Temporal Information Extraction by Predicting Relative Time-lines

Artuur Leeuwenberg and Marie-Francine Moens

Department of Computer Science KU Leuven, Belgium

{tuur.leeuwenberg, sien.moens}@cs.kuleuven.be

Abstract

The current leading paradigm for temporal information extraction from text consists of three phases: (1) recognition of events and temporal expressions, (2) recognition of temporal relations among them, and (3) time-line construction from the temporal relations. In contrast to the first two phases, the last phase, time-line construction, received little attention and is the focus of this work. In this paper, we propose a new method to construct a linear time-line from a set of (extracted) temporal relations. But more importantly, we propose a novel paradigm in which we directly predict start and end-points for events from the text, constituting a time-line without going through the intermediate step of prediction of temporal relations as in earlier work. Within this paradigm, we propose two models that predict in linear complexity, and a new training loss using TimeML-style annotations, yielding promising results.

1 Introduction

The current leading perspective on temporal information extraction regards three phases: (1) a *temporal entity recognition* phase, extracting events (blue boxes in Fig. 1) and their attributes, and extracting temporal expressions (green boxes), and normalizing their values to dates or durations, (2) a *relation extraction* phase, where temporal links (TLinks) among those entities, and between events and the document-creation time (DCT) are found (arrows in Fig. 1, left). And (3), construction of a time-line (Fig. 1, right) from the extracted temporal links, if they are temporally consistent. Much research concentrated on the first two steps, but very little research looks into step 3, time-line construction, which is the focus of this work.

In this paper, we propose a new time-line construction paradigm that evades phase 2, the relation extraction phase, because in the classical paradigm temporal relation extraction comes with many difficulties in training and prediction that arise from the fact that for a text with n temporal entities (events or temporal expressions) there are n^2 possible entity pairs, which makes it likely for annotators to miss relations, and makes inference slow as n^2 pairs need to be considered. Temporal relation extraction models consistently give lower performance than those in the entity recognition phase (UzZaman et al., 2013; Bethard et al., 2016, 2017), introducing errors in the time-line construction pipe-line.

The ultimate goal of our proposed paradigm is to predict from a text in which entities are already detected, for each entity: (1) a probability distribution on the entity's starting point, and (2) another distribution on the entity's duration. The probabilistic aspect is crucial for time-line based decision making. Constructed time-lines allow for further quantitative reasoning with the temporal information, if this would be needed for certain applications.

As a first approach towards this goal, in this paper, we propose several initial time-line models in this paradigm, that directly predict - in a linear fashion - start points and durations for each entity, using text with annotated temporal entities as input (shown in Fig. 1). The predicted start points and durations constitute a relative time-line, i.e. a total order on entity start and end points. The time-line is relative, as start and duration values cannot (yet) be mapped to absolute calender dates or durations expressed in seconds. It represents the relative temporal order and inclusions that temporal entities have with respect to each other by the quantitative start and end values of the entities. Relative time-lines are a first step toward our goal, building models that predict statistical absolute time-lines. To train our relative time-line models, we define novel loss functions that exploit TimeML-style annotations, used in most existing temporal corpora. This work leads to the following contributions:

- A new method to construct a relative time-line from a set of temporal relations (TL2RTL).
- Two new models that, for the first time, directly predict (relative) time-lines in linear complexity from entity-annotated texts without doing a form of temporal relation extraction (S-TLM & C-TLM).
- Three new loss functions based on the mapping between Allen's interval algebra and the end-point algebra to train time-line models from TimeML-style annotations.

In the next sections we will further discuss the related work on temporal information extraction. We will describe the models and training losses in detail, and report on conducted experiments.

2 Related Work

2.1 Temporal Information Extraction

The way temporal information is conveyed in language has been studied for a long time. It can be conveyed directly through verb tense, explicit temporal discourse markers (e.g. during or afterwards) (Derczynski, 2017) or temporal expressions such as dates, times or duration expressions (e.g. 10-05-2010 or yesterday). Temporal information is also captured in text implicitly, through background knowledge about, for example, duration of events mentioned in the text (e.g. even without context, walks are usually shorter than journeys).

Most temporal corpora are annotated with TimeML-style annotations, of which an example is shown in Fig 1, indicating temporal entities, their attributes, and the TLinks among them.

The automatic extraction of TimeML-style temporal information from text using machine learning was first explored by Mani et al. (2006). They proposed a multinomial logistic regression classifier to predict the TLinks between entities. They also noted the problem of missed TLinks by annotators, and experimented with using temporal reasoning (temporal closure) to expand their training data.

Since then, much research focused on further improving the pairwise classification models, by exploring different types of classifiers and features, such as (among others) logistic regression and support vector machines (Bethard, 2013; Lin et al., 2015), and different types of neural network models, such as long short-term memory networks (LSTM) (Tourille et al., 2017; Cheng and Miyao, 2017), and convolutional neural networks (CNN) (Dligach et al., 2017). Moreover, different sievebased approaches were proposed (Chambers et al., 2014; Mirza and Tonelli, 2016), facilitating mixing of rule-based and machine learning components.

Two major issues shared by these existing approaches are: (1) models classify TLinks in a pairwise fashion, often resulting in an inference complexity of $O(n^2)$, and (2) the pair-wise predictions are made independently, possibly resulting in prediction of temporally inconsistent graphs. To address the second, additional temporal reasoning can be used at the cost of computation time, during inference (Chambers and Jurafsky, 2008; Denis and Muller, 2011; Do et al., 2012), or during both training and inference (Yoshikawa et al., 2009; Laokulrat et al., 2015; Ning et al., 2017; Leeuwenberg and Moens, 2017). In this work, we circumvent these issues, as we predict time-lines - in linear time complexity - that are temporally consistent by definition.

2.2 Temporal Reasoning

Temporal reasoning plays a central role in temporal information extraction, and there are roughly two approaches: (1) Reasoning directly with Allen's interval relations (shown in Table 1), by constructing rules like: If event X occurs before Y, and event Y before Z then X should happen before Z (Allen, 1990). Or (2), by first mapping the temporal interval expressions to expressions about interval end-points (start and endings of entities) (Vilain et al., 1990). An example of such mapping is that If event X occurs before Y then the end of X should be before the start of Y. Then reasoning can be done with end-points in a point algebra, which has only three point-wise relations (=,<,>), making reasoning much more efficient compared to reasoning with Allen's thirteen interval relations.

Mapping interval relations to point-wise expressions has been exploited for model inference by Denis and Muller (2011), and for evaluation by UzZaman and Allen (2011). In this work, we ex-

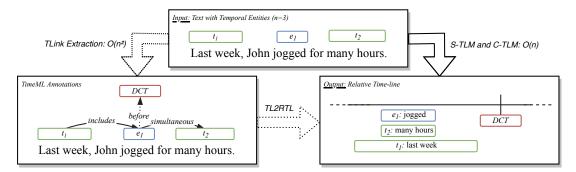


Figure 1: An overview of two paradigms: (1) The indirect approach (dashed arrows), where first TLinks are predicted from which we can build a relative time-line using TL2RTL. And (2), the direct approach (solid arrow), where a relative time-line is predicted directly from the input by S-TLM or C-TLM.

ploit it for the first time for model training, in our loss functions.

3 Models

We propose two model structures for direct time-line construction: (1) a simple context-independent model (S-TLM), and (2) a contextual model (C-TLM). Their structures are shown in Fig. 2. Additionally, we propose a method to construct relative time-lines from a set of (extracted) TLinks (TL2RTL). In this section we first explain the first two direct models S-TLM and C-TLM, and afterwards the indirect method TL2RTL.

3.1 Direct Time-line Models

Word representation

In both S-TLM and C-TLM, words are represented as a concatenation of a word embedding, a POS embedding, and a Boolean feature vector containing entity attributes such as the type, class, aspect, following (Do et al., 2012). Further details on these are given in the experiments section.

Simple Time-line Model (S-TLM)

For the simple context-independent time-line model, each entity is encoded by the word representation of the last word of the entity (generally the most important). From this representation we have a linear projection to the duration d, and the start s. S-TLM is shown by the dotted edges in Fig 2. An advantage of S-TLM is that it has very few parameters, and each entity can be placed on the time-line independently of the others, allowing parallelism during prediction. The downside is that S-TLM is limited in its use of contextual information.

Contextual Time-line Model (C-TLM)

To better exploit the entity context we also propose a contextual time-line model C-TLM (solid edges in Fig 2), that first encodes the full text using two bi-directional recurrent neural networks, one for entity starts (BiRNN $_s$), and one for entity durations (BiRNN $_d$). On top of the encoded text we learn two linear mappings, one from the BiRNN $_d$ output of the last word of the entity mention to its duration d, and similarly for the start time, from the BiRNN $_s$ output to the entity's start s.

Predicting Start, Duration, and End

Both proposed models use linear mappings² to predict the start value s_i and duration d_i for the encoded entity i. By summing start s_i and duration d_i we can calculate the entity's end-point e_i .

$$e_i = s_i + \max(d_i, d_{min}) \tag{1}$$

Predicting durations rather than end-points makes it easy to control that the end-point lies after the start-point by constraining the duration d_i by a constant minimum duration value d_{min} above 0, as shown in Eq. 1.

Modeling Document-Creation Time

Although the DCT is often not found explicitly in the text, it is an entity in TimeML, and has TLinks to other entities. We model it by assigning it a text-independent start $s_{\rm DCT}$ and duration $d_{\rm DCT}$.

Start s_{DCT} is set as a constant (with value 0). This way the model always has the same reference point, and can learn to position the entities w.r.t. the DCT on the time-line.

 $^{^{1}}$ We also experimented with sharing weights among BiRNN $_d$ and BiRNN $_s$. In our experiments, this gave worse performance, so we propose to keep them separate.

²Adding more layers did not improve results.

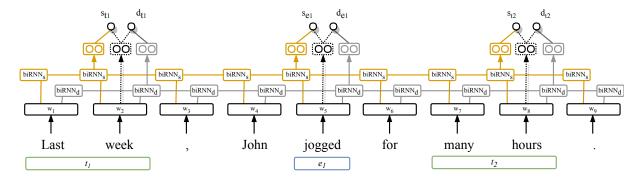


Figure 2: Schematic overview of our two time-line models: C-TLM (solid edges), exploiting entity context, and the simpler S-TLM (dotted edges), which is context independent. The models predict a starting point (s) and duration (d) for each given temporal entity $(t_1, e_1, and t_2)$ in the input.

In contrast, DCT duration $d_{\rm DCT}$ is modeled as a single variable that is learned (initialized with 1). Since multiple entities may be included in the DCT, and entities have a minimum duration d_{min} , a constant $d_{\rm DCT}$ could possibly prevent the model from fitting all entities in the DCT. Modeling $d_{\rm DCT}$ as a variable allows growth of $d_{\rm DCT}$ and averts this issue.³

Training Losses

We propose three loss functions to train time-line models from TimeML-style annotations: a regular time-line loss L_{τ} , and two slightly expanded discriminative time-line losses, $L_{\tau ce}$ and $L_{\tau h}$.

Regular Time-line Loss (L_{τ})

Ground-truth TLinks can be seen as constraints on correct positions of entities on a time-line. The regular time-line loss L_{τ} expresses the degree to which these constraints are met for a predicted time-line. If all TLinks are satisfied in the time-line for a certain text, L_{τ} will be 0 for that text.

As TLinks relate entities (intervals), we first convert the TLinks to expressions that relate the start and end points of entities. How each TLink is translated to its corresponding point-algebraic constraints is given in Table 1, following Allen (1990).

As can be seen in the last column there are only two point-wise operations in the point-algebraic constraints: an order operation (<), and an equality operation (=). To model to what degree each point-wise constraint is met, we employ hinge losses, with a margin m_{τ} , as shown in Eq. 2.

Table 1: Point algebraic interpretation (I_{PA}) of temporal links used to construct the loss function. The start and end points of event X are indicated by s_x and e_x respectively.

Allen Algebra	Temporal Links	Point Algebra	
X precedes Y Y preceded by X	X before Y Y after X	$e_x < s_y$	
X starts Y Y started by X	X begins Y Y begun by X	$s_x = s_y$ $e_x < e_y$	
X finishes Y Y finished by X	X ends Y Y ended by X	$e_x = e_y$ $s_y < s_x$	
X during Y Y includes X	X is included Y Y includes X	$s_y < s_x \\ e_x < e_y$	
X meets Y Y met by X	X immediately before Y Y immediately after X	$e_x = s_y$	
X overlaps Y Y overlapped by X	absent ⁴ absent ⁴	$s_x < s_y$ $s_y < e_x$ $e_x < e_y$	
X equals Y	X simultaneous Y X identity Y	$s_x = s_y$ $e_x = e_y$	

To explain the intuition and notation: If we have a point-wise expression ξ of the form x < y (first case of Eq. 2), then the predicted point \hat{x} should be at least a distance m_{τ} smaller (or earlier on the time-line) than predicted point \hat{y} in order for the loss to be 0. Otherwise, the loss represents the distance \hat{x} or \hat{y} still has to move to make \hat{x} smaller than \hat{y} (and satisfy the constraint). For the second case, if ξ is of the form x = y, then point \hat{x} and \hat{y} should lie very close to each other, i.e. at most a distance m_{τ} away from each other. Any distance further than the margin m_{τ} is counted as loss. Notice that if we set margin m_{τ} to 0, the second case becomes an L1 loss $|\hat{x} - \hat{y}|$. However, we use a small margin m_{τ} to promote some distance between ordered points and prevent con-

 $^{^3}$ Other combinations of modeling $s_{
m DCT}$ and $d_{
m DCT}$ as variable or constant decreased performance.

⁴No TLink for Allen's overlap relation is present in TimeML, also concluded by UzZaman and Allen (2011).

fusion with equality. Fig. 3 visualizes the loss for three TLinks.

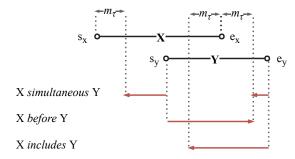


Figure 3: Visualization of the time-line loss L_{τ} with margin m_{τ} , for two events X and Y, and TLinks *simultaneous*, *before*, and *includes*. The red arrows' lengths indicate the loss per relation, i.e. how much the points should be shifted to satisfy each relation.

$$L_p(\xi|t,\theta) = \begin{cases} \max(\hat{x} + m_{\tau} - \hat{y}, 0) & \text{iff } x < y \\ \max(|\hat{x} - \hat{y}| - m_{\tau}, 0) & \text{iff } x = y \end{cases}$$
(2)

The total time-line loss $L_{\tau}(t|\theta)$ of a model with parameters θ on text t with ground-truth TLinks R(t), is the sum of the TLink-level losses of all TLinks $r \in R(t)$. Each TLink-level loss $L_r(r|t,\theta)$ for TLink r is the sum of the pointwise losses $L_p(\xi|t,\theta)$ of the corresponding pointalgebraic constraints $\xi \in I_{PA}(r)$ from Table 1.5

$$L_r(r|t,\theta) = \sum_{\xi \in I_{PA}(r)} L_p(\xi|t,\theta)$$
 (3)

$$L_{\tau}(t,\theta) = \sum_{r \in R(t)} L_r(r|t,\theta)$$
 (4)

Discriminative Time-line Losses

To promote a more explicit difference between the relations on the time-line we introduce two discriminative loss functions, $L_{\tau ce}$ and $L_{\tau h}$, which build on top of L_r . Both discriminative loss functions use an intermediate score $S(r|t,\theta)$ for each TLink r based on the predicted time-line. As scoring function, we use the negative L_r loss, as shown in Eq. 5.

$$S(r|t,\theta) = -L_r(r|t,\theta) \tag{5}$$

Then, a lower time-line loss $L_r(r|t,\theta)$ results in a higher score for relation type r. Notice that the maximum score is 0, as this is the minimum L_r .

Probabilistic Loss ($L_{\tau ce}$)

Our first discriminative loss is a cross-entropy based loss. For this the predicted scores are normalized using a soft-max over the possible relation types (TL). The resulting probabilities are used to calculate a cross-entropy loss, shown in Eq. 6. This way, the loss does not just promote the correct relation type but also distantiates from the other relation types.

$$L_{\tau ce}(t|\theta) = \sum_{r \in R(t)} r \cdot \log\left(\frac{e^{S(r|t,\theta)}}{\sum_{r' \in TL} e^{S(r'|t,\theta)}}\right)$$
(6)

Ranking Loss $(L_{\tau h})$

When interested in discriminating relations on the time-line, we want the correct relation type to have the highest score from all possible relation types TL. To represent this perspective, we also define a ranking loss with a score margin m_h in Eq. 7.

$$L_{\tau h}(t|\theta) = \sum_{r \in R(t)} \sum_{r' \in TL \setminus \{r\}} \max(S(r'|t,\theta) - S(r|t,\theta) + m_h, 0)$$
(7)

Training Procedure

S-TLM and C-TLM are trained by by iterating through the training texts, sampling mini-batches of 32 annotated TLinks. For each batch we (1) perform a forward pass, (2) calculate the total loss (for one of the loss functions), (3) derive gradients using Adam⁶ (Kingma and Ba, 2014), and (4) update the model parameters θ via back-propagation. After each epoch we shuffle the training texts. As stopping criteria we use early stopping (Morgan and Bourlard, 1990), with a patience of 100 epochs and a maximum number of 1000 epochs.

3.2 From TLinks to Time-lines (TL2RTL)

To model the indirect route, we construct a novel method, TL2RTL, that predicts relative time lines from a subset of TLinks, shown in Fig 1. One can choose any method to obtain a set of TLinks R(t) from a text t, serving as input to TL2RTL.

⁵The TLink *during* and its inverse are mapped to *simultaneous*, following the evaluation of TempEval-3.

⁶Using the default parameters from the paper.

TL2RTL constructs a relative time-line, by assigning start and end values to each temporal entity, such that the resulting time-line satisfies the extracted TLinks R(t) by minimizing a loss function that is 0 when the extracted TLinks are satisfied. TL2RTL on itself is a method and not a model. The only variables over which it optimizes the loss are the to be assigned starts and duration values.

In detail, for a text t, with annotated entities E(t), we first extract a set of TLinks R(t). In this work, to extract TLinks, we use the current state-of-the-art structured TLink extraction model by Ning et al. (2017). Secondly, we assign a start variable s_i , and duration variable d_i to each entity $i \in E(t)$. Similar to S-TLM and C-TLM, for each $i \in E(t)$, d_i is bounded by a minimum duration d_{min} to ensure start s_i always lies before end e_i . Also, we model the DCT start s_{DCT} as a constant, and its duration d_{DCT} as a variable. Then we minimize one of the loss functions L_{τ} , $L_{\tau ce}$, or $L_{\tau h}$ on the extracted TLinks R(t), obtaining three TL2RTL variants, one for each loss. If the initially extracted set of TLinks R(t) is consistent, and the loss is minimized sufficiently, all s_i and d_i form a relative time-line that satisfies the TLinks R(t), but from which we can now also derive consistent TLinks for any entity pair, also the pairs that were not in R(t). To minimize the loss we use Adam for 10k epochs until the loss is zero for each document.⁷

4 Experiments

4.1 Evaluation and Data

Because prediction of relative time-lines trained on TimeML-style annotations is new, we cannot compare our model directly to relation extraction or classification models, as the latter do not provide completely temporally consistent TLinks for all possible entity pairs, like the relative time-lines do. Neither can we compare directly to existing absolute time-line prediction models such as Reimers et al. (2018) because they are trained on different data with a very different annotation scheme.

To evaluate the quality of the relative time-line models in a fair way, we use TimeML-style test sets as follows: (1) We predict a time-line for each test-text, and (2) we check for all ground-truth an-

notated TLinks that are present in the data, what would be the derived relation type based on the predicted time-line, which is the relation type that gives the lowest time-line loss L_r . This results in a TLink assignment for each annotated pair in the TimeML-style reference data, and therefor we can use similar metrics. As evaluation metric we employ the temporal awareness metric, used in TempEval-3, which takes into account temporal closure (UzZaman et al., 2013). Notice that although we use the same metric, comparison against relation classification systems would be unfair, as our model assigns consistent labels to all pairs, whereas relation classification systems do not.

For training and evaluation we use two data splits, TE[‡] and TD[‡], exactly following Ning et al. (2017). Some statistics about the data are shown in Table 2.⁸ The splits are constituted from various smaller datasets: the TimeBank (TB) (Pustejovsky et al., 2002), the AQUANT dataset (AQ), and the platinum dataset (PT) all from TempEval-3 (Uz-Zaman et al., 2013). And, the TimeBank Dense (Cassidy et al., 2014), and the Verb-Clause dataset (VC) (Bethard et al., 2007).

4.2 Hyper-parameters and Preprocessing

Hyper-parameters shared in all settings can be found in Table 3. The following hyper-parameters are tuned using grid search on a development set (union of TB and AQ): d_{min} is chosen from $\{1,0.1,0.01\},\ m_{\tau}$ from $\{0,0.025,0.05,0.1\},$ α_d from $\{0,0.1,0.2,0.4,0.8\},$ and α_{rnn} from $\{10,25,50\}.$ We use LSTM (Hochreiter and Schmidhuber, 1997) as RNN units 9 and employ 50-dimensional GloVe word-embeddings pre-trained 10 on 6B words (Wikipedia and NewsCrawl) to initialize the models' word embeddings.

We use very simple tokenization and consider punctuation¹¹ or newline tokens as individual tokens, and split on spaces. Additionally, we lowercase the text and use the Stanford POS Tagger (Toutanova et al., 2003) to obtain POS.

 $^{^7}$ For some documents the extracted TLinks were temporally inconsistent, resulting in a non-zero loss. Nevertheless, > 96% of the extracted TLinks were satisfied.

⁸We explicitly excluded all test documents from training as some corpora annotated the same documents.

⁹We also experimented with GRU as RNN type, obtaining similar results.

¹⁰https://nlp.stanford.edu/projects/glove

¹¹,./\"'=+-;:()!?<>%&\$*|[]{}

Table 2: Dataset splits used for evaluation (indicated with ‡).

Split	Training data	#TLinks	#Documents	Test data	#TLinks	#Documents
TD^{\ddagger}	TD (train+dev)	4.4k	27	TD (test)	1.3k	9
TE3 [‡]	TB, AQ, VC, TD (full)	17.5k	256	PT	0.9k	20

Table 3: Hyper-parameters from the experiments.

Hyper-parameter	Value
Document-creation starting time (s_{DCT})	0
Minimum event duration (d_{min})	0.1
Time-line margin (m_{τ})	0.025
Hinge loss margin (m_h)	0.1
Dropout (α_d)	0.1
Word-level RNN units (α_{rnn})	25
Word-embedding size (α_{wemb})	50
POS-embedding size	10

Table 4: Evaluation of relative time-lines for each model and loss function, where L_* indicates the (unweighted) sum of L_{τ} , $L_{\tau ce}$, and $L_{\tau h}$.

	TE3 [‡]			$\mathbf{T}\mathbf{D}^{\ddagger}$		
Model	P	R	F	P	R	F
Indirect: $O(n^2)$						
TL2RTL (L_{τ})	53.5	51.1	52.3	59.1	61.2	60.1
TL2RTL ($L_{\tau ce}$)	53.9	51.7	52.8	61.2	60.7	60.9
TL2RTL $(L_{\tau h})$	52.8	51.1	51.9	57.9	60.6	59.2
TL2RTL (L_*)	52.6	52.0	52.3	62.3	62.3	62.3
Direct: $O(n)$						
S-TLM (L_{τ})	50.1	50.4	50.2	57.8	59.5	58.6
S-TLM ($L_{\tau ce}$)	50.1	50.0	50.1	53.4	53.5	53.5
S-TLM $(L_{\tau h})$	51.5	51.7	51.6	55.1	56.4	55.7
S-TLM (L_*)	50.9	51.0	51.0	56.5	55.3	55.9
C-TLM (L_{τ})	56.2	56.1	56.1	57.1	59.7	58.4
C-TLM ($L_{\tau ce}$)	54.4	55.4	54.9	52.4	57.3	54.7
C-TLM $(L_{\tau h})$	55.7	55.5	55.6	55.3	54.9	55.1
C-TLM (L_*)	54.0	54.3	54.1	54.6	53.5	54.1

5 Results

We compared our three proposed models for the three loss functions L_{τ} , $L_{\tau ce}$, and $L_{\tau h}$, and their linear (unweighted) combination L_* , on TE3^{\ddagger} and TD^{\ddagger}, for which the results are shown in Table 4.

A trend that can be observed is that overall performance on TD^{\ddagger} is higher than that of $TE3^{\ddagger}$, even though less documents are used for training. We inspected why this is the case, and this is caused by a difference in class balance between both test sets. In $TE3^{\ddagger}$ there are many more TLinks of type *simultaneous* (12% versus 3%), which are very

difficult to predict, resulting in lower scores for TE3[‡] compared to TD[‡]. The difference in performance between the datasets is probably also be related to the dense annotation scheme of TD[‡] compared to the sparser annotations of TE3[‡], as dense annotations give a more complete temporal view of the training texts. For TL2RTL better TLink extraction¹² is also propagated into the final timeline quality.

If we compare loss functions L_{τ} , $L_{\tau ce}$, and $L_{\tau h}$, and combination L_{*} , it can be noticed that, although all loss functions seem to give fairly similar performance, L_{τ} gives the most robust results (never lowest), especially noticeable for the smaller dataset TD^{\ddagger} . This is convenient, because L_{τ} is fastest to compute during training, as it requires no score calculation for each TLink type. L_{τ} is also directly interpretable on the timeline. The combination of losses L_{*} shows mixed results, and has lower performance for S-TLM and C-TLM, but better performance for TL2RTL. However, it is slowest to compute, and less interpretable, as it is a combined loss.

Moreover, we can clearly see that on TE3 ‡ , C-TLM performs better than the indirect models, across all loss functions. This is a very interesting result, as C-TLM is an order of complexity faster in prediction speed compared to the indirect models (O(n)) compared to $O(n^2)$ for a text with n entities). We further explore why this is the case through our error analysis in the next section.

On TD[‡], the indirect models seem to perform slightly better. We suspect that the reason for this is that C-TLM has more parameters (mostly the LSTM weights), and thus requires more data (TD[‡] has much fewer documents than TE3[‡]) compared to the indirect methods. Another result supporting this hypothesis is the fact that the difference between C-TLM and S-TLM is small on the smaller

overall performance decreases consistently with 2-4 points.

 $^{^{12}}$ F1 of 40.3 for TE3[‡] and 48.5 for TD[‡] (Ning et al., 2017) 13 We do not directly compare prediction speed, as it would result in unfair evaluation because of implementation differences. However, currently, C-TLM predicts at \sim 100 w/s incl. POS tagging, and \sim 2000 w/s without. When not using POS,

	В	A	II	I	S
В	24.8%	4.7%	2.8%	1.6%	0.1%
\mathbf{A}	5.0%	15.8%	3.2%	0.5%	0.0%
II	3.2%	3.2%	13.0%	0.6%	0.1%
I	4.0%	1.2%	1.0%	3.2%	0.0%
\mathbf{S}	4.4%	3.0%	2.6%	1.3%	0.4%

	В	A	II	I	S
В	23.0%	8.2%	1.3%	0.9%	0.8%
A	4.7%	17.1%	1.8%	0.3%	0.5%
II	4.3%	4.4%	11.1%	0.4%	0.0%
I	1.6%	5.4%	0.5%	1.3%	0.5%
S	4.3%	4.1%	1.8%	0.6%	0.9%

Figure 4: On the **left**, the confusion matrix of C-TLM (L_{τ}) , and on the **right** of TL2RTL $(L_{\tau ce})$, on TE3[‡] for the top-5 most-frequent TLinks (together 95% of data): BEFORE (B), AFTER (A), IS INCLUDED (II), INCLUDES (I), and SIMULTANEOUS (S). Predictions are shown on the x-axis and ground-truth on the y-axis.

TD[‡], indicating that C-TLM does not yet utilize contextual information from this dataset, whereas, in contrast, on TE3[‡], the larger dataset, C-TLM clearly outperforms S-TLM across all loss functions, showing that when enough data is available C-TLM learns good LSTM weights that exploit context substantially.

6 Error Analysis

We compared predictions of $\text{TL2RTL}(L_{\tau})$ with those of C-TLM (L_{τ}) , the best models of each paradigm. In Table 4, we show the confusion matrices of both systems on TE3^{\ddagger} .

When looking at the overall pattern in errors, both models seem to make similar confusions on both datasets (TD[‡] was excluded for space constraints). Overall, we find that *simultaneous* is the most violated TLink for both models. This can be explained by two reasons: (1) It is the least frequent TLink in both datasets. And (2), simultaneous entities are often co-referring events. Event co-reference resolution is a very difficult task on its own.

We also looked at the average token-distance between arguments of correctly satisfied TLinks by the time-lines of each model. For TL2RTL (L_{τ}) this is 13 tokens, and for C-TLM (L_{τ}) 15. When looking only at the TLinks that C-TLM (L_{τ}) satisfied and TL2RTL (L_{τ}) did not, the average distance is 21. These two observations suggest that the direct C-TLM (L_{τ}) model is better at positioning entities on the time-line that lie further away from each other in the text. An explanation for this can be error propagation of TLink extraction to the time-line construction, as the pairwise TLink extraction of the indirect paradigm extracts TLinks in a contextual window, to prune the $O(n^2)$ number of possible TLink candidates. This

Table 5: Example events from the top-shortest/longest durations and top-earliest/latest start values assigned by the model.

Short d	$\mathbf{Long}\ d$	Early s	Late s
started	going	destroyed	realize
meet	expects	finished	bring
entered	recession	invaded	able
told	war	pronounced	got
arrived	support	created	work
allow	make	took	change
send	think	appeared	start
asked	created	leaving	reenergize

consequently prevents TL2RTL to properly position distant events with respect to each other.

To get more insight in what the model learns we calculated mean durations and mean starts of C-TLM (L_{τ}) predictions. Table 5 contains examples from the top-shortest, and top-longest duration assignments and earliest and latest starting points. We observe that events that generally have more events included are assigned longer duration and vice versa. And, events with low start values are in the past tense and events with high start values are generally in the present (or future) tense.

7 Discussion

A characteristic of our model is that it assumes that all events can be placed on a single timeline, and that it does not assume that unlabeled pairs are temporally unrelated. This has big advantages: it results in fast prediction, and missed annotation do not act as noise to the training, as they do for pairwise models. Ning et al. (2018) argue that actual, negated, hypothesized, expected or opinionated events should possibly be annotated

on separate time-axis. We believe such multi-axis representations can be inferred from the generated single time-lines if hedging information is recognized.

8 Conclusions

This work leads to the following three main contributions¹⁴: (1) Three new loss functions that connect the interval-based TimeML-annotations to points on a time-line, (2) A new method, TL2RTL, to predict relative time-lines from a set of predicted temporal relations. And (3), most importantly, two new models, S-TLM and C-TLM, that - to our knowledge for the first time - predict (relative) time-lines in linear complexity from text, by evading the computationally expensive (often $O(n^2)$) intermediate relation extraction phase in earlier work. From our experiments, we conclude that the proposed loss functions can be used effectively to train direct and indirect relative time-line models, and that, when provided enough data, the - much faster - direct model C-TLM outperforms the indirect method TL2RTL.

As a direction for future work, it would be very interesting to extend the current models, diving further into direct time-line models, and learn to predict absolute time-lines, i.e. making the time-lines directly mappable to calender dates and times, e.g. by exploiting complementary data sources such as the EventTimes Corpus (Reimers et al., 2016) and extending the current loss functions accordingly. The proposed models also provide a good starting point for research into probabilistic time-line models, that additionally model the (un)certainty of the predicted positions and durations of the entities.

Acknowledgments

The authors thank Geert Heyman and the reviewers for their constructive comments which helped us to improve the paper. This work was funded by the KU Leuven C22/15/16 project "MAchine Reading of patient recordS (MARS)", and by the IWT-SBO 150056 project "ACquiring CrUcial Medical information Using LAnguage TEchnology" (ACCUMULATE).

References

- James F Allen. 1990. Maintaining knowledge about temporal intervals. *Readings in Qualitative Reasoning about Physical Systems*, pages 361–372.
- Steven Bethard. 2013. ClearTK-TimeML: A minimalist approach to tempeval 2013. In *Proceedings of SemEval*, volume 2, pages 10–14. ACL.
- Steven Bethard, James H. Martin, and Sara Klingenstein. 2007. Timelines from text: Identification of syntactic temporal relations. In *Proceedings of ICSC*, pages 11–18.
- Steven Bethard, Guergana Savova, Wei-Te Chen, Leon Derczynski, James Pustejovsky, and Marc Verhagen. 2016. Semeval-2016 task 12: Clinical tempeval. In *Proceedings of SemEval*, pages 1052–1062. ACL.
- Steven Bethard, Guergana Savova, Martha Palmer, and James Pustejovsky. 2017. SemEval-2017 Task 12: Clinical TempEval. In *Proceedings of SemEval*, pages 565–572. ACL.
- Taylor Cassidy, Bill McDowell, Nathanael Chambers, and Steven Bethard. 2014. An annotation framework for dense event ordering. In *Proceedings of ACL*, pages 501–506. ACL.
- 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 EMNLP*, pages 698– 706. ACL.
- Fei Cheng and Yusuke Miyao. 2017. Classifying temporal relations by bidirectional LSTM over dependency paths. In *Proceedings of ACL*, volume 2, pages 1–6. ACL.
- Pascal Denis and Philippe Muller. 2011. Predicting globally-coherent temporal structures from texts via endpoint inference and graph decomposition. In *Proceedings of IJCAI*, pages 1788–1793.
- Leon RA Derczynski. 2017. Automatically Ordering Events and Times in Text, volume 677. Springer.
- Dmitriy Dligach, Timothy Miller, Chen Lin, Steven Bethard, and Guergana Savova. 2017. Neural temporal relation extraction. In *Proceedings of EACL*, volume 2, pages 746–751.
- Quang Xuan Do, Wei Lu, and Dan Roth. 2012. Joint inference for event timeline construction. In *Proceedings of EMNLP-CoNLL*, pages 677–687. ACL.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.

¹⁴Code is available at: liir.cs.kuleuven.be/software.php

- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Natsuda Laokulrat, Makoto Miwa, and Yoshimasa Tsuruoka. 2015. Stacking approach to temporal relation classification. *Journal of Natural Language Processing*, 22(3):171–196.
- Artuur Leeuwenberg and Marie-Francine Moens. 2017. Structured learning for temporal relation extraction from clinical records. In *Proceedings of EACL*, volume 1, pages 1150–1158. ACL.
- Chen Lin, Dmitriy Dligach, Timothy A Miller, Steven Bethard, and Guergana K Savova. 2015. Multilayered temporal modeling for the clinical domain. *Journal of the American Medical Informatics Association*, 23(2):387–395.
- Inderjeet Mani, Marc Verhagen, Ben Wellner, Chong Min Lee, and James Pustejovsky. 2006. Machine learning of temporal relations. In *Proceedings* of COLING-ACL, pages 753–760. ACL.
- Paramita Mirza and Sara Tonelli. 2016. CATENA: Causal and temporal relation extraction from natural language texts. *Proceedings of COLING*, pages 64–75.
- Nelson Morgan and Hervé Bourlard. 1990. Generalization and parameter estimation in feedforward nets: Some experiments. In *Advances in Neural Information Processing Systems*.
- Qiang Ning, Zhili Feng, and Dan Roth. 2017. A structured learning approach to temporal relation extraction. *Proceedings of EMNLP*, pages 1038–1048.
- Qiang Ning, Hao Wu, and Dan Roth. 2018. A multiaxis annotation scheme for event temporal relations. In *Proceedings of ACL*, pages 1318–1328. ACL.
- James Pustejovsky, Patrick Hanks, Roser Saurí, Andrew See, Robert Gaizauskas, Andrea Setzer, Beth Sundheim, David Day, Lisa Ferro, and Dragomir. 2002. The TIMEBANK Corpus. *Natural Language Processing and Information Systems*, 4592:647–656.
- Nils Reimers, Nazanin Dehghani, and Iryna Gurevych. 2016. Temporal anchoring of events for the timebank corpus. *Proceedings of ACL*, pages 2195– 2204.
- Nils Reimers, Nazanin Dehghani, and Iryna Gurevych. 2018. Event time extraction with a decision tree of neural classifiers. *Transactions of the Association for Computational Linguistics*, 6:77–89.
- Julien Tourille, 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 ACL*, pages 224–230.

- Kristina Toutanova, Dan Klein, Christopher D Manning, and Yoram Singer. 2003. Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of NAACL-HLT*, pages 173–180. ACL.
- Naushad UzZaman and James F. Allen. 2011. Temporal evaluation. In *Proceedings of ACL*, HLT '11, pages 351–356, Stroudsburg, PA, USA. ACL.
- Naushad UzZaman, Hector Llorens, Leon Derczynski, Marc Verhagen, James Allen, and James Puste-jovsky. 2013. Semeval-2013 task 1: Tempeval-3: Evaluating time expressions, events, and temporal relations. Second joint conference on lexical and computational semantics (* SEM), 2:1–9.
- Marc Vilain, Henry Kautz, and Peter Van Beek. 1990. Constraint propagation algorithms for temporal reasoning: A revised report. In *Readings in Qualitative Reasoning about Physical Systems*, pages 373–381. Elsevier.
- Katsumasa Yoshikawa, Sebastian Riedel, Masayuki Asahara, and Yuji Matsumoto. 2009. Jointly identifying temporal relations with markov logic. In *Proceedings of ACL-IJCNLP*, pages 405–413. ACL.