

Generic Temporal Reasoning with Differential Analysis and Explanation

Yu Feng^{1*} Ben Zhou² Haoyu Wang² Helen Jin² Dan Roth²

¹Duke University ²University of Pennsylvania
yu.feng344@duke.edu {xyzhou, why16gzl, helenjin, danroth}@seas.upenn.edu

Abstract

Temporal reasoning is the task of predicting temporal relations of event pairs with corresponding contexts. While some temporal reasoning models perform reasonably well on in-domain benchmarks, we have little idea of the systems’ generalizability due to existing datasets’ limitations. In this work, we introduce a novel task named TODAY that bridges this gap with **temporal differential analysis**, which as the name suggests, evaluates if systems can correctly understand the effect of incremental changes. Specifically, TODAY makes slight context changes for given event pairs, and systems need to tell how this subtle contextual change will affect temporal relation distributions. To facilitate learning, TODAY also annotates human explanations. We show that existing models, including GPT-3, drop to random guessing on TODAY, suggesting that they heavily rely on spurious information rather than proper reasoning for temporal predictions. On the other hand, we show that TODAY’s supervision style and explanation annotations can be used in joint learning and encourage models to use more appropriate signals during training and outperform across several benchmarks. TODAY can also be used to train models to solicit incidental supervision from noisy sources such as GPT-3 and moves farther towards generic temporal reasoning systems.

1 Introduction

Temporal relation extraction (Pustejovsky et al., 2003; Chambers et al., 2014) is traditionally viewed as an information extraction task, where a model uses explicit temporal signals such as “before” to identify the temporal order of events. While these models have contributed to many downstream pipelines, they are not enough for more complicated tasks such as timeline generation, where most event pairs do not come with explicit clues. These

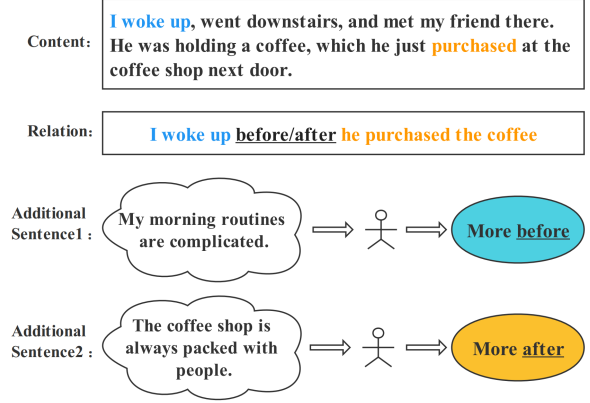


Figure 1: An example of temporal differential analysis, when adding the additional sentence1 to the content, a human will think the relation shifts towards **before** while adding the additional sentence2, he/she will think the relation shifts towards **after**.

implicit temporal relation extractions (Zhou et al., 2021) require temporal reasoning that relies on common sense and semantic understanding of the context. In recent work, a popular approach to address these predictions is to finetune pre-trained language models (PLMs) with annotated supervision data. Unfortunately, existing temporal benchmarks (Pustejovsky et al., 2003; Cassidy et al., 2014; Ning et al., 2018a) only annotate hard labels but ignore the fact that temporal labels based on common sense can be soft and nondeterministic. This allows models to exploit spurious signals and annotation artifacts easily. For example, a model may learn to predict “lunch” before “dinner” regardless of the surrounding context, yet most existing benchmarks will not challenge such beliefs because most “lunch” annotations will happen to be before “dinner.” This means that the current high performances of existing models may be misleading and the community may have a false sense of models’ generalizability.

In this work¹, we bridge this evaluation gap with

*Work partly done when visiting UPenn.

¹We will release data and code upon publication.

a novel benchmark that evaluates whether a temporal reasoning model is getting the correct predictions with the right reasons by properly identifying potential alternatives (e.g., “dinner” can be before “lunch” under certain contexts). Our intuition is to ask models to *explain* temporal relation predictions since the most viable way for humans to demonstrate insights into these problems is by providing satisfactory explanations. While the motivation is sound, automatically evaluating the plausibility of model explanations is extremely difficult. As a result, we use an approximation of such explanations, which we call **temporal differential analysis**. Under this setting, we select event pairs where the temporal relations are not 100% deterministic based on the context, meaning that both before/after relations are possible if additional information in regard to the context is given. Then, we annotate a hypothetical change in the form of an additional sentence added to the beginning of the context. As Figure 1 shows, this hypothetical will shift the event pair’s temporal relation distribution, making it either “*more before*” or “*more after*”. Each hypothetical change is also annotated with human explanations of why the change affects the temporal relation. We collect 2,241 such event pairs with a rigorous human annotation pipeline and call the resulting dataset TODAY (**temporal differential analysis**). If a model is generic enough to provide proper explanations for its temporal decisions, it can also distinguish subtle context changes and understand how each change will affect the distribution of temporal relations.

We find that models that achieve relatively high in-domain test performances are brittle and demonstrate minimal capabilities for differentiating subtle context changes that affect temporal relations. For example, the PatternTime model (Zhou et al., 2021) that achieves 77% binary accuracy on TRACIE (Zhou et al., 2021) - a dataset with similar contexts and events - drops dramatically to 54% on TODAY, which is barely above random guessing. To mitigate this gap, we propose a general technique that uses temporal explanations that TODAY annotates. Specifically, we argue that explanations of temporal relations are a great proxy for understanding temporal reasoning. We show models trained with TODAY’s task formulation and explanation annotation are better at perceiving cross-dataset supervision and achieve superior performances on multiple datasets with a single model. We also find

that while large language models (LLMs) are not good enough for temporal differential analysis, they sometimes produce reasonable explanations for a given temporal relation. We design a pipeline that automatically collects supervision signals based on this finding. The pipeline starts with giving GPT-3 (Brown et al., 2020) an instance from TODAY and a hypothetical temporal relation and then uses GPT-3 to generate several explanations. Finally, we train an explanation verifier based on human annotation, which selects the generated explanations that are more likely to be plausible. We show that adding such explanations from GPT-3 further boosts the performance across our benchmarks.

Our contribution is threefold. 1) We design a novel evaluation framework and collect a new dataset TODAY that uses differential analysis to test whether systems can perform temporal reasoning with the right reasons; 2) We show that the TODAY supervision, especially the use of explanations, contributes towards a generic temporal reasoning model; 3) We use LLMs to generate pseudo explanations and filter them with a novel explanation verification model and show that such distant supervision signals are helpful.

2 Related Work

Temporal Reasoning Models. Significant effort has been devoted to temporal reasoning, a challenging task that requires models to not only recognize the connection between event mentions but their context as well. Several statistical learning models (Mani et al., 2007; Ning et al., 2017, 2018b) have been proposed to characterize events based on features and learn to predict the temporal relations between event pairs. Recently, data-driven temporal reasoning approaches (Trong et al., 2022; Wang et al., 2022; Liu et al., 2021; Mathur et al., 2021; Zhou et al., 2020; Han et al., 2019) have witnessed great improvement over these feature-based models on benchmarks and are generally built upon deep neural models to predict temporal labels in an end-to-end fashion. Nevertheless, lack of interpretability has made these neural models untrustworthy to be deployed in real-world applications (Yin et al., 2022), especially in critical areas such as healthcare, finance, and government. The differential analysis approach first introduced in this paper provides a new paradigm of evaluating the interpretability of temporal reasoning models.

Temporal Relation Datasets. From different perspectives, multiple research projects have focused on constructing temporal reasoning benchmarks. A series of remarkable datasets, TimeBank (Pustejovsky et al., 2003), TempEval 1-3 (Verhagen et al., 2007, 2010; UzZaman et al., 2013), TimeBank-Dense (Cassidy et al., 2014), RED (O’Gorman et al., 2016), MATRES (Ning et al., 2018a) and so forth, are annotated on newswire articles for events and temporal relations between events. TORQUE (Ning et al., 2020) examines models’ capability in temporal reasoning in the format of reading comprehension, whereas a contrast set for MATRES is introduced in (Gardner et al., 2020) to provide a local view of models’ decision boundaries. However, none of these datasets provide reasons for these temporal decisions; thus, current temporal models tend to learn superficial temporal cues given the supervision. In contrast, the newly introduced framework, TODAY, bridges this gap by providing supervision signals under subtle context changes and corresponding explanations in the meantime.

Explanations. The community has been studying explanations and how they can help the reasoning tasks such as question answering. Several models have been proposed (Rajani et al., 2019; Latcinnik and Berant, 2020; Kumar and Talukdar, 2020; Zhou et al., 2022), as well as evaluation benchmarks that aim to test if existing systems can properly utilize explanations (Camburu et al., 2018; Aggarwal et al., 2021). Our work is closely related to this line of effort as we attempt to build a proxy benchmark that can be automatically evaluated for temporal explanations. The recent findings on large pre-trained language models have inspired several works to use them as explanation generators (Wiegreffe et al., 2021; Marasović et al.).

3 Dataset

In this section, we introduce the evaluation framework and the collection process of TODAY.

3.1 Task overview

The TODAY dataset and its overall framework is designed to evaluate systems’ ability to make temporal predictions with plausible reasons. Existing datasets, including MATRES, TORQUE, and TRACIE, annotate only common event pairs that align with human common sense. In other words, If an event pair does not strongly imply a temporal relation (e.g., over 80% confidence), it will not

be annotated and tested on systems. This allows pre-trained language models with millions of parameters to exploit annotation artifacts and certain priors that do not necessarily hold in specific contexts. For example, we know “lunch” is usually before “dinner”, but this also depends on if they are performed by the same subject, at the same location, and/or on the same day. Unfortunately, current models often memorize such relations as immutable facts, leading to prediction errors in instances that are less common in real life. This intuition inspires us to build a framework to evaluate how much spurious information and priors models are using.

Temporal Explanations An ideal method to evaluate if models are doing the right thing when making predictions is to let them explain why a certain prediction is made and evaluate the faithfulness and plausibility of the explanations. However, such an evaluation framework is almost impossible to achieve with current progress in natural language processing, where the two main challenges are 1) it is extremely difficult to collect gold explanations that are sufficient to cover any possible sets of explanations and 2) it is impossible to evaluate system generations using existing summarization metrics automatically.

Temporal Differential Analysis Because of the aforementioned challenges in directly evaluating system explanations, we propose an alternative that is a close proxy to the ideal form, namely temporal differential analysis. The core of temporal differential analysis is to check if models can correctly identify how a subtle change to the context may affect the temporal relations of a given event pair. The intuition behind this choice is two-fold: 1) it is much easier for both annotators and models to produce an explanation if they know which dimension to focus on; 2) this provides a binary evaluation that is deterministic and trustworthy in terms of reflecting how much spurious information models are using.

Specifically, our differential analysis process is defined below. Given an original context \mathcal{C} , event 1 \mathcal{E}_1 and event 2 \mathcal{E}_2 , we assume a gold distribution $\mathbb{D} = \{P_{before}, P_{after}, P_{same}\}$ on the temporal relation between \mathcal{E}_1 and \mathcal{E}_2 concerning \mathcal{C} , where $P_{before}, P_{after}, P_{same}$ are the probabilities of the temporal relation being before, after and simultaneous and they sum to 1. We then annotate two additional sentences \mathcal{AS}_{before} and \mathcal{AS}_{after} , where the temporal relation distribution between

Example

Context \mathcal{C} : Dave wanted to make a great scientific discovery. Dave worked with algae to make electricity. Dave discovered he could make electricity with algae! Dave was awarded for his great discovery.

Additional Sentence 1 (\mathcal{AS}_{before}): Dave was a scientist.

Event 1 (\mathcal{E}_1): Dave applied for a grant for his project.

Event 2 (\mathcal{E}_2): Dave worked with algae to make electricity.

Explanation: The additional sentence implies Dave was a scientist and normally a scientist has to apply for a grant before he starts the project.

Table 1: An example of temporal differential analysis, where \mathcal{AS} shifts the temporal relation between \mathcal{E}_1 and \mathcal{E}_2 to be more “before”. See §3 for more detail.

\mathcal{E}_1 and \mathcal{E}_2 with respect to $\mathcal{AS}_{before} + \mathcal{C}$ has an increased P_{before} , while similarly the distribution using $\mathcal{AS}_{after} + \mathcal{C}$ as the context has a higher P_{after} .

Table 1 shows an example instance of our temporal differential analysis, where an additional sentence \mathcal{AS}_{before} has an effect on the temporal relation between the two events and shifts the label distribution towards “before”. We conduct a pilot human study for this formulation and find that it is easy to annotate and achieve substantial improvement over the explanation quality compared with directly asking for explanations on an event pair. We, therefore, adopt this formulation and create our evaluation dataset TODAY through a multi-stage annotation process as detailed below.

3.2 Dataset Construction

Following the definition of the temporal differential analysis framework above, we collect a dataset to carry out the actual evaluation. Each instance in TODAY contains a context \mathcal{C} , an event pair \mathcal{E}_1 , \mathcal{E}_2 , and an additional sentence of either \mathcal{AS}_{before} or \mathcal{AS}_{after} . In addition, we also annotate a human explanation Exp regarding why the additional sentence affects the temporal relation between the two events. TODAY is constructed in three steps: 1) event pair generation, 2) additional sentence and explanation annotation, and 3) annotation verification and cleaning. We detail this pipeline below.

Generating \mathcal{C} and \mathcal{E} . We randomly sample short stories from the ROCStories dataset (Mostafazadeh et al., 2016) as the context \mathcal{C} . For each story, we use GPT-3 to generate an implicit event phrase based on an explicit event phrase selected by GPT-3 at the same time. An implicit event is a event that is not explicitly mentioned by the given context

but is inferable and relevant. A sample prompt can be referred to in appendix table 5 to construct an event pair. We do this for two main reasons: 1) events that are not explicitly mentioned by the context provide more uncertainty so that the event pair does not come with a deterministic temporal relation decided by the context; 2) this is closer to the format of TRACIE, which we aim to compare system performance changes with.

Crowdsourcing \mathcal{AS} and Exp . After having \mathcal{C} and \mathcal{E} ’s, we use Amazon Turk and ask crowdsourcing annotators to write potential \mathcal{AS}_{before} and \mathcal{AS}_{after} with respect to the provided information. The guideline asks annotators to write additional sentences that can be added to the beginning of the context to prevent models from using text positional information. The annotator is also asked to explain why they write \mathcal{AS} and why it affects the temporal relation distribution. We use this as Exp . We design an annotation interface that is intuitive and filled with examples, and at the same time, we require annotators to pass a rigorous qualification test to demonstrate proper understanding. We list our interfaces and tests in appendix figure 2 and figure 3.

Annotation Verification. We employ an additional verification stage for the human-written instances from the previous step. We provide annotators with the formatted textual entailment instance and ask if the entailment label changes in the expected direction. We collect two individual verifications per instance, and the instances accepted by all annotators will appear in the test set.

3.3 Statistics

We collect 1,000 instances that are agreed upon by both verifications while constructing a silver training set with the rest 1,241 instances.

4 Modeling

In this section, we show how to fully use TODAY’s supervision signals, especially the explanations, to build a more generic temporal reasoning model.

Joint Learning. TODAY only annotates temporal distribution shifts instead of absolute relations. This means that an instance may have a gold label “before” (i.e., the additional sentence \mathcal{AS} makes the relation more “before” compared to the original context), yet the likelihood of “after” can still be higher, and the $argmax$ label will be “after”. As a result, a model cannot sufficiently learn to predict

absolute labels with only supervision signals from TODAY. To mitigate this issue, we propose a joint learning model that requires joint supervision from a dataset that annotates hard labels for temporal relations, such as MATRES or TRACIE.

Modeling. We adopt TRACIE’s formulation (Zhou et al., 2021) to format the temporal reasoning task into textual entailment and use a sequence-to-sequence pre-trained language model as the base of our system. Specifically, the input sequence consists of the premise, which is $\mathcal{AS} + \mathcal{C} + Exp^2$ in our case, as well as the hypothesis, which is \mathcal{E}_1 starts [r] \mathcal{E}_2 . Here r is a hypothetical relation we plug into the hypothesis since systems are unaware of the gold label from the input sequence. The output sequence contains an entailment label, which is answer: positive for entail, and answer: negative for contradiction.

Hard Label Instances. As we note above, a system does not know the gold label when plugging in the hypothetical relation in the hypothesis. As a result, at learning time, we construct two entailment instances for a temporal relation instance with an absolute hard label. The first instance uses a hypothesis that is \mathcal{E}_1 starts before \mathcal{E}_2 . We want the model to learn to output answer: positive for entail if the gold label is also “before”, and answer: negative for contradiction if the gold label is “after”. The second instance uses \mathcal{E}_1 starts after \mathcal{E}_2 as the hypothesis, where the output sequences are reversed compared to the first one. We use the regular cross-entropy loss for optimization and denote the loss as ℓ_{CE} . At test time, we similarly construct two entailment instances for each event pair and conduct a simple probability-based vote to infer a final “before/after” relation.

Relative Label Instances. For instances that do not annotate absolute hard labels, we similarly construct two entailment instances for each event pair in training and evaluation time. However, instead of asking the model to use a cross-entropy loss to learn to output entailment labels, we employ a marginal ranking loss and ask the model to increase the probability of the entailment sequence if the plugged-in relation r is the same as the gold label³

r_g , and vice versa. Specifically, we want⁴

$$\begin{cases} p(\text{ent} | (\mathcal{AS} + \mathcal{C}), r) > p(\text{ent} | \mathcal{C}, r) & r = r_g \\ p(\text{con} | (\mathcal{AS} + \mathcal{C}), r) > p(\text{con} | \mathcal{C}, r) & r = \neg r_g \end{cases} \quad (1)$$

where ent and con represent entailment and contradiction, respectively. The loss function we use can subsequently be written as

$$\begin{aligned} \ell_{MR} &= \max(0, \epsilon + p_{o_g} - p_g) \\ &\quad + \max(0, \epsilon + p_w - p_{o_w}) \\ p_g &= p(\text{ent} | (\mathcal{AS} + \mathcal{C}), r_g) \\ p_o_g &= p(\text{ent} | \mathcal{C}, r_g) \\ p_w &= p(\text{ent} | (\mathcal{AS} + \mathcal{C}), \neg r_g) \\ p_{o_w} &= p(\text{ent} | \mathcal{C}, \neg r_g) \end{aligned} \quad (2)$$

where ϵ is a margin separating the logits. The actual probability of entailment is computed by the word logits in the output sequence of our sequence-to-sequence model.

Aggregated Loss Function. The final loss function we use for training is

$$\ell = \alpha \ell_{CE} + \ell_{MR} \quad (3)$$

where α reduces the two losses into the same scale. As a result, the proposed model is a general-purpose temporal reasoning model that can predict both hard-label temporal relations for an event pair and probability changes for differential analysis as proposed in TODAY.

5 LLM Incidental Supervision

As we both hypothesize and later show in §6, human-annotated explanations benefit generic temporal reasoning models, as they encourage models to learn to use the correct signals. However, it is extremely difficult and expensive to crowdsource such explanations for training purposes since collecting an instance cost \$1 in average. On the other hand, large language models (LLM) can produce a large amount of generated explanations at a much cheaper cost. Yet, these generated explanations are mostly unusable as they are simply model guesses based on textual correlations. In this section, we introduce a knowledge distillation method that combines the benefits of both human annotations and LLM generations by training verification models based on our seed annotation, which is then used to

² \mathcal{AS} and Exp only apply for relative label instances, such as those in TODAY.

³Here “gold label” refers to the direction that \mathcal{AS} shifts the temporal distribution to.

⁴For simplicity, we omit Exp and \mathcal{E} in the condition.

select generations that are more likely to be plausible. Compared to previous work (Wiegreffe et al., 2021), we propose a verification system composed of multiple models that individually verify different aspects of automatically-generated explanations. We detail our pipeline below.

5.1 Temporal Explanations from GPT-3

We adopt the same event pair generation and context selection process as we have detailed in §3. We design a prompt as shown in Table 4 that provides GPT-3 with contexts and event pairs, and ask GPT-3 to generate additional sentences, how these sentences will change the temporal relation, and why. The prompt contains a few examples, which makes this setting few-shot.

5.2 Verification System

Similarity-based Filtering. We filter GPT-3 instances that use exact same sentences from the context as the additional sentence or repeat the event pairs and temporal relations as explanations. We use Sentence-BERT (Reimers and Gurevych, 2019) and a similarity threshold of 0.95 to perform this filtering.

General Explanation Verifier. We use the generic temporal relation model as proposed in §4 trained on TODAY and an additional temporal relation dataset⁵ to verify if the generated additional sentence \mathcal{AS} shifts the temporal relation to the direction that it is supposed to.

Additional Sentence Verifier. The general explanation verifier cannot sufficiently identify if only part of the GPT-3 generation is correct. For example, a generated instance may have a sub-optimal \mathcal{AS} but convincing Exp , which would deceive our temporal relation model. To address this, we train a separate \mathcal{AS} verification model with TODAY that does not use Exp as input. We follow the same training scheme as §4, and similarly, use if \mathcal{AS} shifts the temporal relation as expected as our filtering criteria.

Explanation Sentence Verifier. We also train a binary classification model to individually check the plausibility of Exp . To generate negative Exp instances, for each instance in the TODAY training set, we ask GPT-3 to generate three possible explanation sentences. We use the one that is the least

similar to the human-annotated Exp according to SentenceBert as the negative instance, which we denote as Exp_{neg} . We finetune the base seq-to-seq model with the positive and negative explanations and optimize the loss function as the negative log-likelihood of the positive explanation:

$$\begin{aligned}\ell^E &= -\log \frac{e^{p_{pos}}}{e^{p_{pos}} + e^{p_{neg}}} \\ p_{pos} &= p(e | (\mathcal{AS} + \mathcal{C}, Exp), r_g), \\ p_{neg} &= p(e | (\mathcal{AS} + \mathcal{C}, Exp_{neg}), r_g),\end{aligned}\tag{4}$$

We filter all GPT-3 generated instances whose explanation is deemed as negative by this binary classification model.

6 Experiment

In this section, we conduct experiments to show that 1) existing systems do not truly understand temporal relations, and 2) TODAY and subsequent incidental supervision signals can partially address this issue and contribute to generic temporal reasoning models.

6.1 Datasets, Metrics, and Settings

We evaluate start-time temporal relation predictions with TRACIE (Zhou et al., 2021), MATRES (Ning et al., 2018a), as well as TODAY. Following the settings in (Zhou et al., 2021), we treat MATRES as a binary classification benchmark and use accuracy as the evaluation metric for all three datasets.

We set ϵ in equation 2 to be 0.1. We assign α in equation 3 to be 10. All models and baselines follow a standard TE setup and default parameters. All T5 experiments are trained with the same number of steps and repeated with three seeds.

6.2 Baselines and Systems

We use T5-large implemented by Wolf et al. (2020) as our base temporal reasoning model. We compare our proposed models with a host of baselines, including GPT-3 (Brown et al., 2020) and PatternTime (Zhou et al., 2021). We compare variations of our proposed model based on the same T5-large model including T5(T), where T5 is finetuned with TRACIE training set, T5(T+O), where T5 is finetuned together with TRACIE training set and TODAY training set, T5(T+O+G), where T5 is finetuned together with TRACIE training set, TODAY training set and verifier-filtered GPT-3 generated incidental supervision. We repeat this setting by replacing the TRACIE training set with MATRES

⁵Depending on the task, we choose different temporal relation datasets.

Data	Loss	TRACIE	MATRES	TODAY	TODAY (gen. exp.)	TODAY (gold exp.)	Average
GPT-3	FewShot	52.3	50.1	46.8	-	-	49.7
PatternTime	Distant	77.0	73.0	54.1	59.3	67.7	68.0
T5 (O)	MR	50.6	49.8	52.9	53.7	55.7	51.1
T5 (O+G)	MR	55.4	52.3	55.0	57.8	66.5	54.2
T5 (M)	CE	52.7	81.2	52.5	55.3	57.5	62.1
T5 (M+O)	CE + MR	51.5	81.7	57.4	60.5	82.7	63.5
T5 (M+O+G)	CE + MR	49.9	82.9	61.4	61.9	82.9	64.8
T5 (T)	CE	66.2	63.2	52.3	55.0	56.0	60.7
T5 (T+O)	CE + MR	72.9	69.4	59.9	61.7	81.6	67.4
T5 (T+O+G)	CE + MR	73.5	68.8	62.1	63.1	82.0	68.1
T5 (M+T)	CE	66.2	82.0	52.5	54.7	58.5	66.9
T5 (M+T+O)	CE + MR	73.0	83.5	57.9	60.8	77.8	71.5
T5 (M+T+O+G)	CE + MR	73.3	83.9	63.2	63.1	81.6	73.5
PatternTime (all)	CE + MR	79.9	86.3	62.9	63.4	82.3	76.4

Table 2: System performances under different supervision data and loss function settings across three binary temporal benchmarks. For simplicity, we use T to denote TRACIE training data, and similarly M for MATRES, O for TODAY (ours), and G for GPT-3 generated incidental supervision. TODAY only includes the additional sentence. TODAY (gold exp.) includes the additional sentence and the gold explanation sentence for each instance while TODAY (gen exp.) includes the additional sentence and the explanation sentence generated by GPT-3 after filtering for each instance. Average denotes the average binary accuracy of TRACIE, MATRES and TODAY for each setting. All T5 experiments are trained with the same number of steps and repeated with three seeds.

training set and TRACIE + MATRES combined training set respectively. Note that we only include 1.5k (10%) training instances for MATRES to match the size of other training data. We collect 5000 initial GPT-3 generated incidental supervision and 4811 remained after similarity-based filtering. We apply cross-entropy loss for TRACIE and MATRES training set and margin ranking loss for TODAY training set and GPT-3 generated supervision.

6.3 Inference

For TODAY testing set, given the additional sentence for each instance, we utilize GPT-3 to generate three possible explanation sentences based on the additional sentence for both relation directions of each test instance. We then rely on the explanation sentence verifier to choose the final explanation sentence. Specifically, we adopt the explanation sentence with the highest score under the explanation sentence verifier. To enhance the explanation sentence verifier’s capacity to identify an incorrect explanation sentence given a correct additional sentence, the explanation sentence verifier is further finetuned with GPT-3 generated training set with the same setting.

6.4 Main Results

Table 2 shows system performances under different supervision data and loss function settings across

three binary temporal benchmarks.

The performance of TODAY on existing systems, i.e., GPT-3 and PatternTime is unsatisfactory, revealing there is a gap between current temporal prediction and truly faithful temporal reasoning.

We observe that the average binary accuracy of TRACIE, MATRES and TODAY improves with the increasingly diversified training data and achieves the largest increase from 51.1% to 73.5% under the unified T5 training setting, which indicates that the model is being more generalized. Especially if we apply all the training data to PatternTime, the average binary accuracy increases by 8.4%. The use of explanations contributes to an average increase of 5.6% on the average accuracy compared to merely using the temporal reasoning data, which further verifies the effectiveness of explanations as guidance for models to behave correctly and more like a human towards this task.

We also show that the TODAY supervision contributes towards a better temporal reasoning model, with a 6.7% increase on TRACIE when trained with TRACIE only, 0.5% increase on MATRES when trained with MATRES only, and 6.8% increase on TRACIE and 1.5% increase on MATRES when trained together with TRACIE and MATRES. An increase of average 6% on TODAY without an explanation sentence further proves that the temporal model is drifting towards the right reasoning di-

Data	#GPT	T	M	TODAY	Avg
Ours	1475	73.3	83.9	63.2	73.5
No Exp	1867	73.7	83.5	61.2	72.8
No Addition	2529	70.2	81.4	59.5	70.4
No General	2079	71.0	81.8	59.5	70.8
More #GPT	2483	74.6	84.0	63.2	73.9

Table 3: Ablation study for LLM generated supervision. We test the model performance under different verifier settings. We also test the setting where we include more verifier-filtered GPT-3 data (filtered by three verifiers). #GPT refers to the total number of verifier-filtered GPT-3 data under each setting. T refers to TRACIE, M refers to MATRES, and Avg refers to Average.

rection to focus on the differential highlights that contribute to the shift of temporal relation in the context.

With GPT-3 generated incidental supervision, the model performance further improves on all metrics, with an average increase of 0.5%, 0.8%, 3.8%, 1.3% on MATRES, TRACIE, TODAY and average accuracy respectively. This illustrates that LLM can provide cheap but effective incidental supervision to benefit the model.

We also notice that there is a huge gap between the performance of TODAY without and with gold explanation sentence. This is because a correct explanation sentence can further elaborate and explain the additional sentence, i.e., the differential component. We follow the methods in §6.3 to generate an explanation for TODAY test and further improve over TODAY w/o explanation by approximately 2%, while the performance is still suboptimal compared to including the gold explanation sentence. The reason is that the explanation verifier cannot choose the correct explanation from the possible two explanations of different temporal relations. We leave the research on how to generate and identify a high-quality explanation sentence for future work.

6.5 Ablation Studies and Analysis

We conduct several ablation studies to understand our models’ improvements better. Table 3 demonstrates the results of our model with different settings of verifiers. The results have proved the effectiveness of all the verifiers. The explanation sentence verifier has the least influence. This is expected as we ask GPT-3 to generate an additional sentence followed by an explanation sentence, which largely increases its chance of being

coherent as a single generation. We also utilize similarity-based filtering to drop the explanations that are almost identical to the hypothesis, which alleviates one of the major problems of GPT-3 generated explanations. The additional sentence verifier and the general verifier are more crucial as the quality of incidental supervision heavily relies on if it can first correctly interpret the differences in the context and then draw a corresponding reasonable conclusion.

We also see that including more filter-verified GPT-3 data can further enhance the model performance, suggesting the usefulness of LLMs to generate supervision signals to empower small models. Since the smaller T5 model with LLM distilled knowledge performs much better than the LLM itself, it also directs us to research the trade-off between model scaling and data scaling in temporal reasoning.

7 Conclusion

We introduce a novel differential analysis framework and a dataset named TODAY that aims to interpret and evaluate if a temporal model can make correct predictions instead of using spurious information. We demonstrate that existing temporal models fall short in the performance on TODAY. We further show that training on a temporal relation benchmark together with TODAY leads to a more generic temporal reasoning model, resulting in improved performance on TRACIE, MATRES, and TODAY. Finally, we follow TODAY’s formulation and distill GPT-3 to construct useful incidental supervision for the model by creating a training pipeline that combines GPT-3 with weak explanation verifiers to solicit a large set of cheap and automatic explanations. Despite these advances, the gap in performance on TODAY between using additional sentences only versus including human-annotated gold explanation sentences indicates that TODAY continues to be a challenging task for future work towards generic temporal reasoning.

References

- Shourya Aggarwal, Divyanshu Mandowara, Vishwa-jeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. 2021. Explanations for commonsenseqa: New dataset and models. In *Annual Meeting of the Association for Computational Linguistics*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie

- Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. In *Neural Information Processing Systems*.
- Taylor Cassidy, Bill McDowell, Nathanael Chambers, and Steven Bethard. 2014. [An annotation framework for dense event ordering](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 501–506, Baltimore, Maryland. Association for Computational Linguistics.
- Nathanael Chambers, Taylor Cassidy, Bill McDowell, and Steven Bethard. 2014. [Dense event ordering with a multi-pass architecture](#). *Transactions of the Association for Computational Linguistics*, 2:273–284.
- Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. [Evaluating models’ local decision boundaries via contrast sets](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323, Online. Association for Computational Linguistics.
- Rujun Han, Qiang Ning, and Nanyun Peng. 2019. [Joint event and temporal relation extraction with shared representations and structured prediction](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 434–444, Hong Kong, China. Association for Computational Linguistics.
- Sawan Kumar and Partha Talukdar. 2020. [NILE : Natural language inference with faithful natural language explanations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8730–8742, Online. Association for Computational Linguistics.
- Veronica Latcinnik and Jonathan Berant. 2020. Explaining question answering models through text generation. *ArXiv*, abs/2004.05569.
- Jian Liu, Jinan Xu, Yufeng Chen, and Yujie Zhang. 2021. Discourse-level event temporal ordering with uncertainty-guided graph completion. In *IJCAI*, pages 3871–3877.
- Inderjeet Mani, Ben Wellner, Marc Verhagen, and James Pustejovsky. 2007. Three approaches to learning tlinks in timeml. *Computer Science Department*.
- Ana Marasović, Iz Beltagy, Doug Downey, and Matthew E. Peters. [Few-shot self-rationalization with natural language prompts](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*.
- Puneet Mathur, Rajiv Jain, Franck Dernoncourt, Vlad Morariu, Quan Hung Tran, and Dinesh Manocha. 2021. [TIMERS: Document-level temporal relation extraction](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 524–533, Online. Association for Computational Linguistics.
- 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 *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849, San Diego, California. Association for Computational Linguistics.
- Qiang Ning, Zhili Feng, and Dan Roth. 2017. [A structured learning approach to temporal relation extraction](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1027–1037, Copenhagen, Denmark. Association for Computational Linguistics.
- Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt Gardner, and Dan Roth. 2020. [TORQUE: A reading comprehension dataset of temporal ordering questions](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1158–1172, Online. Association for Computational Linguistics.
- Qiang Ning, Hao Wu, and Dan Roth. 2018a. [A multi-axis annotation scheme for event temporal relations](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1318–1328, Melbourne, Australia. Association for Computational Linguistics.
- Qiang Ning, Ben Zhou, Zhili Feng, Haoruo Peng, and Dan Roth. 2018b. [CogCompTime: A tool for understanding time in natural language](#). In *Proceedings*

- of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 72–77, Brussels, Belgium. Association for Computational Linguistics.
- Tim O’Gorman, Kristin Wright-Bettner, and Martha Palmer. 2016. [Richer event description: Integrating event coreference with temporal, causal and bridging annotation](#). In *Proceedings of the 2nd Workshop on Computing News Storylines (CNS 2016)*, pages 47–56, Austin, Texas. Association for Computational Linguistics.
- James Pustejovsky, Patrick Hanks, Roser Sauri, Andrew See, Robert Gaizauskas, Andrea Setzer, Dragomir Radev, Beth Sundheim, David Day, Lisa Ferro, et al. 2003. The timebank corpus. In *Corpus linguistics*, volume 2003, page 40. Lancaster, UK.
- Nazneen Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning. In *Annual Meeting of the Association for Computational Linguistics*.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Hieu Man Duc Trong, Nghia Ngo Trung, Linh Van Ngo, and Thien Huu Nguyen. 2022. [Selecting optimal context sentences for event-event relation extraction](#). In *AAAI Conference on Artificial Intelligence*.
- Naushad UzZaman, Hector Llorens, Leon Derczynski, James Allen, Marc Verhagen, and James Pustejovsky. 2013. [SemEval-2013 task 1: TempEval-3: Evaluating time expressions, events, and temporal relations](#). In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pages 1–9, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Marc Verhagen, Robert Gaizauskas, Frank Schilder, Mark Hepple, Graham Katz, and James Pustejovsky. 2007. [SemEval-2007 task 15: TempEval temporal relation identification](#). In *Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)*, pages 75–80, Prague, Czech Republic. Association for Computational Linguistics.
- Marc Verhagen, Roser Saurí, Tommaso Caselli, and James Pustejovsky. 2010. [SemEval-2010 task 13: TempEval-2](#). In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 57–62, Uppsala, Sweden. Association for Computational Linguistics.
- Haoyu Wang, Hongming Zhang, Yuqian Deng, Jacob R Gardner, Muhao Chen, and Dan Roth. 2022. Extracting or guessing? improving faithfulness of event temporal relation extraction. *arXiv preprint arXiv:2210.04992*.
- Sarah Wiegreffe, Jack Hessel, Swabha Swayamdipta, Mark O. Riedl, and Yejin Choi. 2021. Reframing human-ai collaboration for generating free-text explanations. In *North American Chapter of the Association for Computational Linguistics*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Fan Yin, Zhouxing Shi, Cho-Jui Hsieh, and Kai-Wei Chang. 2022. [On the sensitivity and stability of model interpretations in NLP](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2631–2647, Dublin, Ireland. Association for Computational Linguistics.
- Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan Roth. 2020. [Temporal common sense acquisition with minimal supervision](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7579–7589, Online. Association for Computational Linguistics.
- Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2021. [Temporal reasoning on implicit events from distant supervision](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1361–1371, Online. Association for Computational Linguistics.
- Ben Zhou, Kyle Richardson, Xiaodong Yu, and Dan Roth. 2022. Learning to decompose: Hypothetical question decomposition based on comparable texts. In *EMNLP*.

<p>Let's add a sentence as the first sentence of the context to let the hypothesis more likely to hold true and explain why.</p> <p>Context: Tara always wanted jewelry. Her birthday was coming up. Test went to the store. He gave her a really nice necklace She adored him for the gift.</p> <p>Hypothesis: Test was being a good friend starts before he give her a really nice necklace</p> <p>Add what sentence as the first sentence of the context and why is the hypothesis more likely to hold true?</p> <p>Test and Tara always hanged out together.</p> <p>This makes the statement true because normally people will only hang out frequently with their friends and friends will send each other gifts on their birthdays.</p> <p>###</p> <p>Context: Tara always wanted jewelry. Her birthday was coming up. Test went to the store. He gave her a really nice necklace She adored him for the gift.</p> <p>Hypothesis: Test was being a good friend starts after he give her a really nice necklace</p> <p>Add what sentence as the first sentence of the context and why is the hypothesis more likely to hold true?</p> <p>Test had always had the biggest crush on his classmate Tara even though she didn't talk to him much.</p> <p>This makes the statement true because it indicates that Test and Tara's relationship wasn't close prior to Test giving Tara the gift.</p> <p>###</p> <p>Context: Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.</p> <p>Hypothesis: Tim scheduled an appointment with his dentist starts after his tooth was hurting like crazy</p> <p>Add what sentence as the first sentence of the context and why is the hypothesis more likely to hold true?</p>
--

Table 4: A sample prompt with an instance for two hypothetical changes to make the event pair's temporal relation "more before" or "more after".

<p>Let's find out an event that is unmentioned but can be inferred from the context and the temporal relation between the two events are not deterministic. The new event should not be longer than ten words and include only one verb.</p> <p>Context: Tara always wanted jewelry. Her birthday was coming up. Test went to the store. He gave her a really nice necklace She adored him for the gift.</p> <p>What is an event that is unmentioned but has some role and can be inferred from the context?</p> <p>Test was being a good friend</p> <p>It can be inferred from She adored him for the gift.</p> <p>###</p> <p>Context: Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.</p> <p>What is an event that is unmentioned but has some role and can be inferred from the context?</p> <p>Tim scheduled an appointment with his dentist</p> <p>It can be inferred from Tim's tooth was hurting like crazy.</p> <p>###</p> <p>Context: Lily went to a nice restaurant. She ordered a steak. To her dismay the steak was rare. Lily was rather upset. She had to send it back.</p> <p>What is an event that is unmentioned but has some role and can be inferred from the context?</p>

Table 5: A sample prompt to generate an implicit event given the context.

Welcome! Please read the paragraph below and the two following statements that use the paragraph for context. For each statement, you are required to: (1) modify the paragraph by adding a new sentence in the front of the paragraph so that the statement will more likely be true and (2) explain why you are adding this sentence.

Note that you should always assume both events mentioned in each statement happened and are inferable and relevant to the paragraph.

[View instructions](#)

Paragraph: Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

Statement1: Tim scheduled an appointment with his dentist **starts before** his tooth was hurting like crazy

1. Use your imagination and **add a sentence** in the **front** of the paragraph so that statement1 will be **more likely** to hold.

The sentence you add **CANNOT directly include the implicit event: Tim scheduled an appointment with his dentist**, i.e. you may not add the same event word for word in the paragraph. **A sample addition** can be seen for reference if you click on **instructions at the beginning**.

Please add a sentence here.

2. Please give an **explanation** for why you added this sentence. **How** does it make statement1 **more likely** to hold true?

Please enter your explanation here...

Figure 2: The interface for differential explanation annotation.

Please read the paragraph below and the two following statements that use the paragraph for context.

Use your imagination and add a sentence in the front of the paragraph so that the statement will be more likely to hold.

The sentence you add CANNOT directly include the implicit event: Tim scheduled an appointment with his dentist

Paragraph: Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

Statement: Tim scheduled an appointment with his dentist starts before his tooth was hurting like crazy.

Question 1.1: Which modified paragraph do you think is the most suitable??

☐ Tim ate a lot of spicy food. Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

☐ Tim didn't schedule an appointment with his dentist. Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

☐ Tim's tooth was usually perfect, so he did not often go to see the dentist. Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

Paragraph: Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

Statement: Tim scheduled an appointment with his dentist starts before his tooth was hurting like crazy.

Question 1.2: Which modified paragraph do you think is the most suitable?

☐ Tim scheduled an appointment with his dentist. Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

☐ Tim was looking for a dentist. Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

☐ Tim always met his dentist regularly. Tim's tooth was hurting like crazy. He could barely eat or drink. His dentist took a look around in his mouth. One of his teeth was rotten. Once the tooth was pulled, Tim felt fine.

Question 2: Do you understand that the additional sentence and the explanation you write down must make the statement more likely to hold true ?

☐ Yes

☐ No

Figure 3: The interface for the qualification test of differential explanation annotation.