

Fine-Grained Temporal Relation Extraction

Siddharth Vashishtha
University of Rochester

Benjamin Van Durme
John Hopkins University

Aaron Steven White
University of Rochester

Abstract

We present a novel semantic framework for modeling temporal relations and event durations that maps pairs of events to real-valued scales for the purpose of constructing document-level event timelines. We use this framework to construct the largest temporal relations dataset to date, covering the entirety of the Universal Dependencies English Web Treebank. We use this dataset to train models for jointly predicting fine-grained temporal relations and event durations. We report strong results on our data and show the efficacy of a transfer-learning approach for predicting standard, categorical TimeML relations.

1 Introduction

Natural languages provide a myriad of formal and lexical devices for conveying the temporal structure of complex events – e.g. tense, aspect, auxiliaries, adverbials, coordinators, subordinators, etc. Yet, these devices are generally insufficient for determining the fine-grained temporal structure of such events. Consider the narrative in (1).

- (1) At 3pm, a boy broke his neighbor’s window. He was running away, when the neighbor rushed out to confront him. His parents were called but couldn’t arrive for two hours because they were still at work.

Most native English speakers would have little difficulty drawing a timeline for these events, likely producing something like that in Figure 1. But how do we know that the breaking, the running away, the confrontation, and the calling were short, while the parents being at work was not? And why should the first four be in sequence, with the last containing the others?

The answers to these questions likely involve a complex interplay between linguistic information, on the one hand, and common sense knowledge

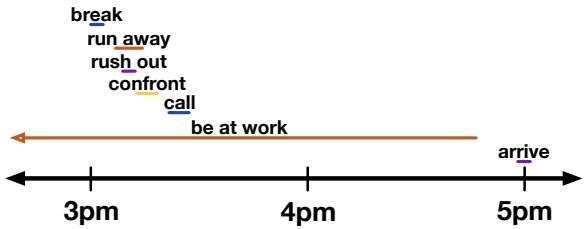


Figure 1: A typical timeline for the narrative in (1). about events and their relationships, on the other (Minsky, 1975; Schank and Abelson, 1975; Lampert, 1978; Allen and Hayes, 1985; Hobbs et al., 1987). But it remains an open question how best to capture this interaction.

A promising line of attack lies in the task of temporal relation extraction. Prior work in this domain has approached this task as a classification problem, labeling pairs of event-referring expressions – e.g. *broke* or *be at work* in (1) – and time-referring expressions – e.g. *3pm* or *two hours* – with categorical temporal relations (Pustejovsky et al., 2003; Styler IV et al., 2014; Minard et al., 2016). The downside of this approach is that we must rely on time-referring expressions to express duration information. But as example (1) highlights, nearly all temporal duration information can be left implicit, meaning it is only explicitly encoded when it is linguistically encoded.

In this paper, we develop a novel framework for temporal relation representation that puts event duration front and center. Like standard approaches using the TimeML standard, we draw inspiration from Allen’s (1983) seminal work on interval representations of time. But instead of annotating text for categorical temporal relations, we map event pairs directly to real-valued relative timeline representations, in addition to mapping events to their likely durations. This change not only supports the goal of giving a more central role to event duration, it also allows us to better reason about the temporal structure of complex events as

described by entire documents.

We begin with a discussion of the literature on temporal relation extraction (§2) and then discuss our own framework and data collection methodology (§3). The resulting Universal Decompositional Semantics Time (UDS-T) dataset is the largest temporal relation dataset to date (available at [decomp.io](#)), covering all of the Universal Dependencies (Silveira et al., 2014; De Marneffe et al., 2014; Nivre et al., 2015) English Web Treebank (Bies et al., 2012). We use this dataset to train a variety of neural models (§4) to jointly predict fine-grained (real-valued) temporal relations and event durations (§5), showing not only that our models obtain strong results on our dataset, the representations they learn can be straightforwardly transferred to the standard categorical relation datasets. (§6).

2 Background

We review prior work on temporal relations frameworks and associated corpora as well as systems for temporal relation extraction.

Corpora Most large datasets capturing temporal relations between events use the TimeML standard (Pustejovsky et al., 2003; Styler IV et al., 2014; Minard et al., 2016). TimeBank is one of the earliest large corpora built using this standard, capturing event pairs that annotators felt were salient (Pustejovsky et al., 2003). The TempEval competitions improved on the number of temporal relations by covering relations between all the events and times in a sentence, but only one of the TempEval tasks covered inter-sentential event relations (Verhagen et al., 2007, 2010; UzZaman et al., 2013, and see Chambers et al. 2014).

Efforts have been made to address the issue of sparsity in event-graphs with corpora such as the TimeBank-Dense (Cassidy et al., 2014) where annotators label all local-edges irrespective of ambiguity. TimeBank-Dense does not capture the complete graph over events and times relations, instead attempting to achieve completeness by capturing all relations within a sentence and the neighboring sentence. We take inspiration from this work for our own annotation protocol.

The Richer Event Description (RED) corpus takes a multi-stage annotation pipeline where various event-event phenomena, including temporal relations and sub-event relations are annotated together in the same datasets (O’Gorman et al.,

2016). Similarly, Hong et al. (2016) build a cross-document event corpora which covers fine-grained event-event relations and roles with more number of event types and sub-types. Another framework called GAF (Fokkens et al., 2013) captures event-identification through both textual and non-textual sources to track events across news articles.

Most of the corpora mentioned above required skilled workers to build the annotations as they follow specific ontologies. We take an alternative approach of capturing temporal relations by designing a protocol that asks simple questions about events which can be answered by any native speaker of English, finding surprisingly high agreement among annotators (see §3).

Models A variety of approaches have been taken to identifying the temporal relations between pairs of events. Early approaches use hand-tagged features modeled with multinomial logistic regression and support vector machines (Mani et al., 2006; Bethard, 2013; Lin et al., 2015). Other approaches use a combination of rule-based and learning-based approaches (D’Souza and Ng, 2013) and sieve-based architectures (Chambers et al., 2014; Mirza and Tonelli, 2016). Ning et al. (2018) jointly model causal and temporal relations using Constrained Conditional Models and formulate the problem as an Integer Linear Programming problem.

We presented a novel joint framework, Temporal and Causal Reasoning (TCR), using CCMs and ILP to the extraction problem of temporal and causal relations between events

In the recent years, neural network-based approaches have used both recurrent (Tourille et al., 2017; Cheng and Miyao, 2017; Leeuwenberg and Moens, 2018) and convolutional architectures (Dligach et al., 2017). Leeuwenberg and Moens (2018) use such models to predict relative timelines constructed from a set of temporal relations. Our annotations allow us to directly predict relative timelines between a pair of events which we then use to create document timelines anchored to some specific event.

The pairwise classification can result in inconsistent temporal graphs, and efforts have been made to avert this issue by employing temporal reasoning (Chambers and Jurafsky, 2008; Yoshikawa et al., 2009; Denis and Muller, 2011; Do et al., 2012; Laokulrat et al., 2016; Ning et al., 2017; Leeuwenberg and Moens, 2017).

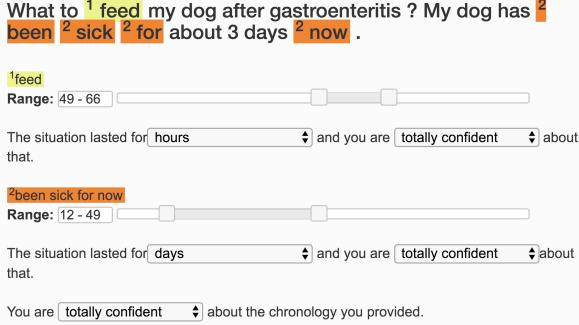


Figure 2: An annotated example from our protocol

People have also worked on modelling event durations from text (Pan et al., 2007; Gusev et al., 2011; Williams and Katz, 2012), but they don’t tie it directly to temporal relations. On the other hand, Filatova and Hovy (2001) assign a time-stamp to every clause in text, but the durations of events are not taken into consideration.

Attention-based models have proven effective in neural machine translation literature (Bahdanau et al., 2014; Luong et al., 2015; Vaswani et al., 2017), but to our knowledge, they have not been explored in identifying temporal relations. We follow up on this work in our models, using a variation of dot-product attention (Luong et al., 2015; Vaswani et al., 2017) to predict the event timelines and durations which is described §4. To cater to temporal reasoning, we treat the document timeline as a hidden representation and build it from the actual pairwise annotations as described in §7.

3 Data Collection

We collect the Universal Decompositional Semantics Time (UDS-T) dataset, which is annotated on top of the Universal Dependencies (Silveira et al., 2014; De Marneffe et al., 2014; Nivre et al., 2015) English Web Treebank (Bies et al., 2012). The main advantages of UD-EWT over other similar corpora are: (i) it covers text from a variety of genres; unlike most other datasets; (ii) it is built upon gold standard Universal Dependency parses; and (iii) it is compatible with various other semantic annotations which use the same predicate extraction standard (White et al., 2016; Zhang et al., 2017; Rudinger et al., 2018). Table 1 compares UDS-T against other temporal relations datasets.

Protocol design Annotators are given two contiguous sentences from a document with two highlighted event-referring expressions (predicates). If the predicate contains a copula, the whole predicate starting from the copula is highlighted. Other-

Dataset	#Events	#Event-Event Relations
TimeBank	7,935	3,481
TempEval 2010	5,688	3,308
TempEval 2013	11,145	5,272
TimeBank-Dense	1,729	8,130
Hong et al. (2016)	863	25,610
UDS-T	32,302	70,368

Table 1: Number of total events, and event-event temporal relations captured in various corpora

wise, only the root of the predicate is highlighted. They are then asked (i) to provide relative timelines on a 0-100 scale for the pair of events referred to by the highlighted predicates; and (ii) to give the likely duration of the event referred to by the predicate from the following list: *instantaneous, seconds, minutes, hours, days, weeks, months, years, decades, centuries, forever*. In addition, annotators were asked to give a confidence ratings for their relation annotation and each of their two duration annotation on the same five-point scale - *not at all confident* (0), *not very confident* (1), *somewhat confident* (2), *very confident* (3), *totally confident* (4).

An example of the annotation instrument is shown in Figure 2. Henceforth, we refer to the situation referred to by the predicate that comes first in linear order (*feed* in Figure 2) as e_1 and the situation referred to by the predicate that comes second in linear order (*sick* in Figure 2) as e_2 .

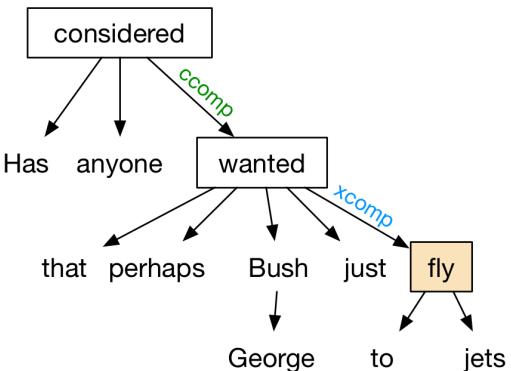


Figure 3: Our heuristic finds *fly* as (the root of) the pivot predicate in *Has anyone considered that perhaps George Bush just wanted to fly jets?*

Predicate extraction We extract predicates from UD-EWT using PredPatt (White et al., 2016; Zhang et al., 2017), which identifies 33,935 predicates from 16,622 sentences. We consider predicates with POS tags in: [ADJ, NOUN, NUM, DET, PROPN, PRON, VERB, AUX].

We concatenate two adjacent sentences to form a combined sentence which allows us to capture

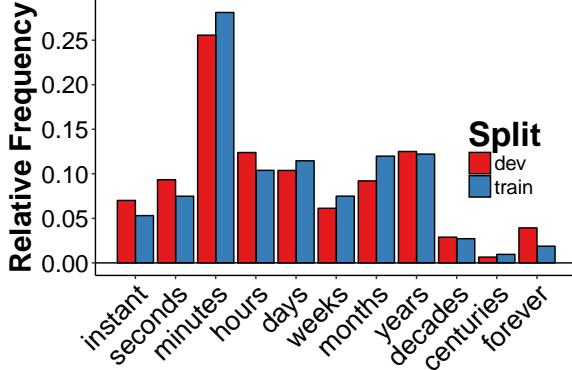


Figure 4: Distribution of event durations in training and development sets.

inter-sentential temporal relations. Considering all possible pairs of events in the combined sentence results into an exploding number of event-event comparisons. Therefore, to reduce the total number of comparisons, we find the *pivot-predicate* of the antecedent of the combined sentence as follows - find the root predicate of the antecedent and if it governs a CCOMP, CSUBJ, or XCOMP, follow that dependency to the next predicate until a predicate is found that doesn't govern a CCOMP, CSUBJ, or XCOMP. We then take all pairs of the antecedent predicates and pair every predicate of the consequent only with the *pivot-predicate*. This results into $\binom{N}{2} + M$ predicates instead of $\binom{N+M}{2}$ per sentence, where N and M are the number of predicates in the antecedent and consequent respectively. This heuristic allows us to find a predicate that loosely denotes the topic being talked about in the sentence. Figure 3 shows an example of finding the pivot predicate.

Annotators We recruited 765 annotators from Amazon Mechanical Turk to annotate predicate pairs in groups of ten. Each predicate pair contained in the UD-EWT training set was annotated by a single annotator, and each predicate in the UD-EWT development and test sets was annotated by three annotators.

Normalization We normalize the slider responses for each event pair by subtracting the minimum slider value from all values, then dividing all such shifted values by the maximum value (after shifting). This ensures that the earliest beginning point for every event pair lies at 0 and that the right-most end-point lies at 1, while preserving the ratio between the durations implied by the sliders.

Summary statistics Figure 4 shows the distribution of duration responses in the training and de-

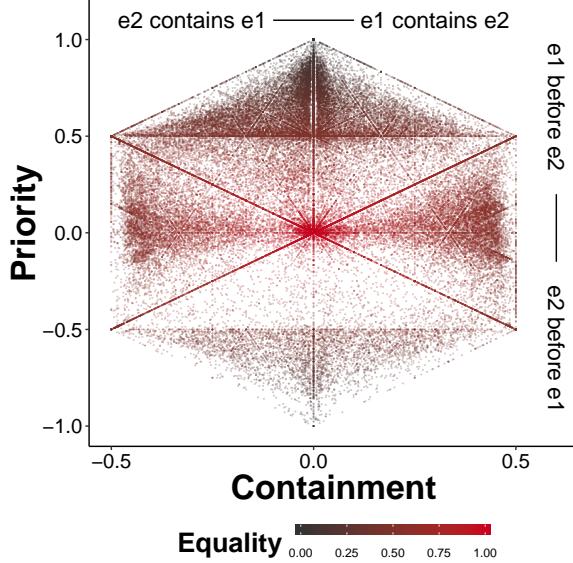


Figure 5: Distribution of event relations in training and development sets.

velopment sets. There is a relatively high density of events lasting *minutes*, with a relatively even distribution across durations of *years* or less and few events lasting *decades* or more.

The raw slider positions themselves are somewhat difficult to directly interpret, and so it is not particularly informative to show their distribution directly. To improve interpretability, we rotate the slider position space to construct four new dimensions: (i) PRIORITY, which is positive when e_1 starts and/or ends earlier than e_2 and most negative when e_2 starts and/or ends earlier than e_1 ; (ii) CONTAINMENT, which is most negative when e_2 contains more of e_1 and most positive when e_1 contains more of e_2 ; (iii) EQUALITY, which is largest when both e_1 and e_2 have the same temporal extents and smallest when they are most unequal; and (iv) SHIFT, which moves the events forward or backward in time. We construct these dimensions by solving for \mathbf{R} in

$$\mathbf{R} \begin{bmatrix} -1 & -1 & 1 & 1 \\ -1 & 1 & 1 & -1 \\ -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = 2\mathbf{S} - 1$$

where $\mathbf{S} \in [0, 1]^{N \times 4}$ contains the slider positions for our N datapoints in the following order: $\text{beg}(e_1), \text{end}(e_1), \text{beg}(e_2), \text{end}(e_2)$.

Figure 5 shows the embedding of the event pairs on the first three of these dimensions of \mathbf{R} . The triangular pattern near the top and bottom of the plot arises because strict priority – i.e. extreme positivity or negativity on the y -axis – precludes any

temporal overlap between the two events, and as we move toward the center of the plot, different priority relations mix with different overlap relations – e.g. the upper-middle left corresponds to event pairs where most of e_1 comes toward the beginning of e_2 , while the upper middle right of the plot corresponds to event pairs where most of e_2 comes toward the end of e_1 .

We see that there is a strong bias for e_1 to start and/or end earlier than e_2 – evidenced by the higher density of points near the upper center of Figure 5 than near the lower center – and a slight bias for e_1 to contain more of e_2 – evidenced by slightly higher density of points near the right center of Figure 5 than near the left center.

Inter-annotator agreement We measure inter-annotator agreement for the temporal relation sliders by calculating the rank (Spearman) correlation between the normalized slider positions for each pair of annotators that annotated a particular group of ten predicate pairs in the development set. Rank correlation is a useful measure in this case because it tells us how much different annotators agree of the relative position of each slider. The average rank correlation between annotators was 0.665 (95% CI=[0.661, 0.669]).

We measure interannotator agreement for the durations by calculating the absolute difference in duration rank between the duration responses for each pair of annotators that annotated a particular group of ten predicate pairs in the development set. On average, annotators disagree by 2.24 scale points (95% CI=[2.21, 2.25]), though there is heavy positive skew ($\gamma_1 = 1.16$, 95% CI=[1.15, 1.18]) – evidenced by the fact that the modal rank difference is 1 (25.3% of the response pairs), with rank difference 0 as the next most likely (24.6%) and rank difference 2 as a distant third (15.4%).

Annotation coherence Annotators were asked to approximate the relative duration of the two events that they were annotating using the distance between the sliders. This means that an annotation is coherent insofar as the ratio of distances between the slider responses for each event matches the ratio of the categorical duration responses. We rejected annotations wherein there was gross mismatch between the categorical responses and the slider responses – i.e. one event is annotated as having a longer duration but is given a shorter slider response – but because this does not guar-

antee that the exact ratios are preserved, we assess that here using a canonical correlation analysis (CCA; Hotelling 1936) between the categorical duration responses and the slider responses.

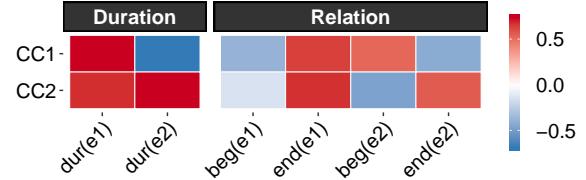


Figure 6: Scores from canonical correlation analysis comparing categorical duration annotations and slider relation annotations.

Figure 6 shows the CCA scores. We find that the first canonical correlation, which captures the ratios between unequal events, is 0.765; and the second, which captures the ratios between roughly unequal events, is 0.427. This preservation of the ratios is quite impressive in light of the fact that our slider scales are bounded; though we hoped for at least a non-linear relationship between the categorical durations and the slider distances, we did not expect such a strong linear relationship.

4 Model

For a given event pair in a sentence, we aim to jointly predict each event’s duration alongside the relative event timelines. We then use these relative timelines to construct timelines for entire documents with a separate model.

Relative timelines The relative timeline model consists of three components: an event model, a duration model, and a relation model. These components use multiple layers of *dot product attention* (Luong et al., 2015) on top of an embedding $\mathbf{H} \in \mathbb{R}^{N \times D}$ for a sentence $s = [w_1, \dots, w_N]$ tuned on the three M -dimensional contextual embeddings produced by ELMo (Peters et al., 2018) for that sentence, concatenated together.¹

$$\mathbf{H} = \tanh(\text{ELMo}(s)\mathbf{W}^{\text{TUNE}} + \mathbf{b}^{\text{TUNE}})$$

where D is the dimension for the tuned embeddings, $\mathbf{W}^{\text{TUNE}} \in \mathbb{R}^{3M \times D}$, and $\mathbf{b}^{\text{TUNE}} \in \mathbb{R}^D$.

Event model We define the model’s representation for the event referred to by predicate k as $\mathbf{g}_{\text{pred}_k} \in \mathbb{R}^D$, where D is the embedding size.

¹We found that correctly aligning BERT’s wordpiece representations with the predicate spans produced by PredPatt was not possible in general. In future work, we aim to introduce tunable intermediary alignment models for this purpose.

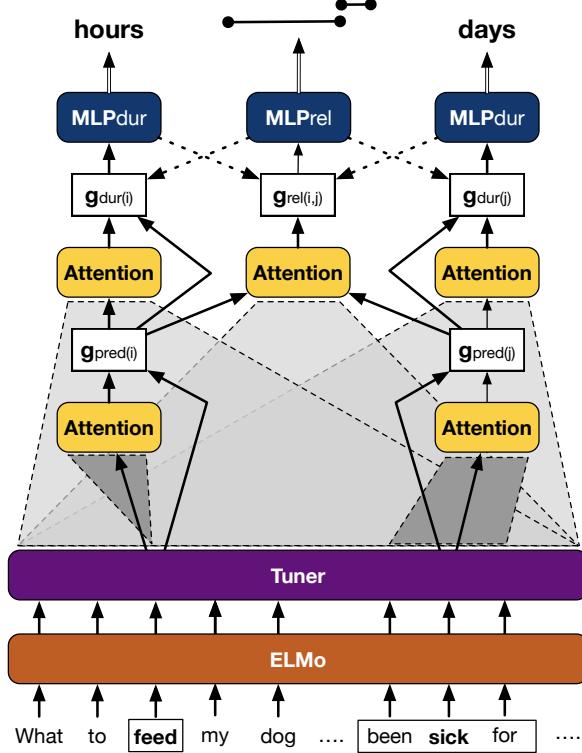


Figure 7: Network diagram for model. Dashed arrows are only included in some models.

We build this representation using a variant of dot-product attention, based on the predicate root.

$$\begin{aligned} \mathbf{a}_{\text{pred}_k}^{\text{SPAN}} &= \tanh(\mathbf{A}_{\text{PRED}}^{\text{SPAN}} \mathbf{h}_{\text{ROOT}(\text{pred}_k)} + \mathbf{b}_{\text{PRED}}^{\text{SPAN}}) \\ \alpha_{\text{pred}_k} &= \text{softmax}(\mathbf{H}_{\text{SPAN}(\text{pred}_k)} \mathbf{a}_{\text{pred}_k}^{\text{SPAN}}) \\ \mathbf{g}_{\text{pred}_k} &= [\mathbf{h}_{\text{ROOT}(\text{pred}_k)}; \alpha_{\text{pred}_k} \mathbf{H}_{\text{SPAN}(\text{pred}_k)}] \end{aligned}$$

where $\mathbf{A}_{\text{PRED}}^{\text{SPAN}} \in \mathbb{R}^{D \times D}$, $\mathbf{b}_{\text{PRED}}^{\text{SPAN}} \in \mathbb{R}^D$; $\mathbf{h}_{\text{ROOT}(\text{pred}_k)}$ is the hidden representation of the k^{th} predicate's root; and $\mathbf{H}_{\text{SPAN}(\text{pred}_k)}$ is obtained by stacking the hidden representations of the entire predicate.

The idea here is that the predicate root itself may be indicative of where within the predicate the relevant temporal information lies. For example, the predicate *been sick for now* in Figure 2 has *sick* as its root, and thus we would take the hidden representation for *sick* as $\mathbf{h}_{\text{ROOT}(\text{pred}_k)}$. Similarly, $\mathbf{H}_{\text{SPAN}(\text{pred}_k)}$ would be equal to taking the hidden-state representations of *been sick for now* and stacking them together. Then, if the model learns that tense information is important, it may weight *been* using the attention mechanism.

Duration model The temporal duration representation $\mathbf{g}_{\text{dur}_k}$ for the event referred to by the k^{th} predicate is defined similarly to the event representation, but instead of stacking the predicate's span, we stack the hidden representations of the entire sentence \mathbf{H} .

$$\mathbf{a}_{\text{dur}_k}^{\text{SENT}} = \tanh(\mathbf{A}_{\text{DUR}}^{\text{SENT}} \mathbf{g}_{\text{pred}_k} + \mathbf{b}_{\text{DUR}}^{\text{SENT}})$$

$$\alpha_{\text{dur}_k} = \text{softmax}(\mathbf{H} \mathbf{a}_{\text{dur}_k}^{\text{SENT}})$$

$$\mathbf{g}_{\text{dur}_k} = [\mathbf{g}_{\text{pred}_k}; \alpha_{\text{dur}_k} \mathbf{H}]$$

where $\mathbf{A}_{\text{DUR}}^{\text{SENT}} \in \mathbb{R}^{D \times \text{size}(\mathbf{g}_{\text{pred}_k})}$ and $\mathbf{b}_{\text{DUR}}^{\text{SENT}} \in \mathbb{R}^D$.

We consider two models of the categorical durations: a softmax model and a binomial model. The main difference is that the binomial model enforces that the probabilities $\mathbf{p}_{\text{dur}_k}$ over the 11 duration values be concave in the duration rank, whereas the softmax model has no such constraint. We employ a cross-entropy loss for both models.

$$\mathbb{L}_{\text{dur}}(d_k; \mathbf{p}) = \log p_{d_k}$$

In the softmax model, we pass the duration representation $\mathbf{g}_{\text{dur}_k}$ for event k through a multilayer perceptron (MLP) with a single hidden layer and ReLU activations, to yield probabilities $\mathbf{p}_{\text{dur}_k}$ over the 11 durations.

$$\begin{aligned} \mathbf{v}_{\text{dur}_k} &= \text{ReLU}(\mathbf{W}_{\text{DUR}}^{(1)} \mathbf{g}_{\text{dur}_k} + \mathbf{b}_{\text{DUR}}^{(1)}) \\ \mathbf{p} &= \text{softmax}(\mathbf{W}_{\text{DUR}}^{(2)} \mathbf{v}_{\text{dur}_k} + \mathbf{b}_{\text{DUR}}^{(2)}) \end{aligned}$$

In the binomial distribution model, we again pass the duration representation through a MLP with a single hidden layer of ReLU activations, but in this case, we yield only a single value π_{dur_k} . With $\mathbf{v}_{\text{dur}_k}$ as defined above:

$$\begin{aligned} \pi &= \sigma(\mathbf{w}_{\text{DUR}}^{(2)} \mathbf{v}_{\text{dur}_k} + \mathbf{b}_{\text{DUR}}^{(2)}) \\ p_c &= \binom{n}{c} \pi^n (1 - \pi)^{(n-c)} \end{aligned}$$

where $c \in \{0, 1, 2, \dots, 10\}$ represents the ranked durations – instant (0), seconds (1), minutes (2), ..., centuries (9), forever (10) – and n is the maximum class rank (10).

Relation model To represent the temporal relation representation between the event referred to by the i^{th} predicate and the event referred to by the j^{th} predicate, we again use a similar attention mechanism.

$$\mathbf{a}_{\text{rel}_{ij}}^{\text{SENT}} = \tanh(\mathbf{A}_{\text{REL}}^{\text{SENT}} [\mathbf{g}_{\text{pred}_i}; \mathbf{g}_{\text{pred}_j}] + \mathbf{b}_{\text{REL}}^{\text{SENT}})$$

$$\alpha_{\text{rel}_{ij}} = \text{softmax}(\mathbf{H} \mathbf{a}_{\text{rel}_{ij}}^{\text{SENT}})$$

$$\mathbf{g}_{\text{rel}_{ij}} = [\mathbf{g}_{\text{pred}_i}; \mathbf{g}_{\text{pred}_j}; \alpha_{\text{rel}_{ij}} \mathbf{H}]$$

where $\mathbf{A}_{\text{REL}}^{\text{SENT}} \in \mathbb{R}^{D \times 2 \text{size}(\mathbf{g}_{\text{pred}_k})}$ and $\mathbf{b}_{\text{REL}}^{\text{SENT}} \in \mathbb{R}^D$.

The main idea behind our temporal model is to map events and states directly to a timeline, which we represent via a *reference interval* $[0, 1]$. For

situation k , we aim to predict the beginning point b_k and end-point $e_k \geq b_k$ of k .

We predict these values by passing $\mathbf{g}_{rel_{ij}}$ through an MLP with one hidden layer of ReLU activations and four real-valued outputs $[\hat{\beta}_i, \hat{\delta}_i, \hat{\beta}_j, \hat{\delta}_j]$, representing the estimated relative beginning points $(\hat{\beta}_i, \hat{\beta}_j)$ and durations $(\hat{\delta}_i, \hat{\delta}_j)$ for events i and j . We then calculate the predicted slider values $\hat{\mathbf{s}}_{ij} = [\hat{b}_i, \hat{e}_i, \hat{b}_j, \hat{e}_j]$

$$[\hat{b}_k, \hat{e}_k] = [\sigma(\hat{\beta}_k), \sigma(\hat{\beta}_k + |\hat{\delta}_k|)]$$

The predicted values $\hat{\mathbf{s}}_{ij}$ are then normalized in the same fashion as the true slider values prior to being entered into the loss. We constrain this normalized $\hat{\mathbf{s}}_{ij}$ using four L1 losses.

$$\begin{aligned} \mathbb{L}_{rel}(\mathbf{s}_{ij}; \hat{\mathbf{s}}_{ij}) = & |(b_i - b_j) - (\hat{b}_i - \hat{b}_j)| + \\ & |(e_i - b_j) - (\hat{e}_i - \hat{b}_j)| + \\ & |(e_j - b_i) - (\hat{e}_j - \hat{b}_i)| + \\ & |(e_i - e_j) - (\hat{e}_i - \hat{e}_j)| \end{aligned}$$

The final loss function is then

$$\mathbb{L} = \frac{\mathbb{L}_{dur} + \epsilon * \mathbb{L}_{rel}}{2}$$

with ϵ set to a fixed value of 2 (see §5).

Duration-relation connections We also experiment with four architectures wherein the duration and relation models are connected to each other in the Dur \rightarrow Rel or Dur \leftarrow Rel directions.

In the first Dur \rightarrow Rel architecture, we modify $\mathbf{g}_{rel_{ij}}$ by additionally concatenating the i^{th} and j^{th} predicate's duration probabilities from the binomial distribution model.

$$\mathbf{g}_{rel_{ij}} = [\mathbf{g}_{pred_i}; \mathbf{g}_{pred_j}; \boldsymbol{\alpha}_{rel_{ij}} \mathbf{H}; \mathbf{p}_i; \mathbf{p}_j]$$

In the second Dur \rightarrow Rel architecture, we do not use the relation representation model at all, just using the i^{th} and j^{th} predicate's duration probabilities from the binomial distribution model.

$$\mathbf{g}_{rel_{ij}} = [\mathbf{p}_i; \mathbf{p}_j]$$

In the first Dur \leftarrow Rel architecture, we modify \mathbf{g}_{dur_k} by concatenating the \hat{b}_k and \hat{e}_k from the relation model.

$$\mathbf{g}_{dur_k} = [\mathbf{g}_{pred_k}; \boldsymbol{\alpha}_{dur_k} \mathbf{H}; \hat{b}_k; \hat{e}_k]$$

In the second Dur \leftarrow Rel architecture, we do not use the duration representation model at all, and instead use the predicted relative duration $\hat{e}_k - \hat{b}_k$ obtained from the relation model, passing it through the binomial distribution model.

$$\pi_{dur_k} = \hat{e}_k - \hat{b}_k$$

Document timelines From the timeline model, we learn the hidden document timelines for UDS-T development set using: (i) actual pairwise slider annotations; (ii) slider values predicted by the best performing model on UDS-T development set. We assume a hidden timeline $\mathbf{T} \in \mathbb{R}_{+}^{n_d \times 2}$, where n_d is the total number of predicates in that document, the two dimensions represent the beginning point and the duration of the predicates. We then construct predicted relative timelines with

$$\tau_{ij} = [t_{i1}, t_{i1} + t_{i2}, t_{j1}, t_{j1} + t_{j2}]$$

$$\hat{\mathbf{s}}_{ij} = \frac{\tau_{ij} - \min(\tau_{ij})}{\max(\tau_{ij} - \min(\tau_{ij}))}$$

We learn \mathbf{T} for each document under the relation loss $\mathbb{L}_{rel}(\mathbf{s}_{ij}, \hat{\mathbf{s}}_{ij})$. We further constrain \mathbf{T} to predict the categorical durations using the binomial distribution model on the durations t_{k2} implied by \mathbf{T} , assuming $\pi_k = \sigma(c \log(t_{k2}))$.

5 Experiments

We implement the neural model and attention in pytorch 1.0. We use the concatenated ELMo layers as word embeddings which are then tuned to a lower dimension of 256. For all experiments, we use stochastic gradient descent to train the ELMo-tuned embeddings, attention, and MLP parameters. The hyperparameter ϵ is set to be 2.0. Both the relation and duration MLP have a single hidden layer with 128 nodes. We weight both \mathbb{L}_{dur} , and \mathbb{L}_{rel} by the ridit-scored confidence ratings of event durations and event relations respectively.

To predict TimeML relations in TempEval3 (Task C - relation only) (UzZaman et al., 2013) and TimeBank-Dense (Cassidy et al., 2014), we use a transfer learning approach. We first use the best-performing model on the UDS-T development set to obtain the relation representation for each pair of annotated predicates in TempEval3 and TimeBank-Dense. We then use this vector as input features to a SVM classifier with a gaussian kernel (sklearn 0.20.0; Pedregosa et al. 2011). to predict the temporal relation on these datasets using the feature vector obtained from our model. We run a hyperparameter grid-search over 4-fold CV with C: (0.1, 1, 10), and gamma: (0.001, 0.01, 0.1, 1). The best performance on cross-validation (C=10 and gamma=0.001) is then evaluated on the test-sets of TempEval3 and TimeBank-Dense.

Model			Duration			Relation		
Duration	Relation	Connection	ρ	rank diff.	R1	Absolute ρ	Relative ρ	R1
softmax	✓	-	32.63	1.86	8.59	77.91	68.00	2.82
binomial	✓	-	37.75	1.75	13.73	77.87	67.68	2.35
-	✓	Dur \leftarrow Rel	22.65	3.08	-51.68	71.65	66.59	-6.09
binomial	-	Dur \rightarrow Rel	36.52	1.76	13.17	77.58	66.36	0.85
binomial	✓	Dur \rightarrow Rel	38.38	1.75	13.85	77.82	67.73	2.58
binomial	✓	Dur \leftarrow Rel	38.12	1.75	13.68	78.12	68.22	2.96

Table 2: Results on test data based on different model representations; ρ denotes the Spearman-correlation coefficient; rank-diff is the duration rank difference. The model highlighted in blue performs best overall on dev-data. The numbers highlighted in **bold** are the best-performing numbers in the respective columns.

Since we require spans of predicates for our model, we pre-process TempEval3 and TimeBank-Dense by removing all xml tags from the sentences and then we pass it through Stanford CoreNLP 3.9.2 (Manning et al., 2014) to get the corresponding conllu format. Roots and spans of predicates are then extracted using Pred-Patt. For our purposes, the *identity* and *simultaneous* relations in TempEval-3 are equivalent when comparing event-event relations. Hence, they are collapsed into one single relation.

Following recent work using continuous labels in event factuality prediction (Lee et al., 2015; Stanovsky et al., 2017; Rudinger et al., 2018; White et al., 2018) and genericity prediction (Govindarajan et al., 2019) we report three metrics for the duration prediction: Spearman correlation (ρ), mean rank difference (*rank diff*), and proportion rank difference explained (R1). We report four metrics for the relation prediction: Spearman correlation between the normalized values of actual beginning and end points and the predicted ones (*absolute ρ*), the Spearman correlation between the actual and predicted values in \mathbb{L}_{rel} (*relative ρ*), and the proportion of MAE explained (R1).

In both cases, the R1 metric corresponds closely to the related R^2 metric, which measures the amount of variance in the data explained by the model, but is defined in terms of mean absolute error (MAE), which assumes an L1 space.

$$R1 = 1 - \frac{\text{MAE}_{\text{model}}}{\text{MAE}_{\text{baseline}}}$$

where $\text{MAE}_{\text{baseline}}$ is always guessing the median. For both ρ and R1, we report the value scaled by 100 for readability.

As Govindarajan et al. (2019) note, these metrics are useful, since ρ tells us how similar the predictions are to the true values, ignoring scale, and

R1 tells us how close the predictions are to the true values, after accounting for variability in the data.

One difficulty that arises in computing metrics for the relation annotations on our test set is that we obtained three annotation each, and taking, e.g., the mean for each slider value in these annotations can result in a qualitatively different temporal relation, with different duration and relation characteristics, than any of the three annotations themselves. So instead of aggregating either the duration or relation annotations, we compute our metrics on all three annotations separately and then aggregate over them. Note that this will result in higher errors than we might see if we aggregate, but we believe it is the fairest way to report.

6 Results

Table 2 shows the results of different model architectures on the UDS-T test set, and Table 3 shows the results of our transfer-learning approach on TempEval-3 and TimeBank-Dense.

Systems	Data	F1	
		Micro	Macro
CAEVO	TD	0.494	-
CATENA	TD	0.519	-
Cheng and Miyao (2017)	TD	0.529	-
This work	TD	0.566	0.327
This work	TE3	0.489	0.208

Table 3: Results of our transfer learning experiment on event-event relations in TimeBank-Dense (TD) and TempEval-3 (TE3) compared against other systems.

UDS-T results The overarching pattern we see is that most of our models are able to predict the relative position of the beginning and ending of events very well (high relation ρ) and the relative duration of events somewhat well (relatively low duration ρ), but they have a lot more trouble predicting relation exactly and relatively less trouble

predicting duration exactly.

Duration model The binomial distribution model outperforms the softmax model for duration prediction by a large margin, though it has effectively no effect on the accuracy of the relation model, with the binomial and softmax models performing comparably. This suggests two things. First, the fact that the duration and relation models share the weights associated with the predicate representation does not affect the models this representation feeds into – i.e. having a bad duration representation does not entail having a bad relation representation, even if they are built upon the same foundation. Second, it seems that enforcing concavity in duration rank on the duration probabilities helps the model better predict durations. Indeed, as an elaboration on the first point, it may not be that the duration representations for the softmax model are worse than for the binomial models, it may just be that the extra constraints from the binomial model are helping.

Connections Connecting the duration and relation model doesn’t improve performance in general. In fact, when the durations are directly predicted from the temporal relation model – i.e. without using the duration representation model – the model’s performance drops by a large margin, with the Spearman correlation down by roughly 15 percentage points. This indicates that constraining the relations model to predict the durations is not enough and that the duration representation is needed to predict durations well.

On the other hand, predicting temporal relations directly from the duration probability distribution – i.e. without using the relation representation model – results in a similar score as that of the top-performing model. This indicates that the duration representation is able to capture most of the relation characteristics of the sentence. Using both duration representation and relation representation separately (model highlighted in blue) results in the best performance overall on the UDS-T development set. This is interesting in light of the fact that, as noted in §3, there is a strong linear relationship between the categorical durations and the durations implied by the relation annotations.

TempEval-3 and TimeBank-Dense We report F1-micro and F1-macro scores on TempEval-3 (TE3) and TimeBank-Dense (TD) in Table 3 and compare our results with some of the other sys-

tems, as reported by Cheng and Miyao (2017).² Our system beats the TD F1-micro scores of all other systems reported in Table 3. The top performing system on TE3 (Mirza and Tonelli, 2016) reports an F1 score of 0.619 over all relations. This indicates that our model is able to achieve competitive performance on other standard temporal classification problems.

Document timelines We apply the document timeline model described in §4 to both the annotations on the development set and the best-performing model’s predictions to obtain timelines for all documents in the development set. Figure 8 shows an example, comparing the two resulting document timelines.

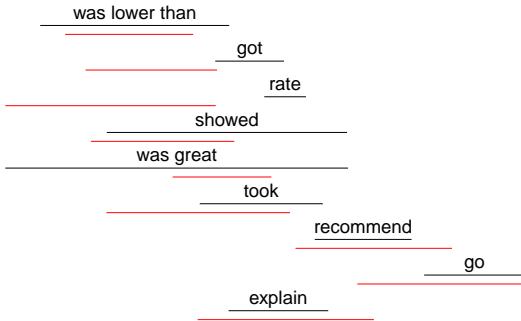


Figure 8: Learned Timeline for the following document based on actual (black) and predicted (red) annotations: “A+. *I would rate Fran pcs an A + because the price was lower than everyone else , i got my computer back the next day , and the professionalism he showed was great . He took the time to explain things to me about my computer , i would recommend you go to him. David*”

For these two timelines, we compare the induced beginning points and durations, obtaining a mean Spearman correlation of 0.28 for beginning points and -0.097 for durations. This suggests that the model agrees to some extent with the annotations about the beginning points of events in most documents but is struggling to find the correct duration spans. One possible reason for poor prediction of durations could be the lack of a direct source of duration information. The model currently tries to identify the duration based only on the slider values, which leads to poor performance in the Dur ← Rel model.

²We do not report the *temporal awareness* scores (F1) of other systems on TE3, since they report their metrics on all relations, including timex-timex, and event-timex relations, and thus they are not directly comparable. For TD, only those systems are reported that report F1-micro scores.

Duration				Relation			
Word	Attention (mean)	Rank (mean)	Freq	Word	Attention (mean)	Rank (mean)	Freq
soldiers	0.911	1.28	69	occupied	0.685	1.33	54
months	0.844	1.38	264	massive	0.522	2.71	66
Nothing	0.777	5.07	114	social	0.510	1.68	57
minutes	0.768	1.33	81	general	0.410	3.52	168
astronauts	0.756	1.37	81	few	0.394	3.07	474
hour	0.749	1.41	84	mathematical	0.393	7.66	132
Palestinians	0.735	1.72	288	are	0.387	3.47	4415
month	0.721	2.03	186	comes	0.339	2.39	51
cartoonists	0.714	1.35	63	or	0.326	3.50	3137
years	0.708	1.94	588	and	0.307	4.86	17615
days	0.635	1.39	84	emerge	0.305	2.67	54
thoughts	0.592	2.90	60	filed	0.303	7.14	66
us	0.557	2.09	483	s	0.298	4.03	1152
week	0.531	2.23	558	were	0.282	3.49	1308
advocates	0.517	2.30	105	gets	0.239	7.36	228

Table 4: The top 15 words in the dev-data which had the highest mean duration-attention and relation-attention weights. For duration, the words highlighted in bold directly correspond to some duration class. For relation, the words in bold are either conjunctions or words containing tense information.

7 Model Analysis and Timelines

We investigate three aspects of the best-performing model on the development set (highlighted in blue in Table 2): what our duration and relation representations attend to, how well we reconstruct the relation space defined in §3, and how well document timelines constructed from the model’s predictions match those constructed from the annotations themselves.

Attention The advantage of using an attention mechanism is that we can often interpret what linguistic information the model is using by analyzing the attention weights. We extract these attention weights for both the duration representation and the relation representation from our best model on the development set. We then compute the mean attention weight for these two attention models for each word type across the corpus. We also compute the mean rank of the attention weight for each word token within a sentence, with rank 1 assigned to the word with highest attention weight. Table 4 shows the top 15 words in the UDS-T development set according to mean attention weight, excluding words with frequency of less than 50 in EWT.

Duration Words that denote some time period – e.g. *month(s)*, *minutes*, *hour*, *years*, *days*, *week* – are among the top words in the duration model, with seven of the top 15 words directly denoting one of the duration classes. This is exactly what one might expect this model to rely heavily on, since time expressions are likely highly

informative for making predictions about duration. It also may suggest that we do not need to directly encode relations between event-referring and time-referring expressions in our framework – as do annotation standards like TimeML – since our models may discover these relations.

The remainder of the top words in the duration model are plurals or mass nouns. This may suggest that the plurality of a predicate’s arguments is an indicator of the likely duration of the event referred to by that predicate. To investigate this possibility, we compute a multinomial regression predicting the attention weights α_s for each sentence s from the K morphological features of each word in that sentence $\mathbf{F}_s \in \{0, 1\}^{\text{length}(s) \times K}$, which are extracted from the UD-EWT features column and binarized. To do this, we optimize coefficients \mathbf{c} in $\arg_{\mathbf{c}} \min \sum_s D(\alpha_s \parallel \text{softmax}(\mathbf{F}_s \mathbf{c}))$, where D is the KL divergence. We find that the five most strongly weighted positive features in \mathbf{c} are all features of nouns – NUMBER=*plur*, CASE=*acc*, PRONTYPE=*prs*, NUMBER=*sing*, GENDER=*masc* – suggesting that good portion of duration information can be gleaned from the arguments of a predicate. We believe this may be because nominal information can be useful in determining whether the clause is about particular events or generic events (Govindarajan et al., 2019). This is corroborated by the fact that the five most strongly weighted negative features in \mathbf{c} tend to be features of function words or predicates: PRONTYPE=*Rel*, DEGREE=*sup*, NUMTYPE=*mult*, VOICE=*pass*, NUMTYPE=*ord*.

Relation A majority of the top words in the relation model are either coordinators – such as *or* and *and* – or bearers of tense information – i.e. lexical verbs and auxiliaries. The first makes sense in light of the fact that, in context, coordinators can carry information about temporal sequencing (Bar-Lev and Palacas, 1980; Carston, 1993; Wilson and Sperber, 1998). The second makes sense in that information about the tense of predicates being compared likely helps the model determine relative ordering of the events they refer to.

To further investigate the role of morphological information, we compute multinomial regression in the same way as for the duration model, using the same morphological featurization. We find that the five most strongly weighted positive features in \mathbf{c} are all features of verbs or auxiliaries – PERSON=1, PERSON=3, TENSE=*pres*, TENSE=*past*, MOOD=*ind*, – suggesting that a majority of the information relevant to relation can be gleaned from the tense-bearing units in a clause. This is corroborated by the fact that the five most strongly weighted negative features in \mathbf{c} tend to be features of nouns or non-coordinator function words: CASE=*acc*, DEGREE=*cmp*, GENDER=*neut*, PRONTYPE=*Rel*, NUMTYPE=*ord*.

Relation space We rotate the predicted slider positions in the relation space defined in §3 and compare it with the rotated space of actual slider positions. We see a Spearman correlation of 0.19 for PRIORITY, 0.23 for CONTAINMENT, and 0.17 for EQUALITY. This suggests that our model is best able to capture CONTAINMENT relations and slightly less good at capturing PRIORITY and EQUALITY relations, though all the numbers are quite low compared to the *absolute* ρ and *relative* ρ metrics reported in Table 2. This may be indicative of the fact that our models do somewhat poorly on predicting more fine-grained aspects of an event relation, and in the future it may be useful to jointly train against the more interpretable PRIORITY, CONTAINMENT, and EQUALITY measures instead of or in conjunction with the slider values.

8 Conclusion

We presented a novel semantic framework for modeling temporal relations and event durations that maps pairs of events to real-valued scales for the purpose of constructing document-level event timelines. We used this framework to construct the largest temporal relations dataset to date – Uni-

versal Decompositional Semantic Time (UDS-T) – covering the entirety of the Universal Dependencies English Web Treebank. We used this dataset to train models for jointly predicting fine-grained temporal relations and event durations, reporting strong results on our data and showing the efficacy of a transfer-learning approach for predicting standard, categorical TimeML relations.

9 Acknowledgment

We thank the FACTS.lab at the University of Rochester for useful comments on framework and protocol design. This research was supported by the University of Rochester, JHU HLTCOE, and DARPA AIDA. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes. The views and conclusions contained in this publication are those of the authors and should not be interpreted as representing official policies or endorsements of DARPA or the U.S. Government.

References

- James F Allen. 1983. Maintaining knowledge about temporal intervals. *Communications of the ACM*, 26(11):832–843.
- James F Allen and Patrick J Hayes. 1985. A common-sense theory of time. In *IJCAI*, volume 85, pages 528–531.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Zev Bar-Lev and Arthur Palacas. 1980. Semantic command over pragmatic priority. *Lingua*, 51(2-3):137–146.
- Steven Bethard. 2013. Cleartk-timeml: A minimalist approach to tempeval 2013. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, volume 2, pages 10–14.
- Ann Bies, Justin Mott, Colin Warner, and Seth Kulick. 2012. English web treebank. *Linguistic Data Consortium, Philadelphia, PA*.
- Robyn Carston. 1993. Conjunction, explanation and relevance. *Lingua*, 90(1-2):27–48.
- Taylor Cassidy, Bill McDowell, Nathaniel Chambers, and Steven Bethard. 2014. An annotation framework for dense event ordering. Technical report, CARNEGIE-MELLON UNIV PITTSBURGH PA.

- 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.
- Nathanael Chambers and Dan Jurafsky. 2008. Jointly combining implicit constraints improves temporal ordering. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 698–706. Association for Computational Linguistics.
- Fei Cheng and Yusuke Miyao. 2017. Classifying temporal relations by bidirectional lstm over dependency paths. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 1–6.
- Marie-Catherine De Marneffe, Timothy Dozat, Natalia Silveira, Katri Haverinen, Filip Ginter, Joakim Nivre, and Christopher D. Manning. 2014. Universal Stanford dependencies: A cross-linguistic typology. In *Proceedings of LREC*, volume 14, pages 4585–4592.
- Pascal Denis and Philippe Muller. 2011. Predicting globally-coherent temporal structures from texts via endpoint inference and graph decomposition. In *IJCAI-11-International Joint Conference on Artificial Intelligence*.
- Dmitriy Dligach, Timothy Miller, Chen Lin, Steven Bethard, and Guergana Savova. 2017. Neural temporal relation extraction. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, volume 2, pages 746–751.
- Quang Xuan Do, Wei Lu, and Dan Roth. 2012. Joint inference for event timeline construction. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 677–687. Association for Computational Linguistics.
- Jennifer D’Souza and Vincent Ng. 2013. Classifying temporal relations with rich linguistic knowledge. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 918–927.
- Elena Filatova and Eduard Hovy. 2001. Assigning time-stamps to event-clauses. In *Proceedings of the workshop on Temporal and spatial information processing-Volume 13*, page 13. Association for Computational Linguistics.
- Antske Fokkens, Marieke van Erp, Piek Vossen, Sara Tonelli, Willem Robert van Hage, BV SynerScope, Luciano Serafini, Rachele Sprugnoli, and Jesper Hoeksema. 2013. Gaf: A grounded annotation framework for events. *NAACL HLT 2013*, page 11.
- Venkata Subrahmanyam Govindarajan, Benjamin Van Durme, and Aaron Steven White. 2019. Decomposing generalization: Models of generic, habitual, and episodic statements. *arXiv preprint arXiv:1901.11429*.
- Andrey Gusev, Nathanael Chambers, Pranav Khaitan, Divye Khilnani, Steven Bethard, and Dan Jurafsky. 2011. Using query patterns to learn the duration of events. In *Proceedings of the ninth international conference on computational semantics*, pages 145–154. Association for Computational Linguistics.
- Jerry R Hobbs, William Croft, Todd Davies, Douglas Edwards, and Kenneth Laws. 1987. Commonsense metaphysics and lexical semantics. *Computational linguistics*, 13(3-4):241–250.
- Yu Hong, Tongtao Zhang, Tim O’Gorman, Sharone Horowitz-Hendler, Heng Ji, and Martha Palmer. 2016. Building a cross-document event-event relation corpus. In *Proceedings of the 10th Linguistic Annotation Workshop held in conjunction with ACL 2016 (LAW-X 2016)*, pages 1–6.
- Harold Hotelling. 1936. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377.
- Leslie Lamport. 1978. Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM*, 21(7):558–565.
- Natsuda Laokulrat, Makoto Miwa, and Yoshimasa Tsuruoka. 2016. Stacking approach to temporal relation classification with temporal inference. *Information and Media Technologies*, 11:53–78.
- Kenton Lee, Yoav Artzi, Yejin Choi, and Luke Zettlemoyer. 2015. Event detection and factuality assessment with non-expert supervision. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1643–1648.
- Artuur Leeuwenberg and Marie-Francine Moens. 2018. Temporal information extraction by predicting relative time-lines. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1237–1246.
- Tuur Leeuwenberg and Marie-Francine Moens. 2017. Structured learning for temporal relation extraction from clinical records. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1150–1158.
- Chen Lin, Dmitriy Dligach, Timothy A Miller, Steven Bethard, and Guergana K Savova. 2015. Multi-layered temporal modeling for the clinical domain. *Journal of the American Medical Informatics Association*, 23(2):387–395.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.

- Inderjeet Mani, Marc Verhagen, Ben Wellner, Chong Min Lee, and James Pustejovsky. 2006. Machine learning of temporal relations. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, pages 753–760. Association for Computational Linguistics.
- Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pages 55–60.
- Anne-Lyse Myriam Minard, Manuela Speranza, Ruben Urizar, Begona Altuna, Marieke van Erp, Anneleen Schoen, and Chantal van Son. 2016. Meantime, the newsreader multilingual event and time corpus. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*. European Language Resources Association (ELRA).
- Marvin Minsky. 1975. A framework for representing knowledge. *The Psychology of Computer Vision*.
- Paramita Mirza and Sara Tonelli. 2016. Catena: Causal and temporal relation extraction from natural language texts. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 64–75.
- 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.
- Qiang Ning, Zhili Feng, Hao Wu, and Dan Roth. 2018. Joint reasoning for temporal and causal relations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 2278–2288.
- Joakim Nivre, Zeljko Agic, Maria Jesus Aranzabe, Masayuki Asahara, Aitziber Atutxa, Miguel Balles-teros, John Bauer, Kepa Bengoetxea, Riyaz Ahmad Bhat, Cristina Bosco, Sam Bowman, Giuseppe G. A. Celano, Miriam Connor, Marie-Catherine de Marneffe, Arantza Diaz de Ilarrazza, Kaja Dobrovolsjc, Timothy Dozat, Toma Erjavec, Richrd Farkas, Jennifer Foster, Daniel Galbraith, Filip Ginter, Iakes Goenaga, Koldo Gojenola, Yoav Goldberg, Berta Gonzales, Bruno Guillaume, Jan Haji, Dag Haug, Radu Ion, Elena Irimia, Anders Jøhannsen, Hiroshi Kanayama, Jenna Kanerva, Simon Krek, Veronika Laippala, Alessandro Lenci, Nikola Ljubei, Teresa Lynn, Christopher Manning, Ctlina Mrnduc, David Mareek, Hector Martnez Alonso, Jan Maek, Yuji Matsumoto, Ryan McDonald, Anna Missil, Verginica Mititelu, Yusuke Miyao, Simonetta Montemagni, Shunsuke Mori, Hanna Nurmi, Petya Osenova, Lilja vrelid, Elena Pascual, Marco Passarotti, Cenel-Augusto Perez, Slav Petrov, Jussi Piitulainen, Barbara Plank, Martin Popel, Prokopis Prokopidis, Sampo Pyysalo, Loganathan Ramasamy, Rudolf Rosa, Shadi Saleh, Sebastian Schuster, Wolfgang Seeker, Mojgan Seraji, Natalia Silveira, Maria Simi, Radu Simionescu, Katalin Simk, Kiril Simov, Aaron Smith, Jan tpnek, Alane Suhr, Zsolt Sznt, Takaaki Tanaka, Reut Tsarfaty, Sumire Uematsu, Larraitz Uria, Viktor Varga, Veronika Vincze, Zdenk abokrtsk, Daniel Zeman, and Hanzhi Zhu. 2015. Universal Dependencies 1.2. <http://universaldependencies.github.io/docs/>.
- 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.
- Feng Pan, Rutu Mulkar-Mehta, and Jerry R Hobbs. 2007. Modeling and learning vague event durations for temporal reasoning. In *Proceedings of the 22nd national conference on Artificial intelligence- Volume 2*, pages 1659–1662. AAAI Press.
- Fabian Pedregosa, Gal Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- 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.
- Rachel Rudinger, Aaron Steven White, and Benjamin Van Durme. 2018. Neural models of factuality. *arXiv preprint arXiv:1804.02472*.
- Roger C Schank and Robert P Abelson. 1975. Scripts, plans, and knowledge. In *IJCAI*, pages 151–157.
- Natalia Silveira, Timothy Dozat, Marie-Catherine De Marneffe, Samuel R Bowman, Miriam Connor, John Bauer, and Christopher D Manning. 2014. A gold standard dependency corpus for english. In *LREC*, pages 2897–2904.
- Gabriel Stanovsky, Judith Eckle-Kohler, Yevgeniy Puzikov, Ido Dagan, and Iryna Gurevych. 2017. Integrating deep linguistic features in factuality prediction over unified datasets. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 352–357.

- William F Styler IV, Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, et al. 2014. Temporal annotation in the clinical domain. *Transactions of the Association for Computational Linguistics*, 2:143.
- Julien Tourville, Olivier Ferret, Aurelie Neveol, and Xavier Tannier. 2017. Neural architecture for temporal relation extraction: A bi-lstm approach for detecting narrative containers. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 224–230.
- 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)*, volume 2, pages 1–9.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.
- 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 4th international workshop on semantic evaluations*, pages 75–80. 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. Association for Computational Linguistics.
- Aaron Steven White, Drew Reisinger, Keisuke Sakaguchi, Tim Vieira, Sheng Zhang, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. 2016. Universal decompositional semantics on universal dependencies. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1713–1723, Austin, TX. Association for Computational Linguistics.
- Aaron Steven White, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. 2018. Lexicosyntactic inference in neural models. *arXiv preprint arXiv:1808.06232*.
- Jennifer Williams and Graham Katz. 2012. Extracting and modeling durations for habits and events from twitter. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers-Volume 2*, pages 223–227. Association for Computational Linguistics.
- Deirdre Wilson and Dan Sperber. 1998. Pragmatics and time. *Pragmatics and Beyond New Series*, pages 1–22.
- Katsumasa Yoshikawa, Sebastian Riedel, Masayuki Asahara, and Yuji Matsumoto. 2009. Jointly identifying temporal relations with markov logic. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*, pages 405–413. Association for Computational Linguistics.
- Sheng Zhang, Rachel Rudinger, and Benjamin Van Durme. 2017. An evaluation of predpatt and open ie via stage 1 semantic role labeling. In *IWCS 201712th International Conference on Computational SemanticsShort papers*.