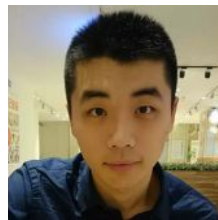
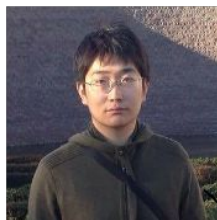




# Joint Constrained Learning for Event-Event Relation Extraction



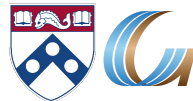
Haoyu Wang<sup>1</sup>, Muhao Chen<sup>1</sup>, Hongming Zhang<sup>2\*</sup>, Dan Roth<sup>1</sup>

<sup>1</sup>Department of Computer and Information Science, UPenn

<sup>2</sup>Department of Computer Science and Engineering, HKUST

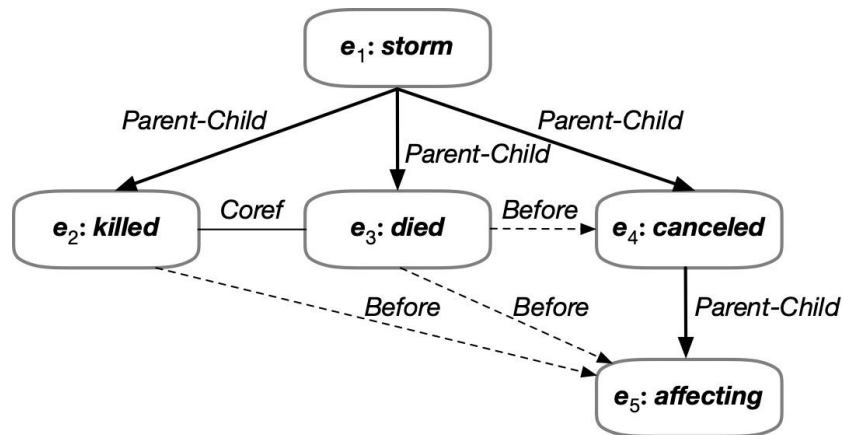
\* This work was done when the author was visiting the University of Pennsylvania

# Event-Event Relation Extraction



- Events are not simple, stand-alone predicates
  - Described at different granularities & may form complex structures
- Event-event relations are needed to induce such “event complex”
  - Temporal Relations
  - Membership
  - Coreference

On Tuesday, there was a typhoon-strength ( $e_1$ :*storm*) in Japan. One man got ( $e_2$ :*killed*) and thousands of people were left stranded. Police said an 81-year-old man ( $e_3$ :*died*) in central Toyama when the wind blew over a shed, trapping him underneath. Later this afternoon, with the agency warning of possible tornadoes, Japan Airlines ( $e_4$ :*canceled*) 230 domestic flights, ( $e_5$ :*affecting*) 31,600 passengers.



# Problem Definition

---



**Input:** document, event triggers

**Output:** event-event relations

## Tasks we study:

- Temporal Relations: Before, After, Equal, Vague
- Membership Relations: Parent-Child, Child-Parent, Coreference, NoRel

## Challenges

- Lack of learning resources
  - No large-scale single resource contains annotation for these interrelated tasks
- Global consistency among interrelated predictions
  - Symmetry constraints
  - Transitivity constraints
  - Conjunction constraints

# Problem Definition



**Input:** document, event triggers

**Output:** event-event relations

$e_2$ : *killed*

$e_3$ : *died*

$e_4$ : *canceled*

**Tasks we study:**

$e_5$ : *affecting*

- Temporal Relations: Before, After, Equal, Vague
- Membership Relations: Parent-Child, Child-Parent, Coreference, NoRel

## Challenges

- Lack of learning resources
  - No large-scale single resource contains annotation for these interrelated tasks
- Global consistency among interrelated predictions
  - Symmetry constraints
  - Transitivity constraints
  - Conjunction constraints

# Problem Definition

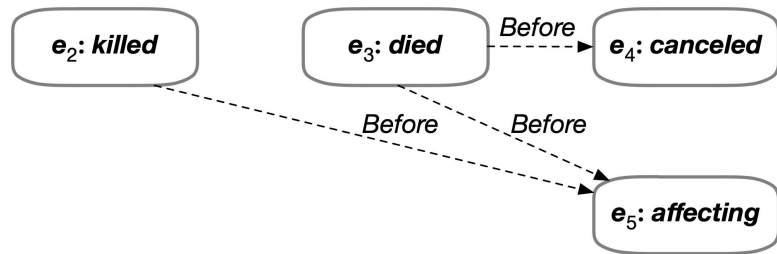


**Input:** document, event triggers

**Output:** event-event relations

**Tasks we study:**

- Temporal Relations: Before, After, Equal, Vague
- Membership Relations: Parent-Child, Child-Parent, Coreference, NoRel



**Challenges**

- Lack of learning resources
  - No large-scale single resource contains annotation for these interrelated tasks
- Global consistency among interrelated predictions
  - Symmetry constraints
  - Transitivity constraints
  - Conjunction constraints

# Problem Definition

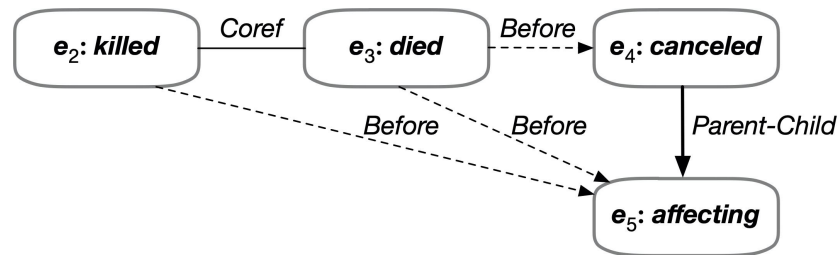


**Input:** document, event triggers

**Output:** event-event relations

**Tasks we study:**

- Temporal Relations: Before, After, Equal, Vague
- Membership Relations: Parent-Child, Child-Parent, Coreference, NoRel



**Challenges**

- Lack of learning resources
  - No large-scale single resource contains annotation for these interrelated tasks
- Global consistency among interrelated predictions
  - Symmetry constraints
  - Transitivity constraints
  - Conjunction constraints

# Constraints



## Symmetry

***e3:died*** is BEFORE ***e4: canceled***  
 => ***e4: canceled*** is AFTER ***e3:died***

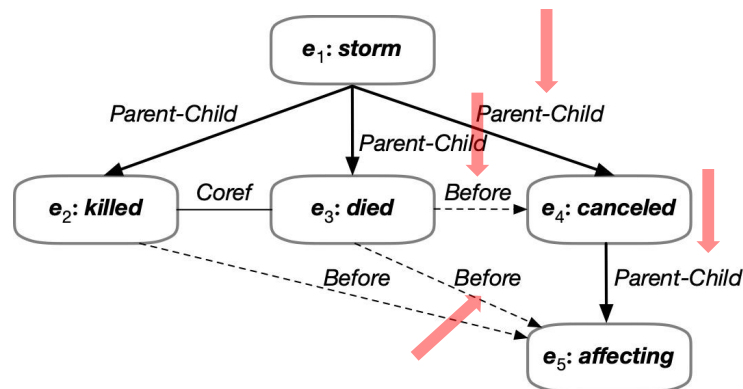
## Conjunction

***e3:died*** is BEFORE ***e4: canceled***  
 + ***e4: canceled*** is a PARENT of ***e5: affecting***  
 => ***e3:died*** BEFORE ***e5: affecting***

## Transitivity

***e1: storm*** is PARENT of ***e4: canceled***  
 + ***e4: canceled*** is a PARENT of ***e5: affecting***  
 => ***e1: storm*** is a PARENT of ***e5: affecting***

$\alpha(e3, e4) = \text{BEFORE}$   
 $\beta(e4, e5) = \text{Parent-Child}$   
 =>  $\gamma(e3, e5) = \text{BEFORE}$



$\alpha \backslash \beta$	PC	CP	CR	NR	BF	AF	EQ	VG
PC	PC, $\neg$ AF	–	PC, $\neg$ AF	$\neg$ CP, $\neg$ CR	BF, $\neg$ CP, $\neg$ CR	–	BF, $\neg$ CP, $\neg$ CR	–
CP	–	CP, $\neg$ BF	CP, $\neg$ BF	$\neg$ PC, $\neg$ CR	–	AF, $\neg$ PC, $\neg$ CR	AF, $\neg$ PC, $\neg$ CR	–
CR	PC, $\neg$ AF	CP, $\neg$ BF	CR, EQ	NR	BF, $\neg$ CP, $\neg$ CR	AF, $\neg$ PC, $\neg$ CR	EQ	VG
NR	$\neg$ CP, $\neg$ CR	$\neg$ PC, $\neg$ CR	NR	–	–	–	–	–
BF	BF, $\neg$ CP, $\neg$ CR	–	BF, $\neg$ CP, $\neg$ CR	–	BF, $\neg$ CP, $\neg$ CR	–	BF, $\neg$ CP, $\neg$ CR	$\neg$ AF, $\neg$ EQ
AF	–	AF, $\neg$ PC, $\neg$ CR	AF, $\neg$ PC, $\neg$ CR	–	–	AF, $\neg$ PC, $\neg$ CR	AF, $\neg$ PC, $\neg$ CR	$\neg$ BF, $\neg$ EQ
EQ	$\neg$ AF	$\neg$ BF	EQ	–	BF, $\neg$ CP, $\neg$ CR	AF, $\neg$ PC, $\neg$ CR	EQ	VG, $\neg$ CR
VG	–	–	VG, $\neg$ CR	–	$\neg$ AF, $\neg$ EQ	$\neg$ BF, $\neg$ EQ	VG	–

**Earlier work:** Incorporate external knowledge as **hard** and **soft** constraints via Constrained Conditional Models (CCMs, Chang et al., 2012); usually post-learning correction

**Constrained learning** (Li et al., EMNLP'19): convert declarative rules to differentiable learning objectives using t-norm

$$\hat{\mathcal{I}} = \operatorname{argmax}_{\mathcal{I}} \sum_{ij \in \mathcal{E}} \sum_{r \in R} (x_r(ij) + \lambda f_r(ij)) \mathcal{I}_r(ij)$$

s.t.  $\sum_r \mathcal{I}_r(ij) = 1$ ,  $\mathcal{I}_r(ij) = \mathcal{I}_{\bar{r}}(ji)$ ,  
 (uniqueness) (symmetry)

$\mathcal{I}_{r_1}(ij) + \mathcal{I}_{r_2}(jk) - \sum_{m=1}^M \mathcal{I}_{r_3^m}(ik) \leq 1$ ,  
 (transitivity)

- $L_A$  Annotation Loss:  $\top \rightarrow r(e_1, e_2) \rightarrow -w_r \log r(e_1, e_2)$
- $L_S$  Symmetric Loss:  $\alpha(e_1, e_2) \leftrightarrow \bar{\alpha}(e_2, e_1) \rightarrow |\log \alpha(e_1, e_2) - \log \bar{\alpha}(e_2, e_1)|$
- $L_C$  Conjunctive Loss:  $\alpha(e_1, e_2) \wedge \beta(e_2, e_3) \rightarrow \gamma(e_1, e_3) \rightarrow \log \alpha(e_1, e_2) + \log \beta(e_2, e_3) - \log \gamma(e_1, e_3)$   
 $\alpha(e_1, e_2) \wedge \beta(e_2, e_3) \rightarrow \neg \delta(e_1, e_3) \rightarrow \log \alpha(e_1, e_2) + \log \beta(e_2, e_3) - \log(1 - \delta(e_1, e_3))$
- Training Objective:  $L = L_A + \lambda_S L_S + \lambda_C L_C$

**Exact Inference** at decision time via CCMs

- Formulate as an ILP problem: Gurobi



- TemProb (Ning et al. 2018a)
  - Source: New York Times 1987-2007 (#Articles~1M)
  - Preprocessing: Semantic Role Labeling & Temporal relation model
  - 51K semantic frames, 80M relations
- ConceptNet
  - “HasSubevent”, “HasFirstSubevent”
  - 30K Positive & 30K Negative

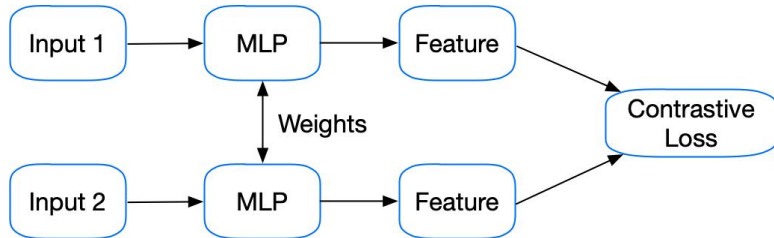
en kill people — HasSubevent → en attack  
Weight: 1.0

en commit murder — HasFirstSubevent → en find a reason  
Weight: 2.0

More than 10 people **died** on their way to the nearest hospital, police said. A suicide car bomb **exploded** on Friday in the middle of a group of men playing volleyball in northwest Pakistan.

Frame1	Frame2	Before	After
concern	protect	92%	8%
conspire	kill	95%	5%
fight	overthrow	92%	8%
accuse	defend	92%	8%
crash	die	97%	3%
<b>explode</b>	<b>die</b>	<b>83%</b>	<b>17%</b>

- MLP Encoders



# Training Process

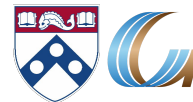


Input sentences

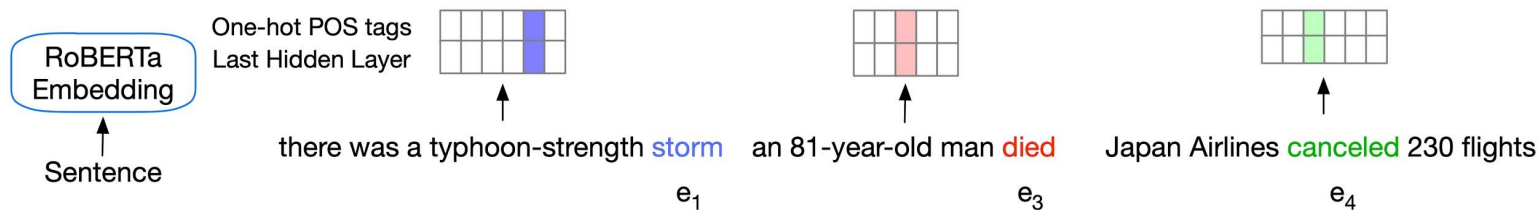
On Tuesday, there was a typhoon-strength ( $e_1$ :*storm*) in Japan. One man got ( $e_2$ :*killed*) and thousands of people were left stranded. Police said an 81-year-old man ( $e_3$ :*died*) in central Toyama when the wind blew over a shed, trapping him underneath. Later this afternoon, with the agency warning of possible tornadoes, Japan Airlines ( $e_4$ :*canceled*) 230 domestic flights, ( $e_5$ :*affecting*) 31,600 passengers.

Sentence	there was a typhoon-strength <b>storm</b>	an 81-year-old man <b>died</b>	Japan Airlines <b>canceled</b> 230 flights
	$e_1$	$e_3$	$e_4$

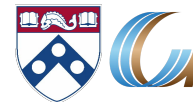
# Training Process



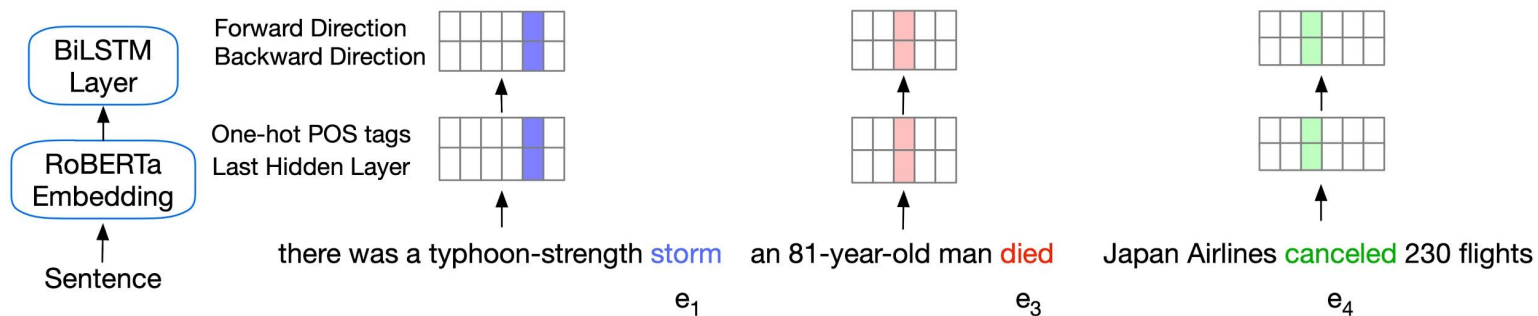
## Contextualized Representation



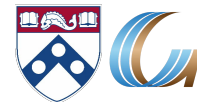
# Training Process



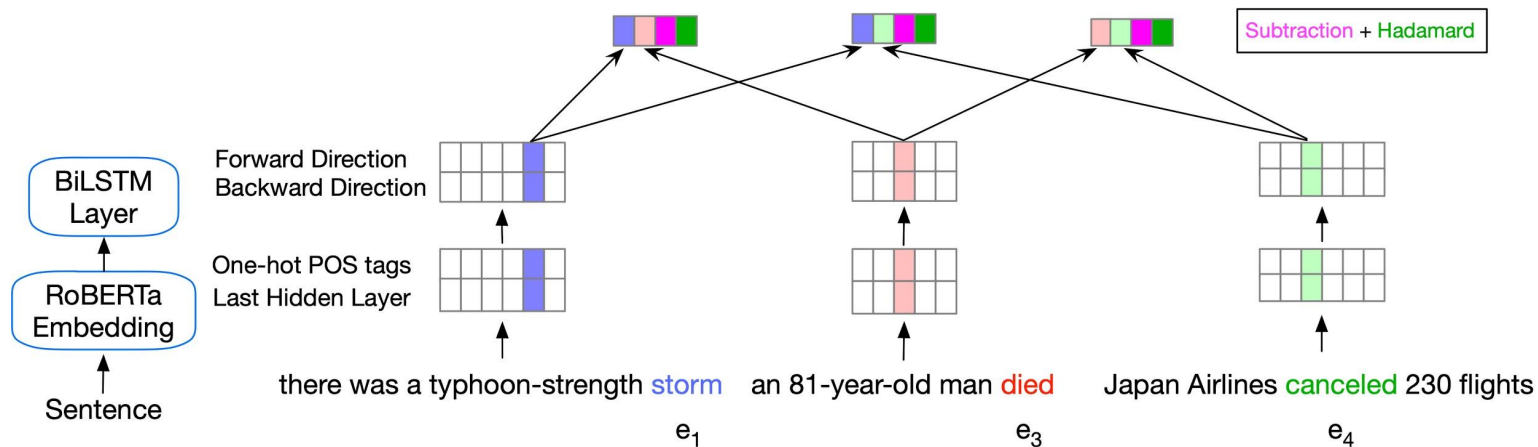
## BiLSTM Encoder



# Training Process



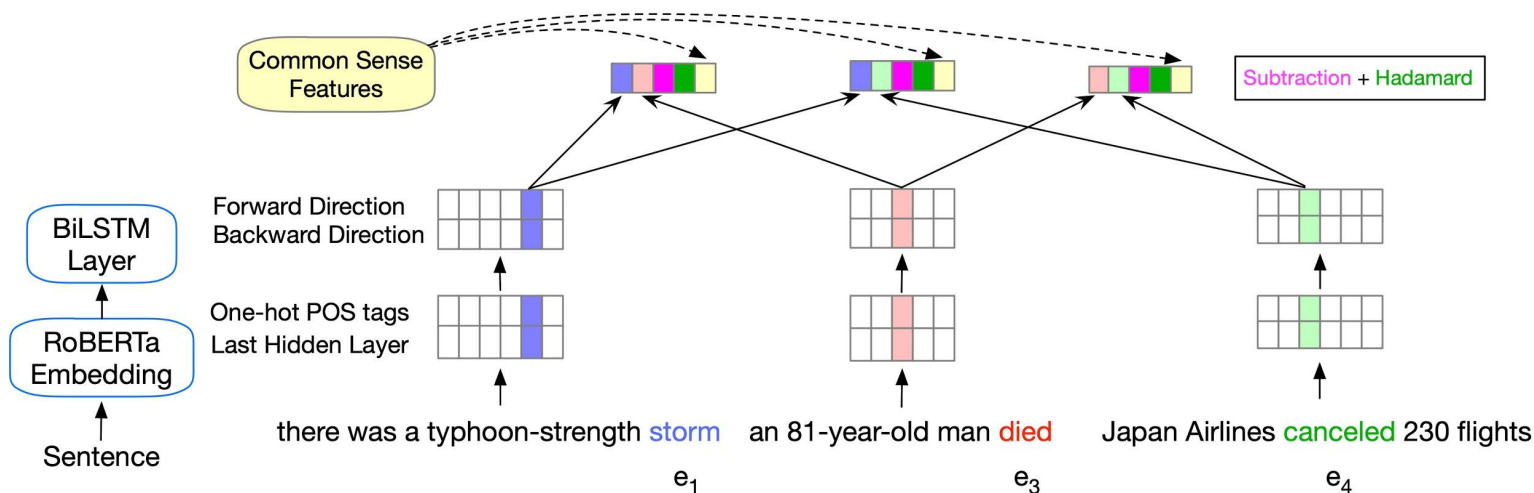
## Event Pair Representation



# Training Process



## Adding Common Sense Features

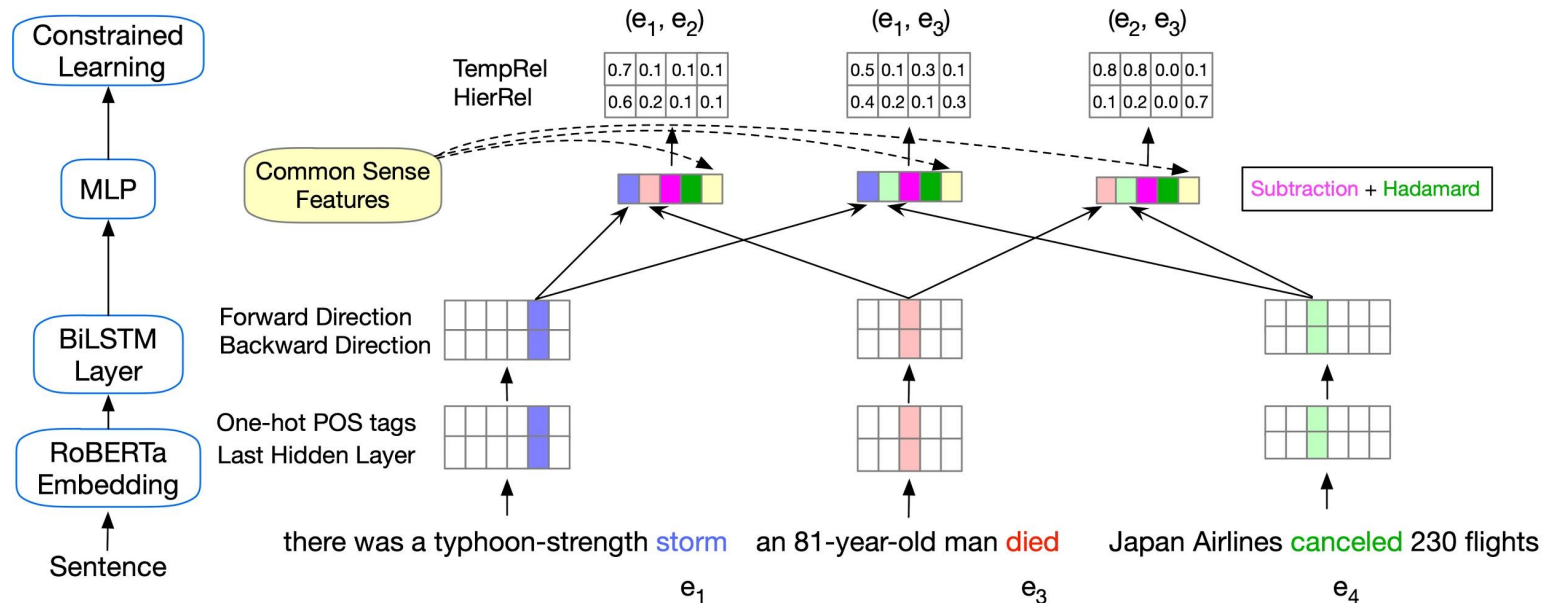


# Training Process



## Constrained Learning

$$\text{Loss Function: } L = L_A + \lambda_S L_S + \lambda_C L_C$$



# Evaluation on Benchmarks



## Benchmark Datasets

- Membership relations: HiEve
- Temporal relations: MATRES
- Case study: RED

## Dataset Statistics

	HiEve	MATRES	RED
# of Documents			
Train	80	183	-
Dev	-	72	-
Test	20	20	35
# of Pairs			
Train	35001	6332	-
Test	7093	827	1718

- The proposed method surpasses the SOTA TempRel extraction method on MATRES by relatively 3.27%  $F_1$ .

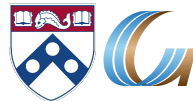
Model	$P$	$R$	$F_1$
CogCompTime (Ning et al., 2018c)	0.616	0.725	0.666
Perceptron (Ning et al., 2018b)	0.660	0.723	0.690
BiLSTM+MAP (Han et al., 2019b)	-	-	0.755
LSTM+CSE+ILP (Ning et al., 2019)	0.713	0.821	0.763
Joint Constrained Learning (ours)	<b>0.734</b>	<b>0.850</b>	<b>0.788</b>

- It also offers promising performance on the HiEve dataset for subevent relation extraction, relatively surpassing previous methods by at least 3.12% in  $F_1$ .

Model	$F_1$ score		
	PC	CP	Avg.
StructLR (Glavaš et al., 2014)	0.522	<b>0.634</b>	0.577
TACOLM (Zhou et al., 2020a)	0.485	0.494	0.489
Joint Constrained Learning (ours)	<b>0.625</b>	0.564	<b>0.595</b>



# Case Study on RED dataset



## Mapping for relations

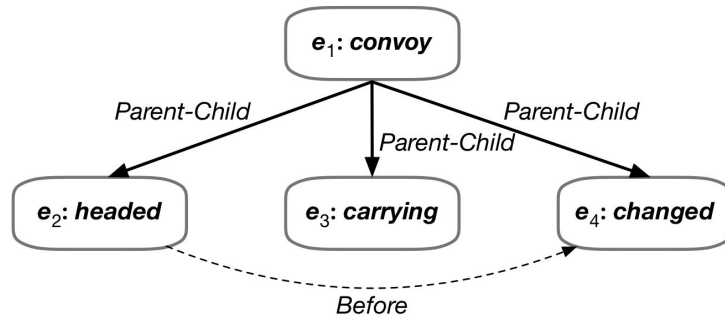
Original labels in RED	Mapped labels
BEFORE, BEFORE/CAUSES, BEFORE/PRECONDITION, ENDS-ON, OVERLAP/PRECONDITION	BEFORE
SIMULTANEOUS	EQUAL
OVERLAP, REINITIATES	VAGUE
CONTAINS, CONTAINS-SUBEVENT	PARENT-CHILD & BEFORE
BEGINS-ON	AFTER

## Performance ( $F_1$ )

Model	TEMPREL	SUBEVENT
Joint Constrained Learning (ours)	<b>0.72</b>	<b>0.54</b>

## Example Output

A (***e1:convoy***) of 280 Russian trucks (***e2:headed***) for Ukraine, which Moscow says is (***e3:carrying***) relief goods for war-weary civilians, has suddenly (***e4:changed***) course, according to a Ukrainian state news agency.



- Single-task Training vs Joint Training
- Constrained Learning
  - Task-specific constraints
  - Cross-task constraints
- Commonsense Knowledge
- Post-learning Correction: ILP

Model	SUBEVENT			TEMPREL		
	<i>P</i>	<i>R</i>	<i>F</i> <sub>1</sub>	<i>P</i>	<i>R</i>	<i>F</i> <sub>1</sub>
Single-task Training	32.5	<b>73.1</b>	45.0	67.7	80.3	73.5
Joint Training	50.4	43.1	46.5	68.4	82.0	74.6
+ Task-specific constrained learning	51.6	59.7	55.4	71.3	82.7	76.6
+ Cross-task constrained learning	51.1	67.0	58.0	72.2	83.8	77.6
+ Commonsense knowledge	56.9	61.6	59.2	73.3	84.2	78.4
+ Global inference (ILP)	<b>57.4</b>	61.7	<b>59.5</b>	<b>73.4</b>	<b>85.0</b>	<b>78.8</b>
All but constrained learning	54.2	41.8	47.2	72.1	80.8	76.2

- A new paradigm for joint constrained learning
  - Bridge temporal and subevent relations with a comprehensive set of logical constraints
  - Convert constraints into differentiable objective functions
- Address the lack of jointly annotated data
  - Complementary supervision signals from two tasks
- Outperform previous methods on benchmark datasets
  - MATRES: brings about 3.27%  $F_1$  improvement
  - HiEve: brings about 3.12%  $F_1$  improvement
- Present promising results on RED
- Generalizable to other relations / constraints

