

Temporal Validity Change Prediction

Anonymous ACL submission

Abstract

Temporal validity is an important property of text that is useful for many downstream applications, such as recommender systems, conversational AI, or story understanding. Existing benchmarking tasks often require models to identify the temporal validity duration of a single statement. However, in many cases, additional contextual information, such as sentences in a story or posts on a social media profile, can be collected from the available text stream. This contextual information may greatly alter the duration for which a statement is expected to be valid. We propose *Temporal Validity Change Prediction*, a natural language processing task benchmarking the capability of machine learning models to detect contextual statements that induce such change. We create a dataset consisting of temporal target statements sourced from Twitter and crowdsource sample context statements. We then benchmark a set of transformer-based language models on our dataset. Finally, we experiment with temporal validity duration prediction as an auxiliary task to improve the performance of the state-of-the-art model.

1 Introduction

In human communication, temporal properties are frequently underspecified when authors assume that the recipient can infer them via commonsense reasoning. For example, when reading “I am moving on Saturday”, a reader is likely to assume the person will be busy for most of the day. On the other hand, when reading “I will make a sandwich on Sunday”, this is likely to only take up a fraction of the author’s day and may not impact other plans. Such reasoning is referred to as *temporal commonsense* (TCS) reasoning (Wenzel and Jatowt, 2023).

Temporal validity (Almquist and Jatowt, 2019; Hosokawa et al., 2023; Lynden et al., 2023) is a property that is vital for the proper understanding of a text. The temporal validity of a statement, i.e.,

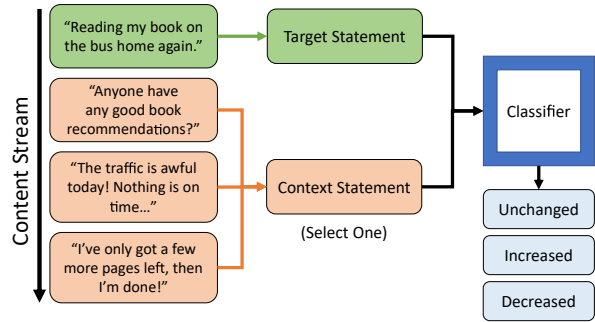


Figure 1: A visualization of the TVCP task

whether it contains valid information at a given time, often requires TCS reasoning to resolve. For example, in determining whether a statement like “I am driving home from work” is still valid after five hours, we may use our prior understanding of the typical duration of related events, such as commuter traffic. While the amount of research into TCS and, to a degree, temporal validity, has risen over the past years (Wenzel and Jatowt, 2023), there are still several properties of temporal validity that have not been considered in previous research. One such property is the impact of *context* on the temporal validity duration of a statement. For example, the sentence “I am driving home from work” may be valid for a longer time when followed by a statement such as “There is a massive traffic jam”.

To model this problem, we propose a new NLP task format called *Temporal Validity Change Prediction* (TVCP), which requires reasoning over whether a context statement changes the temporal validity duration of a target statement. The task is visualized in Figure 1. We propose the following applications for such a system.

Timeline Prioritization: Social media services such as Twitter rely on recommender systems to prioritize the vast amount of content that their users produce. One possible way to improve the prioritization of content is to consider its temporal validity (Takemura and Tajima, 2012; Koul et al., 2022), as

071 users are likely to care about current and relevant in- 119
072 formation over more general, stationary statements. 120
073 TVCP can be used to leverage the stream of social 121
074 media posts by a given user as possible context to 122
075 better estimate the temporal validity duration of a 123
076 previously observed post. 124

077 **User Status Tracking:** Similarly, the content 125
078 of a user’s posts on social media could be utilized 126
079 for other analytical or business purposes, such as 127
080 predicting revenue streams (Asur and Huberman, 128
081 2010; Deng et al., 2011; Lassen et al., 2014; Lu 129
082 et al., 2014) or identifying trends in a community’s 130
083 or an individual user’s behaviour (Li et al., 2018; 131
084 Abe et al., 2018; Shen et al., 2020). TVCP could 132
085 be used to identify posts that refer to previous tem- 133
086 poral information, to detect chains of thought about 134
087 topics that may not be self-contained. 135

088 **Conversational AI:** Foundation models, such 136
089 as CHATGPT (Ouyang et al., 2022) and BARD 137
090 (Manyika, 2023), could use the temporal validity 138
091 of statements provided by the user to keep track 139
092 of knowledge that is still relevant to the conver- 140
093 sation. Using TVCP, new messages by the user 141
094 could be evaluated to adjust the expected temporal 142
095 validity period of previously learned facts. This is 143
096 especially relevant as initial reports indicate that 144
097 foundation models may struggle with TCS reason- 145
098 ing (Bian et al., 2023). 146

099 Our main contributions are the following:

- 100 1. We define a novel NLP task (TVCP). This 147
101 ternary classification task requires models to 148
102 predict the impact of a context statement on a 149
103 target statement’s temporal validity duration. 150
- 104 2. We build a dataset of tuples consisting of time- 151
105 sensitive *target statements*, as well as *follow-* 152
106 *up statements* that act as context for our task. 153
- 107 3. We evaluate the performance of existing pre- 154
108 trained *language models* (LMs) on our dataset, 155
109 including models fine-tuned on other TCS 156
110 tasks as well as CHATGPT. 157
- 111 4. We propose an augmentation to the training 158
112 process that leverages temporal validity dura- 159
113 tion information to help improve the perfor- 160
114 mance of the state-of-the-art classifier. 161

115 2 Related Work 162

116 2.1 Temporal Commonsense Reasoning 163

117 TCS reasoning is often considered one of several 164
118 categories of commonsense reasoning (Storks et al., 165

119 2019a; Bhargava and Ng, 2022). A major driver 120
of research specifically into TCS appears to have 121
been the transformer architecture (Vaswani et al., 122
2017) and resulting LMs. In recent years, several 123
datasets that specifically aim to benchmark TCS 124
understanding have been published (Zhou et al., 125
2019; Ning et al., 2020; Zhang et al., 2020; Qin 126
et al., 2021; Zhou et al., 2021), while ROCSTO- 127
RIES (Mostafazadeh et al., 2016) appears to be 128
the only dataset focussing on this type of reason- 129
ing before the publication of the transformer archi- 130
tecture. Small adjustments to transformer-based 131
LMs are often proposed as state-of-the-art solu- 132
tions for these datasets (Pereira et al., 2020; Yang 133
et al., 2020; Zhou et al., 2020; Pereira et al., 2021; 134
Kimura et al., 2021; Zhou et al., 2021, 2022; Cai 135
et al., 2022; Yu et al., 2022). Similarly, temporal- 136
ized transformer models are popular solutions for 137
tasks such as document dating or semantic change 138
detection (Rosin and Radinsky, 2022; Rosin et al., 139
2022; Wang et al., 2023).

140 The TCS taxonomy defined by Zhou et al. (2019) 141
is frequently referenced. It contains the five di- 142
mensions of *duration* (how long an event takes), 143
temporal ordering (typical order of events), *typi-* 144
cal time (when an event happens), *frequency* (how 145
often an event occurs) and *stationarity* (whether a 146
state holds for a very long time or indefinitely).

147 2.2 Temporal Validity 148

149 Compared to TCS reasoning, temporal validity in 150
text is a less well-researched field. It effectively 151
combines three dimensions of the taxonomy by 152
Zhou et al. (2019): *Stationarity*, to reason about 153
whether a statement contains temporal information, 154
typical time, to reason about when the temporal 155
information occurs, and *duration*, to reason about 156
how long the temporal information takes to resolve.

157 Takemura and Tajima (2012) classify the lifetime 158
duration of tweets, i.e., the informational value of 159
a tweet over time. They use handcrafted, domain- 160
specific features to train a *support vector classifier* 161
(SVC) on supervised samples. 162

163 Almquist and Jatowt (2019) similarly design fea- 164
tures to estimate the temporal validity duration of 165
sentences collected from news, blog posts, and 166
Wikipedia using SVCs. Their features contain gen- 167
eral properties such as the word- or sentence length, 168
but also more complex ones, such as latent seman- 169
tic analysis. 170

Method	Task	Data Source	Duration Bias	Model	# Samples
Takemura and Tajima (2012)	TV _d	Twitter	N/A	SVC	9,890
Almquist and Jatowt (2019)	TV _d	Blogs, News, Wikipedia	years	SVC	1,762
Hosokawa et al. (2023)	TNLI	Image Captions	seconds ¹	LM	10,659
Lynden et al. (2023)	TV _d	WikiHow	hours	LM	339,184
Ours	TVCP	Twitter	hours	LM	5,055

Table 1: Summary of related work

Hosokawa et al. (2023) define the *Temporal Natural Language Inference* (TNLI) task. The goal of TNLI is to determine whether the temporal validity of a given hypothesis sentence is supported by a premise sentence.

Lynden et al. (2023) build a large dataset of human annotations specifying the duration required to perform various actions on WikiHow as well as their respective temporal validity durations.

2.3 Comparison with Related Work

Table 1 shows the most closely related research. As noted, our dataset is based on the proposed TVCP task, whereas previous work was based on the TV_d and TNLI tasks. All tasks are described in more detail in Section 3.

Another prominent distinctive attribute is the text source and the resulting temporal validity duration bias. For example, sentences sourced from news or Wikipedia articles often appear to be valid for years or longer. On the other hand, image captions may only be valid for a few seconds or minutes. We decided to source our sentences from Twitter due to its alignment with our downstream use cases. Similar to Lynden et al. (2023), our collected temporal information tends to be valid for a few hours.

We follow recent research by evaluating our dataset using transformer-based LMs, whereas earlier approaches relied on methods such as SVCs.

Except for the CoTAK dataset (Lynden et al., 2023), the datasets tend to be relatively small. As crowdsourcing is used in all datasets referenced in Table 1 to annotate text spans with common-sense information, the costs of dataset creation can quickly escalate. In addition, we ask participants to create examples of follow-up statements that cause temporal validity change. This approach further restricts the overall size of our dataset due to the relative difficulty of the task.

¹Based on analysis of a sample. TV_d labels are not available for the full dataset.

3 Task

3.1 Defining Temporal Validity

Temporal validity, in essence, is simply the time-dependent validity of a text. As shown in Equation 1, the temporal validity of a statement s at a time t is a binary value that determines whether the information in s is valid at the given time.

$$\text{TV}(s, t) = \begin{cases} \text{True} & \text{if information in } s \text{ is valid at } t, \\ \text{False} & \text{otherwise} \end{cases} \quad (1)$$

In some previous research (Hosokawa et al., 2023; Lynden et al., 2023), the scope of evaluated temporal information is limited to actions, such as “I am *baking bread*”. However, we note that other types of temporal information exist, such as events (e.g., in the sentence “*Job interview tomorrow*”) or temporary states (e.g., in the sentence “*It is nice out today*”). In an analysis of a subset of our collected statements, shown in Figure 2, we find that these alternative types of temporal information constitute a significant portion (28%) of samples. Additionally, one-third of sampled statements contained at least two distinct pieces of temporal information with differing temporal validity spans. This indicates that the true scope of determining the temporal validity of a text may exceed what current datasets are benchmarking.

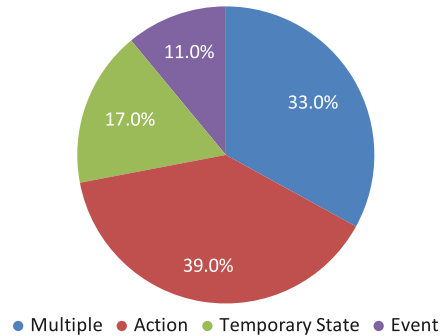


Figure 2: Distribution of different types of temporal information in a sample of our dataset

We assume that the temporal validity of stationary statements is constant for any given timestamp t . A stationary statement may be continuously true (e.g., “Japan lies in Asia”), or continuously false (e.g., “Japan lies in Europe”). This includes information that is fully contained in the past (e.g., “I went to the bank yesterday”). In general, we do not expect the validity of such a statement to change.

For contemporary or future information, we assume the statement is valid from the moment of sentence conception until the information is no longer ongoing. We include the duration of the information, rather than just its occurrence time, as humans are likely to still consider durative information relevant while it is ongoing. For example, we may reason that the statement “I will take a shower at 8 p.m.” still has informational value at 8:05 p.m., as it allows us to infer the current action of the author.

3.2 Formalizing Existing Tasks

3.2.1 Temporal Validity Duration Estimation

Temporal Validity Duration Estimation (TV_d) is the primary task that is evaluated in temporal validity research (Takemura and Tajima, 2012; Almquist and Jatowt, 2019; Lynden et al., 2023). The goal is to estimate the duration for which the statement is valid, starting at the statement creation time. We formalize this task in Equation 2, where t_s is the timestamp at which the statement s is created.

$$TV_d(s) = \max_{t \geq t_s} \{t \mid TV(s, t) = \text{True}\} \quad (2)$$

The TV_d task is useful in downstream applications such as social media, where information on the posting time of a statement is readily available and can be used to infer the span during which the statement is valid.

3.2.2 Temporal Natural Language Inference

The goal of TNLI (Hosokawa et al., 2023) is to infer whether a statement is temporally valid, given additional context, using typical NLI terminology (MacCartney, 2009; Storks et al., 2019b). TNLI requires a *hypothesis statement* (that we call *target statement*, or s_t) and a *premise sentence* (that we call *follow-up statement*, or s_f). Implicitly, the inference takes place at t_{s_f} , that is, the posting time of the follow-up statement, but no explicit duration information is required to solve this task. Formally, we define TNLI in Equation 3 (SUO = supported, INV = invalidated, UNK = unknown),

where $TV^c(s, t)$ is the temporal validity of a statement s at a time t given context c . The UNK class is assigned in cases where $TV^{s_f}(s_t, t_{s_f})$ is neither clearly supported nor invalidated by the context.

$$TNLI(s_t, s_f) = \begin{cases} \text{SUO} & TV^{s_f}(s_t, t_{s_f}) = \text{True} \\ \text{INV} & TV^{s_f}(s_t, t_{s_f}) = \text{False} \\ \text{UNK} & TV^{s_f}(s_t, t_{s_f}) = \text{Unclear} \end{cases} \quad (3)$$

Unlike TV_d , this task format lends itself to downstream applications such as story understanding, wherein a larger text stream of individual statements is provided with no clear explicit notion of time passing between each sentence (e.g., in a book).

3.3 Temporal Validity Change Prediction

We propose *Temporal Validity Change Prediction* (TVCP), which combines ideas from both the inference- and duration-based tasks. Like TNLI, we require s_t and s_f for classification, and determine a ternary label that provides information about the impact of s_f on s_t . Unlike TNLI, our goal is to predict a *change* in the temporal validity *duration* of s_t .

We consider TVCP a necessary step in accurately determining a statement’s temporal validity. Simply estimating the duration of the statement alone may not yield very precise results when it is, as in many use cases, extracted from a rich context, such as a book, a story, a news article, a step-by-step guide, or a social media profile. In these cases, surrounding information may provide additional context that could lead us to a different TV_d estimate. Simply concatenating s_t and s_f may lead to the classification of temporal information within s_f , which is undesired. Our segmentation of s_t and s_f into different semantic roles, similar to TNLI, prevents this issue.

Formally, we define TVCP in Equation 4 (DEC = decreased, UNC = unchanged, INC = increased), where $TV_d^c(s)$ is the temporal validity duration of a statement s given context c . Figure 3 shows a concrete comparative example of the goal of all three tasks.

$$TVCP(s_t, s_f) = \begin{cases} \text{DEC} & TV_d(s_t) > TV_d^{s_f}(s_t) \\ \text{UNC} & TV_d(s_t) = TV_d^{s_f}(s_t) \\ \text{INC} & TV_d(s_t) < TV_d^{s_f}(s_t) \end{cases} \quad (4)$$

Since TVCP is a signal measuring the difference between TV_d with- and without s_f , respectively, a

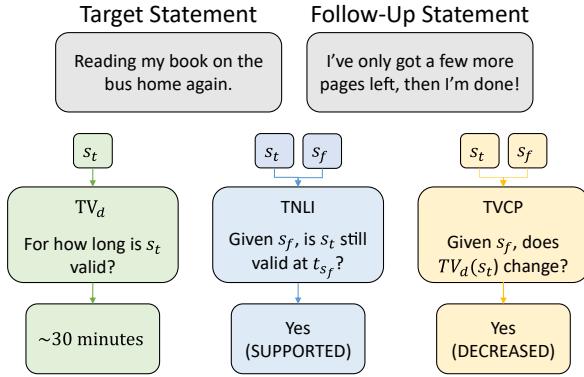


Figure 3: An example of TV_d , TNLI and TVCP

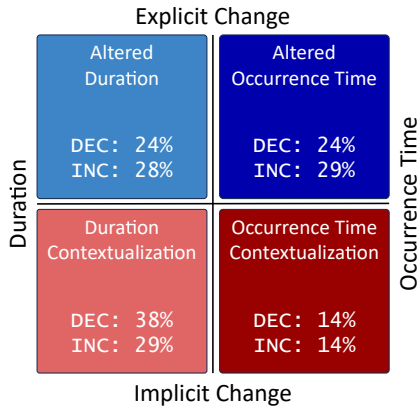


Figure 4: Dimensions of temporal validity change. The frequency of each category for DEC and INC classes in our sample is appended.

more fine-grained TV_d classification increases the number of TVCP instances that can be detected. On the other hand, evaluating TV_d on a very fine-grained scale may be more difficult for both models and humans (Honda et al., 2022), and the resulting uncertainty and inaccuracies could lead to a degradation of the system as a whole.

In our sample analysis, we find that temporal validity change generally occurs along two axes, shown in Figure 4. The first dimension is *implicit* versus *explicit* change. For example, an appointment may be postponed, which is an explicit change. On the other hand, the author may note in a follow-up statement that the appointment is in a sleep laboratory, which causes us to re-evaluate for how long the original statement is valid, although the information itself has not changed.

The second dimension is a change to the *occurrence time* versus the *duration* of the information. For example, a flight may be delayed, in which case the occurrence time changes. Alternatively, the flight might have to be re-routed mid-air due to bad weather, in which case the duration changes.

In our sample, we find that all four categories are present to a reasonable degree in both the DEC and INC classes. Generally, changes to the duration tend to be slightly more frequent than changes to the occurrence time. This makes sense, as the occurrence time is a dimension that is only present when the information occurs in the future, whereas the duration of temporal information can change irrespective of the occurrence time.

4 Dataset

We create a dataset for training and benchmarking TVCP, where each sample is a quintuple $\langle s_t, s_f, TV_d(s_t), TV_d^{s_f}(s_t), TVCP(s_t, s_f) \rangle$.

s_t is collected by querying the Twitter API for tweets with no external context (e.g., no tweets that are retweets or replies, or tweets containing media). We apply several pre-processing steps to remove tweets whose content may not be self-contained. We aim to minimize spam and offensive content by applying publicly available LMs and word-list-based filters. To decrease the number of stationary statements, we employ an ensemble classifier based on the ALMQUIST2019 (Almquist and Jatowt, 2019) and COTAK datasets and select the most likely statements to contain temporal information. Finally, crowdworkers can tag any remaining stationary samples during the annotation process. A summary of our pre-processing pipeline is shown in Figure 5. Our code, including all preprocessing steps, is published under the Apache 2.0 licence.

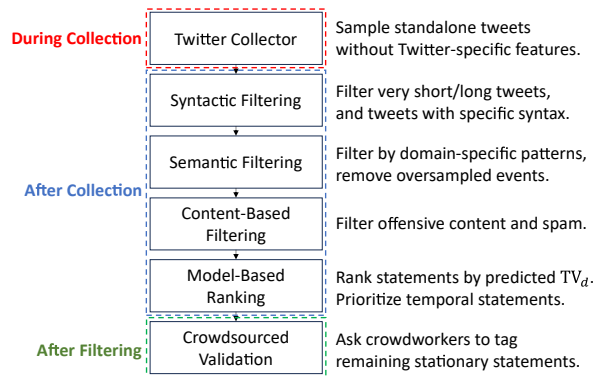


Figure 5: A summary of our tweet collection pipeline

For each target statement, we ask two crowdworkers to estimate $TV_d(s_t)$ from the logarithmic class design shown in Equation 5, which is modelled after human timeline understanding (Jatowt and Au Yeung, 2011; Varshney and Sun, 2013; Howard, 2018). If the annotators disagreed, we supplied a third vote. We discarded any tweets that

were annotated as *less than one minute*, *more than one month*, or *no time-sensitive information* (i.e., stationary), as well as tweets where no majority agreement could be reached. Of 2,996 annotated target tweets, 571 were discarded without a third annotation, 867 were added without a third annotation, 546 were discarded after providing a third vote, and 1,012 were added after providing a third vote. The distribution of resulting TV_d labels before temporal validity change is shown in Figure 6.

$$t \in \{< 1 \text{ minute}, 1\text{-}5 \text{ minutes}, 5\text{-}15 \text{ minutes}, 15\text{-}45 \text{ minutes}, 45 \text{ minutes}\text{-}2 \text{ hours}, 2\text{-}6 \text{ hours}, \text{more than } 6 \text{ hours}, 1\text{-}3 \text{ days}, 3\text{-}7 \text{ days}, 1\text{-}4 \text{ weeks}, \text{more than } 1 \text{ month}\} \quad (5)$$

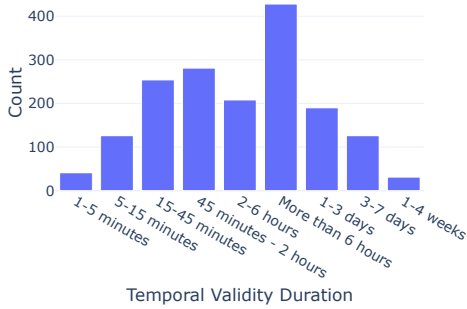


Figure 6: Distribution of TV_d labels (before temporal validity change) in our dataset

Both s_f and $TV_d^{sf}(s_t)$ were provided by a separate set of crowdworkers, given s_t and $TV_d(s_t)$ as an input. In total, we collected 5,055 samples from 1,685 target statements. In Figure 7, we plot the *temporal validity change delta*, which is the class distance between the original and the updated TV_d estimate. We find that, in most cases, the temporal validity duration of a target statement is shifted only by one class.

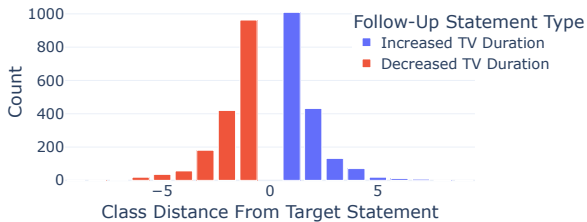


Figure 7: Temporal validity change delta distribution

All crowdsourcing tasks were set up on Amazon Mechanical Turk, using qualification tests, partic-

ipation criteria, and manual verification of results to ensure high-quality samples (see Appendix A). We publish the resulting dataset for public use under the CC BY 4.0 licence. In accordance with the Twitter developer policy², we only publish the Tweet IDs of sourced statements. This also means original tweet authors retain the ability to delete their content, effectively removing it from the dataset.

5 Experiments

5.1 Language Models

We evaluate a set of transformer-based LMs on our dataset. We test four different archetypes in total:

- **TRANSFORMERCLASSIFIER**: Builds a hidden representation from the sentence-embedding token of the concatenation of s_t and s_f .
- **SIAMESECLASSIFIER**: Builds a hidden representation from the concatenated embeddings $[h_{s_t}, h_{s_f}, h_{s_t} - h_{s_f}, h_{s_t} \otimes h_{s_f}]$, where h_{s_t} and h_{s_f} are the sentence-embedding tokens of the target- and follow-up statement, respectively (Bromley et al., 1993; Nandy et al., 2020).
- **SELFEXPLAIN** (Sun et al., 2020): Builds a hidden representation from the embeddings of spans between arbitrary tokens in either s_t or s_f , selected by the model.
- **CHATGPT**: A chain-of-thought (Wei et al., 2022) reasoning prompt based on few-shot learning (one sample per TVCP class), passed to the gpt-3.5-turbo model via the OpenAI API.³

For the TRANSFORMERCLASSIFIER and SIAMESECLASSIFIER pipelines, we evaluate BERT-BASE-UNCASED (Kenton and Toutanova, 2019; 110M parameters) and ROBERTA-BASE (Liu et al., 2019; 125M parameters) embeddings. For SELFEXPLAIN, we only test the original implementation with ROBERTA-BASE embeddings. To evaluate transfer learning from other TCS tasks, we test the TRANSFORMERCLASSIFIER pipeline on regular BERT-BASE-UNCASED pre-training weights as well as two variants TACOLM (Zhou et al., 2020) and COTAK (Lynden et al., 2023),

²<https://developer.twitter.com/en/developer-terms/policy>, accessed 12.10.2023

³This call uses the most recent GPT3.5 model. We collected CHATGPT responses in July 2023.

which have the same underlying architecture, but use weights fine-tuned on existing TCS datasets.

We use the ADAMW optimizer (Loshchilov and Hutter, 2018) with $\varepsilon = 1e-8$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\text{weight_decay} = 0.01$. We optimize for cross-entropy loss. SELFEXPLAIN adds an additional loss parameter in the form of squared span-attention weights, to encourage the model to more sharply choose which spans should be used to build the hidden representation.

We set the dropout probabilities and learning rates as defined in Table 2 as a result of our hyperparameter optimization (see Appendix C). For all models, the hidden embedding size is 768. For some ROBERTA-based models, we freeze embedding layers (i.e., only fine-tune intermediate and classification weights), as training all parameters leads to poor performance.

Model	Dropout	LR	Frozen
TF - BERT	0.25	1e-4	False
S - BERT	0.25	1e-4	False
TF - ROBERTA	0.25	1e-3	True
S - ROBERTA	0.10	1e-4	True
SELFEXPLAIN	0.00	2e-5	False

Table 2: Hyperparameter settings for different models. TF = TRANSFORMERCLASSIFIER, S = SIAMESECLASSIFIER

5.2 Multitask Implementation

For all archetypes except CHATGPT, we provide a second implementation, in which we add two regression layers that aim to respectively predict $\text{TV}_d(s_t)$ and $\text{TV}_d^{sf}(s_t)$ from the same hidden representation. For these layers, we calculate the mean squared error between a single output neuron and a linear mapping of the TV_d class index to the range $[0, 1]$. Our intuition is that embeddings with an understanding of TV_d may be better suited for TVCP. Inspiration for this approach are models that utilize the interplay between temporal dimensions to improve the TCS reasoning performance in LMs, such as SYMTIME (Zhou et al., 2021) or SLEER (Cai et al., 2022).

5.3 Evaluation

We evaluate two metrics, accuracy and *exact match* (EM). Accuracy is simply the fraction of correctly classified samples. EM is the fraction of *target statements* for which all three samples were correctly classified. This metric punishes inconsistency in the model more strictly, thus providing a

better view of the true performance and task understanding of each model (Wenzel and Jatowt, 2023), while disincentivizing shallow reasoning behaviours commonly seen in transformer models (Helwe et al., 2021; Tan et al., 2023).

We report the mean EM and accuracy across a five-fold cross-validation split. Each evaluation consists of 70% training data, 10% validation data, and 20% test data. If the validation EM does not exceed the best previously observed value for 5 consecutive epochs, we stop training. The model with the best validation EM is used for evaluation on the test set. The results are shown in Table 3.

Model	Acc (+ MT)	EM (+ MT)
TF - ROBERTA	64.0 (+1.5)	21.2 (+2.5)
CHATGPT	66.3 (N/A)	29.3 (N/A)
S - ROBERTA	78.7 (+1.1)	48.2 (+2.1)
TF - CoTAK	83.2 (+0.6)	58.2 (+1.4)
S - BERT	83.8 (-0.3)	59.1 (-1.5)
TF - TACOLM	83.5 (+1.4)	59.1 (+2.9)
TF - BERT	84.8 (-0.2)	61.2 (+0.9)
SELFEXPLAIN	88.5 (+1.1)	69.8 (+2.8)

Table 3: Model evaluation results, sorted by mean EM score. TF = TRANSFORMERCLASSIFIER, S = SIAMESECLASSIFIER, MT = Multitask Implementation

We note a positive impact on the EM score from implementing multitasking in all models except the Siamese architecture with BERT-based embeddings. We use bootstrapping to calculate the statistical significance of implementing multitask learning on the best-performing model, SELFEXPLAIN, evaluating the number of bootstrap samples in which the multitask implementation outperforms regular SELFEXPLAIN. We find $p = 0.0027$ for accuracy, with a 95% confidence interval of $[0.0036, 0.0192]$. For EM, $p = 0.0025$, with a 95% confidence interval of $[0.0089, 0.0487]$.

The use of weights from other TCS tasks does not seem to have a positive impact on the performance of the TRANSFORMERCLASSIFIER pipeline. It is possible that, although the resulting embeddings are more aligned with temporal properties (Zhou et al., 2020), other important information in the embeddings is lost, leading to an overall decreased performance.

Due to some ROBERTA-based models having frozen embedding layers, the baseline performance by ROBERTA is much worse, but it improves much more when switching to the SIAMESECLASSIFIER implementation. We hypothesize that ROBERTA’s sentence embedding token, $\langle s \rangle$, may contain less

information about the full sequence than BERT’s [SEP] token due to the lack of a next-sentence-prediction task during pre-training.

CHATGPT ranks among the lower-performing models, which is consistent with other studies on TCS understanding (Bian et al., 2023). Its shortcomings may be due to the few-shot learning approach and a lack of knowledge about dataset specifics traits, which a trained classifier could leverage. Additionally, we do not specify our class design in the CHATGPT prompt, which could make it harder for CHATGPT to isolate the UNC class.

To evaluate the impact of training data quantity on classifier performance, we train our best-performing classifier (SELFEXPLAIN with multi-tasking, which we dub MULTITASK) on a single train-val-test split (80%/10%/10%) with different amounts of training data. The results can be seen in Figure 8. We find that performance increases as more data is used for training, but this effect starts to diminish as we approach 100% of our training data.



Figure 8: Training data vs. performance metrics in MULTITASK

In testing SELFEXPLAIN and MULTITASK on various temporal validity change deltas (Figure 9), we find they perform comparably on the UNC class, but MULTITASK slightly outperforms SELFEXPLAIN on all delta values greater than zero. While CHATGPT’s subpar performance in the UNC class can partially be attributed to prompt design, it continues to lag far behind other models in the DEC and INC classes.



Figure 9: Temporal validity change delta vs. accuracy in MULTITASK, SELFEXPLAIN and CHATGPT

All models were trained and evaluated on an MSI GeForce RTX 3080 GAMING X TRIO 10G GPU using CUDA 11.7. Training and evaluation of all final models as well as hyperparameter tests took around 15 GPU hours.

6 Conclusion and Future Work

In this work, we have introduced TVCP, an NLP task, to aid in the accurate determination of the temporal validity duration of text by incorporating surrounding context. We create a benchmark dataset for our task and provide a set of baseline evaluation results for our dataset. We find that the performance of most classifiers can be improved by explicitly incorporating the temporal validity duration as a loss signal during training to improve the resulting embeddings. Despite the impressive feats performed by foundation models, we report, similar to previous work (Bian et al., 2023), poor performance in the TCS domain. These findings show that users should carefully evaluate whether a model like CHATGPT properly understands a given task before choosing it over smaller, fine-tuned LMs. We hypothesize that the performance of all models could further increase with additional training data.

Possible future work includes using the provided dataset and classifiers to collect a larger number of TVCP samples and annotating them with an updated temporal validity duration. A comparison of context-aware TV_d classifiers with prior models, like those by Almquist and Jatowt (2019), would shed light on the significance of accurate semantic segmentation between target and context. Similarly, the use of our dataset for generative approaches could be explored, for example, in the context of generative adversarial networks. For our multi-tasking implementation, directions for future work could be changes to hyperparameters such as the weight of the auxiliary loss, changes to the definition of the auxiliary task (e.g., log-scaled regression or ordinal classification), or even entirely new auxiliary tasks. Finally, current methods face limitations due to the effort of manual removal of stationary samples (Almquist and Jatowt, 2019; ours) or altering task definitions to avoid them (Hosokawa et al., 2023; Lynden et al., 2023). Research into models differentiating temporal and stationary information could enhance the development and definition of future TCS reasoning tasks.

607 Limitations

608 Although we focus on creating a reproducible
609 training- and evaluation environment, some vari-
610 ables are out of our control. For example, bit-wise
611 reproducibility is only guaranteed on the same
612 CUDA toolkit version and when executed on a
613 GPU with the same architecture and the same num-
614 ber of streaming multiprocessors. This means that
615 an exact reproduction of the models discussed in
616 this article may not be possible. Nevertheless, we
617 expect trends to remain the same across GPU archi-
618 tectures.

619 The use of CHATGPT as an example of foun-
620 dation model performance may be limiting due to
621 its black box design. In the future, open-source
622 models, such as LLAMA 2 (Touvron et al., 2023),
623 could be evaluated to improve the reproducibil-
624 ity of foundation model performance claims. We
625 chose to benchmark CHATGPT due to its common
626 use as a baseline and in end-user scenarios, but
627 the evaluation results may not be transferrable to
628 other foundation models or even other versions of
629 CHATGPT.

630 One of the major limitations of our approach
631 is likely the dataset size. Although a relatively
632 small dataset size is common in TCS reasoning,
633 we find that our model performance still increases
634 with the amount of training data used. The existing
635 synthesized context statements in our dataset could
636 be used to bootstrap an approach for automatically
637 extracting additional samples from social media to
638 alleviate this issue.

639 The data we collect is not personal in nature.
640 However, the possibility of latent demographic bi-
641 ases in our data exists, for example, with respect
642 to certain language structures or expressions used
643 in the creation of follow-up statements. This could
644 lead to the propagation of any such bias when the
645 dataset is used to bootstrap further data collection,
646 which should be considered in future work.

647 Our external validity is mainly threatened by
648 two factors. First, our context statements are crowd-
649 sourced. While we apply several steps to ensure the
650 produced context is sensible, it is unclear whether
651 downstream context, such as on social media plat-
652 forms, manifests in similar structures as in our
653 dataset, with respect to traits such as sentence
654 length, grammaticality, and phrasing.

655 Second, similar to how pre-training weights
656 from other TCS tasks do not seem to improve the
657 classifier performance on our dataset, the weights

658 generated as part of our training process are likely
659 very task-specific, and may not generalize well to
660 other tasks or text sources.

661 Overall, we recommend the use of the TVCP
662 dataset and classifiers for bootstrapping further re-
663 search into combining the duration- and inference-
664 based temporal validity tasks, as well as research
665 into directly predicting updated temporal validity
666 durations and improving the generalizability to dif-
667 ferent text sources, rather than for a direct down-
668 stream task application.

References 669

- 670 Shun Abe, Masumi Shirakawa, Tatsuya Nakamura,
671 Takahiro Hara, Kazushi Ikeda, and Keiichiro Hoashi.
672 2018. Predicting the occurrence of life events from
673 user’s tweet history. In *2018 IEEE 12th International
674 Conference on Semantic Computing (ICSC)*, pages
675 219–226. IEEE.
- 676 Axel Almqvist and Adam Jatowt. 2019. Towards con-
677 tent expiry date determination: Predicting validity
678 periods of sentences. In *European Conference on
679 Information Retrieval*, pages 86–101. Springer.
- 680 Sitaram Asur and Bernardo A Huberman. 2010. Pre-
681 dicting the future with social media. In *2010
682 IEEE/WIC/ACM international conference on web in-
683 telligence and intelligent agent technology*, volume 1,
684 pages 492–499. IEEE.
- 685 Prajwal Bhargava and Vincent Ng. 2022. Common-
686 sense knowledge reasoning and generation with pre-
687 trained language models: a survey. In *Proceedings
688 of the AAAI Conference on Artificial Intelligence*,
689 volume 36, pages 12317–12325.
- 690 Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie
691 Lu, and Ben He. 2023. Chatgpt is a knowledgeable
692 but inexperienced solver: An investigation of com-
693 mon-sense problem in large language models. *arXiv
694 preprint arXiv:2303.16421*.
- 695 Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard
696 Säckinger, and Roopak Shah. 1993. Signature verifi-
697 cation using a " siamese" time delay neural network.
698 *Advances in neural information processing systems*,
699 6.
- 700 Bibo Cai, Xiao Ding, Bowen Chen, Li Du, and Ting Liu.
701 2022. Mitigating reporting bias in semi-supervised
702 temporal commonsense inference with probabilistic
703 soft logic. In *Proceedings of the AAAI Conference
704 on Artificial Intelligence*, volume 36, pages 10454–
705 10462.
- 706 Shangkun Deng, Takashi Mitsubuchi, Kei Shioda, Tat-
707 suro Shimada, and Akito Sakurai. 2011. Combining
708 technical analysis with sentiment analysis for stock
709 price prediction. In *2011 IEEE ninth international
710 conference on dependable, autonomous and secure
711 computing*, pages 800–807. IEEE.

712	Chadi Helwe, Chloé Clavel, and Fabian M Suchanek.	2014. Integrating predictive analytics and social media. In <i>2014 IEEE Conference on Visual Analytics Science and Technology (VAST)</i> , pages 193–202. IEEE.	765
713	2021. Reasoning with transformer-based models: Deep learning, but shallow reasoning. In <i>3rd Conference on Automated Knowledge Base Construction</i> .		766
714			767
715			768
716	Hidehito Honda, Rina Kagawa, and Masaru Shirasuna.	Steven Lynden, Mehari Heilemariam, Kyoung-Sook Kim, Adam Jatowt, Akiyoshi Matono, Hai-Tao Yu, Xin Liu, and Yijun Duan. 2023. Commonsense temporal action knowledge (cotak) dataset. In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM 2023)</i> .	769
717	2022. On the round number bias and wisdom of crowds in different response formats for numerical estimation. <i>Scientific Reports</i> , 12(1):1–18.		770
718			771
719			772
720	Taishi Hosokawa, Adam Jatowt, Masatoshi Yoshikawa, and Kazunari Sugiyama. 2023. Temporal natural language inference: Evidence-based evaluation of temporal text validity. <i>Proceedings of the 45th European Conference on Information Retrieval (ECIR 2023)</i> , Springer LNCS.	Bill MacCartney. 2009. <i>Natural language inference</i> . Stanford University.	774
721			775
722			776
723			777
724		James Manyika. 2023. An overview of bard: an early experiment with generative ai.	778
725			779
726	Marc W Howard. 2018. Memory as perception of the past: compressed time in mind and brain. <i>Trends in cognitive sciences</i> , 22(2):124–136.	Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 839–849.	780
727			781
728			782
729	Adam Jatowt and Ching-man Au Yeung. 2011. Extracting collective expectations about the future from large text collections. In <i>Proceedings of the 20th ACM international conference on Information and knowledge management</i> , pages 1259–1264.		783
730			784
731			785
732			786
733			787
734	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of naacL-HLT</i> , volume 1, page 2.	Abhilash Nandy, Sushovan Haldar, Subhashis Banerjee, and Sushmita Mitra. 2020. A survey on applications of siamese neural networks in computer vision. In <i>2020 International Conference for Emerging Technology (INCET)</i> , pages 1–5. IEEE.	788
735			789
736			790
737			791
738	Mayuko Kimura, Lis Kanashiro Pereira, and Ichiro Kobayashi. 2021. Towards a language model for temporal commonsense reasoning. In <i>Proceedings of the Student Research Workshop Associated with RANLP 2021</i> , pages 78–84.	Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt Gardner, and Dan Roth. 2020. Torque: A reading comprehension dataset of temporal ordering questions. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1158–1172.	792
739			793
740			794
741			795
742			796
743	Yashasvi Koul, Kanishk Mamgain, and Ankit Gupta. 2022. Lifetime of tweets: a statistical analysis. <i>Social Network Analysis and Mining</i> , 12(1):101.		797
744			798
745		Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	799
746	Niels Buus Lassen, Rene Madsen, and Ravi Vatrpu. 2014. Predicting iphone sales from iphone tweets. In <i>2014 IEEE 18th International Enterprise Distributed Object Computing Conference</i> , pages 81–90. IEEE.		800
747			801
748			802
749			803
750	Pengfei Li, Hua Lu, Nattiya Kanhabua, Sha Zhao, and Gang Pan. 2018. Location inference for non-geotagged tweets in user timelines. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 31(6):1150–1165.	Lis Pereira, Fei Cheng, Masayuki Asahara, and Ichiro Kobayashi. 2021. Alice++: Adversarial training for robust and effective temporal reasoning. In <i>Proceedings of the 35th Pacific Asia Conference on Language, Information and Computation</i> , pages 373–382.	804
751			805
752			806
753			807
754			808
755	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	Lis Pereira, Xiaodong Liu, Fei Cheng, Masayuki Asahara, and Ichiro Kobayashi. 2020. Adversarial training for commonsense inference. In <i>Proceedings of the 5th Workshop on Representation Learning for NLP (RepLANLP-2020)</i> , pages 55–60. Association for Computational Linguistics.	809
756			810
757			811
758			812
759			813
760	Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In <i>International Conference on Learning Representations</i> .	Lianhui Qin, Aditya Gupta, Shyam Upadhyay, Luheng He, Yejin Choi, and Manaal Faruqui. 2021. Time-dial: Temporal commonsense reasoning in dialog. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the</i>	814
761			815
762			816
763	Yafeng Lu, Robert Krüger, Dennis Thom, Feng Wang, Steffen Koch, Thomas Ertl, and Ross Maciejewski.		817
764			818
			819
			820

821			
822			
823			
824	Guy Rosin and Kira Radinsky. 2022. Temporal attention for language models. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 1498–1508.		
825			
826			
827			
828	Guy D Rosin, Ido Guy, and Kira Radinsky. 2022. Time masking for temporal language models. In <i>Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining</i> , pages 833–841.		
829			
830			
831			
832	Lingfeng Shen, Zhuoming Liu, and Xiongtao Zhou. 2020. Forecasting people’s action via social media data. In <i>2020 IEEE International Conference on Big Data (Big Data)</i> , pages 5254–5259. IEEE.		
833			
834			
835			
836	Shane Storks, Qiaozi Gao, and Joyce Y Chai. 2019a. Commonsense reasoning for natural language understanding: A survey of benchmarks, resources, and approaches. <i>arXiv preprint arXiv:1904.01172</i> , pages 1–60.		
837			
838			
839			
840			
841	Shane Storks, Qiaozi Gao, and Joyce Y Chai. 2019b. Recent advances in natural language inference: A survey of benchmarks, resources, and approaches. <i>arXiv preprint arXiv:1904.01172</i> .		
842			
843			
844			
845	Zijun Sun, Chun Fan, Qinghong Han, Xiaofei Sun, Yuxian Meng, Fei Wu, and Jiwei Li. 2020. Self-explaining structures improve nlp models. <i>arXiv preprint arXiv:2012.01786</i> .		
846			
847			
848			
849	Hikaru Takemura and Keishi Tajima. 2012. Tweet classification based on their lifetime duration. In <i>Proceedings of the 21st ACM international conference on Information and knowledge management</i> , pages 2367–2370.		
850			
851			
852			
853			
854	Qingyu Tan, Hwee Tou Ng, and Lidong Bing. 2023. Towards benchmarking and improving the temporal reasoning capability of large language models. <i>arXiv preprint arXiv:2306.08952</i> .		
855			
856			
857			
858	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .		
859			
860			
861			
862			
863			
864	Lav R Varshney and John Z Sun. 2013. Why do we perceive logarithmically? <i>Significance</i> , 10(1):28–31.		
865			
866	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.		
867			
868			
869			
870			
871	Jiexin Wang, Adam Jatowt, Masatoshi Yoshikawa, and Yi Cai. 2023. Bitimebert: Extending pre-trained language representations with bi-temporal information.		
872			
873			
		In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 812–821.	874
			875
			876
		Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	877
			878
			879
			880
			881
		Georg Wenzel and Adam Jatowt. 2023. An overview of temporal commonsense reasoning and acquisition. <i>arXiv preprint arXiv:2308.00002</i> .	882
			883
			884
		Zonglin Yang, Xinya Du, Alexander M Rush, and Claire Cardie. 2020. Improving event duration prediction via time-aware pre-training. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3370–3378.	885
			886
			887
			888
			889
		Changlong Yu, Hongming Zhang, Yangqiu Song, and Wilfred Ng. 2022. Cocolm: Complex commonsense enhanced language model with discourse relations. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 1175–1187.	890
			891
			892
			893
			894
		Li Zhang, Qing Lyu, and Chris Callison-Burch. 2020. Reasoning about goals, steps, and temporal ordering with wikihow. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4630–4639.	895
			896
			897
			898
			899
		Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. “going on a vacation” takes longer than “going for a walk”: A study of temporal commonsense understanding. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> . Association for Computational Linguistics.	900
			901
			902
			903
			904
			905
			906
			907
		Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan Roth. 2020. Temporal common sense acquisition with minimal supervision. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7579–7589.	908
			909
			910
			911
			912
		Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2021. Temporal reasoning on implicit events from distant supervision. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1361–1371.	913
			914
			915
			916
			917
			918
			919
		Bo Zhou, Yubo Chen, Kang Liu, Jun Zhao, Jiexin Xu, Xiaojian Jiang, and Qiuxia Li. 2022. Generating temporally-ordered event sequences via event optimal transport. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 1875–1884.	920
			921
			922
			923
			924
			925

926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974

A Crowdsourcing Definitions

In this section, we provide details on the crowdsourcing implementation. As noted, we use Amazon Mechanical Turk to collect crowdsourced data from participants.

A.1 Temporal Validity Duration Estimation

We assume the average layman is not familiar with the term *temporal validity*. Thus, we define the task as “determining how long the information within the tweet remains relevant after its publication”, i.e., for how long the user would consider the tweet timely and relevant. We provide the option *no time-sensitive information*, which should be selected when tweets do not contain any information, when information is not expected to change over time, or when it is fully contained in the past.

The task is otherwise a relatively straightforward classification task. We split our dataset into batches of 10 samples that are grouped into a single *human intelligence task* (HIT). For each HIT, we offer a compensation of USD0.25, based on an estimated 6-9 seconds of processing time per individual statement (i.e., 60-90 seconds per HIT). Figures 10 to 13 show the crowdsourcing layout.

A.2 Follow-Up Content Generation

Compared to the temporal validity duration estimation task, the follow-up content generation task requires a much more robust understanding of the overall concept of temporal validity and the respective semantic roles of the target- and follow-up statements. Hence, we focus on providing a more detailed explanation of the task. Figures 14 to 16 show the crowdsourcing setup. The detailed instructions tab is not listed due to its length, but contains instructions that can also be found in the public repository.

Notably, we labelled the target statement as *context tweet* rather than *target tweet* in this crowdsourcing task to emphasize that participants should not alter this statement directly, as this was a problem that occurred somewhat frequently during pilot tests. This contrasts with our formal definition of TVCP, where providing context is the role of the follow-up statement.

Each HIT requires participants to provide three follow-up statements, one for each TVCP class (DEC, UNC, INC). For each HIT, we offer a compensation of USD0.35. We base our compensation on an estimated 30–40 seconds of processing time

per follow-up statement (i.e., 90–120 seconds per HIT) due to the creative writing involved .

A.3 Discouraging Dishonest Activity

In initial pilot runs, we find that many submissions are the result of spam, dishonest activity, or a complete lack of task understanding, with many provided annotations being inexplicable by common sense in any possible interpretation of the statement.

To increase the quality of work on both tasks, we introduced three measures. First, we required participants to have an overall approval rate of 90% on the platform, as well as 1,000 approved HITs. Without these requirements, the amount of blatant spam (e.g., copy-pasted content) increases significantly.

Second, we devised qualification tests for both tasks. Participants had to determine the temporal validity durations for a set of sample statements to work on the temporal validity duration estimation task, and determine the correctness of follow-up statements and their updated duration labels to work on the follow-up content generation task.

Finally, we vet all participants’ responses individually up to a certain threshold. For each task, we manually verify the first 20 submissions of each annotator on their quality. We provide feedback and manually adapt submissions when they are partially incorrect. If submission quality is appropriate by the time a participant reaches 20 submitted HITs, we consider them as trusted, and only spot-check every 5th submission thereafter. If submission quality does not sufficiently improve at this point, we prohibit the participant from further working on the task.

Despite these efforts, the follow-up content generation task specifically still received several low-quality submissions that had to be manually filtered out and corrected. In future work, a preferable approach may be to replace the qualification test with an unpaid qualification HIT, in which a feedback loop between participants and requesters can be established on data that will not be included in the final dataset, and participants can manually be assigned a qualification once their quality of work is sufficient.

B ChatGPT Setup

We provide the following system prompt to the CHATGPT API:

975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023

Task Description

For the tweets below, select for how long you would consider information within them to be relevant. (i.e., the timespan for which each tweet is likely to contain relevant and timely information after its publication). If multiple options seem plausible, choose the most likely one. Please **follow the provided instructions carefully**. The task remains identical for each tweet.

[View Instructions](#)

Tweet 1: "\${tweet1}"

For how long does this tweet contain relevant information after being posted?

- This tweet contains no time-sensitive information.
- Less than one minute
- 1-5 minutes
- 5-15 minutes
- 15-45 minutes
- 45 minutes - 2 hours
- 2-6 hours
- More than 6 hours
- 1-3 days
- 3-7 days
- 1-4 weeks
- More than one month

Tweet 2: "\${tweet2}"

Figure 10: The interface of the temporal validity duration estimation task

Instructions

Summary

Detailed Instructions

Examples

Task

For each tweet, your task is to **determine how long the information within the tweet remains relevant after its publication**. First, read the tweet carefully and consider what information it is trying to convey. Then, classify the lifetime of information in the tweet from the time of its publication. In other words, imagine you are a user interested in the tweet's information. The lifetime of information is the period during which you would consider the tweet timely and relevant.

Guidelines

A tweet can be considered to have "no time-sensitive information" when its information is expected to always remain true (i.e., we do not expect the information to change over time, or it is fully contained in the past).

Further guidelines:

- Do not use real-world (contextual) knowledge to reason about when information becomes outdated if this information is not included in the tweet itself.
- Assume the content of the tweet is truthful and accurate.

Figure 11: The summary section of the temporal validity duration estimation task guidelines

Instructions

Summary

Detailed Instructions

Examples

Task

The goal of this task is to gather commonsense judgments about the duration of relevance for actions and events commonly discussed on social media. For each tweet, your task is to **determine how long the information within the tweet remains relevant after its publication**. First, read the tweet carefully and consider what information it is trying to convey. Then, classify the lifetime of information in the tweet from the time of its publication. In other words, imagine you are a user interested in the tweet's information. The lifetime of information is the period during which you would consider the tweet timely and relevant. For example, a tweet like "Check out the circus, coming to town this weekend only!" would have a lifespan of "3-7 days" (specifically, until the end of the week). If someone were to read this tweet a few weeks after it was posted, the information would have lost its value.

Guidelines

A tweet can be considered to have "no time-sensitive information" in the following cases:

- The tweet contains information that is not expected to change over time (e.g., "My name is Georg." or "Japan lies in Asia.").
- The tweet contains no information at all (e.g., "Dartsssss" or "Endless.").
- The tweet contains information that is fully contained in the past (e.g., "I slept for 10 hours."). This **also applies** if such a statement is tied to a temporal expression (e.g., "I slept for 10 hours yesterday."). In this case, despite the statement being tied to the current day due to the expression "yesterday", since the actual information is fully contained in the past (and the action is already fully completed), the sentence is considered to have no time-sensitive information. This is because a statement about past actions or events is considered to always remain true.

Further guidelines:

- Do not use real-world knowledge (i.e., contextual knowledge about entities or events that is not stated in the tweet itself) to reason about when information becomes outdated. For example, for the sentence "The world cup finals are coming up.", do not consider the actual date of the next world cup finals, but rather consider how far before the finals of any given world cup someone would be expected to post this tweet.
- Assume the content of the tweet is truthful. For example, for the sentence "I am going to meet the queen.", do not consider the actual likelihood of this event occurring or real-life circumstances which cause this particular event to be impossible, but instead, assume that information in the tweet holds true and that events are expected to occur.

Figure 12: The detailed description of the temporal validity duration estimation task guidelines

Instructions

Summary

Detailed Instructions

Examples

Good examples

Tweet: *This breakfast was pretty bad, but at least I'm going out for dinner tonight.*

Classification: More than 6 hours

Comment: As the user mentions breakfast, we can assume this tweet was written early in the day. Without this context, "2-6 hours" would also be acceptable.

Tweet: *I hate Thursdays.*

Classification: No time-sensitive information.

Comment: The tweet is phrased in a way that implies it is a recurring feeling and not limited to the current week. Thus, we do not expect this sentiment to change.

Tweet: *I just want to finish getting all of my tattoos so badly, but I have more important things to spend money on right now.*

Classification: More than one month

Comment: Note the user's intent to finish getting their tattoos, which indicates that the tweet contains time-sensitive information. However, the tweet indicates that this change is not expected to occur soon.

Bad examples

Tweet: *This day is awful! I don't even know how it could get any worse.*

Classification: 1-3 days

Comment: Since the tweet is only relevant on the current day, the correct classification is "More than 6 hours". "1-3 days" should only be used as a classification when the tweet is relevant until at least the next day.

Tweet: *This year all my family is getting coal and a hug.*

Classification: Less than one minute

Comment: The target action (giving family coal and a hug) may take less than one minute. However, unless we expect the action to take place immediately, this is not equal to the duration of relevance of the tweet.

Tweet: *Idk if I wanna go to dc today or tomorrow*

Classification: No time-sensitive information

Comment: Even though there is no concrete action specified, the intents of the user are focused on a specific duration. The correct classification is "1-3 days".

Figure 13: The examples section of the temporal validity duration estimation task guidelines

For the "Context Tweet" shown below, assume that its content is relevant for the duration of the "Expected Lifetime" annotation. Propose some follow-up tweets that the original author could write after the context tweet, respectively. Each follow-up tweet should affect the expected information lifetime of the context tweet in a certain way. Additionally, after writing each follow-up tweet that changes the information lifetime, specify the new expected lifetime of the context tweet by choosing from the corresponding dropdown menus. (The expected lifetime should now be different due to the follow-up tweet.)

[Important - Help Us Avoid Rejections](#)

The results of this task are important for our research. On the other hand, we understand the impact of rejections on a worker's account. Therefore, we ask workers to **follow the task description carefully** to facilitate a positive collaboration. Note especially the following guidelines:

- The updated expected lifetime estimates must be **shorter** or **longer** than the original expected lifetime.
You may not specify the same value as the original expected lifetime!
- The follow-up tweets must appropriately alter the information lifetime of the **context tweet**.
This is explained in detail in the instructions! The updated information lifetime refers to information in the **context tweet only!**

If you are unsure about your understanding of the task, please read the instructions carefully, work on a small number of HITs at first (3-5), and wait for our feedback. We will **not reject** single submissions that do not fit the task description completely (as long as an effort was made) and will instead provide **individual feedback**. However, if larger quantities of incorrect work are submitted, **we may have to reject such batches** to ensure an appropriate sample size for our research. Therefore, please do not work on larger quantities of HITs unless several of your submissions have been **accepted without feedback**. It is also possible that your qualification may be revoked if provided feedback is ignored.

[View Instructions](#)

Context Tweet: "\${text}"

Expected Lifetime: \${expected}

Follow-up tweet to decrease the expected lifetime:

Your follow-up tweet here.

For how long does the context tweet contain relevant information when considering your follow-up tweet?

Follow-up tweet with unchanged lifetime:

Your follow-up tweet here.

Follow-up tweet to increase the expected lifetime:

Your follow-up tweet here.

For how long does the context tweet contain relevant information when considering your follow-up tweet?

Figure 14: The interface of the follow-up content generation task

Instructions

Summary

Detailed Instructions

Examples

Task

In this crowdsourcing task, you are given a context tweet with an "expected lifetime" that indicates how long the information in the tweet will be relevant. Your task is to write three follow-up tweets:

- One where the expected lifetime of information in the context tweet **decreases**.
- One where the expected lifetime of information in the context tweet **remains unchanged**.
- One where the expected lifetime of information in the context tweet **increases**.

For the follow-up tweets that change the expected lifetime, you must also provide an updated lifetime estimate for the context tweet. Note that this new estimate is a period starting at the creation of the context tweet, **not** the follow-up tweet. Additionally, it must be a different class than the original expected lifetime (at least the adjacent shorter/longer class).

Guidelines

- Do not change the context tweet itself. Write follow-up tweets instead.
- You may not specify the same value as the original expected lifetime for your updated lifetime estimates.
- Focus on changing the lifetime of information in the context tweet.
- Give your best effort when the context tweet is unclear.
- Try to avoid using explicit temporal expressions.
- Be creative and come up with varied scenarios that change the expected information lifetime.

Possible Reasons for Rejection

We appreciate your contributions to our crowdsourcing tasks and strive to avoid rejecting work. However, in cases where the work submitted does not meet the requirements of the provided task, we may be unable to issue payment. **Work may be rejected if you submit a large number of HITs that do not follow the task description.** Most notably, some reasons for rejection may be:

- The work submitted does not adhere to the task description, especially the guidelines highlighted within the HIT interface and the instruction summary.
- The work appears to be "low-effort" (e.g., simply stating that an action will take a certain amount of time without providing further context).
- The work is written in poor English. While perfect grammar is not required, the level of English should at least match that of the context tweet.
- The provided updated lifetime estimates do not follow the task description. For instance, if the objective is to increase the lifetime of information, the work may be rejected if the updated lifetime estimate is not longer than the original estimate.

Figure 15: The summary section of the follow-up content generation task guidelines

Instructions

Summary

Detailed Instructions

Examples

Good examples	Bad examples
<p>Context Tweet: "Going to the gym after work today!" Expected Lifetime: More than 6 hours</p> <p>Follow-up to decrease the expected lifetime: "Actually, I'll get a quick workout in during my lunch break at the gym next door." New expected lifetime: 2-6 hours</p> <p>Why? The main information in the context tweet (going to the gym / working out) remains valid, but the action will now occur within a more immediate timeframe.</p>	<p>Context Tweet: "Going to the gym after work today!" Expected Lifetime: More than 6 hours</p> <p>Follow-up to decrease the expected lifetime: "I'm so sore from yesterday's workout that I can barely move. Skipping the gym today." New expected lifetime: Less than one minute</p> <p>Why? Consider the difference between an action that does not occur at all, and an action that occurs very quickly. In this example, the context tweet's information does not apply.</p>
<p>Follow-up with unchanged lifetime: "I think I'll try out a new HIIT workout."</p> <p>Why? The follow-up tweet relates to the context tweet, but does not change the expected lifetime of information.</p>	<p>Follow-up with unchanged lifetime: "I think I'll get pizza for dinner tonight."</p> <p>Why? In the unchanged lifetime task, there should still be some topical connection between the follow-up and the context tweet.</p>
<p>Follow-up to increase the expected lifetime: "Have to work overtime today. The gym will have to wait until tomorrow." New expected lifetime: 1-3 days</p> <p>Why? As the author confirms that the planned action will still take place, we consider the information lifetime in the context tweet as continuously valid until this new date.</p>	<p>Follow-up to increase the expected lifetime: "The gym is closed today due to a maintenance issue. Guess I'm not going." New expected lifetime: 1-3 days</p> <p>Why? A follow-up tweet cancelling plans can only be considered an appropriate follow-up when it is clear the action will still take place at a later date, which is not the case in this example.</p>

Figure 16: The examples section of the follow-up content generation task guidelines

1024 “You are a language model specializ-
1025 ing in temporal commonsense reason-
1026 ing. Each prompt contains Sentence
1027 A and Sentence B. You should deter-
1028 mine whether Sentence B changes the ex-
1029 pected temporal validity duration of Sen-
1030 tence A, i.e., the duration for which the
1031 information in Sentence A is expected to
1032 be relevant to a reader.

1033 To achieve this, in your responses, first,
1034 estimate for how long the average reader
1035 may expect Sentence A to be relevant on
1036 its own. Then, consider if the informa-
1037 tion introduced in Sentence B increases
1038 or decreases this duration. Surround this
1039 explanation in triple backticks (```).

1040 After your explanation, respond with one
1041 of the three possible classes correspond-
1042 ing to your explanation: Decreased, Neu-
1043 tral, or Increased.”

1044 After this system prompt, we provide three sam-
1045 ple responses, one for each of the classes. These
1046 sample responses are shown in Figure 17.

1047 Similar to the crowdsourcing task setup, we use
1048 the concept of the *expected relevance duration* of
1049 the target statement (called Statement A in the
1050 CHATGPT prompt) to explain statement-level tem-
1051 poral validity. Additionally, instead of prompting
1052 the model to classify the sample directly, we ask
1053 it to provide an explanation for its decision based
1054 on chain-of-thought reasoning. Wei et al. (2022)
1055 show that chain-of-thought prompting significantly
1056 increases several types of reasoning capabilities,
1057 including commonsense, in LLMs. We prompt
1058 CHATGPT to first estimate the temporal validity
1059 duration of the target statement on its own. In a
1060 second step, the model should then determine if the
1061 information introduced in the follow-up statement
1062 changes this temporal validity duration. After giv-
1063 ing its explanation, the model should respond with
1064 one of the three target classes.

1065 C Hyperparameters

1066 We perform hyperparameter testing regarding
1067 dropout probability before the classification layer
1068 (0.1, 0.25, 0.5), the base learning rate (1e-2, 1e-3,
1069 1e-4), and whether to freeze embedding layers (i.e.,
1070 training only intermediary and classification lay-
1071 ers). For both BERT and ROBERTA in the frozen
1072 and unfrozen setting, we perform grid-search over

Sentence A: I’m ready to go to the beach

Sentence B: I forgot all the beach towels are still
in the dryer, but I’ll be ready to go as soon as the
dryer’s done running.

Target Class: Increased

Sample Explanation: Going to the beach may
take a few minutes to an hour, depending on the
distance. However, if the author first needs to wait
on the dryer to finish in order to retrieve their beach
towels, this may take an additional 30-60 minutes.

Sentence A: taking bad thoughts out of my mind
thru grinding my assignments

Sentence B: I just have to get through a short
math homework assignment and memorize a few
spelling words so it shouldn’t take long.

Target Class: Decreased

Sample Explanation: Grinding through assign-
ments may take several hours, depending on the
number of assignments to complete. In Sentence B,
the author states they only have a few short assign-
ments remaining, so they may only take an hour or
less to finish them.

Sentence A: Slide to my dm guys, come on

Sentence B: Instagram DMs are such a fun way to
communicate.

Target Class: Neutral

Sample Explanation: The author encourages peo-
ple to direct message them, which may be relevant
for several minutes to a few hours. Sentence B does
not change the duration for which Sentence A is
expected to be relevant.

Figure 17: Sample items, target classes, and explana-
tions provided to CHATGPT for few-shot reasoning

the learning rate and dropout probability. For these
benchmarks, we use a predefined train-val-test split
(80%/10%/10%). The remaining setup is the same
as in Section 5.

Table 4 shows the three best-performing config-
urations for BERT and ROBERTA in the freeze
and nofreeze settings, respectively, on the TRANS-
FORMERCLASSIFIER pipeline. Table 5 shows
the same results for the SIAMESECLASSIFIER
pipeline.

The most notable finding appears to be that

Model	DO	LR	#Epochs	EM
BERT-nofreeze	0.25	1e-4	5	0.613
BERT-nofreeze	0.10	1e-4	6	0.548
BERT-nofreeze	0.50	1e-4	4	0.548
BERT	0.25	1e-4	17	0.321
BERT	0.10	1e-4	8	0.315
BERT	0.10	1e-3	10	0.304
ROBERTA	0.25	1e-3	14	0.262
ROBERTA	0.10	1e-4	16	0.256
ROBERTA	0.50	1e-3	15	0.238
ROBERTA-nofreeze	0.25	1e-3	1	0.000
ROBERTA-nofreeze	0.50	1e-3	1	0.000
ROBERTA-nofreeze	0.10	1e-4	1	0.000

Table 4: Best three models for each of the proposed configurations in the TRANSFORMERCLASSIFIER pipeline

Model	DO	LR	#Epoch	EM
BERT-nofreeze	0.25	1e-4	7	0.589
BERT-nofreeze	0.10	1e-4	4	0.577
BERT-nofreeze	0.50	1e-4	2	0.565
ROBERTA	0.10	1e-4	21	0.548
ROBERTA	0.50	1e-4	13	0.518
ROBERTA	0.25	1e-4	17	0.512
BERT	0.50	1e-4	9	0.387
BERT	0.25	1e-3	8	0.357
BERT	0.25	1e-4	5	0.339
ROBERTA-nofreeze	0.25	1e-3	1	0.000
ROBERTA-nofreeze	0.50	1e-3	1	0.000
ROBERTA-nofreeze	0.10	1e-4	1	0.000

Table 5: Best three models for each of the proposed configurations in the SIAMESECLASSIFIER pipeline

1084 ROBERTA gets stuck in a false minimum of pre-
1085 dicting a constant class when embedding layers are
1086 unfrozen, leading to an accuracy of 0.33 and an
1087 EM of 0. Hence, we freeze embedding layers for
1088 these model types in our main evaluation. As noted
1089 in Section 5, a possible reason for this behaviour
1090 could be differences in the embeddings contained
1091 within BERT’s [CLS] and ROBERTA’s <s> token.