 | 60 | 60 | 58 | 60 | 59 | 60 |

Figure 4: Temporal misalignment in finetuning affects task performance (§3.1). In all cases, higher scores are better. The heatmap is shaded per column, i.e., the darkest shade of orange in a cell means the cell has the highest score in that column. Mismatch between the the training and evaluation data can result in massive performance drop; the degree varies by task. For example, YELPCLS shows minimal degradation. In contrast, POLIAFF shows major deterioration over time. Additional tables for remaining tasks can be found in Appendix D.

| Domain | Task (metric) | TD | $r$ |
|---|---|---|---|
| Twitter | POLIAFF (F1) | 7.72 | 0.98 |
|  | TWiERC (F1) | 0.96 | 0.74 |
| Science | SCIERC (F1) | 0.67 | 0.93 |
|  | AIC (F1) | 1.79 | 0.93 |
| News | PUBCLS (F1) | 5.46 | 0.85 |
|  | NEWSUM (Rouge-L) | 1.38 | 0.91 |
|  | MFC (F1) | 0.98 | 0.86 |
| Reviews | YELPCLS (F1) | 0.26 | 0.30 |

Table 2: Finetuned models' temporal degradation summary scores (TD; §2.3; details in Figure 4). These scores estimate how fast a model degrades as the time period of training and evaluation data diverge (higher scores imply faster degradation). We note that since we normalize by the overall time range of a task, the temporal partitions we used do have an effect on the TD scores. For example, AIC spans ten years, even though there are only four partitions. We also show the correlation coefficient, $r$, that measures the strength of a linear relationship (0 meaning no correlation, 1 being perfectly correlated). In all cases but Yelp, the degree of degradation has a moderate correlation with the distance between the training and evaluation years ($r > 0.5$, $p < 0.05$). We use the Wald test with the null hypothesis that the slope is 0.

sults, in terms of TD scores (§2.3) is in Table 2 which indicates the speed of temporal degradation, for every year that the training and evaluation data diverges. Recall that this score (applied to task performance measures) summarizes the strength of the effect of temporal misalignment on the score, using evidence from across experiments.

($Q_1$) **Temporal misalignment degrades task performance substantially.** Fig. 4, similar to earlier work (Röttger and Pierrehumbert, 2021), shows that models trained on data from the same time period as the test data perform far better than those from the past. The performance drop is most severe for POLIAFF (TD=7.72) and PUBCLS (TD=5.45).

($Q_1$) **Temporal misalignment has a measurable effect on most tasks.** Half of our tasks see an average loss of at least 1 point for each time period that the training data diverges from the test data. For datasets like SCIERC that make use of data from three decades or more, this effect could add up.

Moreover, 1 point of difference can be meaningful, especially for the summarization task where we measure Rouge-L. According to the leaderboard,[10] the best three performing models are within a point of each other in Rouge-L (Shi et al., 2019, 2021; Mendes et al., 2019). The task has a TD score of 1.38. On average, a time period of temporal misalignment results has a larger effect on performance than changing between the three best models.

($Q_1$) **Performance loss from temporal misalignment occurs in both directions.** Another observation in Table 4 is that degradation happens in both directions (past and future). While most of the emphasis on temporal misalignment is on how to adapt our stale models/data to the present time (Dhingra et al., 2021; Lazaridou et al., 2021; Röttger and Pierrehumbert, 2021), our experiments also show that models trained on newer data can be misaligned from the past, as well. Weak performance in older texts has been noted in NLP for historical documents (Yang and Eisenstein, 2016; Han and Eisenstein, 2019). However, our findings indicate deterioration can occur sooner—just a few years rather than decades or centuries.

($Q_2$) **Tasks, even in the same domain, are affected differently.** Consider the two tasks of POLIAFF and TWIERC (both in the Twitter domain), with TD scores of 7.72 and 0.96, respectively. Of our 8 tasks, TWIERC, MFC, and YELPCLS are the most robust to temporal misalignment (TD scores of 0.96, 0.98 and 0.26, respectively). The high levels of variation show that temporal misalignment affects performance through labeled datasets, not just unlabeled pretraining data.

### 3.2 Temporal Misalignment in LMs

As LMs are widely used in modern NLP systems, it is important to inspect how robust they are to temporal misalignment. We seek to understand how temporal misalignment affects the language modeling task in our four domains and if temporal domain adaptation helps in downstream tasks.

**Vocabulary Shift** We first consider an extremely simple measurement of language shift: how do vocabularies change across time periods?[11] We use a similar procedure to the one Gururangan et al. (2020) used for analyzing domain similarity. Fixing a domain, we compare the (unigram) vocabularies of each pair of training sets. The vocabularies are built using the 10K most frequent terms from each time period. We note that vocabulary overlap is higher between two time periods the closer they are. Most domains see a sizeable amount of shift; however, Yelp is relatively stagnant. Fig. 5 visualizes the overlap measurement. Table 6 in Appendix D shows the correlation between model performance and the word overlap.

**Temporal Domain Adaptation** Researchers have studied the broader problem of distributional shift (Shimodaira, 2000; Zhang et al., 2013). The NLP community has historically tackled these problems via domain adaptation (Jiang and Zhai, 2007; Daumé III, 2007; Gururangan et al., 2020). Taking inspiration from these approaches, we next apply DAPT to GPT2, treating each time period as a domain: for each time period, we continue pretraining and then evaluate perplexity. We consider how the perplexity varies with the (mis)alignment between the DAPT training data and the evaluation data. We measure the TD score, which summarizes how much performance is affected by temporal misalignment (now applied to perplexity). The results of temporal domain adaptation are in Fig. 6.

[11]This can be understood as a model-free way to measure covariate shift for NLP tasks that take text as input.

($Q_3$) **Domains are a major driver of temporal misalignment in LMs.** Consistent with Lazaridou et al. (2021), Fig. 6 shows degradation of LM due to temporal misalignment; it further shows considerable variation by text domain. Twitter changes most rapidly, and food reviews are much slower. This observation is consistent with past work on language change in social media (Stewart and Eisenstein, 2018; Eisenstein et al., 2014). To the extent that a LM's practical usefulness is associated with its fit to new data, researchers and practitioners should understand the temporal dynamics of their target text domains and plan LM updates accordingly.

**Joint Effects of Temporal Adaptation and Finetuning** As discussed in §2, continued pretraining of an LM on in-domain text has been shown to improve task performance. Our prior results show that both downstream tasks and language modeling are affected by temporal misalignment. Can temporal domain adaptation help mitigate the effects of misalignment in downstream tasks?

Here we consider how the time period of the data selected for continued pretraining affects task performance. For each task's evaluation set, we apply DAPT twice: once with the earliest available time period's unannotated data and once with the latest's. We then finetune and evaluate on data from the same time periods as in the earlier experiment.

($Q_4$) **Temporal adaptation does <u>not</u> overcome degradation from temporally misaligned labeled data.** In Table 3, we see small performance gains from temporal domain adaptation on LMs, and in some cases it is harmful. These observations underscore the importance of the labeled data; adjustments to the LM alone do not yet appear sufficient to mitigate the effects of temporal misalignment. In contrast to temporal domain adaptation, which does not mitigate temporal misalignment's effects, finetuning on temporally-updated labeled data is more effective.

This can be observed in each task-specific subtable of in Table 3: the top-left and bottom-right quadrants (fine-tuning on time-stamp that is used for evaluation) generally lead to higher scores.

## 4 Limitations and Future Work

We provided a well-controlled suite of experiments to study the effects of temporal misalignment on model performance. However, the setup has some
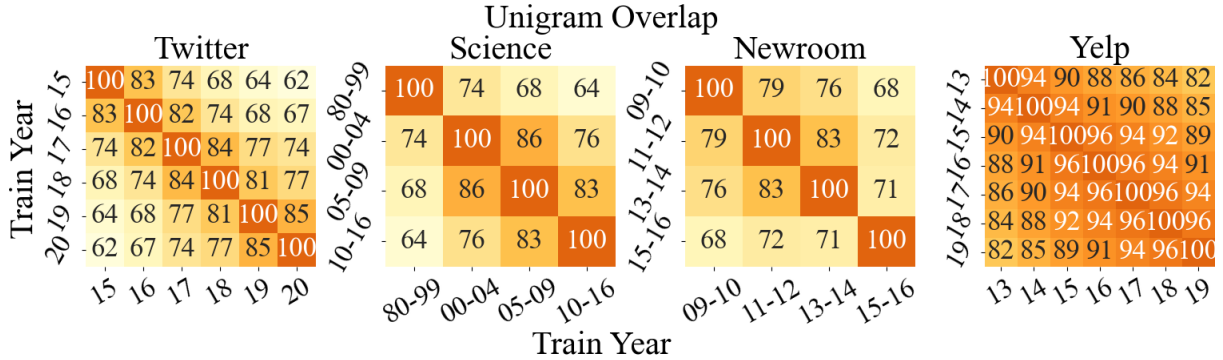
Figure 5: Vocabulary overlap between time periods, over a subset of our tasks' datasets. Each cell shows the % overlap between the vocabularies of two time periods.

| Domain (Task) ↓ | Finetune Year ↓ | Evaluation → Pretrain ↓ | 2015 | 2020 |
|---|---|---|---|---|
| Twitter (PoliAff) *F1* | 2015 | Default | 91.4 | 48.4 |
| | | Default → 2015 | 92.2 | 47.5 |
| | | Default → 2020 | 90.9 | 50.8 |
| | 2020 | Default | 45.8 | 78.0 |
| | | Default → 2015 | 47.2 | 76.9 |
| | | Default → 2020 | 44.2 | 78.3 |

| Domain (Task) ↓ | Finetune Year ↓ | Evaluation → Pretrain ↓ | 1980-1999 | 2010-2016 |
|---|---|---|---|---|
| Scientific (SciERC) *F1* | 1980-1999 | Default | 67.9 | 57.2 |
| | | Default → 1980-1999 | 73.2 | 66.4 |
| | | Default → 2010-2016 | 73.7 | 66.8 |
| | 2010-2016 | Default | 60.3 | 72.5 |
| | | Default → 1980-1999 | 63.4 | 75.0 |
| | | Default → 2010-2016 | 64.8 | 76.0 |

| Domain (Task) ↓ | Finetune Year ↓ | Evaluation → Pretrain ↓ | 2009-2010 | 2015-2016 |
|---|---|---|---|---|
| News (NewsSum) *Rouge-L* | 2009-2010 | Default | 36.4 | 29.0 |
| | | Default → 2009-2010 | 36.4 | 29.1 |
| | | Default → 2015-2016 | 36.1 | 28.9 |
| | 2015-2016 | Default | 27.8 | 31.8 |
| | | Default → 2009-2010 | 28.2 | 31.8 |
| | | Default → 2015-2016 | 27.8 | 31.6 |

| Domain (Task) ↓ | Finetune Year ↓ | Evaluation → Pretrain ↓ | 2014 | 2019 |
|---|---|---|---|---|
| Food Reviews (Yelp) *F1* | 2009-2010 | Default | 58.6 | 58.3 |
| | | Default → 2014 | 63.3 | 60.1 |
| | | Default → 2019 | 60.2 | 62.3 |
| | 2015-2016 | Default | 58.3 | 58.3 |
| | | Default → 2014 | 60.2 | 62.3 |
| | | Default → 2019 | 60.8 | 62.3 |

Table 3: Combination of temporal adaptation and finetuning (§3.2) on our tasks. The row labeled "Default" corresponds to a model that has not been adapted (uses the default pretraining). The models with temporal domain adaptation are shown in rows labeled "Default → $y$" and each is comparable to the "Default" row above it. The color coding is proportional to the magnitude of the performances of each task (darker shade of orange indicates higher scores). It can be observed that temporal finetuning has a greater impact than temporal pretraining. Each quadrant of 3 for each task, indicating the same finetune and evaluation years, but different pretraining conditions, are mostly uniform. In contrast, we notice a sharper difference in performance when varying the finetuning year (comparing the quadrants vertically).

drawbacks. For example, we expect that models trained on data accumulated across multiple time periods would perform well (Lazaridou et al., 2021; Röttger and Pierrehumbert, 2021; Jin et al., 2021).

We chose the time periods in our study so that they would each have sufficient and consistent training data sizes. However, amounts of data in a particular domain or task will fluctuate over time. Moreover, the rate of language use change may not be uniform. Time periods should be selected with these two considerations in mind.

Our findings indicate that temporal misalignment's effects depend heavily on the task. Though not studied here, the same issues may arise in annotation efforts; consider, for example, recent work on controversy (Zhang et al., 2018) and social norms (Xu et al., 2021; Zhou et al., 2021) likely hinges on constructs that may be time sensitive. Annotations that are temporally misaligned with the original data being annotated may be anachronistic.

An opportunity for future exploration is in the context of real-world events with sudden changes such as COVID-19 pandemic (Cao et al., 2021) or political changes, which influence tasks such as question answering (Dhingra et al., 2021; Zhang and Choi, 2021).

Extensive work has been done on modeling and detecting lexical semantic change, or how words evolve in meaning (Hamilton et al., 2016; Rudolph and Blei, 2018; Gonen et al., 2020). Techniques
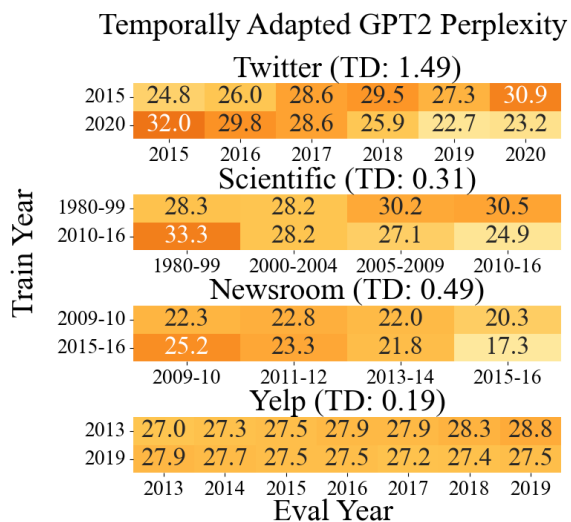
Figure 6: Perplexity of GPT2 after adaptive pretraining on temporally-selected data in different domains (lower is better). The TD score (in parentheses) estimates the expected perplexity rise (i.e., degradation) for every time period of misalignment between evaluation and training times. Degradation follows the expected pattern, but the magnitude varies by domain.

and intuition from this body of work may be useful in finding solutions to mitigate degradation due to misalignment. We believe that this phenomenon is an important aspect of temporal misalignment, but leave disentangling semantic shifts from other, perhaps task-related factors, for future work.

Continual learning, which allows models to learn from a continuous stream of data, could also be one way to mitigate temporal misalignment. Most prior work in this space has focused on continual learning in LMs (Jin et al., 2021) or learning disparate tasks (de Masson d'Autume et al., 2019; Huang et al., 2021). Future work may investigate continual learning algorithms for tasks that change over time.

Our results showed that straightforward domain adaptation was unable to mitigate the effects of temporal misalignment. Recent work in language modeling has elevated the importance of domains by using a mixture of domains (Gururangan et al., 2021) or giving domains a hierarchical structure (Chronopoulou et al., 2021). More sophisticated approach to domains, in line with these works, could lead to temporally robust models.

While we found that task-specific finetuning is more effective than temporal adaptation, new labeled data can be expensive. Ways to characterize or detect changes in a task could be helpful in efficiently updating datasets (Lu et al., 2019; Webb

et al., 2018). Future work can also treat dataset maintenance as an optimization problem between the cost and gains of annotating new data (Bai et al., 2021).

## 5 Conclusion

Changes in language use over time, and how language relates to other quantities of interest in NLP applications, has clear effects on the performance of those applications. We have explored how temporal misalignment between training data—both data used to train LMs and annotated data used to finetune them—affects performance across a range of NLP tasks and domains, taking advantage of datasets where timestamps are available. We compile these datasets as a benchmark for future research as well. We also introduced a summary metric, TD score, that makes it easier to compare models in terms of their temporal misalignment.

Our experiments revealed considerable variation in temporal degradation accross tasks, more so than found in previous studies (Röttger and Pierrehumbert, 2021). These findings motivate continued study of temporal misalignment across applications of NLP, its consideration in benchmark evaluations,[12] and vigilance on the part of practitioners able to monitor live system performance over time.

Notably, we observed that continued training of LMs on temporally aligned data does not have much effect, motivating further research to find effective temporal adaptation methods that are less costly than ongoing collection of annotated/labeled datasets over time.

---

[12]Indeed, for benchmarks where training and testing data *are* aligned, our findings suggest that measures of performance may be in some cases inflated.

# References

Eduardo G Altmann, Janet B Pierrehumbert, and Adilson E Motter. 2009. Beyond word frequency: Bursts, lulls, and scaling in the temporal distributions of words. *PLOS one*, 4(11):e7678.

Waleed Ammar, Dirk Groeneveld, Chandra Bhagavatula, Iz Beltagy, Miles Crawford, Doug Downey, Jason Dunkelberger, Ahmed Elgohary, Sergey Feldman, Vu Ha, et al. 2018. Construction of the literature graph in semantic scholar. In *NAACL*.

Fan Bai, Alan Ritter, and Wei Xu. 2021. Pre-train or annotate? domain adaptation with a constrained budget. In *EMNLP*, pages 5002–5015.

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. Scibert: A pretrained language model for scientific text. In *EMNLP*.

Andrew Eliot Borthwick. 1999. *A maximum entropy approach to named entity recognition*. New York University.

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. *arXiv preprint arXiv:2005.14165*.

Ivy Cao, Zizhou Liu, Giannis Karamanolakis, Daniel Hsu, and Luis Gravano. 2021. Quantifying the effects of COVID-19 on restaurant reviews. In *Proceedings of the International Workshop on Natural Language Processing for Social Media*, Online. Association for Computational Linguistics.

Dallas Card, Amber E. Boydstun, Justin H. Gross, Philip Resnik, and Noah A. Smith. 2015. The media frames corpus: Annotations of frames across issues. In *ACL*.

Alexandra Chronopoulou, Matthew E. Peters, and Jesse Dodge. 2021. Efficient hierarchical domain adaptation for pretrained language models. *ArXiv*, abs/2112.08786.

Hal Daumé III. 2007. Frustratingly easy domain adaptation. In *ACL*, pages 256–263.

Kushal Dave, Steve Lawrence, and David M. Pennock. 2003. Mining the peanut gallery: opinion extraction and semantic classification of product reviews. In *WWW*.

Cyprien de Masson d'Autume, Sebastian Ruder, Lingpeng Kong, and Dani Yogatama. 2019. Episodic memory in lifelong language learning. In *nips*.

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. *CoRR*, abs/2106.15110.

Jesse Dodge, Maarten Sap, Ana Marasovic, William Agnew, Gabriel Ilharco, Dirk Groeneveld, and Matt Gardner. 2021. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. In *EMNLP*.

Jacob Eisenstein, Brendan O'Connor, Noah A Smith, and Eric P Xing. 2014. Diffusion of lexical change in social media. *PloS one*, 9(11).

Robert M. Entman. 1983. Framing: Toward clarification of a fractured paradigm. *Journal of Communications*.

Hila Gonen, Ganesh Jawahar, Djamé Seddah, and Yoav Goldberg. 2020. Simple, interpretable and stable method for detecting words with usage change across corpora. In *ACL*.

Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies. In *NAACL*.

Suchin Gururangan, Mike Lewis, Ari Holtzman, Noah A. Smith, and Luke Zettlemoyer. 2021. Demix layers: Disentangling domains for modular language modeling. *CoRR*, abs/2108.05036.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *ACL*.

William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic word embeddings reveal statistical laws of semantic change. In *ACL*.

Xiaochuang Han and Jacob Eisenstein. 2019. Unsupervised domain adaptation of contextualized embeddings for sequence labeling. In *EMNLP*.

Yufan Huang, Yanzhe Zhang, Jiaao Chen, Xuezhi Wang, and Diyi Yang. 2021. Continual learning for text classification with information disentanglement based regularization. In *ACL*.

Mohit Iyyer, Peter Enns, Jordan Boyd-Graber, and Philip Resnik. 2014. Political ideology detection using recursive neural networks. In *ACL*.

Kokil Jaidka, Niyati Chhaya, and Lyle Ungar. 2018. Diachronic degradation of language models: Insights from social media. In *ACL*.

Jing Jiang and ChengXiang Zhai. 2007. Instance weighting for domain adaptation in nlp. In *ACL*.

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*.

William Labov. 2011. *Principles of linguistic change, volume 3: Cognitive and cultural factors*, volume 36. John Wiley & Sons.

Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d'Autume, Sebastian Ruder, Dani Yogatama, et al. 2021. Pitfalls of static language modelling. *arXiv preprint arXiv:2102.01951*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Proc. of Text Summarization Branches Out*.

Wei-Hao Lin, Eric P. Xing, and Alexander Hauptmann. 2008. A joint topic and perspective model for ideological discourse. In *ECML/PKDD*.

Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel S Weld. 2020. S2orc: The semantic scholar open research corpus. In *ACL*.

Jie Lu, Anjin Liu, Fan Dong, Feng Gu, João Gama, and Guangquan Zhang. 2019. Learning under concept drift: A review. *IEEE Transactions on Knowledge and Data Engineering*.

Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction. In *EMNLP*.

Afonso Mendes, Shashi Narayan, Sebastião Miranda, Zita Marinho, André F. T. Martins, and Shay B. Cohen. 2019. Jointly extracting and compressing documents with summary state representations. In *NAACL*.

Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? sentiment classification using machine learning techniques. In *EMNLP*.

Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *NAACL*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 21:1–67.

Shruti Rijhwani and Daniel Preoţiuc-Pietro. 2020. Temporally-informed analysis of named entity recognition. In *ACL*.

Paul Röttger and Janet Pierrehumbert. 2021. Temporal adaptation of bert and performance on downstream document classification: Insights from social media. In *Findings of EMNLP*, pages 2400–2412.

Maja R. Rudolph and David M. Blei. 2018. Dynamic embeddings for language evolution. *WWW*.

Tian Shi, Yaser Keneshloo, Naren Ramakrishnan, and Chandan K. Reddy. 2021. Neural abstractive text summarization with sequence-to-sequence models. *ACM Transactions on Data Science*.

Tian Shi, Ping Wang, and Chandan K. Reddy. 2019. LeafNATS: An open-source toolkit and live demo system for neural abstractive text summarization. In *NAACL*.

Hidetoshi Shimodaira. 2000. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*.

Ian Stewart and Jacob Eisenstein. 2018. Making "fetch" happen: The influence of social and linguistic context on nonstandard word growth and decline. In *EMNLP*.

Shivashankar Subramanian, Daniel King, Doug Downey, and Sergey Feldman. 2021. S2and: A benchmark and evaluation system for author name disambiguation. In *ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, pages 170–179. IEEE.

Nadine Tamburrini, Marco Cinnirella, Vincent AA Jansen, and John Bryden. 2015. Twitter users change word usage according to conversation-partner social identity. *Social Networks*, 40:84–89.

David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. In *EMNLP*.

Yunli Wang and Cyril Goutte. 2017. Detecting changes in twitter streams using temporal clusters of hashtags. In *Proceedings of the Events and Stories in the News Workshop*.

Geoffrey I. Webb, Loong Kuan Lee, Bart Goethals, and François Petitjean. 2018. Analyzing concept drift and shift from sample data. *Data Mining and Knowledge Discovery*.

Albert Xu, Eshaan Pathak, Eric Wallace, Suchin Gururangan, Maarten Sap, and Dan Klein. 2021. Detoxifying language models risks marginalizing minority voices. In *NAACL*.

Yi Yang and Jacob Eisenstein. 2016. Part-of-speech tagging for historical english. In *NAACL*.

Justine Zhang, Jonathan Chang, Cristian Danescu-Niculescu-Mizil, Lucas Dixon, Yiqing Hua, Dario Taraborelli, and Nithum Thain. 2018. Conversations gone awry: Detecting early signs of conversational failure. In *ACL*.

Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. 2013. Domain adaptation under target and conditional shift. In *ICML*.

Michael J.Q. Zhang and Eunsol Choi. 2021. SituatedQA: Incorporating extra-linguistic contexts into QA. *EMNLP*.

Yi Zhang, Zachary Ives, and Dan Roth. 2020. "who said it, and why?" provenance for natural language claims. In *ACL*.

Xuhui Zhou, Maarten Sap, Swabha Swayamdipta, Yejin Choi, and Noah A. Smith. 2021. Challenges in automated debiasing for toxic language detection. In *EACL*.

# Supplementary Material

## A   A Metric for Temporal Degradation

Let $t$ be the time period of the training data and $t'$ the time period of the evaluation data.[13] We aim to summarize the general effect of temporal misalignment (the difference between $t$ and $t'$) on task performance, in an interpretable way that is comparable across tasks.

Let $S_{t' \to t}$ indicate the performance a model trained on timestamp $t'$ data and evaluated on the timestamp $t$. Let

$$D(t' \to t) = - (S_{t' \to t} - S_{t \to t}) \times \text{sign}(t' - t),$$

In other words, $D(t' \to t)$ is a modified difference in performance between a aligned and misaligned models. The modification ensures that, as performance deteriorates, $D$ increases, regardless of the direction of time between $t$ and $t'$.

Our temporal degradation (TD) score for a fixed evaluation timestamp $t$ for models trained on a set of timestamps $\mathcal{T}$ is defined as:

$$\text{TD}(\mathcal{T} \to t') = \left| \frac{\sum_{t \in \mathcal{T}} \left( D(t' \to t) - \bar{D} \right) (t - \bar{t})}{\sum_{t \in \mathcal{T}} (t - \bar{t})^2} \right|,$$

where $\bar{t} = \text{avg}_{t \in \mathcal{T}} t'$ and $\bar{D} = \text{avg}_{t \in \mathcal{T}} D(t' \to t)$. This metric is the *slope* of a line fitting the the performance change of models trained on a variety of timestamps, when evaluated on a fixed timestamp. It can be interpreted as the average rate of performance deterioration per time period.

Fig. 7 shows three examples of TD scores from POLIAFF (the first) and YELPCLS (the latter two). These illustrate cases with and without temporal sensitivity. In practice, most examples with deterioration showed a linear trend and thus the rate of degradation was suitible to be approximated by a line. The final TD score is averaged over all evaluation years $\mathcal{T}'$.

$$\text{TD} = \frac{\sum_{t \in \mathcal{T}'} \text{TD}(\mathcal{T} \to t)}{|\mathcal{T}'|}$$

## B   Details of Model Development

**Training Details for Temporal Adaptation**   We train GPT2 over each domain and timestamp for $k$ steps using Huggingface's implementation of GPT2. Hyperparameter details can be seen in Table 4.

---

[13]See examples in Fig. 4.

| Hyperparameter | DAPT Assignment |
|---|---|
| Number of steps | 10k |
| Batch size | 32 |
| Maximum learning rate | 5e-05 |
| Adam Epsilon | 1e-08 |
| Adam Beta | 0.9. 0.999 |
| Block size | 1024 |

Table 4: Hyperparameters for temporal adaptation across the four domains.

| Hyperparameter | Cls. Assign | Summ. Assign |
|---|---|---|
| Number of Epochs | 50 | 10 |
| Batch size | 32 | 8 |
| Max learning rate | 2e-05 | 2e-05 |
| Adam Epsilon | 1e-08 | 1e-08 |
| Adam Beta | 0.9. 0.999 | 0.9. 0.999 |
| top p (sampling) | - | 0.05 |
| top k | - | 20 |
| temperature | - | 1 |
| max length | - | 512 |

Table 5: Hyperparameters for temporal finetuning across the eight tasks.

**Training Details for Temporal Finetuning**   We use Huggingface's implementation of GPT2 for finetuning for both the classification and summarization tasks. We train on Quadro RTX 800 GPUs. See Table 5 for details.

## C   Data Collection

We describe the postprocessing and data collection in greater detail. All data released is intended for non-commercial use.

**POLIAFF**   We acquire a list of U.S. politician names and Twitter handles.[14] One of the authors manually annotated whether each politician was a Republican or Democrat. In addition, one volunteer double checked to ensure correctness. We discard any politician who changed parties between 2015 and 2020, any independents, and anyone suspended by Twitter (e.g., RealDonaldTrump).

---

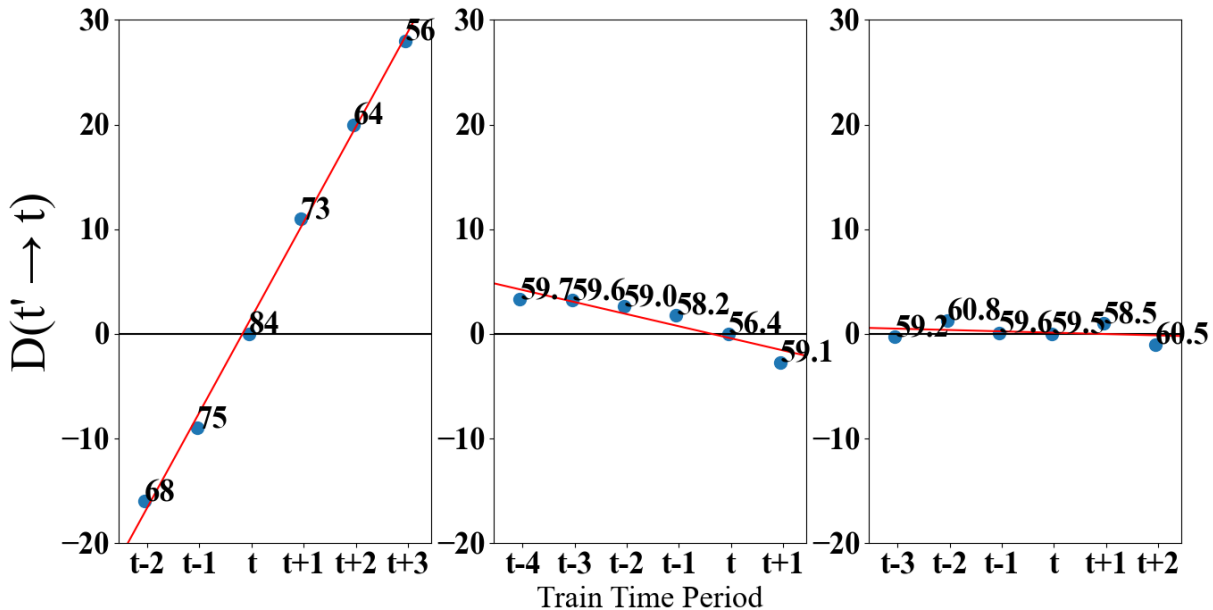[14]*https://files.pushshift.io/twitter/US_PoliticalTweets.tar.gz*

Figure 7: Three example calculations of the TD score (left from POLIAFF and the center and right from YELP-CLS). The annotated numbers are the raw evaluation scores $S_{t' \to t}$ and the plotted markers represent the modified differences $D(t' \to t)$ discussed in Section 2.3. For a particular plot, the red line is the line of best fit and its slope is the TD$(t)$ score for evaluation timestep $t$. The final TD score is averaged between all evaluation timesteps for the particular task.

**AIC** We randomly sample science documents in Semantic Scholar's corpus.[15] Of those, we only keep documents that (1) are published in ICML or AAAI, (2) are classified as 'computer science' documents, and (3) have an abstract of at least 50 tokens.

**Newsroom** The following applies to the postprocessing and data selection for both supervised temporal finetuning and unsupervised temporal adaptation of PUBCLS and NEWSUM. We use the Newsroom dataset.[16] We only keep articles where (1) the year in the metadata also appears in the main text and (2) no future year is mentioned in the main text.

**PUBCLS** We carry out additional postprocessing and ensure that each of the three labels (Fox News, New York Times, and Washington Post) have an equal distribution across years. We do so by uniform-random downsampling.

## D   Extended Results

We provide further results from our experiments described in Section 3.

| Domain | Task (metric) | Pearson's $r$ |
|--------|---------------|---------------|
| Twitter | POLIAFF ($F_1$) | 0.84 |
| | TWIERC ($F_1$) | 0.51 |
| Science | SCIERC ($F_1$) | 0.72 |
| | AIC ($F_1$) | 0.79 |
| News | PUBCLS ($F_1$) | 0.65 |
| | NEWSUM (Rouge-L) | 0.72 |
| | MFC ($F_1$) | 0.80 |
| Reviews | YELPCLS ($F_1$) | 0.14 |

Table 6: Pearson $r$ correlation coeffecients between the word overlap and performance of each task.

**Word Overlap Correlation with Performance** In addition to measuring vocabularies' change over time in Section 3.2, we find correlations between the word overlap and model performance of each task in Table 6.

**Finetuning Results** We provide the full results from our fientuning experiments in Section 3.1 in Fig. 8. These results are for downstream tasks with no domain adaptation.

**Finetuning with Temporal Domain Adaptation** We provide the full results from our finetuning with temporal domain adaptation in Section 3.2 in Fig. 7.

---

[15] *https://api.semanticscholar.org/corpus/*; licensed under an ODC-BY

[16] *https://lil.nlp.cornell.edu/newsroom/*

## PoliAff (F1)

| Train Year \ Eval Year | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|
| 20 | 91 | 77 | 65 | 56 | 57 | 48 |
| 19 | 81 | 83 | 72 | 62 | 57 | 49 |
| 18 | 68 | 75 | 84 | 73 | 64 | 56 |
| 17 | 61 | 66 | 77 | 79 | 69 | 64 |
| 16 | 52 | 58 | 69 | 73 | 80 | 72 |
| 15 | 46 | 53 | 65 | 70 | 76 | 78 |

## SciERC (F1)

| | 80-99 | 00-04 | 05-09 | 10-16 |
|---|---|---|---|---|
| 80-99 | 68 | 61 | 60 | 57 |
| 00-04 | 64 | 70 | 66 | 67 |
| 05-09 | 65 | 69 | 76 | 69 |
| 10-16 | 60 | 62 | 65 | 73 |

## NewSum (R-L)

| | 09-10 | 11-12 | 13-14 | 15-16 |
|---|---|---|---|---|
| 09-10 | 36 | 39 | 33 | 29 |
| 11-12 | 31 | 43 | 35 | 26 |
| 13-14 | 29 | 39 | 36 | 27 |
| 15-16 | 28 | 32 | 31 | 32 |

## YelpCls (F1)

| | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|
| 19 | 60 | 61 | 62 | 59 | 60 | 61 | 58 |
| 18 | 58 | 60 | 62 | 58 | 59 | 60 | 59 |
| 17 | 58 | 60 | 61 | 58 | 61 | 60 | 61 |
| 16 | 58 | 59 | 60 | 62 | 60 | 59 | 59 |
| 15 | 56 | 58 | 59 | 57 | 60 | 58 | 59 |
| 14 | 55 | 58 | 56 | 56 | 58 | 56 | 59 |
| 13 | 58 | 60 | 60 | 58 | 60 | 59 | 60 |

## TwiERC (F1)

| | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|
| 19 | 76 | 77 | 76 | 72 | 69 | 69 |
| 18 | 72 | 74 | 77 | 72 | 69 | 68 |
| 17 | 72 | 74 | 78 | 71 | 69 | 69 |
| 16 | 74 | 77 | 79 | 76 | 73 | 71 |
| 15 | 72 | 76 | 79 | 71 | 74 | 73 |
| 14 | 71 | 72 | 77 | 71 | 72 | 73 |

## AIC (F1)

| | 09-11 | 12-14 | 15-17 | 18-20 |
|---|---|---|---|---|
| 09-11 | 86 | 79 | 71 | 66 |
| 12-14 | 83 | 86 | 74 | 63 |
| 15-17 | 82 | 85 | 83 | 84 |
| 18-20 | 72 | 79 | 78 | 85 |

## MFC (F1)

| | 09-10 | 11-12 | 13-14 | 15-16 |
|---|---|---|---|---|
| 09-10 | 27 | 25 | 25 | 26 |
| 11-12 | 24 | 28 | 24 | 27 |
| 13-14 | 22 | 24 | 26 | 26 |
| 15-16 | 24 | 26 | 25 | 33 |

## PubCls (F1)

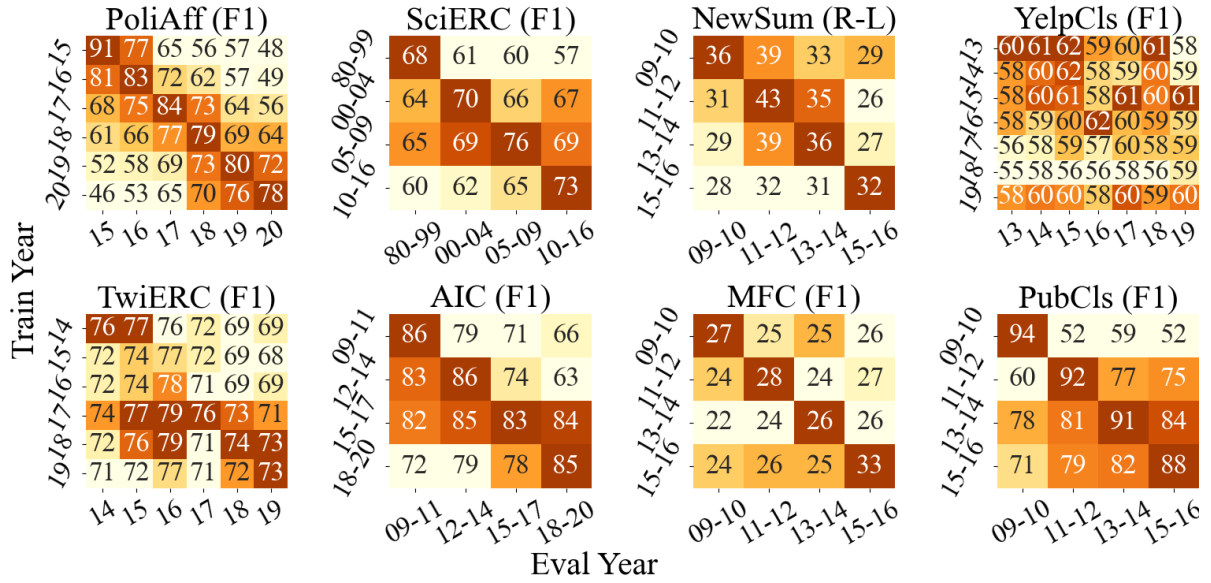| | 09-10 | 11-12 | 13-14 | 15-16 |
|---|---|---|---|---|
| 09-10 | 94 | 52 | 59 | 52 |
| 11-12 | 60 | 92 | 77 | 75 |
| 13-14 | 78 | 81 | 91 | 84 |
| 15-16 | 71 | 79 | 82 | 88 |

Figure 8: Temporal misalignment in finetuning affects task performance (§3.1). In all cases, higher scores are better. The heatmap is shaded per column, i.e., the darkest shade of orange in a cell means the cell has the highest score in that column. Mismatch between the the training and evaluation data can result in massive performance drop; degree varies by task. For example, YELPCLS, MFC, and TWIERC show minimal degradation. In contrast, POLIAFF and NEWSUM major deterioration over time.

| Domain (Task) | Finetune Year | Evaluation → Pretrain ↓ | 2015 | 2020 | Domain (Task) | Finetune Year | Evaluation → Pretrain ↓ | 2014 | 2019 |
|---|---|---|---|---|---|---|---|---|---|
| Twitter (PoliAff) *F1* | 2015 | Default | 91.4 | 48.4 | Twitter (TwiERC) *F1* | 2014 | Default | 74.3 | 68.9 |
| | | Default → 2015 | 92.2 | 47.5 | | | Default → 2014 | 76.1 | 69.6 |
| | | Default → 2020 | 90.9 | 50.8 | | | Default → 2019 | 74.1 | 68.9 |
| | 2020 | Default | 45.8 | 78.0 | | 2019 | Default | 71.0 | 74.6 |
| | | Default → 2015 | 47.2 | 76.9 | | | Default → 2014 | 73.1 | 75.2 |
| | | Default → 2020 | 44.2 | 78.3 | | | Default → 2019 | 73.7 | 75.8 |

| Domain (Task) | Finetune Year | Evaluation → Pretrain ↓ | 2009-11 | 2018-20 | Domain (Task) | Finetune Year | Evaluation → Pretrain ↓ | 1980-1999 | 2010-2016 |
|---|---|---|---|---|---|---|---|---|---|
| Scienctific (AIC) *F1* | 2009-2011 | Default | 79.0 | 72.0 | Scientific (SciERC) *F1* | 1980-1999 | Default | 67.9 | 57.2 |
| | | Default → 2009-2011 | 94.5 | 68.8 | | | Default → 1980-1999 | 73.2 | 66.4 |
| | | Default → 2018-2020 | 88.4 | 86.0 | | | Default → 2010-2016 | 73.7 | 66.8 |
| | 2018-2020 | Default | 72.0 | 85.0 | | 2010-2016 | Default | 60.3 | 72.5 |
| | | Default → 2009-2011 | 87.2 | 65.2 | | | Default → 1980-1999 | 63.4 | 75.0 |
| | | Default → 2018-2020 | 86.8 | 79.4 | | | Default → 2010-2016 | 64.8 | 76.0 |

| Domain (Task) | Finetune Year | Evaluation → Pretrain ↓ | 2009-2010 | 2015-2016 | Domain (Task) | Finetune Year | Evaluation → Pretrain ↓ | 2009-2010 | 2015-2016 |
|---|---|---|---|---|---|---|---|---|---|
| News (MFC) *F1* | 2009-2010 | Default | 27.0 | 26.0 | News (PubCls) *F1* | 2009-2010 | Default | 94.1 | 52.4 |
| | | Default → 2009-2010 | 30.6 | 31.8 | | | Default → 2009-2010 | 95.4 | 54.0 |
| | | Default → 2015-2016 | 29.8 | 30.0 | | | Default → 2015-2016 | 95.4 | 53.5 |
| | 2015-2016 | Default | 23.8 | 33.4 | | 2015-2016 | Default | 71.3 | 88.2 |
| | | Default → 2009-2010 | 29.7 | 41.6 | | | Default → 2009-2010 | 80.4 | 90.7 |
| | | Default → 2015-2016 | 32.7 | 41.9 | | | Default → 2015-2016 | 78.7 | 91.1 |

| Domain (Task) | Finetune Year | Evaluation → Pretrain ↓ | 2009-2010 | 2015-2016 | Domain (Task) ↓ | Finetune Year ↓ | Evaluation → Pretrain ↓ | 2014 | 2019 |
|---|---|---|---|---|---|---|---|---|---|
| News (NewSum) *Rouge-L* | 2009-2010 | Default | 36.4 | 29.0 | Food Reviews (Yelp) *F1* | 2013 | Default | 58.6 | 58.3 |
| | | Default → 2009-2010 | 36.4 | 29.1 | | | Default → 2013 | 63.3 | 60.1 |
| | | Default → 2015-2016 | 36.1 | 28.9 | | | Default → 2019 | 60.2 | 62.3 |
| | 2015-2016 | Default | 27.8 | 31.8 | | 2019 | Default | 58.3 | 58.3 |
| | | Default → 2009-2010 | 28.2 | 31.8 | | | Default → 2013 | 60.2 | 62.3 |
| | | Default → 2015-2016 | 27.8 | 31.6 | | | Default → 2019 | 60.8 | 62.3 |

Table 7: Combination of temporal adaptation and finetuning (§3.2) on our tasks. The row labeled "Default" corresponds to a model that has not been adapted (uses the default pretraining). The color coding is proportional to the magnitude of the performances of each task (darker shade of orange indicates higher scores). We see that models that were finetuned on similar time periods performed similarly, no matter how their DAPT conditions differed.