

September 25, 2018  
DRAFT

*Thesis Proposal*  
**Diving Deep into Event Semantics**

Zhengzhong Liu

October 2018

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

**Thesis Committee:**

Teruko Mitamura (Chair)  
Eduard Hovy  
Taylor Berg-Kirkpatrick  
Vicent Ng

*Submitted in partial fulfillment of the requirements  
for the degree of Doctor of Philosophy.*

September 25, 2018  
DRAFT

**Keywords:** event,script,event schema,semantics,machine learning

# Abstract

Among various discourse elements (DEs) in natural language discourse, events are important due to their rich structures. For example, event mentions are often connected with other DEs, such as event participants, time and location. Multiple events can further form larger structures such as coreference clusters and scripts. The structural perspective makes events perfect bridges across DEs. Prior work exploiting the event structures normally makes an identity assumption: DEs filling in the same role in a discourse are assumed to be identical. However, we often observe violations.

We argue that only partial information of DEs are intended for communication in text, thus analysis based on the full information will encounter problems. To address them, we propose a linguistic framework featuring a **facet-based** representation. Facets are smaller semantic units that constitute the DEs. We propose to solve the problem by focusing on the **active facets** — the facets that are relevant to the context.

The facet-based framework allows us to study hypotheses from two perspectives: static and dynamic. Statically, we hypothesize that DEs are written to highlight the active facets and that linking on the facet level is more consistent. Dynamically, we hypothesize that events are functions that operate on the active facets of the entity mentions, which can help analyze the states within a script.

Our computation contributions are in two-fold. We first introduce our early **supervised** attempts on event semantics, including an event detector and coreference engine, an event sequencing system, and an event ellipsis system. However, all the datasets of these tasks suffer in both scale and domain, which prevents us from exploring rich linguistic phenomenon. To this end, we attempt to obtain **indirect supervision** from the data. In Part II, we present some successful examples on this line. Finally, in Part III, we propose to validate the two hypotheses using indirect supervision on important NLP problems such as entity coreference and semantic role labeling.

September 25, 2018  
DRAFT

# Contents

<b>Abstract</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Event Semantics . . . . .	2
1.1.1 Event Definition . . . . .	2
1.1.2 Event Structure . . . . .	2
1.1.3 Basic Terminology . . . . .	2
1.2 Learning Semantics . . . . .	2
1.3 Thesis Outline . . . . .	2
<b>I Learning with Supervised Approaches</b>	<b>3</b>
<b>2 Event Detection</b>	<b>5</b>
2.1 Introduction . . . . .	5
2.2 Event Detection and Type Classification . . . . .	5
2.3 Realis . . . . .	7
2.4 Adapting to a Chinese Event Detection system . . . . .	7
2.4.1 A trick to deal with Chinese data . . . . .	7
2.5 Evaluation Results . . . . .	8
2.6 Discussion . . . . .	9
2.6.1 Chinese Data Annotation . . . . .	9
2.6.2 Towards Unsupervised Methods . . . . .	11
<b>3 Event Coreference</b>	<b>13</b>
3.1 Introduction . . . . .	13
3.2 Related Work . . . . .	13
3.2.1 Problem Definition . . . . .	13
3.2.2 Dataset and Settings . . . . .	15
3.2.3 Gold Standard Annotations . . . . .	16
3.2.4 Availability . . . . .	16
3.3 Corpus . . . . .	16
3.4 System description . . . . .	16
3.4.1 Procedure . . . . .	16

3.4.2	Features . . . . .	17
3.4.3	Clustering . . . . .	18
3.5	Evaluation . . . . .	18
3.5.1	Evaluation Metrics . . . . .	18
3.5.2	Experiments and Results . . . . .	18
3.6	Discussion . . . . .	19
3.7	Conclusion and Future Work . . . . .	20
<b>4</b>	<b>Event Sequencing</b>	<b>23</b>
4.1	Introduction . . . . .	23
4.2	Related Work . . . . .	25
4.3	Model . . . . .	25
4.3.1	Graph-Based Decoding Model . . . . .	25
4.3.2	Features . . . . .	29
4.4	Experiments . . . . .	30
4.4.1	Dataset . . . . .	30
4.4.2	Baselines and Benchmarks . . . . .	30
4.4.3	Evaluation Metrics . . . . .	31
4.4.4	Evaluation Results for Event Coreference . . . . .	31
4.4.5	Evaluation Results for Event Sequencing . . . . .	32
4.5	Discussion . . . . .	33
4.5.1	Event Coreference Challenges . . . . .	33
4.5.2	Event Sequencing Challenges . . . . .	34
4.6	Conclusion . . . . .	35
<b>5</b>	<b>Verb Phrase Ellipsis</b>	<b>37</b>
5.1	Introduction . . . . .	37
5.2	Related Work . . . . .	38
5.3	Approaches . . . . .	38
5.3.1	Target Detection . . . . .	38
5.3.2	Antecedent Head Resolution . . . . .	39
5.3.3	Antecedent Boundary Determination . . . . .	40
5.4	Joint Modeling . . . . .	41
5.5	Experiments . . . . .	43
5.5.1	Datasets . . . . .	43
5.5.2	Evaluation . . . . .	43
5.5.3	Baselines and Benchmarks . . . . .	44
5.6	Results . . . . .	44
5.6.1	Target Detection . . . . .	44
5.6.2	Antecedent Head Resolution . . . . .	45
5.6.3	Antecedent Boundary Determination . . . . .	45
5.6.4	End-to-End Evaluation . . . . .	47
5.7	Discussion . . . . .	47

## II Learning with Indirect Supervision 49

<b>6</b>	<b>Event Saliency</b>	<b>51</b>
6.1	Introduction . . . . .	51
6.2	Related Work . . . . .	53
6.3	The Event Saliency Corpus . . . . .	53
6.4	Feature-Based Event Saliency Model . . . . .	54
6.4.1	Features . . . . .	54
6.4.2	Model . . . . .	55
6.5	Neural Event Saliency Model . . . . .	56
6.5.1	Kernel-based Centrality Estimation . . . . .	56
6.5.2	Integrating Entities into KCE . . . . .	56
6.6	Experimental Methodology . . . . .	57
6.6.1	Event Saliency Detection . . . . .	57
6.6.2	The Event Intrusion Test: A Study . . . . .	58
6.7	Evaluation Results . . . . .	59
6.7.1	Event Saliency Performance . . . . .	59
6.7.2	Intrusion Test Results . . . . .	61
6.8	Conclusion . . . . .	62

## III Proposed Work 65

<b>7</b>	<b>Proposed Theoretical Framework</b>	<b>67</b>
7.1	Introduction . . . . .	67
7.1.1	The Exact Frame Identity Assumption and Its Failures . . . . .	67
7.2	The Complexity about Identity . . . . .	69
7.2.1	Evidence from Corpus Annotation . . . . .	70
7.3	The Facet View of Discourse Elements . . . . .	71
7.3.1	Facets for Static Identity Analysis . . . . .	72
7.3.2	Facets and the Dynamic Interactions . . . . .	74
7.3.3	Facets for Events . . . . .	74
7.4	Relations to Prior Research . . . . .	75
7.4.1	Quasi-Identity . . . . .	76
7.4.2	Script Modeling . . . . .	76
7.4.3	Area Similar to Quasi-Identity . . . . .	76
7.4.4	Context Dependent Modeling . . . . .	77
<b>8</b>	<b>Proposed Experiments</b>	<b>79</b>
8.1	Facets Understanding from a Static View . . . . .	79
8.1.1	Facet Representation . . . . .	79
8.1.2	Facet Identification . . . . .	80
8.2	Understanding States and Facets from a Dynamic View . . . . .	81
8.2.1	State Modeling . . . . .	81

<b>9</b>	<b>Timeline</b>	<b>83</b>
	<b>Bibliography</b>	<b>85</b>



# List of Figures

4.1	Example of Event Coreference and Sequence relations. Red lines are coreference links; solid blue arrows represent Subevent relations; dotted green arrows represent After relations. . . . .	24
4.2	Latent Tree Model (left): tree structure formed by undirected links. Latent Graph Model (right): a DAG form by directed links. Dashed red links highlight the discrepancy between prediction and gold standard. The dotted yellow link (bottom right) can be inferred from other links. . . . .	26
6.1	Examples annotations. Underlying words are annotated event triggers; the red bold ones are annotated as salient. . . . .	52
6.2	Learned Kernel Weights of $KCE$ . . . . .	62
6.3	Intruder study results. X-axis shows the percentage of intruders inserted. Y-axis is the AUC score scale. The left and right figures are results from salient and non-salient intruders respectively. The blue bar is AUC. The orange shaded bar is SA-AUC. The line shows the SA-AUC of the frequency baseline. . . . .	63
7.1	While the ideal script model (Left) assumes the shared frame elements to be exactly identical, actual text (Right) demonstrates more complexities. . . . .	68
7.2	Example of active facets for an entity mention of “paper”. . . . .	73
7.3	A Venn Graph showing the relations between different sets of facets. . . . .	75

September 25, 2018  
DRAFT

# List of Tables

2.1	Event Nugget score over 5-fold validation on training data for LTI1. . . . .	8
2.2	Official Event Nugget score for LTI1. . . . .	8
2.3	5-fold validation for Realis detection on training data with gold span and types for LTI1 . . . . .	9
2.4	Performance on English Nugget Detection . . . . .	9
2.5	Performance on Chinese Nugget Detection . . . . .	9
2.6	Top 5 double character mentions in ACE . . . . .	10
2.7	Top 5 double character mentions in ERE . . . . .	10
3.1	Recent computational approaches to event coreference resolution. . . . .	14
3.2	Corpus Statistics . . . . .	17
3.3	Evaluation results and comparisons . . . . .	19
3.4	List of features (with counts) in the pairwise model. Entity coreference are from the Stanford Entity Coreference Engine [79]. . . . .	21
4.1	Coreference Features. Parsing is done using Stanford CoreNLP [93]; frame names are produced by Semafor [43]. . . . .	29
4.2	Test Results for Event Coreference with the Singleton and Matching baselines.	32
4.3	Ablation study for Event Coreference. . . . .	32
4.4	Test results for event sequencing. The Oracle Cluster+Temporal system is running CAEVO on the Oracle Clusters. . . . .	33
4.5	Ablation Study for Event Sequencing. . . . .	33
5.1	Antecedent Features . . . . .	42
5.2	Corpus statistics . . . . .	43
5.3	Results for Target Detection . . . . .	46
5.4	Results for Antecedent Head Resolution . . . . .	46
5.5	Soft results for Antecedent Boundary Determination . . . . .	46
5.6	Soft results for Antecedent Boundary Determination with non-gold heads . . . .	48
5.7	Soft end-to-end results . . . . .	48
6.1	Dataset Statistics. . . . .	54
6.2	Event Saliency Features. . . . .	55

6.3	Event salience performance. (-E) and (-F) marks removing Features and Entity information from the full KCM model. The relative performance differences are computed against <code>Frequency</code> . <code>W/T/L</code> are the number of documents a method wins, ties, and loses compared to <code>Frequency</code> . † and ‡ mark the statistically significant improvements over <code>Frequency</code> †, <code>LeToR</code> ‡ respectively. . . . .	59
6.4	Event Salience Feature Ablation Results. The + sign indicates adding feature groups to <code>Frequency</code> . <code>SL</code> is the sentence location feature. <code>Event</code> is the event voting feature. <code>Entity</code> is the entity voting feature. <code>Local</code> is the local entity voting feature. † marks the statistically significant improvements over + <code>SL</code> . . . .	60
6.5	Examples of pairs of Events/Entities in the kernels. The <b>Word2vec</b> column shows the cosine similarity using pre-trained word vectors. The <b>Kernel</b> column shows the closest kernel they belong after training. Items marked with (E) are entities. .	61

# Chapter 1

## Introduction

Modern NLP systems excel particularly at recognizing and extracting information content from text documents. Yet many challenges exist in further understanding of semantics. Examples include creating structures to connect the information pieces in discourse (e.g. Semantic Role Labeling, Discourse Parsing), and inferring implicit information from the surface (e.g. Textual entailment). Many of these inference problems involve understanding the structures among discourse elements (DEs). One important type of DE is event, which are linguistic expressions that describe state change.

In an underlying world (real or imaginary), we consider the static configuration of entities and their properties as *states*. For example, a vase is placed on top a table. If the vase falls down from the table, its physical location has then changed, which is a change of state. And we may name it the “fall” event.

Events are important building blocks of documents and play a key role in the process of document understanding. Conceptually, a single event happens at a location, in a time interval, participated by several entities and can interact with other events. Events may also collectively form larger events. The structure capturing such interaction thus forms the backbone of a document, and can be viewed as a representation of the document. A famous formulation is Schank’s script theory [137], which suggests that the information is centered around the event sequences, which enables language understanding and inference.

In text documents, events are normally realized as spans of text, which are often referred as **Event Mentions** or **Event Nuggets**. Analogy to the entity coreference problem, the same event can also be represented by multiple text spans, one thus needs to resolve **coreference** to recover events from mentions. On the basis of events, various other types of relations can then be established between events. Many of them are related to Schank’s script. For example, the “subevent” relation proposed in [74] organize events into script clusters. The TempEval tasks [147] considers temporal ordering of events. Following these directions, a new task named **Event Sequencing** (ES) is proposed for TAC KBP 2017, which aims at ordering events from text documents that belong to the same script. Both clustering events into the same script and event ordering are required in the task.

In this thesis, we studied various aspects of event semantics and showcased several useful application of such semantic knowledge. Textual realization of events and entities are the main medium connecting to the underlying world. Accurate parsing of them are very important towards

document understanding.

## **1.1 Event Semantics**

### **1.1.1 Event Definition**

Define the problems and terms in our context, such as scripts, events, mentions

### **1.1.2 Event Structure**

### **1.1.3 Basic Terminology**

## **1.2 Learning Semantics**

## **1.3 Thesis Outline**

# **Part I**

## **Learning with Supervised Approaches**

September 25, 2018  
DRAFT



# Chapter 2

## Event Detection

### 2.1 Introduction

A short introduction goes here.

### 2.2 Event Detection and Type Classification

LTI1 deploys a Conditional Random Field (CRF) model to detect mention span and event type, and a linear Support Vector Machine (SVM) model to determine mention realis status. The CRF model is trained with the structured perceptron [35], which is outlined in Algorithm 1. The decoding step is done using standard Viterbi algorithm. Our final system always make use of the average weight variation as described in Collins [35].

A number of “Double tagging” mentions are annotated in the corpus, in which a mention might have one or more event types. In LTI1 we simply combine multiple labels for each mention into a single label<sup>1</sup>. The intuition behind is that some surface forms are usually associated with some fixed mention types, for instance, “KILL” is normally associated with “Life.Die” and “Conflict attack”.

<sup>1</sup>Multi-label classification can be solved with sophisticated methods when simple concatenation creates too many classes. In the KBP event dataset, the number of extended classes is small (56 in the training data).

---

**Algorithm 1** Structured perceptron.

---

**Input:** training examples  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$

**Input:** number of iterations  $T$

**Output:** weight vector  $\mathbf{w}$

```

1:  $\mathbf{w} \leftarrow \mathbf{0}$  ▷ Initialization.
2: for  $t \leftarrow 1..T$  do
3:   for  $i \leftarrow 1..N$  do
4:      $\hat{y}^{(i)} = \arg \max_{y \in \mathcal{Y}(x^{(i)})} \mathbf{w} \cdot \Phi(x^{(i)}, y)$ 
5:     if  $\hat{y}^{(i)} \neq y^{(i)}$  then
6:        $\mathbf{w} \leftarrow \mathbf{w} + \Phi(x^{(i)}, y^{(i)}) - \Phi(x^{(i)}, \hat{y}^{(i)})$ 
return  $\mathbf{w}$ 

```

---

An event mention is normally composed by its mention trigger and the arguments. To get a list of arguments for the event mention. We run two Semantic Role Labeling system, the PropBank style Fanse Parser [43] and the FrameNet style Semafor Parser [144]. In addition, to reduce sparsity, LTI1 also employs a few external data resources, including WordNet [55] and a Brown Clustering trained on newswire data [143]. We also use the Stanford CoreNLP system to obtain lemma, part-of-speech and parsing information [93].

Using these resources, LTI1 employ regular linguistic features for mention type detection, which are summarized as followed:

1. Part-of-Speech, lemma, named entity tag of words in the 2-word window of the trigger (both side), the trigger word itself and the direct dependent words of the trigger.
2. Brown clusters, WordNet Synonym and derivative forms of the trigger.
3. Whether the words in the 5 word window match some selected WordNet senses, including “Leader”, “Worker”, “Body Part”, “Monetary System”, “Possession”, “Government”, “Crime” and “Pathological State”.
4. Closest named entity type.
5. Dependency features, including lemma, dependency type and part-of-speech of the child dependencies and head dependencies.
6. Semantic role related features includes the frame name and the argument role, named entity tag, argument head word lemma and WordNet sense (selected from the above list as well) of the arguments.

The WordNet related features are selected following the intuition that certain category of words are likely to imply the existence of certain events. For example, “Leader” are normally associated with “Personnel” type. We are currently working on methods that can automatically select such senses instead of using explicit human intervention.

In the current implementation, all the raw features are simply concatenated with the mention type. However, one can further make use of the hierarchy of mention to extract features on the higher ontology. For example, when extracting features for “Personnel.End-Position”, one can also add a feature about “Personnel”. We didn’t finish the implementation of this variation and will leave it to our future work.

## 2.3 Realis

The realis model is a simple SVM model trained using LIBLINEAR<sup>2</sup>. We use similar window features and dependency features as in mention type detection. We also add frame argument role names in to the model (i.e. Attacker). However, we didn't include most of the lexicalized features used in type detection to avoid overfitting. We design a specific feature to capture whether the phrase containing the event mention is quoted (if the whole sentence is quoted, we do not fire this feature). In order to test the effect of our realis model, we evaluate its performance given gold standard mention span and types in the training data. We report the 5-fold validation result in Table 2.3<sup>3</sup>.

## 2.4 Adapting to a Chinese Event Detection system

To extend our system to handle Chinese documents, we develop similar features for Chinese. Most of the features can be reused without changes in the Chinese system, which includes: window based features<sup>4</sup>, syntactic based features, entity features, head word features and SRL features. We also use the Brown clustering features with clusters induced from Chinese Gigaword 3<sup>5</sup>.

In addition, Chinese tokens normally contain internal structure and each single character in the token may convey useful semantic information. We add the following character related features:

1. Whether the token contains a character.
2. The contained character and its character level Part-of-Speech.
3. The first character of the token.
4. The last character of the token.
5. Base verb structure feature as described in [82]: we use a feature to represent one of the base verb structure. In addition to the 6 main structures proposed by Li et al. [82], we added 3 structures for completeness: 1. No verb character found 2. The verb character is found after 2 characters and 3. Other: any cases that are not defined above.

### 2.4.1 A trick to deal with Chinese data

During the system development, we observe that our Chinese system suffers from serious low recall despite all the features we added in. By following the training procedures, we hypothesize that the annotated Chinese data is not complete (see § 2.6.1 for more discussion). As a result, our learning algorithm will be biased by the missed events and learn incorrect negative signals. The final model thus will be very conservative in making predictions, leading to a low recall.

The problem can only be solved by manually polishing the data, which is too expensive for us. We mitigate the problem by ignoring all training sentences that do not contain an event mention, which reduce the probability of missed annotations. On our development experiments, we found

<sup>2</sup><https://www.csie.ntu.edu.tw/~cjlin/liblinear/>

<sup>3</sup>Because the span detection is given, the precision and recall are the same all the time.

<sup>4</sup>However since the discussion forum training data are quite noisy, we restrict the POS window to 1 instead.

<sup>5</sup><http://www.cs.brandeis.edu/~clp/conll15st/dataset.html>

	Precision	Recall	F1
Span	74.36	55.72	63.62
Type	67.08	50.25	57.38
Realis	51.79	38.75	44.27
All	46.29	34.63	39.56

Table 2.1: Event Nugget score over 5-fold validation on training data for LTI1.

	Precision	Recall	F1
Span	82.46	50.30	62.49
Type	73.68	44.94	55.83
Realis	62.09	37.87	47.05
All	55.12	33.62	41.77

Table 2.2: Official Event Nugget score for LTI1.

that this simple trick can directly raise our nugget detection performance by about 3%. The performance improvement also support our hypothesis that the Chinese dataset is indeed not fully annotated.

## 2.5 Evaluation Results

We report our system results in the official TAC KBP evaluation workshop.

We report our score averaged on the 5-fold validation in the training data in Table 2.1, and our official score in Table 2.2. LTI1 ranks 2nd among all the participants in terms of the final event type + realis status (All in the table) F1 score. LTI1 also rank high according the other metrics. Notably, we found that the system is very robust against across datasets. Although we have different precision, recall trade off between the training and test set, the F-score is relatively stable.

Both our English and Chinese event detection and coreference systems produce competitive results. Our Chinese system is the first place based on all the event nugget attributes. Our English system is the second in mention type detection<sup>6</sup>. To our surprise, our English nugget detection performance drops about 13% (span and type) comparing to last year. However, our relative ranking is almost unchanged. We leave the investigation of the problem to future work. Since

<sup>6</sup>Due to the Realis component bug our combined ranking drops.

Fold	Precision	Recall	F1
1	71.68	71.63	71.66
2	64.06	64.06	64.06
3	62.07	62.07	62.07
4	72.66	72.66	72.66
5	62.21	62.21	62.21

Table 2.3: 5-fold validation for Realis detection on training data with gold span and types for LTI1

the event coreference component and coreference evaluation relies highly on the performance of nugget detection, we also observe a big drop in coreference performance comparing to last year.

	Prec.	Recall	F1
Span	69.82	39.54	50.49
Type	61.69	34.94	44.61
Realis	45.78	25.93	33.11
All	40.19	22.76	29.06

Table 2.4: Performance on English Nugget Detection

	Prec.	Recall	F1
Span	56.46	39.55	46.52
Type	50.72	35.53	41.79
Realis	42.7	29.92	35.18
All	38.91	27.26	32.06

Table 2.5: Performance on Chinese Nugget Detection

## 2.6 Discussion

### 2.6.1 Chinese Data Annotation

We hypothesize that the Chinese datasets are not fully annotated. We take a closer look in the data and found a number of missed event nuggets. Here we list a couple examples:

(2.1) 支持香港同胞争取[Personnel.Elect 选举]与被[Personnel.Elect 选举]权！

(2.2) 司务长都是骑着二八去[TransferOwnership 买]菜去。

(2.3) 海豹行动是绝密，塔利班竟然可以预先得知？用个火箭就可以[Conflict.Attack打]下来，这个难度也实在是太高了吧。

In the above examples, we show several event nuggets. However, mentions annotated in red are not actually annotated in the Rich ERE datasets. Especially, in example 2.1, the first 选举 is annotated but the second one is not. Such inconsistencies happen a lot across the dataset. When training with such data, the classifier will likely to be quite conservative on making event nugget predictions. We conduct a very simple quantitative analysis by comparing the ACE 2005 Chinese annotation against the Rich ERE Chinese annotation. Table 2.6 and Table 2.7 summarize the top 5 double-character tokens annotated in ACE and RichERE. For the most popular event mentions, Rich ERE annotated only a smaller percentage comparing to ACE.

In addition, we find that the most popular event nuggets are mostly single character in the Rich ERE datasets, such as 打(170), 说(148), 死(131), 杀(118). In fact, in top 20 most popular event nuggets of Rich ERE, there are 17 single-character nuggets, this number is only 6 in ACE. These single character tokens are more ambiguous comparing to a double character mention (for example, 打 can represent the action of "calling someone" or "attacking someone", which corresponds to very different event type. This is because language in discuss forum posts are normally not formal. This actually challenges our event nugget systems to deal with deeper semantic problems.

Token	Annotated	Total	%
冲突	100	119	84.03%
访问	64	90	71.11%
受伤	53	59	89.83%
死亡	46	50	92.00%
前往	44	52	84.62%

Table 2.6: Top 5 double character mentions in ACE

Token	Annotated	Total	%
战争	96	223	43.05%
死亡	24	33	72.73%
暗杀	22	40	55.00%
入侵	18	22	81.82%
自杀	17	33	51.52%

Table 2.7: Top 5 double character mentions in ERE

### **2.6.2 Towards Unsupervised Methods**

Another major limitation of our current systems is that we relies highly on annotated dataset. However, creating high quality event nugget and coreference dataset is very expensive, especially on different languages. Furthermore, the datasets are normally small in size and are narrow in terms of the event types.

September 25, 2018  
DRAFT



# Chapter 3

## Event Coreference

### 3.1 Introduction

Coreference resolution, the task of linking surface mentions to their underlying discourse entities, is an important task in natural language processing. Most of the early work focused on coreference of entity mentions. Recently, event coreference has attracted attention on both theoretical and computational aspects. However, most event coreference work is preliminary and applied in quite different circumstances, making comparisons difficult or impossible.

In this paper, we first provide an overview of all the relevant literature to identify the ways each experiment differs from the others. Our claim here is that the comparisons to related work in prior papers are not really appropriate due to these differences. This makes future research difficult. We then present a supervised approach to event coreference, and describe a method for propagating information between events and their arguments that can improve results. In our method, different argument types support different methods of propagation. For these experiments, we annotate and use a corpus of 65 documents in the Intelligence Community (IC) domain that contains a rich set of within-document coreference links [74].

### 3.2 Related Work

Table 3.1 summarizes recent work on event coreference resolution. For the reasons below, only one supervised system [1] and two unsupervised [11, 38] on within-document event coreference are suitable as a basis for ongoing comparison.

#### 3.2.1 Problem Definition

Different approaches use different definitions of the problem (see Compatible Definition column). However, as discussed in recent linguistic studies [74, 128], the existence of different types and degrees of coreference makes it necessary to agree on the definition of coreference before performance can be compared. The lack of clarity about what coreference should encompass rules out several systems for comparison. OntoNotes created restricted event coreference [122], linking only some nominalized events and some verbs, without reporting event-specific results. Both

	Gold standard used	Cross/within Document		Compatible definition	Corpus
		Within	Cross		
Lee et al. [80]			✓	✓	ECB
Sangeetha and Arock [136]	ACE mention, arguments and attributes	✓		✓	ACE
Cybulska and Vossen [38]	Mention head word	✓		✓	IC
McConky et al. [95]	ACE mention, arguments and attributes	✓		✓	ACE
Li et al. [83]	Human entity and event mention detection	✓		✓	Unavailable
Bejan and Harabagiu [11]		✓ (ACE, ECB)	✓ (ECB)	✓	ECB, ACE
Chen and Ji [27]	ACE mention, arguments and attributes	✓		✓	ACE
Elkhlifi and Faiz [52]		✓			Unavailable
Naughton [104]		✓			IBC, ACE
Pradhan et al. [122]		✓			OntoNotes
Ahn [1]	ACE mention, arguments and attributes	✓		✓	ACE
Bagga and Baldwin [6]			✓		Unavailable

Table 3.1: Recent computational approaches to event coreference resolution.

Naughton [104] and Elkhilfi and Faiz [52] worked on sentence-level coreference, which is closer to the definition of Danlos [41]. However it is unclear when one sentence contains multiple event mentions, and hence these are not comparable to systems that process more specific coreference units.

### 3.2.2 Dataset and Settings

Early work by Bagga and Baldwin [6] conduct experiments only on cross-document coreference. Recent advanced work on event coreference is by Bejan and Harabagiu [11] and Lee et al. [80] use the ECB corpus<sup>1</sup> (or a refined version<sup>2</sup>) to evaluate performance, which is annotated mainly for cross-document coreference. In this corpus, within-document coreference is only very partially annotated; most difficult coreference instances are not marked.

- (3.1)
1. Indian naval forces came to the rescue (E1) of a merchant vessel under attack (E3) by pirates in the Gulf of Ade on Saturday, capturing (E2) 23 of the raiders, India said (E4).
  2. Indian commandos boarded the larger pirate boat, **seizing** 12 Somali and 11 Yemeni nationals as well as arms and equipment, the statement said.
- (3.2)
1. The Indian navy captured (E2) 23 piracy suspects who tried (E5) to take over (E3) a merchant vessel in the Gulf of Aden, between the Horn of Africa and the Arabian Peninsula, Indian officials said (E4).
  2. In addition to the 12 Somali and 11 Yemeni suspects, the Indian navy seized two small boats and “a substantial cache of arms and equipment”, the military said in a statement.

We will discuss the argument difference problem now and start promoting the facet based understanding.

The examples sentences are extracted from two documents from the ECB. In both documents, event mentions appear in the first sentence are annotated once, but not in the subsequent sentences. In example 3.1, we find in one of the subsequent sentences the event mention “seizing” which should actually marked as coreferent with “capturing (E2)”. In example 2, we find a more tricky case: the mention “seized”, which has semantics similar to “captured” but this pair is not marked as coreference due to different patients. In cross-document settings, we also find discrepancies between the definitions. In ECB, “attack (E3)” in example 3.1 is annotated as coreferent with “take over (E3)” in example 1, which we believe is wrong: at best, the attack is only a part of the attempt to take over the merchant vessel.

Goyal et al. [64] use a distributional semantic approach on event coreference. However, they didn’t adopt a conventional evaluation setting. They draw from the IC corpus an equal number of positive and negative testing examples, which is different from the natural data distribution.

<sup>1</sup><http://adi.bejan.ro/data/ECB1.0.tar.gz>

<sup>2</sup><http://nlp.stanford.edu/pubs/jcoref-corpus.zip>

### 3.2.3 Gold Standard Annotations

Recent work using the ACE 2005 corpus<sup>3</sup> Ahn [1], Chen and Ji [27], Chen et al. [28], McConky et al. [95], Sangeetha and Arock [136] agrees with our definition of coreference. However, the ACE corpus annotations, in addition to event mentions, also include argument structures, entity ids, and time-stamps. Most coreference systems on the ACE corpus make use of this additional information. This makes them impossible to compare to systems that do not make this simplifying assumption. It also makes results achieved on ACE hard to compare to results on corpora without this additional information. Among these work, only Ahn [1] reported some results using system generating arguments, we compare our system against it.

### 3.2.4 Availability

Li et al. [83] use a hand-annotated web corpus, which is not publicly available for comparison.

In summary, anyone wanting to work on within-document event coreference has to obtain a corpus that is fully annotated, that does not include additional facilitating information, whose definition of coreference respects the theoretical considerations of partial coreference, and that has other systems freely available for comparison. Meeting these criteria is not easy. The closest work we find is by Cybulska and Vossen [38] and Bejan and Harabagiu [11], both adopt unsupervised methods for event coreference. Ahn [1] also reported results on ACE by swapping gold standard annotations with system results. We compare our system to their results on their corresponding corpus.

## 3.3 Corpus

Our system is trained and evaluated on the IC domain corpus, which annotates several different event relations. Table 3.2 summarizes the corpus level statistics and the average over documents. In this work, we focus on full coreference relations. The inter-annotator agreement among 3 annotators for full coreference is 0.614 in terms of Fleiss’ kappa [60]. For detailed definition for the corpus, we refer readers to Hovy et al. [74]. To facilitate future research, We also report our system results on the ACE 2005 training dataset, which contains 599 documents.

## 3.4 System description

Our system is almost end-to-end, except that we start with a minimal gold standard head word annotations in order to focus on the core coreference problem. This approach is the same as Cybulska and Vossen [38] and Bejan and Harabagiu [11]<sup>4</sup>.

<sup>3</sup><http://www.itl.nist.gov/iad/mig/tests/ace/2005/doc/>

<sup>4</sup>Although Bejan and Harabagiu [11] use automatic mention detection to extend the mention set for training, they only use true mentions of the ACE dataset at evaluation time.

	Total	Avg.
Event Mention	2678	41.2
Non-elliptical Domain Event Mention	1998	30.7
Reporting Event Mention	669	10.29
Full coreference relations	1253	21.6
Subevent relations (parent-child)	455	8
Membership relations (parent-child)	161	2.9

Table 3.2: Corpus Statistics

### 3.4.1 Procedure

Similar to Chen et al. [28], we approach the problem first with a conventional pairwise model:

1. Supervised classification that determines the probability whether two mentions co-refer. The classifier used in the experiment is Random Forest [14], implemented in Weka [67].
2. Clustering that processes all the pairwise scores to output the final clusters of pairs.
3. In addition, we added a third step after clustering, information is propagated between event mentions to enrich the original feature set.

The last step tries to enrich the event representations during clustering. Typically, the information carried from one event to its coreferent mention is about the participants (agent, patient, etc.). When an event has been enriched by receiving information from another, it may in turn now be linkable to a third event. The system repeats this process until no more information can be propagated. Currently, the propagation includes two parts: 1. if one mention has missing arguments, they will be copied over from the co-referred counterpart; 2. if both arguments are present, information not presented in one will be copied from another.

Similarly, Lee et al. [80] show that jointly modeling references to events and entities can boost the performance on both. We hold a similar assumption. But by focusing on events and their arguments, we can perform propagation specific to each type of argument, for instance, geographical reasoning as described below.

### 3.4.2 Features

In addition to typical lexical and discourse features, we also model an event mention with its surface form and its arguments, including agent, patient<sup>5</sup>, and location. We use a rich set of 105 semantic features, described in table 3.3.

<sup>5</sup>Specifically, these are defined as ARG0, ARG1 in PropBank. They could be more-specific variants roles such as *experiencer*, but we prefer a smaller set for simplicity.

## Agent, patient extraction and propagation

We use the semantic parser *Fanse* [144] to annotate the predicate arguments defined in PropBank. For nominal events, we extract agent and patient using heuristics such as finding the token attached to the event mention with specific words (such as “by”) and modifiers as agent (e.g., HAMAS in HAMAS’s attack). During the propagation step, information not present in one entity can be copied from another.

## Location extraction and propagation

In contrast to agent and patient, the propagation of location information employs external information to gain additional power. We use the *Stanford Entity Recognition* [59] engine to identify location mentions. *DBpedia Spotlight* [96] is run to disambiguate location entities. DBpedia [4] information, such as cities, country, and alternative names, are then injected. When the location is not found in DBpedia, we search the mention string in *GeoNames*<sup>6</sup> and use the first result with highest Dice coefficient with the surface string. This world knowledge enriches annotation. For example, we can now match the mention “Istanbul” with the country name “Turkey”.

### 3.4.3 Clustering

We conduct experiments with two simple clustering methods. The first is a pure transitive closure that links all pairs mentions that the classification engine judges as positive. The second is the Best-Link algorithm of Ng and Cardie [105], which links each mention to its antecedent with the highest likelihood when the classifier judges as positive.

## 3.5 Evaluation

### 3.5.1 Evaluation Metrics

Coreference evaluation metrics have been discussed by the community for years. To enable comparison, we report most metrics used by the CoNLL 2012 shared task [121], including MUC [31], B-Cubed [5], entity-based CEAF [92], and BLANC [126]. Pairwise scores are used to provide a direct view on performance.

### 3.5.2 Experiments and Results

We split the documents in IC corpus randomly into 40 documents for training and development, and 25 for testing. Parameters such as the probability threshold to determine coreference are tuned on the 40 documents using five-fold cross validation. Optimization is not done separately for each metric. We simply use a universal classifier threshold optimized for pairwise case. During experiment, the propagation step is actually performed for only one iteration, since no further

<sup>6</sup><http://www.geonames.org/>

information is propagated. On the ACE corpus, we simply apply the best model configuration from IC corpus and train on 90% of the documents (539) for training and 10% for testing (60).

Table 3.3 summarizes the overall average results obtained by BestLink on both ACE and IC corpus (BestLink consistently outperforms naively full transitive closure). We also attach three other reported results at the end. Note that these results are not directly comparable: Cybulska and Vossen [38] and Bejan and Harabagiu [11] use unsupervised methods, thus their reported results are evaluated on the whole corpus; Ahn [1] also use a 9:1 train-test split, but the split might be different with ours. A simple comparison shows that our results outperform these systems in all metrics, which is notable because all these metrics are designed to capture the performance from different aspects.

To interpret the results, it should also be noted that because of the existence of large number of singleton clusters, some measures such as  $B^3$  seem to be high even using the most naive feature set. By looking at the pairwise performance, however, we see that current best F-score is only about 50%. There are still many challenges in event coreference.

	Pairwise			MUC			$B^3$			CEAF-e			BLANC		
	R	P	F	R	P	F	R	P	F	R	P	F	R	P	F
<b>IC corpus</b>															
Discourse + Lexical	32.69	25.11	28.40	41.7	33.58	37.2	79.46	74.06	76.67	66.89	73.95	70.24	59.77	61.2	60.43
+ Syntactic	47.12	35.15	40.26	52.6	47.63	50.0	82.24	81.46	81.85	76.91	80.21	78.53	64.76	68.59	66.42
+ Semantic (no arguments)	51.15	42.22	46.26	54.5	49.1	51.68	82.12	82.08	82.1	74.93	78.31	76.58	65.41	69.98	67.35
+ Arguments	55.96	47.86	51.60	56.87	<b>55.81</b>	56.33	83.38	<b>85.58</b>	<b>84.46</b>	<b>88.13</b>	80.73	<b>80.43</b>	68.77	<b>75.21</b>	71.46
+ Propagation	<b>59.04</b>	<b>48.27</b>	<b>53.11</b>	<b>68.72</b>	55.5	<b>61.44</b>	<b>89.28</b>	79.89	84.33	75.14	<b>82.9</b>	78.83	<b>82.28</b>	70.77	<b>75.06</b>
Cybulska and Vossen [38]	—	—	—	—	—	—	81.0	71.0	76.0	—	—	—	—	—	—
<b>ACE corpus</b>															
This work	55.86	40.52	<b>46.97</b>	53.42	48.75	50.98	89.9	88.86	<b>89.38</b>	85.54	87.42	<b>86.47</b>	70.88	70.01	70.43
Bejan and Harabagiu [11]	43.3	47.1	45.1	—	—	—	83.4	84.2	83.8	76.9	76.5	76.7	—	—	—
Ahn [1]	—	—	43.3	—	—	—	—	—	—	—	—	—	—	—	—

Table 3.3: Evaluation results and comparisons

### 3.6 Discussion

The evaluation results show that almost all types of features help to improve the performance over all metrics rather consistently. However, preliminary error analysis shows that some events are still clustered incorrectly even when arguments match. We argue that limitations in argument extraction and entity coreference prevent these features from contributing directly to correct coreference decisions. On the other hand, the results of propagation show that new information helps to find more links but inevitably comes with a drop in precision. We consider that modeling event and arguments holistically like Lee et al. [80] would help guide the propagation. By inspecting the data, we hypothesize that the main benefits brought by the propagation scheme is to match arguments of two coreferent events. If the arguments are nominal events, they will be then marked as coreferent due to the feature “Event as Entity” (See Semantic features in table 3.4). In the following example, if the two event mentions “planning” are marked as coreference, then the corresponding argument “attack” will be also marked as coreference.

- (3.3) A member of the Islamic militant movement HAMAS suspected of **planning** a suicide **attack** against Israel surrendered to Palestinian police here after a six-hour shootout on Friday. HAMAS’s military wing, was on the run from both Palestinian and Israeli police for **planning** anti-Israeli **attacks**.

This hypothesis is also in line with our observation that propagation can only be performed for one round, because the nominal event themselves are unlikely to have other nominal events as arguments. Such interactions between event mentions also remind us that conference can be possibly improved by other types of event relations, such as subevent relations.

Furthermore, the system tends to merge clusters where the event mention head words are the same because the head word feature receives a high weight in the model, even when this is not appropriate. More work should be performed on disambiguating such difficult cases.

### 3.7 Conclusion and Future Work

In this paper we first describe why most previous work on coreference is not directly comparable to one another, for a variety of reasons. In particular, reports of high coreference performance on one corpus do not really transfer over to other corpora or other definitions of coreference. Event coreference is not a solved problem.

We then present a simple supervised pairwise event coreference system. We show that rich linguistic features, especially event arguments, can improve event coreference performance. Argument specific information propagation further help finding new relations. In the future, we propose to implement propagation based on temporal and other types of event relations.



Type (counts)	Feature Name	Description
Discourse (5)	Sentence Distance	Number of sentences between two events mentions.
	Event Distance	Number of event mentions between two event mentions.
	Position	One event is in the title, or the first sentence.
Lexical (12)	Surface Similarity	Several string similarity measures computed between event mention headwords: Dice coefficient, edit distance, Jaro coefficient, lemma match and exact phrase match.
	Modifier Similarity	Dice coefficient similarity of the modifiers of the event mentions.
Syntactic (38)	Part Of Speech	Binary features for plurality, tense, noun or verb for the event mention head words.
	Dependency	The dependency label connecting the two event mentions.
	Negation	Whether two mention head words are both negated.
	Determiner	Whether the event mentions are modified by determiners
Semantic (16)	Coreference	Whether the predicates are in the same entity coreference cluster (only for nominal events).
	WordNet Similarity	Wu-Palmer similarity [112] of the headword pair.
	Senna Embeddings	Cosine similarity of event mention head word embeddings (Senna embeddings [36]).
	Distributional	Distributional similarity between predicates in Goyal et al. [64].
	Verb Ocean	Predicate word relations in “Verb Ocean” [32].
	Semantic Frame	Whether two predicates trigger the same semantic frame (extracted by Semafor [42]).
	Mention Type	Predicate word type (generated by IBM Sire [61]) match.
Arguments (34)	Surface	Dice coefficient and Wu-Palmer similarity between argument pairs;
	Coreference	Entity coreference between argument pairs; Numeric word match between the argument pairs.
	Existence	Whether each argument slot is instantiated.
	Location	Containment and alternatives name match between the location arguments based on geographical resources such as DBpedia [4] and GeoNames.

Table 3.4: List of features (with counts) in the pairwise model. Entity coreference are from the Stanford Entity Coreference Engine [79].

September 25, 2018  
DRAFT

# Chapter 4

## Event Sequencing

### 4.1 Introduction

Events are important building blocks of documents. They play a key role in document understanding tasks, such as information extraction [22], news summarization [149], story understanding [103]. Conceptually, **events** correspond to state changes and normally include a location, a time interval, and several entities/participants. In a text, events are realized as text spans, normally as verbs and nouns that indicate state changes [148]. The text spans are often referred to as **event mentions** or **event nuggets** (We use the term event mention in this paper.). The textual mentions of events have rich relations among them, and collectively convey the meaning of one or more related documents. In this paper, we study two different types of relation: **Event Hopper Coreference (EH)** and **Event Sequencing (ES)**.

**Event Coreference:** There is a rich literature on the Event Coreference problem [2, 39, 89, 90, 91, 115]. By analogy to entity coreference, the “same” conceptual event may be realized by multiple text spans (event mentions). The coreference problem aims at identifying these relations to recover events from the text spans. The **Event Hopper** Coreference task in the TAC-KBP evaluation campaign defines coreference links as follows [99]: Two event mentions are considered coreferent if they refer to the conceptually same underlying event, even if their arguments are not strictly identical. For example, mentions that share similar temporal and location scope, though not necessarily the same expression, are considered to be coreferent (*Attack in Baghdad on Thursday* vs. *Bombing in the Green Zone last week*). This means that the event arguments of coreferential events mentions can be non-coreferential (18 killed vs. dozens killed), as long as they refer to the same event, judging from the available evidence.

**Event Sequencing:** The coreference relations build up events from scattered mentions. On the basis of events, various other types of relations can then be established between them. The Event Sequencing task studies one such relation. The task is motivated by Schank’s *scripts* [137], which suggests that human organize information through procedural data structures, reassembling sequences of events. For example, the list of verbs *order*, *eat*, *pay*, *leave* may trigger the restaurant script. A human can conduct reasoning with a typical ordering of these events based on common sense (e.g., *order* should be the first event, *leave* should be the last event).

The ES task studies how to group and order events from text documents belonging to the same

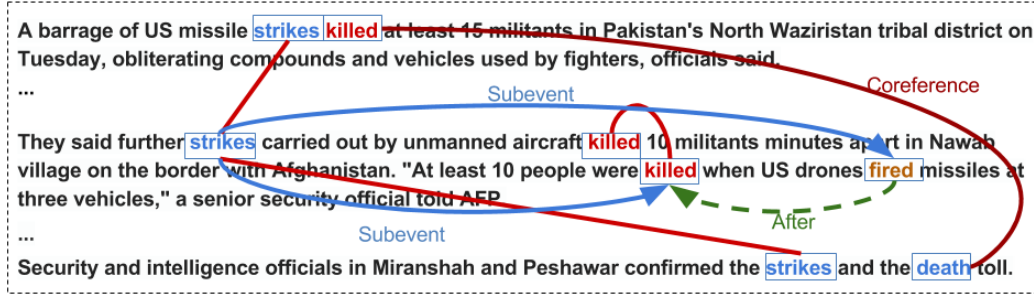


Figure 4.1: Example of Event Coreference and Sequence relations. Red lines are coreference links; solid blue arrows represent Subevent relations; dotted green arrows represent After relations.

script. Figure ?? shows some annotation examples. Conceptually, event sequencing relations hold between the events, while coreference relations hold between textual event mentions. Given a document, the ES task requires systems to identify events within the same script and classify their inter-relations. These relations can be represented as labeled Directed Acyclic Graphs (DAGs). There are two types of relations<sup>1</sup>: **After** relations connect events following script orders (e.g. *order* followed by *eating*); **Subevent** relations connect events to a larger event that contains them. In this paper, we focus only on the **After** relations.

Since script-based understanding is built in the ES task, it has some unique properties comparing to pure temporal ordering: 1. event sequences from different scripts provide separate logical divisions of text, while temporal ordering considers all events to lie on a single timeline; 2. temporal relations for events occurring at similar time points may be complicated. Script-based relations may alleviate the problem. For example, if a *bombing* kills some people, the temporal relation of the *bombing* and *kill* may be “inclusion” or “after”. This is considered an **After** relation in ES because *bombing* causes the *killing*.

For structure prediction, decoding — recovering the complex structure from local decisions — is one of the core problems. The most successful decoding algorithm for coreference nowadays is mention ranking based [12, 51, 81]. These models rank the antecedents (mentions that appear earlier in discourse) and recover the full coreference clusters from local decisions. However, unlike coreference relations, sequencing relations are directed. Coreference decoding algorithms cannot be directly applied to such relations (§4.3.1). To solve this problem, we propose a unified graph-based framework that tackles both event coreference and event sequencing. Our method achieves state-of-the-art results on the event coreference task (§4.4.4) and beats an informed baseline on the event sequencing task (§4.4.5). Finally, we analyze the results and discuss the difficult challenges for both tasks (§4.5). Detailed definitions of these tasks can be found in the corresponding task documents<sup>2</sup>.

<sup>1</sup>Detailed definition of relations can be found in <http://cairo.lti.cs.cmu.edu/kbp/2016/after/>

<sup>2</sup><http://cairo.lti.cs.cmu.edu/kbp/2017/event/documents>

## 4.2 Related Work

Many researchers have worked on event coreference tasks since Humphreys et al. [75]. Recent advances in event coreference have been promoted by the availability of annotated corpora. However, due to the complex nature of events, approaches to event coreference adopt quite different assumptions and definitions. Most of event coreference researches are conducted on the popular ACE corpus [25, 26, 27, 28, 136]. Unlike the TAC KBP setting, the definition of event coreference in the ACE corpus requires strict argument matching. Work on the Intelligence Community (IC) Corpus [3, 38, 74, 89] considers event relations on a restricted domain (i.e., terrorist events). Works on the ECB corpus [39, 80] focuses on both within-document and cross-document coreference.

Our work follows the line of work promoted by the TAC-KBP event nugget tasks [98]. There is a small but growing amount of work on conducting event coreference on the TAC-KBP datasets [90, 91, 115]. The TAC dataset uses a relaxed coreference definition comparing to other corpora, requiring two event mentions to intuitively refer to the same real-world event despite differences of their participants.

For event sequencing, there are few supervised methods on script-like relation classification due to the lack of data. To the best of our knowledge, the only work in this direction is by Araki et al. [3]. This work focuses on the other type of relations in the event sequencing task: **Subevent** relations. There is also a rich literature on unsupervised script induction [19, 30, 58, 118, 133] that extracts scripts as a type of common-sense knowledge from raw documents. The focus of this work is to make use of massive collections of text documents to mine event co-occurrence patterns. In contrast, our work focuses on parsing the detailed relations between event mentions in each document.

Another line of work closely related to event sequencing is to detect other temporal relations between events. Recent computational approaches for temporal detection are mainly conducted on the TimeBank corpus [124]. There have been several studies on building automatic temporal reasoning systems [18, 47, 146]. In comparison, the Event Sequencing task is motivated by the Script theory, which places more emphasis on common-sense knowledge about event chronology.

## 4.3 Model

### 4.3.1 Graph-Based Decoding Model

In the Latent Antecedent Tree (LAT) model popularly used for entity coreference decoding [12, 56], each node represents an event mention and each arc a coreference relation, and new mentions are connected to some past mention considered most similar. Thus the LAT model represents the decoding structure as a tree. This can represent any coreference cluster, because coreference relations are by definition equivalence relations<sup>3</sup>.

In contrast, tree structures cannot always fully cover an Event Sequence relation graph, because 1. the After links are directed, not symmetric, and 2. multiple event nodes can link to one node, resulting in multiple parents.

<sup>3</sup>An equivalence relation is reflexive, symmetric and transitive.

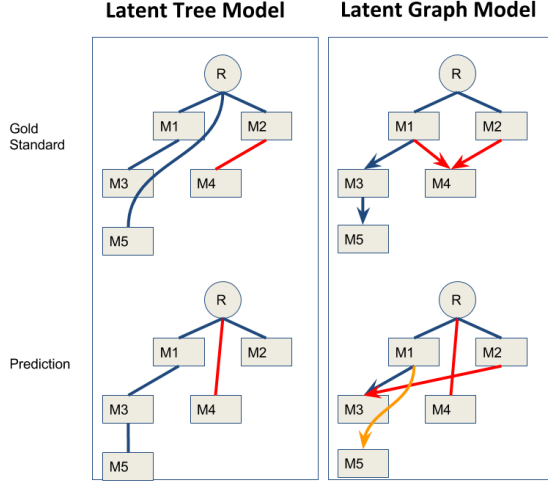


Figure 4.2: Latent Tree Model (left): tree structure formed by undirected links. Latent Graph Model (right): a DAG form by directed links. Dashed red links highlight the discrepancy between prediction and gold standard. The dotted yellow link (bottom right) can be inferred from other links.

To solve this problem, we extend the LAT model and propose its graph version, namely the Latent Antecedent Graph (LAG) model. Figure 4.2 contrast LAT and LAG with decoding examples. The left box shows two example decoded trees in LAT, where each node has one single parent. The right box shows two example decoded trees in LAG, where each node can be linked to multiple parents.

Formally, we define the series of (pre-extracted) event mentions of the document as  $M = \{m_0, m_1, \dots, m_n\}$ , following their discourse order.  $m_0$  is an artificial root node preceding all mentions. For each mention  $m_j$ , let  $A_j = \{m_0, m_1, \dots, m_{j-1}\}$  be the set of its potential antecedents: Let  $\mathcal{A}$  denotes the set of antecedents for all the mentions in the sequence  $\{A_0, A_1, \dots, A_n\}$ . The two tasks in question can be considered as finding the appropriate antecedent(s) from  $\mathcal{A}$ . Similarly, we define the gold antecedent set  $\tilde{\mathcal{A}} = \{\tilde{A}_0, \tilde{A}_1, \dots, \tilde{A}_n\}$ , where  $\tilde{A}_i$  represent the set of antecedents of  $m_i$  allowed by the gold standard. In the coreference task,  $\tilde{A}_i$  contains all antecedents that are coreferent with  $m_i$ . In the sequencing task,  $\tilde{A}_i$  contains all antecedents that have an *After* relation to  $m_i$ .

We can now describe the decoding process. We represent each arc as  $\langle m_i, m_j, r \rangle (i < j)$ , where  $r$  is the relation name. The relation direction can be specified in the relation name  $r$  (e.g.  $r$  can be *after.forward* or *after.backward*). Further, an arc from the root node  $m_0$  to node  $m_j$  represents that  $m_j$  does not have any antecedent. The score of the arc is the dot product between the weight parameter  $\vec{w}$  and a feature vector  $\Phi(\langle m_i, m_j, r \rangle)$ , where  $\Phi$  is an arc-wise feature function. The decoded graph  $z$  can be determined by a set of binary variables  $\vec{z}$ , where  $\vec{z}_{ijr} = 1$  if there is an arc  $\langle m_i, m_j, r \rangle$  or 0 otherwise. The final score of  $z$  is the sum of scores of all arcs:

$$score(z) = \sum_{i,j,r} \vec{z}_{ijr} \vec{w} \cdot \Phi(\langle m_i, m_j, r \rangle) \quad (4.1)$$

The decoding step is to find the output  $\hat{z}$  that maximizes the scoring function:

$$\hat{z} = \arg \max_{z \in \mathcal{Z}(\mathcal{A})} \text{score}(z) \quad (4.2)$$

where  $\mathcal{Z}(\mathcal{A})$  denotes all possible decoding structures given the antecedent sets  $\mathcal{A}$ . It is useful to note that the decoding step can be applied in the same way to the gold antecedent set  $\tilde{\mathcal{A}}$ .

Algorithm 2 shows the Passive-Aggressive training algorithm [37] used in our decoding framework. Line 7 decodes the maximum scored structure from all possible gold standard structures using the current parameters  $\vec{w}$ . Intuitively, this step tries to find the **“easiest” correct graph** — the correct graph with the highest score — for the current model. Several important components remain unspecified in algorithm 2: (1) the decoding step (line 5, 7); (2) the match criteria: whether to consider the system decoding structure as correct (line 6); (3) feature delta: computation of feature difference (line 8); (4) loss computation (line 9). We detail the actual implementation of these steps in §4.3.1.

---

**Algorithm 2** PA algorithm for training

---

**Input:** Training data  $D$ , number of iterations  $T$

```

1:
Output: Weight vector  $\vec{w}$ 
2:
3:  $\vec{w} = \vec{0}$ ;
4:  $\langle \mathcal{A}, \tilde{\mathcal{A}} \rangle \in D$ ;
5: for  $t \leftarrow 1..T$  do  $\hat{z} = \arg \max_{z \in \mathcal{Z}(\mathcal{A})} \text{score}(z)$ 
6:   if  $\neg \text{Match}(\hat{z}, \tilde{\mathcal{A}})$  then
7:      $\tilde{z} = \arg \max_{z \in \mathcal{Z}(\tilde{\mathcal{A}})} \text{score}(z)$ 
8:      $\Delta = \text{FeatureDelta}(\tilde{z}, \hat{z})$ 
9:      $\tau = \frac{\text{loss}(\tilde{z}, \hat{z})}{\|\Delta\|^2}$ 
10:     $w = w + \tau \Delta$ 
return  $w$ 
```

---

### Minimum Decoding Structure

Similar to the LAT model, there may be many decoding structures representing the same configuration. In LAT, since there is exactly one link per node, the number of links in different decoding structures is the same, hence comparable. In LAG, however, one node is allowed to link to multiple antecedents, creating a potential problem for decoding. For example, consider the sequence  $m_1 \xrightarrow{\text{after}} m_2 \xrightarrow{\text{after}} m_3$ , both of the following structures are correct:

1.  $\langle m_1, m_2, \text{after} \rangle, \langle m_2, m_3, \text{after} \rangle$
2.  $\langle m_1, m_2, \text{after} \rangle, \langle m_2, m_3, \text{after} \rangle, \langle m_1, m_3, \text{after} \rangle$

However, the last relation in the second decoding structure can actually be inferred via transitivity. We do not intend to spend the modeling power on such cases. We empirically avoid such redundant cases by using the **transitive reduction graph** for each structure. For a directed acyclic graph, a transitive reduction graph contains the fewest possible edges that

have the same reachability relation as the original graph. In the example above, structure 1 is a transitive reduction graph for structure 2. We call the decoding structures that corresponding to the reduction graphs as *minimum decoding structures*. For LAG, we further restrict  $\mathcal{Z}(\mathcal{A})$  to contain only minimum decoding structures.

### Training Details in Latent Antecedent Graph

In this section, we describe the decoding details for LAG. Note that if we enforce a single antecedent for each node (as in our coreference model), it falls back to the LAT model [12].

**Decoding:** We use a greedy **best-first decoder** [105], which makes a left-to-right pass over the mentions. The decoding step is the same for line 5 and 7. The only difference is that we will use gold antecedent set ( $\tilde{\mathcal{A}}$ ) at line 7. For each node  $m_j$ , we keep all links that score higher than the root link  $\langle 0, m_j, r \rangle$ .

**Cycle and Structure Check:** Incremental decoding a DAG may introduce cycles to the graph, or violate the minimum decoding structure criterion. To solve this, we maintain a set  $R(m_i)$  that is reachable from  $m_i$  during the decoding process. We reject a new link  $\langle m_j, m_i \rangle$  if  $m_j \in R(m_i)$  to avoid cycles. We also reject a redundant link  $\langle m_i, m_j \rangle$  if  $m_j \in R(m_i)$  to keep a minimum decoding structure. Our current implementation is greedy, we leave investigations of search or global inference based algorithms to future work.

**Selecting the Latent Event Mention Graph:** Note that sequence relations are on the event level. Given a unique event graph, it may still correspond to multiple mention graphs. In our implementation, we use a minimum set of event mentions to represent the full event graph by taking one single mention from each event. Following the “easiest” intuition, we select the single mention that will result in the highest score given the current feature weight  $w$ .

**Match Criteria:** We consider two graphs to match when their inferred graphs are the same. The inferred graph is defined by taking the transitive closure of the graph and propagate the links through the coreference relations. For example, in Figure 4.1, the mention `fired` will be linked to two `killed` mentions after propagation.

**Feature Delta:** In structural perceptron training [35], the weights are updated directly by the feature delta. For all the features  $\tilde{f}$  of the gold standard graph  $\tilde{z}$  and features  $\hat{f}$  of a decoded graph  $\hat{z}$ , the feature delta is simply:  $\Delta = \tilde{f} - \hat{f}$ . However, a decoded graph may contain links that are not directly presented but inferable from the gold standard graph. For example, in Figure 4.2, the prediction graph has a link from  $M5$  to  $M1$  (the orange arc), which is absent but inferable from the gold standard tree. If we keep these links when computing  $\Delta$ , the model does not converge well. We thus remove the features on the inferable links from  $\hat{f}$  when computing  $\Delta$ .

**Loss:** We define the loss to be the number of different edges in two graphs. Following Björkelund and Kuhn [12], we further penalize erroneous root attachment: an incorrect link to the root  $m_0$  adds the loss by 2. For example, in Figure 4.2 the prediction graph (bottom right) incorrectly links  $m_4$  to Root and misses a link to  $m_3$ , which cause a total loss of 3. In addition, to be consistent with the feature delta computation, we do not compute loss for predicted links that are inferable from the gold standard.



Head	Headword token and lemma pair, and whether they are the same.
Type	The pair of event types, and whether they are the same.
Realis	The pair of realis types and whether they are the same.
POS	POS pair of the two mentions and whether they are the same.
Exact Match	Whether the 5-word windows of the two mentions matches exactly.
Distance	Sentence distance between the two mentions.
Frame	Frame name pair of the two mentions and whether they are the same.
Syntactic	Whether a mention is the syntactic ancestor of another.

Table 4.1: Coreference Features. Parsing is done using Stanford CoreNLP [93]; frame names are produced by Semafor [43].

## 4.3.2 Features

### Event Coreference Features

For event coreference, we design a simple feature set to capture syntactic and semantic similarity of arcs. The main features are summarized in Table 4.1. In the TAC KBP 2015 coreference task setting, the event mentions are annotated with two attributes. There are 38 event types and subtype pairs (e.g., *Business.Merge-Org*, *Conflict.Attack*). There also 3 realis type: events that actually occurred are marked as *Actual*; events that are not specific are marked as *Generic*; other events such as future events are marked as *Other*. For these two attributes, we use the gold annotations in our feature sets.

### Event Sequencing Features

An event sequencing system needs to determine whether the events are in the same script and order them. We design separate feature sets to capture these aspects: the Script Compatibility set considers whether mentions should belong to the same script; the Event Ordering set determines the relative ordering of the mentions. Our final features are the cross products of features from the following 3 sets.

1. **Surface-Based Script Compatibility:** these features capture whether two mentions are script compatible based on the surface information, including:
  - Mention headword pair.
  - Event type pair.
  - Whether two event mentions appear in the same cluster in Chambers’s event schema database [21].

- Whether the two event mentions share arguments, and the semantic frame name of the shared argument (produced by the Semafor parser [43]).
2. **Discourse-Based Script Compatibility:** these features capture whether two event mentions are related given the discourse context.
    - Dependency path between the two mentions.
    - Function words (words other than Noun, Verb, Adjective and Adverb) in between the two mentions.
    - The types of other event mentions between the two mentions.
    - The sentence distance of two event mentions.
    - Whether there are temporal expressions in the sentences of the two mentions, extracted from the AGM-TMP slot using a PropBank parser [144]
  3. **Event Ordering:** this feature set tries to capture the ordering of events. We use the discourse ordering of two mentions (forward: the antecedent is the parent; backward: the antecedent is the child), and temporal ordering produced by CAEVO [18].

Taking the *after* arc from *fired* to *killed* in Figure 4.1 as an example, a feature after the cross product is: Event type pair is *Conflict.Attack* and *Life.Die*, discourse ordering is *backward*, and sentence distance is 0.

## 4.4 Experiments

### 4.4.1 Dataset

We conduct experiments on the dataset released in Text Analysis Coreference (TAC-KBP) 2017 Event Sequencing task (released by LDC under the catalog name LDC2016E130). This dataset contains rich event relation annotations, with event mentions and coreference annotated in TAC-KBP 2015, and additional annotations on Event Sequencing<sup>4</sup>. There are 158 documents in the training set and 202 in the test set, selected from general news articles and forum discussion threads. The event mentions are annotated with 38 type-subtype and 3 realis status (Actual, Generic, Other). Event Hopper, After, and Subevent links are annotated between event mentions. For all experiments, we develop our system and conduct ablation studies using 5-fold cross-validation on the training set, and report performance on the test set.

### 4.4.2 Baselines and Benchmarks

**Coreference:** we compare our event coreference system against the top performing systems from TAC-KBP 2015 (LCC, UI-CCG, and LTI). In addition, we also compare the results against two official baselines [98]: the Singleton baseline that put each event mention in its own cluster and the Match baseline that creates clusters based on mention type and realis status match.

**Sequencing:** This work is an initial attempt to this problem, so there is currently no comparable

<sup>4</sup><http://cairo.lti.cs.cmu.edu/kbp/2016/after/>

prior work on the same task. We instead compare with a baseline using event temporal ordering systems. We use a state-of-the-art temporal system named *Caevo* [18]. To make a fair comparison, we feed the gold standard event mentions to the system along with mentions predicted by *Caevo*<sup>5</sup>. However, since the script-style After links are only connected between mentions in the same script, directly using the output of *Caevo* produces very low precision. Instead, we run a stronger baseline: we take the gold standard script clusters and then only ask *Caevo* to predict links within these clusters (Oracle Cluster + Temporal).

### 4.4.3 Evaluation Metrics

**Evaluating Event Coreference:** We evaluate our results using the official scorer provided by TAC-KBP, which uses 4 coreference metrics: *BLANC* [126], *MUC* [31], *B<sup>3</sup>* [5] and *CEAF-E* [92]. Following the TAC KBP task, systems are ranked using the average of these 4 metrics.

**Evaluating Event Sequencing:** The TAC KBP scorer evaluates event sequencing using the metric of the TempEval task [145, 147]. The TempEval metric calculates special precision and recall values based on the closure and reduction graphs:

$$Precision = \frac{|Response^- \cap Reference^+|}{|Response^-|} \quad Recall = \frac{|Reference^- \cap Response^+|}{|Reference^-|}$$

where *Response* represents the After link graph from the system response and *Reference* represents the After link graph from the gold standard.  $G^+$  represents the graph closure for graph  $G$  and  $G^-$  represents the graph reduction for graph  $G$ . As preprocessing, relations are automatically propagated through coreference clusters (currently using gold standard clusters). The final score is the standard F-score: geometric mean of the precision and recall values.

### 4.4.4 Evaluation Results for Event Coreference

The test performance on Event Coreference is summarized in Table 4.2. Comparing to the top 3 coreference systems in TAC-KBP 2015, we outperform the best system by about 2 points absolute F-score on average. Our system is also competitive on individual metrics. Our model performs the best based on *B<sup>3</sup>* and *CEAF-E*, and is comparable to the top performing systems on *MUC* and *BLANC*.

Note that while the *Matching* baseline only links event mentions based on event type and realis status, it is very competitive and performs close to the top systems. This is not surprising since these two attributes are based on the gold standard. To take a closer look, we conduct an ablation study by removing the simple match features one by one. The results are summarized in Table 4.3. We observe that some features produce mixed results on different metrics: they provide improvements on some metrics but not all. This is partially caused by the different characteristics of different metrics. On the other hand, these features (parsing and frames) are automatically predicted, which make them less stable. Furthermore, the Frame features contain duplicate information to event types, which makes it less useful in this setting.

Besides the presented features, we have also designed features using event argument. However, we do not report the results since the argument features decrease the performance on all metrics.

<sup>5</sup>We keep the mentions predicted by *Caevo* because its inference may be affected by these mentions.

	$B^3$	CEAF-E	MUC	BLANC	AVG.
Singleton	78.10	68.98	0.00	48.88	52.01
Matching	78.40	65.82	<b>69.83</b>	76.29	71.94
LCC	82.85	74.66	68.50	<b>77.61</b>	75.69
UI-CCG	83.75	75.81	63.78	73.99	74.28
LTI	82.27	75.15	60.93	71.57	72.60
This work	<b>85.59</b>	<b>79.65</b>	67.81	77.37	<b>77.61</b>

Table 4.2: Test Results for Event Coreference with the Singleton and Matching baselines.

	$B^3$	CEAF-E	MUC	BLANC	AVG.
ALL	81.97	74.80	76.33	76.07	77.29
-Distance	81.92	74.48	76.02	77.55	77.50
-Frame	82.14	75.01	76.28	77.74	77.79
-Syntactic	81.87	74.89	75.79	76.22	77.19

Table 4.3: Ablation study for Event Coreference.

#### 4.4.5 Evaluation Results for Event Sequencing

The evaluation results on Event Sequencing is summarized in Table 4.4. Because the baseline system has access to the oracle script clusters, it produces high precision. However, the low recall value shows that it fails to produce enough After links. Our analysis shows that a lot of After relations are not indicated by clear temporal clues, but can only be solved with script knowledge. In Example 4.3, the baseline system is able to identify “fled” is after “ousted” from explicit marker “after”. However, it fails to identify that “extradited” is after “arrested”, which requires knowledge about prototypical event sequences.

- (4.3) Eight months after the [transport fled] Ivory Coast when Gbagbo, the former president, was [End.Position ousted] by the French military. Blé Goudé was subsequently [Jail arrested] in Ghana and [transport extradited] Megrahi, [Jail jailed] for [Attack killing] 270 people in 1988.<sup>6</sup>

In our error analysis, we noticed that our system produces a large number of relations due to coreference propagation. One single wrong prediction can cause the error to propagate.

Besides memorizing the mention pairs, our model also tries to capture script compatibility through discourse signals. To further understand how much these signals help, we conduct an ablation study of the features in the discoursed based compatibility features (see §4.3.2). Similarly, we remove each feature group from the full feature set one by one and observe the performance change.

<sup>6</sup>The small red text indicates the event type for each mention.

	Prec.	Recall	F-Score
Oracle Cluster+Temporal	<b>46.21</b>	8.72	14.68
Our Model	18.28	<b>16.91</b>	<b>17.57</b>

Table 4.4: Test results for event sequencing. The Oracle Cluster+Temporal system is running CAEVO on the Oracle Clusters.

	Prec.	Recall	F-Score	$\Delta$
Full	37.92	36.79	36.36	
- Mention Type	32.78	29.81	30.07	6.29
- Sentence	33.90	30.75	31.00	5.36
- Temporal	37.21	36.53	35.81	0.55
- Dependency	38.18	36.44	36.23	0.13
- Function words	38.08	36.51	36.18	0.18

Table 4.5: Ablation Study for Event Sequencing.

The results are reported in Table 4.5. While most of the features only affect the performance by less than 1 absolute F1 score, the feature sets after removing *mention* or *sentences* show a significant drop in both precision and recall. This shows that discourse proximity is the most significant ones among these features. In addition, the *mention* feature set captures the following *explain away* intuition: the event mentions A and B are less likely to be related if there are similar mentions in between. One such example can be seen in Figure 4.1, the event mention *fired* is more likely to relate to the closest *killed*, instead of the other *killed* in the first paragraph.

In addition, our performance on the development set is higher than the test set. Further analysis reveals two causes: 1. the coreference propagation step causes the scores to be very unstable, 2. our model only learns limited common sense ordering based on lexical pairs, which can easily overfit to the small training corpus. Since the annotation is difficult to scale, it is important to use methods to harvest script common sense knowledge automatically, as in the script induction work [19].

## 4.5 Discussion

### 4.5.1 Event Coreference Challenges

Although we have achieved good performance on event coreference, upon closer investigation we found that most of the coreference decisions are still made based on simple word/lemma matching

(note that the type and realis baseline is as high as 0.72 F1 score). The system exploits little semantic information to resolve difficult event coreference problems. A major challenge is that our system is not capable of utilizing event arguments: in fact, Hasler and Orasan [72] found that only around 20% of the arguments in the same event slot are actually coreferent for coreferential event pairs in the ACE 2005 corpus. Furthermore, the TAC-KBP corpus uses a relaxed participant identity requirement for event coreference, which makes argument-based matching more difficult.

### 4.5.2 Event Sequencing Challenges

Our event sequencing performance is still low despite the introduction of many features. This task is inherently difficult because it requires a system to solve both the script clustering and event ordering tasks. The former task requires both common-sense knowledge and discourse reasoning. Reasoning is more important for long-term links since there are no explicit clues like prepositions and dependencies to be exploited. The ablation study shows that discourse features like sentence distance are more effective, which indicates that our model mainly relies on surface clues and has limited reasoning power.

Furthermore, we observe a strong locality property of After links by skimming the training data: most After link relations are found in a small local region. Since reasoning and coreference based propagation will accumulate local decisions, a system must be accurate on them.

#### The Ambiguous Boundary of a Script

Besides the above-mentioned challenges, a more fundamental problem is to define the boundary of scripts. Since the definition of scripts is only prototypical event sequences, the boundaries between them are not clear. In Example 4.3, the event `jailed` is considered to belong to a “Judicial Process” script and `killing` is considered to belong to an “Attack” script<sup>7</sup>. No link is annotated between these two mentions since they are considered to belong to different clusters, even though the “jailed” event is to punish the “killing”. Therefore essentially, the current Event Sequencing task simply requires the system to fit these human defined boundaries. In principle, the “Judicial Process” script and the “Attack” script can form a larger script structure, on a higher hierarchical level.

While it is possible to manually define scripts and what kind of events they may contain specifically in a controlled domain, it is difficult to generalize the relations. Most previous work on script induction [19, 30, 58, 118, 133] treats scripts as statistical models where probabilities can be assigned, thereby avoiding the boundary problem. While the script boundaries may be application dependent, a possible solution may rely on the “Goals” in Schank’s script theory. The Goal of a script is the final state expected (by the script protagonist) from the sequence of events. Goal oriented scripts may be able to help us explain whether `killing` and `jailed` should be separate: if we take the “killer” as the protagonist, the goal of “kill” is achieved at the point of the victim dying. We leave the investigation on proper theoretical justification to future work.

<sup>7</sup>Script names are taken from the annotation guideline: <http://cairo.lti.cs.cmu.edu/kbp/2016/after/annotation>

## 4.6 Conclusion

In this paper, we presented a unified graph framework to conduct event coreference and sequencing. We have achieved state-of-the-art results on event coreference and report the first attempt at event sequencing. While we only studied two types of relations, we believe the method can be adopted in broader contexts. In the future, we plan to build a joint model to allow the tasks to mutually improve each other.

In general, analyzing event structure can bring new aspects of knowledge from text. For instance, Event Coreference systems can help group scattered information together. Understanding Event Sequencing can help clarify the discourse structure, which can be useful in other NLP applications, such as solving entity coreference problems [114]. However, in our investigation, we find that the linguistic theory and definitions for events are not adequate for the computational setting. For example, proper theoretical justification is needed to define event coreference, which should explain the problems, such as argument mismatches. In addition, we also need a theoretical basis for script boundaries. In the future, we will devote our effort to understanding the theoretical and computational aspects of events relations, and utilizing them for other NLP tasks.

September 25, 2018  
DRAFT



# Chapter 5

## Verb Phrase Ellipsis

### 5.1 Introduction

Verb Phrase Ellipsis (VPE) is the anaphoric process where a verbal constituent is partially or totally unexpressed, but can be resolved through an antecedent in the context, as in the following examples:

(5.1) His wife also [antecedent *works for the paper*], as **did** his father.

(5.2) In particular, Mr. Coxon says, businesses are [antecedent *paying out a smaller percentage of their profits and cash flow in the form of dividends*] than they **have** historically.

In example 5.1, a light verb **did** is used to represent the verb phrase *works for the paper*; example 5.2 shows a much longer antecedent phrase, which in addition differs in tense from the elided one. Following Dalrymple et al. [40], we refer to the full verb expression as the “antecedent”, and to the anaphor as the “target”.

VPE resolution is necessary for deeper Natural Language Understanding, and can be beneficial for instance in dialogue systems or Information Extraction applications.

Computationally, VPE resolution can be modeled as a pipeline process: first detect the VPE targets, then identify their antecedents. Prior work on this topic [69, 110] has used this pipeline approach but without analysis of the interaction of the different steps.

In this paper, we analyze the steps needed to resolve VPE. We preserve the target identification task, but propose a decomposition of the antecedent selection step in two subtasks. We use learning-based models to address each task separately, and also explore the combination of contiguous steps. Although the features used in our system are relatively simple, our models yield state-of-the-art results on the overall task. We also observe a small performance improvement from our decomposition modeling of the tasks.

There are only a few small datasets that include manual VPE annotations. While Bos and Spenader [13] provide publicly available VPE annotations for Wall Street Journal (WSJ) news documents, the annotations created by Nielsen [110] include a more diverse set of genres (e.g., articles and plays) from the British National Corpus (BNC).

We semi-automatically transform these latter annotations into the same format used by the former. The unified format allows better benchmarking and will facilitate more meaningful

comparisons in the future. We evaluate our methods on both datasets, making our results directly comparable to those published by Nielsen [110].

## 5.2 Related Work

Considerable work has been done on VPE in the field of theoretical linguistics: e.g., [40, 140]; yet there is much less work on computational approaches to resolving VPE.

Hardt [1992, 1997] presents, to our knowledge, the first computational approach to VPE. His system applies a set of linguistically motivated rules to select an antecedent given an elliptical target. Hardt [71] uses Transformation-Based Learning to replace the manually developed rules. However, in Hardt’s work, the targets are selected from the corpus by searching for “empty verb phrases” (constructions with an auxiliary verb only) in the gold standard parse trees.

Nielsen [2005] presents the first end-to-end system that resolves VPE from raw text input. He describes several heuristic and learning-based approaches for target detection and antecedent identification. He also discusses a post-processing substitution step in which the target is replaced by a transformed version of the antecedent (to match the context). We do not address this task here because other VPE datasets do not contain relevant substitution annotations. Similar techniques are also described in Nielsen [2003, 2004, 2004].

Results from this prior work are relatively difficult to reproduce because the annotations on which they rely are inaccessible. The annotations used by Hardt [70] have not been made available, and those used by Nielsen [110] are not easily reusable since they rely on some particular tokenization and parser. Bos and Spenader [13] address this problem by annotating a new corpus of VPE on top of the WSJ section of the Penn Treebank, and propose it as a standard evaluation benchmark for the task. Still it is desirable to use Nielsen’s annotations on the BNC which contain more diverse text genres with more frequent VPE.

## 5.3 Approaches

We focus on the problems of target detection and antecedent identification as proposed by Nielsen [110]. We propose a refinement of these two tasks, splitting them into these three:

1. **Target Detection (T)**, where the subset of VPE targets is identified.
2. **Antecedent Head Resolution (H)**, where each target is linked to the head of its antecedent.
3. **Antecedent Boundary Determination (B)**, where the exact boundaries of the antecedent are determined from its head.

The following sections describe each of the steps in detail.

### 5.3.1 Target Detection

Since the VPE target is annotated as a single word in the corpus<sup>1</sup>, we model their detection as a binary classification problem. We only consider modal or light verbs (*be*, *do*, *have*) as candidates,

<sup>1</sup>All targets in the corpus of Bos and Spenader [13] are single-word by their annotation guideline.

and train a logistic regression classifier ( $\mathbf{Log}^T$ ) with the following set of binary features:

1. The POS tag, lemma, and dependency label of the verb, its dependency parent, and the immediately preceding and succeeding words.
2. The POS tags, lemmas and dependency labels of the words in the dependency subtree of the verb, in the 3-word window, and in the same-size window after (as bags of words).
3. Whether the subject of the verb appears to its right (i.e., there is subject-verb inversion).

### 5.3.2 Antecedent Head Resolution

For each detected target, we consider as potential antecedent heads all verbs (including modals and auxiliaries) in the three immediately preceding sentences of the target word<sup>2</sup> as well as the sentence including the target word (up to the target<sup>3</sup>). This follows Hardt [69] and Nielsen [110].

We perform experiments using a logistic regression classifier ( $\mathbf{Log}^H$ ), trained to distinguish correct antecedents from all other possible candidates. The set of features are shared with the Antecedent Boundary Determination task, and are described in detail in Section 5.3.3.

However, a more natural view of the resolution task is that of a ranking problem. The gold annotation can be seen as a partial ordering of the candidates, where, for a given target, the correct antecedent ranks above all other candidates, but there is no ordering among the remaining candidates. To handle this specific setting, we adopt a ranking model with domination loss [45].

Formally, for each potential target  $t$  in the determined set of targets  $T$ , we consider its set of candidates  $C_t$ , and denote whether a candidate  $c \in C_t$  is the antecedent for  $t$  using a binary variable  $a_{ct}$ . We express the ranking problem as a bipartite graph  $\mathcal{G} = (V^+, V^-, E)$  where vertices represent antecedent candidates:

$$\begin{aligned} V^+ &= \{(t, c) \mid t \in T, c \in C_t, a_{ct} = 1\} \\ V^- &= \{(t, c) \mid t \in T, c \in C_t, a_{ct} = 0\} \end{aligned}$$

and the edges link the correct antecedents to the rest of the candidates for the same target<sup>4</sup>:

$$E = \{((t, c^+), (t, c^-)) \mid (t, c^+) \in V^+, (t, c^-) \in V^-\}$$

We associate each vertex  $i$  with a feature vector  $\mathbf{x}_i$ , and compute its score  $s_i$  as a parametric function of the features  $s_i = g(\mathbf{w}, \mathbf{x}_i)$ . The training objective is to learn parameters  $\mathbf{w}$  such that each positive vertex  $i \in V^+$  has a higher score than the negative vertices  $j$  it is connected to,  $V_i^- = \{j \mid j \in V^-, (i, j) \in E\}$ .

The combinatorial domination loss for a vertex  $i \in V^+$  is 1 if there exists any vertex  $j \in V_i^-$  with a higher score. A convex relaxation of the loss for the graph is given by [45]:

$$f(w) = \frac{1}{|V^+|} \sum_{i \in V^+} \log(1 + \sum_{j \in V_i^-} \exp(s_j - s_i + \Delta))$$

<sup>2</sup>Only 1 of the targets in the corpus of Bos and Spenader [13], has an antecedent beyond that window.

<sup>3</sup>Only 1% of the targets in the corpus are cataphoric.

<sup>4</sup>During training, there is always 1 correct antecedent for each gold standard target, with several incorrect ones.

---

**Algorithm 3** Candidate generation

---

**Input:**  $a$ , the antecedent head

**Input:**  $t$ , the target

**Output:**  $B$ , the set of possible antecedent boundaries (begin, end)

```

1:  $a_s \leftarrow \text{SemanticHeadVerb}(a)$ 
2:  $E \leftarrow \{a_s\}$  // the set of ending positions
3: for  $ch \in \text{RightChildren}(a_s)$  do
4:    $e \leftarrow \text{RightMostNode}(ch)$ 
5:   if  $e < t \wedge \text{ValidEnding}(e)$  then
6:      $E \leftarrow E \cup \{e\}$ 
7:  $B \leftarrow \emptyset$ 
8: for  $e \in E$  do
9:    $B \leftarrow B \cup \{(a, e)\}$ 
10:  if  $a == \text{"be"}$  then
11:    if  $\text{IsVerb}(a + 1)$  then
12:       $A \leftarrow A \cup \{(a + 1, e)\}$ 
13:    for  $s \in \{a + 1, a + 2 \dots e - 1\}$  do
14:      if  $\text{IsAdverb}(s) \wedge \text{IsVerb}(s + 1)$  then
15:         $B \leftarrow B \cup \{(s + 1, e)\}$ 
return  $B$ 

```

---

Taking  $\Delta = 0$ , and choosing  $\mathbf{g}$  to be a linear feature scoring function  $s_i = \mathbf{w} \cdot \mathbf{x}_i$ , the loss becomes:

$$f(w) = \frac{1}{|V^+|} \sum_{i \in V^+} \log \sum_{j \in V_i^-} \exp(\mathbf{w} \cdot \mathbf{x}_j) - \mathbf{w} \cdot x_i$$

The loss over the whole graph can then be minimized using stochastic gradient descent. We will denote the ranker learned with this approach as **Rank<sup>H</sup>**.

### 5.3.3 Antecedent Boundary Determination

From a given antecedent head, the set of potential boundaries for the antecedent, which is a complete or partial verb phrase, is constructed using Algorithm 3.

Informally, the algorithm tries to generate different valid verb phrase structures by varying the amount of information encoded in the phrase. To do so, it accesses the semantic head verb  $a_s$  of the antecedent head  $a$  (e.g., *paying* for *are* in Example 5.2), and considers the rightmost node of each right child. If the node is a valid ending (punctuation and quotation are excluded), it is added to the potential set of endings  $E$ . The set of valid boundaries  $B$  contains the cross-product of the starting position  $S = \{a\}$  with  $E$ .

For instance, from Example 5.2, the following boundary candidates are generated for *are*:

- are paying

- are paying out
- are paying out a smaller percentage of their profits and cash flow
- are paying out a smaller percentage of their profits and cash flow in the form of dividends

We experiment with both logistic regression ( $\mathbf{Log}^B$ ) and ranking ( $\mathbf{Rank}^B$ ) models for this task. The set of features is shared with the previous task, and is described in the following section.

### Antecedent Features

The features used for antecedent head resolution and/or boundary determination try to capture aspects of both tasks. We summarize the features in Table 5.1. The features are roughly grouped by their type. **Labels** features make use of the parsing labels of the antecedent and target; **Tree** features are intended to capture the dependency relations between the antecedent and target; **Distance** features describe distance between them; **Match** features test whether the context of the antecedent and target are similar; **Semantic** features capture shallow semantic similarity; finally, there are a few **Other** features which are not categorized.

On the last column of the feature table, we indicate the design purpose of the feature: head selection (H), boundary detection (B) or both (B&H). However, we use the full feature set for all three tasks.

## 5.4 Joint Modeling

Here we consider the possibility that antecedent head resolution and target detection should be modeled jointly (they are typically separate). The hypothesis is that if a suitable antecedent for a target cannot be found, the target itself might have been incorrectly detected. Similarly, the suitability of a candidate as antecedent head can depend on the possible boundaries of the antecedents that can be generated from it.

We also consider the possibility that antecedent head resolution and antecedent boundary determination should be modeled independently (though they are typically combined). We hypothesize that these two steps actually focus on different perspectives: the antecedent head resolution (**H**) focuses on finding the correct antecedent position; the boundary detection step (**B**) focuses on constructing a well-formed verb phrase. We are also aware that **B** might be helpful to **H**, for instance, a correct antecedent boundary will give us correct context words, that can be useful in determining the antecedent position.

We examine the joint interactions by combining adjacent steps in our pipeline. For the combination of antecedent head resolution and antecedent boundary determination (**H+B**), we consider simultaneously as candidates for each target the set of all potential boundaries for all potential heads. Here too, a logistic regression model ( $\mathbf{Log}^{H+B}$ ) can be used to distinguish correct (target, antecedent start, antecedent end) triplets; or a ranking model ( $\mathbf{Rank}^{H+B}$ ) can be trained to rank the correct one above the other ones for the same target.

The combination of target detection with antecedent head resolution (**T+H**) requires identifying the targets. This is not straightforward when using a ranking model since scores are only comparable for the same target. To get around this problem, we add a “null” antecedent head. For a given target candidate, the null antecedent should be ranked higher than all other candidates if it

Type	Feature Description	Purpose
Labels	The POS tag and dependency label of the antecedent head	H
	The POS tag and dependency label of the antecedent’s last word	B
	The POS tag and lemma of the antecedent parent	H
	The POS tag, lemma and dependency label of within a 3 word around around the antecedent	B
	The pair of the POS tags of the antecedent head and the target, and of their auxiliary verbs	H
	The pair of the lemmas of the auxiliary verbs of the antecedent head and the target	H
Tree	Whether the antecedent and the target form a comparative construction connecting by <i>so</i> , <i>as</i> or <i>than</i>	H&B
	The dependency labels of the shared lemmas between the parse tree of the antecedent and the target	H
	Label of the dependency between the antecedent and target (if exists)	H
	Whether the antecedent contains any descendant with the same lemma and dependency label as a descendant of the target.	H
	Whether antecedent and target are dependent ancestor of each other	H
	Whether antecedent and target share prepositions in their dependency tree	H
Distance	The distance in sentences between the antecedent and the target (clipped to 2)	H
	The number of verb phrases between the antecedent and the target (clipped to 5)	H
Match	Whether the lemmas of the heads, and words in the the window (=2) before the antecedent and the target match respectively	H
	Whether the lemmas of the $i$ th word before the antecedent and $i - 1$ th word before the target match respectively (for $i \in \{1, 2, 3\}$ , with the 0th word of the target being the target itself)	H&B
Semantic	Whether the subject of the antecedent and the target are coreferent	H
Other	Whether the lemma of the head of the antecedent is <i>be</i> and that of the target is <i>do</i> (be-do match, used by Hardt and Nielsen)	H
	Whether the antecedent is in quotes and the target is not, or vice versa	H&B

Table 5.1: Antecedent Features

is not actually a target. Since this produces many examples where the null antecedent should be selected, random subsampling is used to reduce the training data imbalance. The “null” hypothesis

	Documents		VPE Instances	
	Train	Test	Train	Test
WSJ	1999	500	435	119
BNC	12	2	641	204

Table 5.2: Corpus statistics

approach is used previously in ranking-based coreference systems [50, 125].

Most of the features presented in the previous section will not trigger for the null instance, and an additional feature to mark this case is added.

The combination of the three tasks (**T+H+B**) only differs from the previous case in that all antecedent boundaries are considered as candidates for a target, in addition to the potential antecedent heads.

## 5.5 Experiments

### 5.5.1 Datasets

We conduct our experiments on two datasets (see Table 5.2 for corpus counts). The first one is the corpus of Bos and Spenader [13], which provides VPE annotation on the WSJ section of the Penn Treebank. Bos and Spenader [13] propose a train-test split that we follow<sup>5</sup>.

To facilitate more meaningful comparison, we converted the sections of the British National Corpus annotated by Nielsen [110] into the format used by Bos and Spenader [13], and manually fixed conversion errors introduced during the process<sup>6</sup> (Our version of the dataset is publicly available for research<sup>7</sup>.) We use a train-test split similar to Nielsen [110]<sup>8</sup>.

### 5.5.2 Evaluation

We evaluate and compare our models following the metrics used by Bos and Spenader [13].

VPE target detection is a per-word binary classification problem, which can be evaluated using the conventional precision (Prec), recall (Rec) and F1 scores.

Bos and Spenader [13] propose a token-based evaluation metric for antecedent selection. The antecedent scores are computed over the correctly identified tokens per antecedent: precision is the number of correctly identified tokens divided by the number of predicted tokens, and recall is

<sup>5</sup>Section 20 to 24 are used as test data.

<sup>6</sup>We also found 3 annotation instances that could be deemed errors, but decided to preserve the annotations as they were.

<sup>7</sup><https://github.com/hunterhector/VerbPhraseEllipsis>

<sup>8</sup>Training set is CS6, A2U, J25, FU6, H7F, HA3, A19, A0P, G1A, EWC, FNS, C8T; test set is EDJ, FR3

the number of correctly identified tokens divided by the number of gold standard tokens. Averaged scores refer to a “macro”-average over all antecedents.

Finally, in order to assess the performance of antecedent head resolution, we compute precision, recall and F1 where credit is given if the proposed head is included inside the golden antecedent boundaries.

### 5.5.3 Baselines and Benchmarks

We begin with simple, linguistically motivated baseline approaches for the three subtasks. For target detection, we re-implement the heuristic baseline used by Nielsen [110]: take all auxiliaries as possible candidates and eliminate them using part-of-speech context rules (we refer to this as **Pos<sup>T</sup>**). For antecedent head resolution, we take the first non-auxiliary verb preceding the target verb. For antecedent boundary detection, we expand the verb into a phrase by taking the largest subtree of the verb such that it does not overlap with the target. These two baselines are also used in Nielsen [110] (and we refer to them as **Prev<sup>H</sup>** and **Max<sup>B</sup>**, respectively).

To upper-bound our results, we include an oracle for the three subtasks, which selects the highest scoring candidate among all those considered. We denote these as **Ora<sup>T</sup>**, **Ora<sup>H</sup>**, **Ora<sup>B</sup>**.

We also compare to the current state-of-the-art target detection results as reported in Nielsen [110] on the BNC dataset (**Nielsen<sup>T</sup>**)<sup>9</sup>.

## 5.6 Results

The results for each one of the three subtasks in isolation are presented first, followed by those of the end-to-end evaluation. We have not attempted to tune classification thresholds to maximize F1.

### 5.6.1 Target Detection

Table 5.3 shows the performance of the compared approaches on the Target Detection task. The logistic regression model **Log<sup>T</sup>** gives relatively high precision compared to recall, probably because there are so many more negative training examples than positive ones. Despite a simple set of features, the F1 results are significantly better than Nielsen’s baseline **Pos<sup>T</sup>**.

Notice also how the oracle **Ora<sup>T</sup>** does not achieve 100% recall, since not all the targets in the gold data are captured by our candidate generation strategy. The loss is around 7% for both corpora.

The results obtained by the joint models are low on this task. In particular, the ranking models **Rank<sup>T+H</sup>** and **Rank<sup>T+H+B</sup>** fail to predict any target in the WSJ corpus, since the null antecedent is always preferred. This happens because joint modeling further exaggerates the class imbalance: the ranker is asked to consider many incorrect targets coupled with all sorts of hypothesis antecedents, and ultimately learns just to select the null target. Our initial attempts at subsampling the negative examples did not improve the situation. The logistic regression models

<sup>9</sup>The differences in the setup make the results on antecedent resolution not directly comparable.



$\mathbf{Log}^{T+H}$  and  $\mathbf{Log}^{T+H+B}$  are most robust, but still their performance is far below that of the pure classifier  $\mathbf{Log}^T$ .

### 5.6.2 Antecedent Head Resolution

Table 5.4 contains the performance of the compared approaches on the Antecedent Head Resolution task, assuming oracle targets ( $\mathbf{Ora}^T$ ).

First, we observe that even the oracle  $\mathbf{Ora}^H$  has low scores on the BNC corpus. This suggests that some phenomena beyond the scope of those observed in the WSJ data appear in the more general corpus (we developed our system using the WSJ annotations and then simply evaluated on the BNC test data).

Second, the ranking-based model  $\mathbf{Rank}^H$  consistently outperforms the logistic regression model  $\mathbf{Log}^H$  and the baseline  $\mathbf{Prev}^H$ . The ranking model’s advantage is small in the WSJ, but much more pronounced in the BNC data. These improvements suggest that indeed, ranking is a more natural modeling choice than classification for antecedent head resolution.

Finally, the joint resolution models  $\mathbf{Rank}^{H+B}$  and  $\mathbf{Log}^{H+B}$  give poorer results than their single-task counterparts, though  $\mathbf{Rank}^{H+B}$  is not far behind  $\mathbf{Rank}^H$ . Joint modeling requires more training data and we may not have enough to reflect the benefit of a more powerful model.

### 5.6.3 Antecedent Boundary Determination

Table 5.5 shows the performance of the compared approaches on the Antecedent Boundary Determination task, using the soft evaluation scores (the results for the strict scores are omitted for brevity, but in general look quite similar). The systems use the output of the oracle targets ( $\mathbf{Ora}^T$ ) and antecedent heads ( $\mathbf{Ora}^H$ ).

Regarding boundary detection alone, the logistic regression model  $\mathbf{Log}^B$  outperforms the ranking model  $\mathbf{Rank}^B$ . This suggests that boundary determination is more a problem of determining the compatibility between target and antecedent extent than one of ranking alternative boundaries. However, the next experiments suggest this advantage is diminished when gold targets and antecedent heads are replaced by system predictions.

#### Non-Gold Antecedent Heads

Table 5.6 contains Antecedent Boundary Determination results for systems which use oracle targets, but system antecedent heads. When  $\mathbf{Rank}^H$  or  $\mathbf{Log}^H$  are used for head resolution, the difference between  $\mathbf{Log}^B$  and  $\mathbf{Rank}^B$  diminishes, and it is even better to use the latter in the BNC corpus. The models were trained with gold annotations rather than system outputs, and the ranking model is somewhat more robust to noisier inputs.

On the other hand, the results for the joint resolution model  $\mathbf{Rank}^{H+B}$  are better in this case than the combination of  $\mathbf{Rank}^H + \mathbf{Rank}^B$ , whereas  $\mathbf{Log}^{H+B}$  performs worse than any 2-step combination. The benefits of using a ranking model for antecedent head resolution seem thus to outperform those of using classification to determine its boundaries.

	WSJ			BNC		
	Prec	Rec	F1	Prec	Rec	F1
<b>Ora<sup>T</sup></b>	100.00	93.28	96.52	100.00	92.65	96.18
<b>Log<sup>T</sup></b>	80.22	61.34	69.52	80.90	70.59	75.39
<b>Pos<sup>T</sup></b>	42.62	43.7	43.15	35.47	35.29	35.38
<b>Log<sup>T+H</sup></b>	23.36	26.89	25.00	12.52	38.24	18.86
<b>Rank<sup>T+H</sup></b>	0.00	0.00	0.00	15.79	5.88	8.57
<b>Log<sup>T+H+B</sup></b>	25.61	17.65	20.90	21.50	32.35	25.83
<b>Rank<sup>T+H+B</sup></b>	0.00	0.00	0.00	16.67	11.27	13.45
<b>Nielsen<sup>T</sup></b>	—	—	—	72.50	72.86	72.68

Table 5.3: Results for Target Detection

	WSJ			BNC		
	Prec	Rec	F1	Prec	Rec	F1
<b>Ora<sup>H</sup></b>	94.59	88.24	91.30	79.89	74.02	76.84
<b>Rank<sup>H</sup></b>	70.27	65.55	67.83	52.91	49.02	50.89
<b>Prev<sup>H</sup></b>	67.57	63.03	65.22	39.68	36.76	38.17
<b>Log<sup>H</sup></b>	59.46	55.46	57.39	38.62	35.78	37.15
<b>Rank<sup>H+B</sup></b>	68.47	63.87	66.09	51.85	48.04	49.87
<b>Log<sup>H+B</sup></b>	39.64	36.97	38.26	30.16	27.94	29.01

Table 5.4: Results for Antecedent Head Resolution

	WSJ			BNC		
	Prec	Rec	F1	Prec	Rec	F1
<b>Ora<sup>B</sup></b>	95.06	88.67	91.76	85.79	79.49	82.52
<b>Log<sup>B</sup></b>	89.47	83.46	86.36	81.10	75.13	78.00
<b>Rank<sup>B</sup></b>	83.96	78.32	81.04	75.68	70.12	72.79
<b>Max<sup>B</sup></b>	78.97	73.66	76.22	73.70	68.28	70.88

Table 5.5: Soft results for Antecedent Boundary Determination

### 5.6.4 End-to-End Evaluation

Table 5.7 contains the end-to-end performance of different approaches, using the soft evaluation scores.

The trends we observed with gold targets are preserved: approaches using the **Rank**<sup>H</sup> maintain an advantage over **Log**<sup>H</sup>, but the improvement of **Log**<sup>B</sup> over **Rank**<sup>B</sup> for boundary determination is diminished with non-gold heads. Also, the 3-step approaches seem to perform slightly better than the 2-step ones. Together with the fact that the smaller problems are easier to train, this appears to validate our decomposition choice.

## 5.7 Discussion

In this chapter we have explored a decomposition of Verb Phrase Ellipsis resolution into subtasks, which splits antecedent selection in two distinct steps. By modeling these two subtasks separately with two different learning paradigms, we can achieve better performance than doing them jointly, suggesting they are indeed of different underlying nature.

Our experiments show that a logistic regression classification model works better for target detection and antecedent boundary determination, while a ranking-based model is more suitable for selecting the antecedent head of a given target. However, the benefits of the classification model for boundary determination are reduced for non-gold targets and heads. On the other hand, by separating the two steps, we lose the potential joint interaction of them. It might be possible to explore whether we can bring the benefits of the two side: use separate models on each step, but learn them jointly. We leave further investigation of this to future work.

We have also explored jointly training a target detection and antecedent resolution model, but have not been successful in dealing with the class imbalance inherent to the problem.

Our current model adopts a simple feature set, which is composed mostly by simple syntax and lexical features. It may be interesting to explore more semantic and discourse-level features in our system. We leave these to future investigation.

All our experiments have been run on publicly available datasets, to which we add our manually aligned version of the VPE annotations on the BNC corpus. We hope our experiments, analysis, and more easily processed data can further the development of new computational approaches to the problem of Verb Phrase Ellipsis resolution.

	WSJ			BNC		
	Prec	Rec	F1	Prec	Rec	F1
<b>Ora<sup>H</sup>+Ora<sup>B</sup></b>	95.06	88.67	91.76	85.79	79.49	82.52
<b>Rank<sup>H</sup>+Log<sup>B</sup></b>	64.11	59.8	61.88	47.04	43.58	45.24
<b>Rank<sup>H</sup>+Rank<sup>B</sup></b>	63.90	59.6	61.67	49.11	45.5	47.24
<b>Log<sup>H</sup>+Log<sup>B</sup></b>	53.49	49.89	51.63	34.77	32.21	33.44
<b>Log<sup>H</sup>+Rank<sup>B</sup></b>	53.27	49.69	51.42	36.26	33.59	34.88
<b>Rank<sup>H+B</sup></b>	67.55	63.01	65.20	50.68	46.95	48.74
<b>Log<sup>H+B</sup></b>	40.96	38.20	39.53	30.00	27.79	28.85

Table 5.6: Soft results for Antecedent Boundary Determination with non-gold heads

	WSJ			BNC		
	Prec	Rec	F1	Prec	Rec	F1
<b>Ora<sup>T</sup>+Ora<sup>H</sup>+Ora<sup>B</sup></b>	95.06	88.67	91.76	85.79	79.49	82.52
<b>Log<sup>T</sup>+Rank<sup>H</sup>+Rank<sup>B</sup></b>	52.68	40.28	45.65	43.03	37.54	40.10
<b>Log<sup>T</sup>+Rank<sup>H</sup>+Log<sup>B</sup></b>	52.82	40.40	45.78	40.21	35.08	37.47
<b>Log<sup>T</sup>+Log<sup>H</sup>+Rank<sup>B</sup></b>	49.45	37.82	42.86	33.12	28.90	30.86
<b>Log<sup>T</sup>+Log<sup>H</sup>+Log<sup>B</sup></b>	49.41	37.79	42.83	31.32	27.33	29.19
<b>Pos<sup>T</sup>+Prev<sup>H</sup>+Max<sup>B</sup></b>	19.04	19.52	19.27	12.81	12.75	12.78
<b>Log<sup>T</sup>+Rank<sup>H+B</sup></b>	54.82	41.92	47.51	41.86	36.52	39.01
<b>Log<sup>T</sup>+Log<sup>H+B</sup></b>	38.85	29.71	33.67	26.11	22.78	24.33

Table 5.7: Soft end-to-end results

## **Part II**

# **Learning with Indirect Supervision**

September 25, 2018  
DRAFT

# Chapter 6

## Event Saliency

### 6.1 Introduction

Automatic extraction of prominent information from text has always been a core problem in language research. While traditional methods mostly concentrate on the word level, researchers start to analyze higher-level discourse units in text, such as entities [49] and events [34].

Events are important discourse units that form the backbone of our communication. They play various roles in documents. Some are more central in discourse: connecting other entities and events, or providing key information of a story. Others are less relevant, but not easily identifiable by NLP systems. Hence it is important to be able to quantify the “importance” of events. For example, Figure 6.1 is a news excerpt describing a debate around a jurisdiction process: “*trial*” is central as the main discussing topic, while “*war*” is not.

Researchers are aware of the need to identify central events in applications like detecting salient relations [155], and identifying climax in storyline [149]. Generally, the saliency of discourse units is important for language understanding tasks, such as document analysis [10], information retrieval [152], and semantic role labeling [29]. Thus, proper models for finding important events are desired.

In this work, we study the task of **event saliency detection**, to find events that are most relevant to the main content of documents. To build a saliency detection model, one core observation is that salient discourse units are forming discourse relations. In Figure 6.1, the “*trial*” event is connected to many other events: “*charge*” is pressed before “*trial*”; “*trial*” is being “*delayed*”.

We present two saliency detection systems based on the observations. First is a feature based learning to rank model. Beyond basic features like frequency and discourse location, we design features using cosine similarities among events and entities, to estimate the *content organization* [65]: how lexical meaning of elements relates to each other. Similarities from within-sentence or across the whole document are used to capture interactions on both local and global aspects (§6.4). The model significantly outperforms a strong “Frequency” baseline in our experiments.

However, there are other discourse relations beyond lexical similarity. Figure 6.1 showcases some: the script relation [137]<sup>1</sup> between “*charge*” and “*trial*”, and the frame relation [7] between

<sup>1</sup>Scripts are prototypical sequences of events: a *restaurant* script normally contains events like “order”, “eat” and

Federal prosecutors **urged** a **trial judge** today to **deny** defense **requests** to **delay** the **trial** of Zacarias Moussaoui and suggested that Mr. Moussaoui, the only person **charged** in the Sept. 11 **attacks**, was to **blame** for many of the **delays** so far. The **attacks** "were volleys in a **declared war** against the United States and were more than just **acts** of terror," the prosecutors said in a **filing** to the Federal District Court in Alexandria, Va. "Thus, the **victims**' and the nation's interest in a fair and speedy **trial** is beyond **dispute**." Last week, court-appointed defense lawyers **asked** that the **starting** date of the **trial**, now set for Sept. 30, be **delayed** by at least two months to allow them to **wade** through volumes of evidence that prosecutors have presented to them, including more than 1,300 computer discs.

Figure 6.1: Examples annotations. Underlying words are annotated event triggers; the red bold ones are annotated as salient.

“attacks” and “trial” (“attacks” fills the “charges” role of “trial”). Since it is unclear which ones contribute more to salience, we design a Kernel based Centrality Estimation (KCE) model (§6.5) to capture salient specific interactions between discourse units automatically.

In KCE, discourse units are projected to embeddings, which are trained end-to-end towards the salience task to capture rich semantic information. A set of soft-count kernels are trained to weigh salient specific latent relations between discourse units. With the capacity to model richer relations, KCE outperforms the feature-based model by a large margin (§6.7.1). Our analysis shows that KCE can identify several relations between discourse units: including script and frames (Table 6.5). To further understand the nature of KCE, we conduct an *intrusion test* (§6.6.2), which requires a model to identify events from another document. The test shows salient events form tightly related groups with relations captured by KCE.

The notion of salience is subjective and may vary from person to person. We follow the empirical approaches used in entity salience research [49]. We consider the *summarization test*: an event is considered salient if a summary written by a human is likely to include it, since events about the main content are more likely to appear in a summary. This approach allows us to create a large-scale corpus (§6.3).

In this paper, we make three main contributions. First, we present two event salience detection systems, which capture rich relations among discourse units. Second, we observe interesting connections between salience and various discourse relations (§6.7.1 and Table 6.5), implying potential research on these areas. Finally, we construct a large scale event salience corpus, providing a testbed for future research. Our dataset and models are publicly available<sup>2</sup>.

“pay”.

<sup>2</sup>Omit for blind review.



## 6.2 Related Work

Events have been studied on many aspects due to their importance in language. To name a few: event detection [84, 106, 115], coreference [89, 90], temporal analysis [18, 47], sequencing [3], script induction [8, 19, 118, 133].

However, studies on event salience are premature. Some previous work attempts to approximate event salience with word frequency or discourse position [149, 155]. Parallel to ours, Choubey et al. [34] propose a task to find the most dominant event in news articles. They draw connections between event coreference and importance, on hundreds of closed-domain documents, using several oracle event attributes. In contrast, our proposed models are fully learned and applied on more general domains and at a larger scale. We also do not restrict to a single most important event per document.

There is a small but growing line of work on entity salience [48, 49, 120, 152]. In this work, we study the case for events.

Text relations have been studied in tasks like text summarization, which mainly focused on cohesion [68]. Grammatical cohesion methods make use of document level structures such as anaphora relations [9] and discourse parse trees [94]. Lexical cohesion based methods focus on repetitions and synonyms on the lexical level [53, 102, 141]. Though sharing similar intuitions, our proposed models are designed to learn richer semantic relations in the embedding space.

Comparing to the traditional summarization task, we focus on events, which are at a different granularity. Our experiments also unveil interesting phenomena among events and other discourse units.

## 6.3 The Event Salience Corpus

This section introduces our approach to construct a large-scale event salience corpus, including methods for finding event mentions and obtaining saliency labels. The studies are based on the Annotated New York Times corpus [135], a newswire corpus with expert-written abstracts.

**Event Mention Annotation:** Despite many annotation attempts on events [15, 124], automatic labeling of them in general domain remains an open problem. Most of the previous work follows empirical approaches. For example, Chambers and Jurafsky [19] consider all verbs together with their subject and object as events. Do et al. [46] additionally include nominal predicates, using the nominal form of verbs and lexical items under the *Event* frame in FrameNet [7].

There are two main challenges in labeling event mentions. First, we need to decide which lexical items are event triggers. Second, we have to disambiguate the word sense to correctly identify events. For example, the word “phone” can refer to an entity (a physical phone) or an event (a phone call event). We use FrameNet to solve these problems. We first use a FrameNet based parser: Semafor [43], to find and disambiguate triggers into frame classes. We then use the FrameNet ontology to select event mentions.

Our frame based selection method follows the Vendler classes [148], a four way classification of eventuality: *states*, *activities*, *accomplishments* and *achievements*. The last three classes involve state change, and are normally considered as events. Following this, we create an “event-evoking frame” list using the following procedure:

	Train	Dev	Test
# Documents	526126	64000	63589
Avg. # Word	794.12	790.27	798.68
Avg. # Events	61.96	60.65	61.34
Avg. # Salience	8.77	8.79	8.90

Table 6.1: Dataset Statistics.

1. We keep frames that are subframes of *Event* and *Process* in the FrameNet ontology.
2. We discard frames that are subframes of state, entity and attribute frames, such as *Entity*, *Attributes*, *Locale*, etc.
3. We manually inspect frames that are not subframes of the above-mentioned ones (around 200) to keep event related ones (including subframes), such as *Arson*, *Delivery*, etc.

This gives us a total of 569 frames. We parse the documents with Semafor and consider predicates that trigger a frame in the list as candidates. We finish the process by removing the light verbs<sup>3</sup> and reporting events<sup>4</sup> from the candidates, similar to previous research [129].

**Salience Labeling:** For all articles (around 664,911) with a human written abstract in the New York Times Annotated Corpus, we extract event mentions. We then label an event mention as salient if we can find its lemma in the corresponding abstract (Mitamura et al. [98] showed that lemma matching is a strong baseline for event coreference.). For example, in Figure 6.1, event mentions in bold and red are found in the abstract, thus labeled as salient. Data split is detailed in Table 6.1 and §6.6.

## 6.4 Feature-Based Event Salience Model

This section presents the feature-based model, including the features and the learning process.

### 6.4.1 Features

Our features are summarized in Table 6.2.

**Basic Discourse Features:** We first use two basic features similar to Dunietz and Gillick [49]: `Frequency` and `Sentence Location`. `Frequency` is the lemma count of the mention’s syntactic head word [93]. `Sentence Location` is the sentence index of the mention, since the first few sentences are normally more important. These two features are often used to estimate salience [10, 149].

<sup>3</sup>Light verbs carry little semantic information: “appear”, “be”, “become”, “do”, “have”, “seem”, “do”, “get”, “give”, “go”, “have”, “keep”, “make”, “put”, “set”, “take”.

<sup>4</sup>Reporting verbs are normally associated with the narrator: “argue”, “claim”, “say”, “suggest”, “tell”.

Name	Description
Frequency	The frequency of the event lemma in document.
Sentence Location	The location of the first sentence that contains the event.
Event Voting	Average cosine similarity with other events in document.
Entity Voting	Average cosine similarity with other entities in document.
Local Entity Voting	Average cosine similarity with entities in the sentence.

Table 6.2: Event Saliency Features.

**Content Features:** We then design several lexical similarity features, to reflect Grimes’ content relatedness [65]. In addition to events, the relations between events and entities are also important. For example, Figure 6.1 shows some related entities in the legal domain, such as “*prosecutors*” and “*court*”. Ideally, they should help promote the saliency status for event “*trial*”.

Lexical relations can be found both within-sentence (local) or across sentence (global) [68]. We compute the local part by averaging similarity scores from other units in the same sentence. The global part is computed by averaging similarity scores from other units in the document. All similarity scores are computed using cosine similarities on pre-trained embeddings [97].

These lead to 3 content features: `Event Voting`, the average similarity to other events in the document; `Entity Voting`, the average similarity to entities in the document; `Local Entity Voting`, the average similarity to entities in the same sentence. Local event voting is not used as a sentence often contains only 1 event.

## 6.4.2 Model

A Learning to Rank (LeToR) model [88] is used to combine the features. Let  $ev_i$  denote the  $i$ th event in a document  $d$ . Its saliency score is computed as:

$$f(ev_i, d) = W_f \cdot F(ev_i, d) + b \quad (6.1)$$

where  $F(ev_i, d)$  is the features for  $ev_i$  in  $d$  (Table 6.2);  $W_f$  and  $b$  are the parameters to learn.

The model is trained with pairwise loss:

$$\sum_{ev^+, ev^- \in d} \max(0, 1 - f(ev^+, d) + f(ev^-, d)), \quad (6.2)$$

w.r.t.  $y(ev^+, d) = +1$  &  $y(ev^-, d) = -1$ .

$$y(e_i, d) = \begin{cases} +1, & \text{if } e_i \text{ is a salient entity in } d, \\ -1, & \text{otherwise.} \end{cases}$$

where  $ev^+$  and  $ev^-$  represent the salient and non-salient events;  $y$  is the gold standard function. Learning can be done by standard gradient methods.

## 6.5 Neural Event Saliency Model

As discussed in §6.1, the saliency of discourse units is reflected by rich relations beyond lexical similarities, for example, script (“*charge*” and “*trial*”) and frame (a “*trial*” of “*attacks*”). In this section, we present a neural model to capture these relations for event saliency estimation.

### 6.5.1 Kernel-based Centrality Estimation

Inspired by the kernel ranking model [153], we propose Kernel-based Centrality Estimation (KCE), to find and weight semantic relations of interests, in order to better estimate saliency.

Formally, given a document  $d$ , the set of annotated events  $\mathbb{V} = \{ev_1, \dots, ev_i, \dots, ev_n\}$ , KCE first embed an event into vector space:  $ev_i \xrightarrow{Emb} \vec{ev}_i$ . It then extract  $K$  features for each  $ev_i$ :

$$\Phi_K(ev_i, \mathbb{V}) = \{\phi_1(\vec{ev}_i, \mathbb{V}), \dots, \phi_k(\vec{ev}_i, \mathbb{V}), \dots, \phi_K(\vec{ev}_i, \mathbb{V})\}, \quad (6.3)$$

$$\phi_k(\vec{ev}_i, \mathbb{V}) = \sum_{ev_j \in \mathbb{V}} \exp \left( -\frac{(\cos(\vec{ev}_i, \vec{ev}_j) - \mu_k)^2}{2\sigma_k^2} \right). \quad (6.4)$$

$\phi_k(\vec{ev}_i, \mathbb{V})$  is the  $k$ -th Gaussian kernel with mean  $\mu_k$  and variance  $\sigma_k^2$ . It models the interactions between events in its kernel range defined by  $\mu_k$  and  $\sigma_k$ .  $\Phi_K(ev_i, \mathbb{V})$  captures multi-level interactions among events — relations that contribute similarly to saliency are expected to be grouped into the same kernels. Kernels are effective in modeling interactions in Information Retrieval [153], but have not yet been used in discourse analysis.

The final saliency score is computed as:

$$f(ev_i, d) = W_v \cdot \Phi_K(ev_i, \mathbb{V}) + b, \quad (6.5)$$

where  $W_v$  is learned to weight the contribution of the certain relations captured by each kernel.

We then use the exact same learning objective as in equation (6.2). The pairwise loss is first back-propagated through the network to update the kernel weights  $W_v$ . Then the kernels use the gradients thus to update the embeddings thus to capture the meaningful discourse relations for saliency.

Since the features and KCE capture different aspects, combining them may give superior performance. This can be done by combining the two vectors in the final linear layer:

$$f(ev_i, d) = W_v \cdot \Phi_K(ev_i, \mathbb{V}) + W_f \cdot F(ev_i, d) + b \quad (6.6)$$

## 6.5.2 Integrating Entities into KCE

KCE is also used to model the relations between events and entities. For example, in Figure 6.1, the entity “*court*” is a frame element of the event “*trial*”; “*United States*” is a frame element of the event “*war*”. It is not clear which pair contributes more to salience. We again let KCE to learn it.

Formally, let  $\mathbb{E}$  be the list of entities in the document, i.e.  $\mathbb{E} = \{en_1, \dots, en_i, \dots, en_n\}$ , where  $en_i$  is the  $i$ th entity in document  $d$ . KCE extracts the kernel features about entity-event relations as follows:

$$\Phi_K(ev_i, \mathbb{E}) = \{\phi_1(\vec{ev_i}, \mathbb{E}), \dots, \phi_k(\vec{ev_i}, \mathbb{E}), \dots, \phi_K(\vec{ev_i}, \mathbb{E})\}, \quad (6.7)$$

$$\phi_k(\vec{ev_i}, \mathbb{E}) = \sum_{en_j \in \mathbb{E}} \exp \left( -\frac{(\cos(\vec{ev_i}, \vec{en_j}) - \mu_k)^2}{2\sigma_k^2} \right) \quad (6.8)$$

similarly,  $en_i$  is embedded by:  $en_i \xrightarrow{Emb} \vec{en_i}$ .

We reach the full KCE model by combining all the vectors using a linear layer:

$$\begin{aligned} f(ev_i, d) = & W_e \cdot \Phi_K(ev_i, \mathbb{E}) + W_v \cdot \Phi_K(ev_i, \mathbb{V}) \\ & + W_f \cdot F(ev_i, d) + b \end{aligned} \quad (6.9)$$

The model is again trained by equation (6.2).

## 6.6 Experimental Methodology

This section describes our experiment settings.

### 6.6.1 Event Salience Detection

**Dataset:** We conduct our experiments on the salience corpus described in §6.3. Among the 664,911 articles with abstracts, we sample 10% of the data as the test set and then randomly leave out another 10% documents for development. Overall, there are 4359 distinct event lexical items, at a similar scale with previous work [19, 46]. The corpus statistics are summarized in Table 6.1.

**Input:** The inputs to models are the documents and the extracted events. The models are required to rank the events from the most to least salience.

**Baselines:** Three methods from previous researches are used as baselines: *Frequency*, *Location* and *PageRank*. The first two are often used to simulate saliency [10, 149]. The *Frequency* baseline ranks events based on the count of the headword lemma; the *Location* baseline ranks events using the order of their appearances in discourse. Ties are broken randomly.

Similar to entity salience ranking with PageRank scores [152], our *PageRank* baseline runs PageRank on a fully connected graph whose nodes are the events in documents. The edges are weighted by the embedding similarities between event pairs. We conduct supervised PageRank on this graph, using the same pairwise loss setup as in KCE. We report the best performance obtained by linearly combining *Frequency* with the scores obtained after a one-step random walk.

**Evaluation Metric:** Since the importance of events is on a continuous scale, the boundary between “important” and “not important” is vague. Hence we evaluate it as a ranking problem. The metrics are the precision and recall value at 1, 5 and 10 respectively. It is adequate to stop at 10 since there are less than 9 salient events per document on average (Table 6.1). We also report Area Under Curve (AUC). Statistical significance values are computed by permutation (randomization) test with  $p < 0.05$ .

**Implementation Details:** We pre-trained a set of 128-dimensional embeddings on the whole Annotated New York Times corpus using Word2Vec [97]. Entities are extracted using the TagMe entity linking toolkit [57]. Words or entities that appear only once in training are replaced with special “unknown” tokens.

The hyper-parameters of the KCE kernels follow previous literature [153]. There is one exact match kernel ( $\mu = 1, \sigma = 1e^{-3}$ ) and ten soft-match kernels evenly distributed between  $(-1, 1)$ , i.e.  $\mu \in \{-0.9, -0.7, \dots, 0.9\}$ , with the same  $\sigma = 0.1$ .

The parameters of the models are optimized by Adam [78], with batch size 128. The vectors of entities are initialized by the pre-trained embeddings. Event embeddings are initialized by their headword embedding.

### 6.6.2 The Event Intrusion Test: A Study

KCE is designed to estimate salience by modeling relations between discourse units. To better understand its behavior, we design the following **event intrusion test**, following the word intrusion test used to assess topic model quality [24].

**Event Intrusion Test:** The test will present to a model a set of events, including: the **origins**, all events from one document; the **intruders**, some events from another document. Intuitively, if events inside a document are organized around the core content, a model capturing their relations well should easily identify the intruder(s).

Specifically, we take a bag of unordered events  $\{O_1, O_2, \dots, O_p\}$ , from a document  $O$ , as the origins. We insert into it intruders, events drawn from another document,  $I$ :  $\{I_1, I_2, \dots, I_q\}$ . We ask a model to rank the mixed event set  $M = \{O_1, I_1, O_2, I_2, \dots\}$ . A good model should rank the intruders  $I_i$  below the origins  $O_i$ .

**Intrusion Instances:** From the development set, we randomly sample 15,000 origin and intruding document pairs. To simplify the analysis, we only take documents with at least 5 salient events. The intruder events, together with the entities in the same sentences, are added to the origin document.

**Metrics:** AUC is used to quantify ranking quality, where events in  $O$  are positive and events in  $I$  are negative. To observe the ranking among the salient origins, we compute a separate AUC score between the intruders and the salient origins, denoted as SA-AUC. SA-AUC is the AUC score on the list with non-salient origins removed.

**Experiments Details:** We take the best-performing KCE model to compute salient scores for events in the mixed event set  $M$ , which are directly used for ranking. Frequency is recounted. All other features (Table 6.2) are set to 0 to emphasize the relational aspects,

We experiment with two settings: 1. adding only the salient intruders. 2. adding only the non-salient intruders. Under both settings, the intruders are added one by one, allowing us to observe

Method	P@01		P@05		P@10		AUC	
Location	0.3555	–	0.3077	–	0.2505	–	0.5226	–
PageRank	0.3628	–	0.3438	–	0.3007	–	0.5866	–
Frequency	0.4542	–	0.4024	–	0.3445	–	0.5732	–
LeToR	0.4753 <sup>†</sup>	+4.64%	0.4099 <sup>†</sup>	+1.87%	0.3517 <sup>†</sup>	+2.10%	0.6373 <sup>†</sup>	+11.19%
KCE (–EF)	0.4420	–2.69%	0.4038	+0.34%	0.3464 <sup>†</sup>	+0.54%	0.6089 <sup>†</sup>	+6.23%
KCE (–E)	0.4861 <sup>†‡</sup>	+7.01%	0.4227 <sup>†‡</sup>	+5.04%	0.3603 <sup>†‡</sup>	+4.58%	0.6541 <sup>†‡</sup>	+14.12%
KCE	0.5049 <sup>†‡</sup>	+11.14%	0.4277 <sup>†‡</sup>	+6.29%	0.3638 <sup>†‡</sup>	+5.61%	0.6557 <sup>†‡</sup>	+14.41%

Method	R@01		R@05		R@10		W/T/L
Location	0.0807	–	0.2671	–	0.3792	–	–/–/–
PageRank	0.0758	–	0.2760	–	0.4163	–	–/–/–
Frequency	0.0792	–	0.2846	–	0.4270	–	–/–/–
LeToR	0.0836 <sup>†</sup>	+5.61%	0.2980 <sup>†</sup>	+4.70%	0.4454 <sup>†</sup>	+4.31%	8037 / 48493 / 6770
KCE (–EF)	0.0714	–9.77%	0.2812	–1.18%	0.4321 <sup>†</sup>	+1.20%	6936 / 48811 / 7553
KCE (–E)	0.0925 <sup>†‡</sup>	+16.78%	0.3172 <sup>†‡</sup>	+11.46%	0.4672 <sup>†‡</sup>	+9.41%	11676 / 43294 / 8330
KCE	0.0946 <sup>†‡</sup>	+19.44%	0.3215 <sup>†‡</sup>	+12.96%	0.4719 <sup>†‡</sup>	+10.51%	12554 / 41461 / 9285

Table 6.3: Event salience performance. (–E) and (–F) marks removing Features and Entity information from the full KCM model. The relative performance differences are computed against Frequency. W/T/L are the number of documents a method wins, ties, and loses compared to Frequency. <sup>†</sup> and <sup>‡</sup> mark the statistically significant improvements over Frequency<sup>†</sup>, LeToR<sup>‡</sup> respectively.

the scores regarding the number of intruders added. For comparison, we add a Frequency baseline, that directly ranks events by the Frequency feature.

## 6.7 Evaluation Results

This section presents the evaluations and analyses.

### 6.7.1 Event Salience Performance

We summarize the main results in Table 6.3.

**Baselines:** Frequency is the best performing baseline. Its precision at 1 and 5 are higher than 40%. PageRank performs worse than Frequency on all the precision and recall metrics. Location performs the worst.

**Feature Based:** LeToR outperforms the baselines significantly on all metrics. Particularly, its P@1 value outperforms the Frequency baseline the most (4.64%), indicating a much better estimation on the most salient event. In terms of AUC, LeToR outperforms Frequency by a large margin (11.19% relative gain).

Feature Groups	P@1	P@5	P@10	R@1	R@5	R@10	AUC
SL	0.3548	0.3069	0.2497	0.0807	0.2671	0.3792	0.5226
Frequency	0.4536	0.4018	0.3440	0.0792	0.2846	0.4270	0.5732
+ SL	0.4734	0.4097	0.3513	0.0835	0.2976	0.4436	0.6354
+ SL + Event	0.4726	0.4101 <sup>†</sup>	0.3516	0.0831	0.2969	0.4431	0.6365 <sup>†</sup>
+ SL + Entity	0.4739	0.4100	0.3518	0.0812	0.2955	0.4418	0.6374
+ SL + Entity + Event	0.4739	0.4100	0.3518 <sup>†</sup>	0.0832	0.2974	0.4452 <sup>†</sup>	0.6374 <sup>†</sup>
+ SL + Entity + Event + Local	0.4754 <sup>†</sup>	0.4100	0.3517 <sup>†</sup>	0.0837	0.2981	0.4454 <sup>†</sup>	0.6373 <sup>†</sup>

Table 6.4: Event Saliency Feature Ablation Results. The + sign indicates adding feature groups to Frequency. SL is the sentence location feature. Event is the event voting feature. Entity is the entity voting feature. Local is the local entity voting feature. † marks the statistically significant improvements over +SL.

**Feature Ablation:** To understand the contribution of individual features, we conduct an ablation study of various feature settings in Table 6.4. We gradually add feature groups to the Frequency baseline. The combination of Location (sentence location) and Frequency almost sets the performance for the whole model. Adding each voting feature individually produces mixed results. However, adding all voting features improves all metrics. Though the margin is small, 4 of them are statistically significant over Frequency+Location.

**Kernel Centrality Estimation:** The KCE model further beats LeToR significantly on all metrics, by around 5% on AUC and precision values, and by around 10% on the recalls. Notably, the P@1 score is much higher, reaching 50%. The large relative gain on all the recall metrics and the high performance on precision show that KCE works really well on the top of the rank list.

**Kernel Ablation:** To understand the source of performance gain of KCE, we conduct an ablation study by removing its components: -E removes of entity kernels; -EF removes the entity kernels and the features. We observe a performance drop in both cases. Without entities and features, the model only using event information still performs similarly to Frequency. The drops are also a reflection of the small number of events ( $\approx 60$  per document) comparing to entities ( $\approx 200$  per document). The study indicates that the relational signals and features contain different but both important information.

**Discussion:** The superior results of KCE demonstrate its effectiveness in predicting saliency. So what additional information does it capture? We revisit the changes made by KCE: 1. it adjusts the embeddings during training. 2. it introduces weighted soft count kernels. However, the PageRank baseline also does embedding tuning but produces poor results, thus the second change should be crucial. We plot the learned kernel weights of KCE in Figure 6.2. Surprisingly, the salient decisions are not linearly related, nor even positively correlated to the weights. In fact, besides the “Exact Match” bin, the highest absolute weights actually appear at 0.3 and -0.3. This implies that embedding similarities do not directly imply saliency, breaking some assumptions of the feature based model and PageRank.



		Word2Vec	Kernel
attack	kill	0.69	0.3
arrest	charge	0.53	0.3
USA (E)	war	0.46	0.3
911 attack (E)	attack	0.72	0.3
attack	trade	0.42	0.9
hotel (E)	travel	0.49	0.9
charge	murder	0.49	0.7
business(E)	increase	0.43	0.7
attack	walk	0.44	-0.3
people (E)	work	0.40	-0.3

Table 6.5: Examples of pairs of Events/Entities in the kernels. The **Word2vec** column shows the cosine similarity using pre-trained word vectors. The **Kernel** column shows the closest kernel they belong after training. Items marked with (E) are entities.

**Case Study:** We inspect some pairs of events and entities in different kernels and list some examples in Table 6.5. The pre-trained embeddings are changed a lot. Pairs of units with different raw similarity values are now placed in the same bin. The pairs in Table 3 exhibit interesting types of relations: e.g., “*arrest-charge*” and “*attack-kill*” form script-like chains; “911 attack” forms a quasi-identity relation [127] with “attack”; “business” and “increase” are candidates as frame-argument structure. While these pairs have different raw cosine similarities, they are all useful in predicting salience. KCE learns to gather these relations into bins assigned with higher weights, which is not achieved by pure embedding based methods. This partially explains why the raw voting features and PageRank are not as effective.

## 6.7.2 Intrusion Test Results

Figure 6.3 plots results of the intrusion test . The left figure shows the results of setting 1: adding non-salient intruders. The right one shows the results of setting 2: adding salient intruders. The AUC is 0.493 and the SA-AUC is 0.753 if all intruders are added.

The left figure shows that KCE successfully finds the non-salient intruders. The SA-AUC is higher than 0.8. Yet the AUC scores, which include the rankings of non-salience events, are rather close to random. This shows that the salient events in the origin documents form a more cohesive group, making them more robust against the intruders; the non-salient ones are not as cohesive.

In both settings, KCE produces higher SA-AUC than *Frequency* at the first 30%. However, in setting 2, KCE starts to produce lower SA-AUC than *Frequency* after 30%, then gradually drops to 0.5 (random). This trend shows that salient intruders form groups connected by meaningful relations that confuse KCE. Again, this phenomenon is observed only on the salient intruders, which again confirms the cohesive relations among salient events.

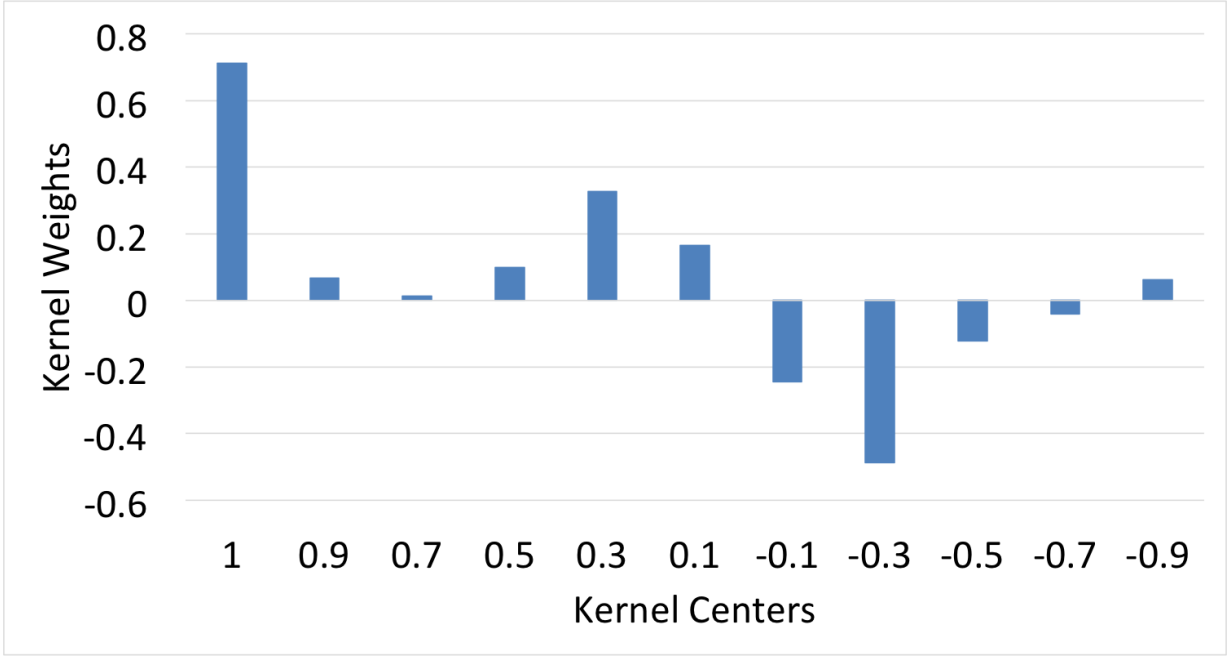


Figure 6.2: Learned Kernel Weights of KCE

In conclusion, we observe that the salient events form tight groups connected by discourse relations while the non-salient events are not as related. The observations imply that the main scripts in documents are mostly anchored by small groups of salient events (such as the “Trial” script in Example 6.1). Other events may serve as “backgrounds” [30]. Similarly, Choubey et al. [34] find that relations like event coreference and sequence are important for saliency.

## 6.8 Conclusion

We propose two salient detection models, based on lexical relatedness and semantic relations. The feature-based model with lexical similarities is effective, but cannot capture semantic relations like scripts and frames. The KCE model uses kernels and embeddings to capture these relations, thus outperforms the baselines and feature-based models significantly. All the results are tested on our newly created large-scale event salience dataset.

Our case study shows that the salience model finds and utilize a variety of discourse relations: script chain (*attack* and *kill*), frame argument relation (*business* and *increase*), quasi-identity (*911 attack* and *attack*). Such complex relations are not as prominent in the raw word embedding space.

In the intrusion test, we observe that the small number of salient events are forming tight connected groups. While KCE captures these relations quite effectively, it can be confused by salient intrusion events. The phenomenon indicates that the salient events are tightly connected, which form the main scripts of documents.

This paper empirically reveals many interesting connections between discourse phenomena

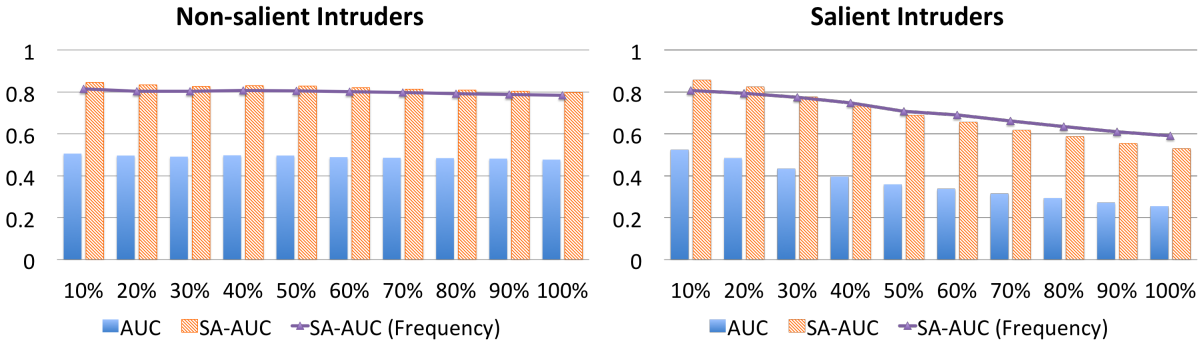


Figure 6.3: Intruder study results. X-axis shows the percentage of intruders inserted. Y-axis is the AUC score scale. The left and right figures are results from salient and non-salient intruders respectively. The blue bar is AUC. The orange shaded bar is SA-AUC. The line shows the SA-AUC of the frequency baseline.

and salience. One important implication is that, core script information may reside mostly in the salient events. In the future, we plan to study the promising direction of large-scale semantic relation discovery, for example, frames and scripts, with a focus on salient discourse units.

September 25, 2018  
DRAFT

# **Part III**

## **Proposed Work**

September 25, 2018  
DRAFT

# Chapter 7

## Proposed Theoretical Framework

### 7.1 Introduction

The information that can be convey via natural language is limited by its expressiveness. Fauconnier [54] has pointed out that language is a “superficial manifestation of hidden, highly abstract, cognitive constructions”, and that “a natural language sentence is a completely different kind of thing from a sentence in a logical calculus”. Natural language sentences can not provide precise and unambiguous information. Instead, people suggest that language only has the potential to convey meaning, and that meaning is only created after projection to the mental space. Language acts as a succinct medium for communication that triggers the process of understanding.

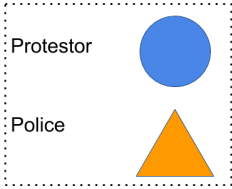
Textual mentions of events and entities are no exception. They trigger particular understanding of language users. Since these discourse elements (DEs) can have rich meaning, the information would be ambiguous without proper restriction and inference. Prior theoretical work [128] has shown that such ambiguity is apparent when modeling the coreference between DEs. In this chapter, we will discuss that the same type of ambiguity also hinders interactions between DEs.

Across this thesis, we are trying to model the structures and interactions between events and entities. There are two different level. The **mention level** is about the relations of elements within an event mention (e.g. predicate argument relations), and the **script level** considers the relations of elements across events (e.g. whether the arguments in different mentions are filling the same frame slot). It is beneficial if one can explore the mutual benefits across these levels.

#### 7.1.1 The Exact Frame Identity Assumption and Its Failures

Script modeling (or event sequencing) normally requires studying both levels. A popular approach of script modeling uses the arguments to connect up the script [16, 114, 119], in a form of  $\langle \text{event}, \text{role} \rangle$  chains<sup>1</sup>. For example, a shared criminal suspect could connect a chain of events like:  $\langle \text{search}, \text{Arg1} \rangle$ ,  $\langle \text{arrest}, \text{Arg1} \rangle$ ,  $\langle \text{plead}, \text{Arg0} \rangle$ ,  $\langle \text{convict}, \text{Arg1} \rangle$ ,  $\langle \text{sentence}, \text{Arg1} \rangle$ . Such chains can then be modelled with a language modeling-like approach. Prior work normally find them by entity coreference chains, which in fact use the following assumption.

<sup>1</sup>Some approaches also use multiple arguments instead of one, but our discussion still hold in that case



(b) The actual script model.

Figure 7.1: While the ideal script model (Left) assumes the shared frame elements to be exactly identical, actual text (Right) demonstrates more complexities.

**Assumption 7.1.1** (Exact Frame Identity Assumption). *The shared frames elements of the event mentions in a script should be identical.*

As a special case of this assumption, if a pair of event mentions are coreferent, then the fillers in their respective argument slots (e.g. the agent for both events) should be coreferent.

At the first glance, this assumption is reasonable. However, it does not always hold, especially when we use the current widely-accepted definition of coreference (we will give a formal definition in the next section). This can be demonstrated by analyzing the scripts in the following paragraph:

(7.1) The police warned the crowd to disperse, but the protesters refused to leave Taksim Square and chanted anti-government slogans. Dozens of police officers then moved toward the crowd and began spraying the protesters with water cannons. . . . said as she held red carnations in her hand. “I wanted to hand these to them, but instead they pushed me away with their shields and said our right to protest was over.”

The paragraph describes a protest script with several events and their participants. Following the Exact Frame Identity Assumption, the fillers in the shared arguments slots should be identical. One chain of shared frame elements are the “police” ( $\langle warn, Arg0 \rangle$ ,  $\langle refuse, Arg1 \rangle$ ,  $\langle spray, Arg0 \rangle$ ,  $\langle push, Arg0 \rangle$ ). Another chain is about the “protestor” ( $\langle warn, Arg1 \rangle$ ,  $\langle refuse, Arg0 \rangle$ ,  $\langle spray, Arg1 \rangle$ ,  $\langle push, Arg1 \rangle$ ). Such ideal setting is shown in Figure 7.1a.

However, the actual analysis should be described as in Figure 7.1b. While one chain is still about the “police”, it is unclear whether “Dozens of police officers” should be identical to “The Police” — we do not have enough information to determine whether the two sets contain exactly the same members. In the “protestors” chain, it is also inappropriate to treat “she”, a single protestor, to be identical to the group of protestors.

Similar violations are also observed our prior experiments, in both event coreference (Chapter 3) and script modeling (Chapter 4). These violations actually cause problems on many



perspectives. Essentially, the defect in this assumption hurts the effectiveness to model the interaction between events and entities. Corpus annotation based on it will also suffer from many inconsistencies, which makes supervised study difficult.

In the remaining of this chapter, we first analyze and define the problems of this assumption. We then introduce our proposed theoretical framework targeting at solving it. The framework is inspired by the quasi-identity theory [128]. We further discuss that this framework have implications to boarder event semantics, especially on event state modeling.

## 7.2 The Complexity about Identity

Another type of failure of the Exact Frame Identity Assumption can be observed in event coreference. Example 7.2 is taken from the TAC-KBP event coreference data, the two `arrested` event mentions are considered as coreference, although their corresponding locations are not exactly the same. The two locations are in the same scope, but described with different level of granularity.

- (7.2)
1. Man Linked by U.S. to Hezbollah is `arrested` in Brazil.
  2. The suspect was `arrested` Thursday in the city of Curitiba in southern Brazil.

Similar to Example 7.1, it also demonstrates the gap between entity coreference and frame sharing, which prevents us to conduct proper inference across DEs. In this section, we argue that the problems of assumption 7.1.1 arise from the definition of coreference: the identity of discourse elements. To be exact, the gaps are inevitable under the widely accepted notion of *exact coreference*, which is phrased by Recasens et al. [128] as:

**Definition 7.2.1** (Exact Coreference). Coreference is a relation holding between two (or more) linguistic expressions that refer to the same entity in the real world.<sup>2</sup>

Recasens et al. [128] pointed out the complexity of identity and argue that identity should be treated as a continuous value instead of a boolean decision. This is formally defined as followed:

**Definition 7.2.2** (Quasi Coreference). Coreference is a scalar relation holding between two (or more) linguistic expressions that refer to discourse entities considered to be at the same granularity level relevant to the linguistic and pragmatic context.

The Quasi Coreference definition differs Exact Coreference from several aspects. First, it considers coreference as a scalar relation. Second, it states that coreference happens in *the discourse model* instead of *the real world*: the discourse model is a projected world, a partial mental replica of the actual world built by language users (see [128, 150] for more details). Third, the coreference criteria depends on the context. Here we briefly introduce this new coreference definition with the following examples (taken from [128]):

- (7.3) Last night in Tel Aviv, *Jews* attacked a restaurant that employs Palestinians. “We want war,” *the crowd* chanted.
- (7.4) The *plant* colonized the South of France, from where *it* entered Catalonia in the 80s, spreading quickly.

<sup>2</sup>Their discussion on coreference are general and can be applied to other discourse units, such as events.

In 7.3, the mention *Jews* is a conceptual set that contain the participants of the attacks. The mention *we* is a generic set of all the members involved in this incident. The mention *the crowd* only refer to members who are “chanting”. It is not likely that we can establish a full set equality between these groups. However, Recasens et al. [128] argue that the sets themselves lose their distinctive features given the purpose of the discourse, thus quasi-coreference can be established between them. In 7.4, the mention *plant* are described twice with different locations (the ones in the South of France and the ones in Catalonia). However, quasi-coreference can be established because the purpose of the discourse is trying to emphasize that the same plant colonized and then entered and spread.

All the changes in the new definition have implications to our problem. The most prominent one here is the differences between a boolean relation and a scalar relation. Apparently, it is ambiguous to use boolean values to approximate scalar values. The shared frame elements should be better modelled as scalar relations, but are approximated as boolean values by the Exact Frame Identity Assumption. This happens in the coreference Example 7.2 and the script Example 7.1. We can summarize the problem as followed:

*Problem 7.2.1* (Frame Element Inconsistency). The same frame element may be instantiated as multiple mentions across a script. These mentions may not be coreferent under the exact coreference definition.

Based on the theory of quasi-identity, the solution is apparent, we can incorporate quasi-identity to Assumption 7.1.1 to reach Assumption 7.2.1.

**Assumption 7.2.1** (Quasi Frame Identity Assumption). *The shared frames elements in a script have non-zero quasi-coreferent value.*

In the coreference Example 7.2, the arguments can be classified as quasi-coreferent by unifying the granularity differences. Similarly, in the script Example 7.1, we can establish quasi Set.Set and Part.Whole<sup>3</sup> relations between the conflicting fillers. Readers who are interested in the details of determining quasi-identity types should refer to Recasens et al. [128].

*Frame Element Inconsistency* creates several problems. Computationally, it creates difficulties for models to exploit the mutual benefits between entity and event relations. For example, it may have reduced the effectiveness of entity coreference for within-document event coreference [33]. Script mining based on entity coreference will be sparser than expected [154]. Or that script information is not helpful to general entity coreference cases [114]. Besides these computational evidences, we will also discuss how the problems have affected corpus annotation on events and entities.

## 7.2.1 Evidence from Corpus Annotation

The frame inconsistency phenomenon can be frequently observed in datasets under two popular annotation schemes for events and entities: RichERE [87] and ACE [86]. These schemes take different approaches to handle the issue.

The RichERE annotation scheme defines a new relation named *Event Hopper*, which is a relaxed version for Exact Coreference, particularly for dealing with granularity variation of event and entity mentions [87]. The guideline makes several relaxations, such as:

<sup>3</sup>These cases belongs to different types of Meronymy quasi-identity as classified by Recasens et al. [127].

1. Trigger granularity can be different (assaulting 32 people vs. wielded a knife)
2. Event arguments may be non-coreferential or conflicting (18 killed vs. dozens killed)

The former rule relaxes the identity decisions between event predicates, which actually is a evidence for the continuum on event identity. The latter rule relaxes the identity decision between the entity pairs in the event arguments. However, given the relaxation rules, the guideline is under-specified in terms of coreference definition: *Event hopsers contain mentions of events that “feel” coreferential to the annotator even if they do not meet the earlier<sup>4</sup> strict event identity requirement.* The definition is ad-hoc and may result in some arbitrary annotation decisions.

Unlike RichERE, the ACE corpus [86] takes a strict approach to both event and entity coreference. Generally, the ACE annotation guideline discourages marking any ambiguous coreference links. Such approach reduces the ambiguity during annotation, but also lessen the connections between events and entities. Furthermore, even under the strict constraints, the same frame inconsistency problem can still occur. Hasler and Orasan [72] have done a thorough annotation study on the NP4E corpus based on the ACE ontology, following the guideline of ACE. Their study still reports that the arguments of coreferential event mentions are not coreferent in many cases. To be specific, they found that indirect anaphoric relations (which is defined similar to quasi-identity) between arguments exist in 70% of the event chains annotated. They advocate that indirect referential relations among arguments should be taken into account in building event coreference chains, which is one of the goals we are pursuing in this thesis.

### 7.3 The Facet View of Discourse Elements

We find prior work on quasi-identity inspiring. The continuum definition of identity is theoretically sound and elegant. However, the continuous scalar value termed in definition 7.2.2 is difficult to be obtained in practice. There is a need for a more **feasible computational view** than the full continuum view.

One approach along this line taken by Recasens et al. [127] is taxonomy-based. They defined 4 types and 15 subtypes for quasi-identity. For example, Example 7.3 shows a quasi-identical link with the Set.Set subtype under Meronymy. However, while the taxonomy provides guidance on which entity pairs can be considered as quasi-identity, the taxonomy type itself is inadequate to demonstrate the nature of the relation: two DEs fitting in the taxonomy can be interpreted as near identical only when both are functionally equivalent given the context. For example, not all meronymy relations are quasi-identical. We propose to seek for an approach with more **explainable power**.

In this section, we are proposing an interpretation of quasi-coreference, based on *facets*, which also provides explanations for coreference decisions. The proposed facet language not only provides a new way to describe and model quasi-coreference, but also has further implications on modeling event entity interactions. In §7.3.1 we will describe how they can help us describe the static properties of entities, with applications on identity/coreference problems. In §7.3.2 we will discuss how the new view may have interesting extensions for modelling events as dynamic functions.

<sup>4</sup>The LightERE guideline, see Song et al. [142].

### 7.3.1 Facets for Static Identity Analysis

The traditional fully identical view considers a DE as a non-dividable whole, and the continuum view of quasi-coreference considers a DE to be fully continuous. There exists, obviously, a middle ground in between these two. A DE can be viewed as a collection of discrete units, which we will refer to as *facets*<sup>5</sup> in this thesis:

**Definition 7.3.1** (Facet). The facets of an discourse element are semantic units corresponding to the possible interpretations given the textual description in a discourse model.

For example, the mention of a person can be interpreted as his/her social role or family role. The mention of a organization can be interpreted as a collection of members, or a collection of facilities, or its social function. These are all different facets described by the same mention. In natural language, the mentions of DEs serve the purpose for communication, and generally, only partial facets of the mention is to be communicated. These facets should possibly be identified given the linguistic and pragmatic context, which we call *active facets*:

**Definition 7.3.2** (Active Facet). The active facets of a mention of a DE are the facets that are considered to be relevant to the language user given the linguistic and pragmatic context.

In other words, the *active facets* are the facets that the language users are trying to communicate. Note that similar to quasi-coreference, the interpretations may subject to the language user, they happen in the discourse model [150] projected by the user instead of the real world. Figure 7.2 provides one example, the mention “paper” may have facets like “burnable”, “foldable”, but in this sentence, only some of them are relevant, such as “cuttable”.

From the facet perspective, we can now describe the identity between discourse elements: If two DEs refer to exactly the same set of facets, they are fully identical. If they are only identical on part of the facets, we enter the realm of quasi-identity and more analysis are needed. Since the active facets are the ones related to the context, a quasi-identity relation should be established when they are they same. We can rephrased the quasi-coreference definition a little using the facet language:

**Definition 7.3.3** (Facet Based Quasi Coreference). Coreference is a scalar relation holding between two (or more) linguistic expressions when the active facets of the discourse elements are considered to be the same.

We can explain the quasi-coreference examples using the facet language. In 7.4, the *active facets* of the *plant* are the organization perspective, instead of its particular facilities. This can be induced from the uses of the verbs: colonize, enter and spread, since these verbs are not compatible with physical facilities in terms of selectional preference, but are compatible to organizations. Since the active facets of *plant* and *it* are the organization, which are the same one, thus we can establish quasi-coreference between them.

The facet language can be used to model the interaction between events and entities. We can rephrase Assumption 7.2.1 as followed:

<sup>5</sup>Some work uses the word “facet” to represent the taxonomic features for DEs belonging to more than one taxonomy, we overload “facet” to represent various features for DEs, including the taxonomic ones

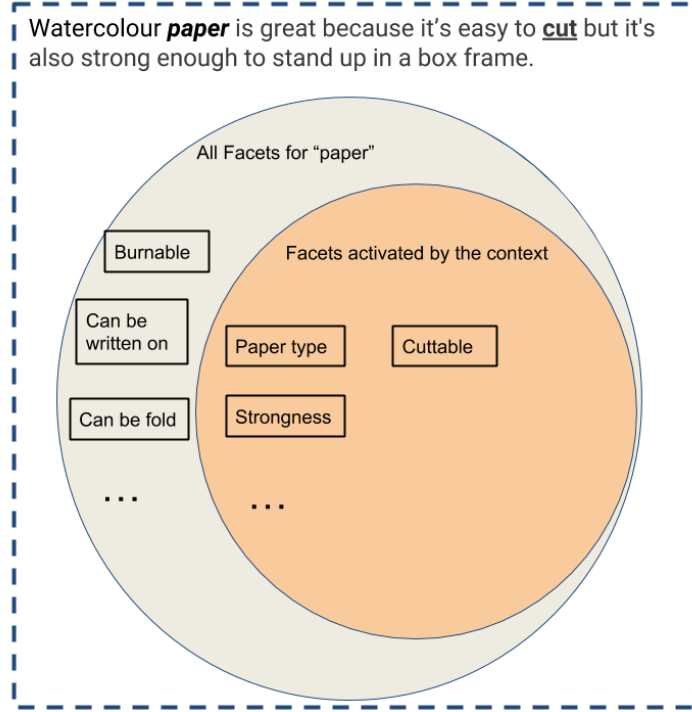


Figure 7.2: Example of active facets for an entity mention of “paper”.

**Assumption 7.3.1** (Facet based Quasi Frame Identity Assumption). *The shared frame elements in a script share active facets given the context of the script.*

In this new assumption, we have replaced “non-zero coreference value” to “sharing of active facets”, which we believe has more explainable power. For example, in 7.1, the script emphasizes the roles of the two parties. Thus the shared active facets for both set of frame elements are their role in protest (e.g. police force and protesters). In 7.2, the purpose of the script can be met without emphasizing the granularity differences between a country and a city. Thus the active facets of the two locations may not share the granularity hence can be considered identical.

Note that we emphasize “context”, but we restrict ourselves to the script context only (i.e. event related). There are other context to analyze these facets. For example, the “18 killed vs. dozens killed” example illustrated in the RichERE guideline are due to reporting variations. The active facets of the mention “dozens” can be considered compatible to those of the mention “18” because “dozens” is used as an approximate expression here. We are aware of these complex interpretations for mentions, but studying all of them are out of the scope of this thesis. Our work is mostly under the umbrella of script-based interpretations.

So far, we focus our discussions on the static properties of the entities and events, such as the role of the entities. However, one of the main aspects of events is that they are about the changes of states (such as entity properties). These are the dynamic perspective of events. In the next section, we will discuss how facets are also relevant from the dynamic perspective.

### 7.3.2 Facets and the Dynamic Interactions

A particular configuration of the entity properties in the underlying world (the real world or a discourse model) are called *states*. In the state view, entities are objects with properties, and events are functions that change the properties of these objects<sup>6</sup>. In our prior discussion, the analyses mainly concern the static part of the entities: the event context is used to help identify the prominent facets of the entity mentions. However, we have not modelled the important dynamic perspective of events, that they represent the change of states. In this section, we discuss that facets also play very important roles in modeling state change.

Similar to entity mentions, event mentions represent the events in the underlying world. We can consider events as functions that take a state as input, and produce a new state. Similarly, textual mentions of events are the triggers for understanding. In the discussion on entity mention understanding, we argue that not all facets of the entities are activated. Similarly, when analyzing the event functions, we do not need to take the whole entities as arguments. The inference procedure should be more efficient if we focus only on the facets relevant to the event functions.

Conveniently, the active facets are relevant to the context by definition. Since events are part of the context, the facets needed to model the event functions should be a subset of the active facets. The set relations can be illustrated by the Venn Graph in Figure 7.3. A few of the “paper” facets are activated by the context, for example, the “strong-ness” facet is triggered by the adjective “strong” and event mention “stand”; The “cuttable” facet is triggered by the event mention “cut”. The “cuttable” facet should be the relevant one when analyzing the “cut” event.

By considering the partial facets of event arguments, we can focus our modeling effort on state changes. In the previous example, the state changes triggered by the event mention “cut” can be reduced to applying “cut” to an “cuttable” object. This may make it easy to generalize the learned knowledge of state changes: if we can understand the states after applying “cut” to “paper”, we may be able to generalize it to other “cuttable” objects.

As the core for event understanding, state modeling is a very important process. In terms of script, state modeling enables us to trace the change of some particular facets of a frame element throughout the discourse. For example, in 7.5, understanding the state change indicated by *crashed* and *broke* allow us to infer that part of the *plane* become the *pieces*: the state of it change from a whole to a collection of parts.

(7.5) Federal investigators say a stolen Horizon Air turboprop *plane* *broke* into *many pieces* when it *crashed* into an area of thick brush on a small island in the Puget Sound.

Interestingly, the state analysis helps to establish a quasi-identity relations between *plane* and *pieces*, and the different facets between them are the ones being changed due to the event. The dynamic and static analyses now meet each other: the static analyses help highlight the relevant facets, and the dynamic analyses also help establish relations.

### 7.3.3 Facets for Events

We haven’t discussed much about the facets for events, which is a not that obvious comparing to entities. One type of multi-facet for events is that they have affected states of multiple objects. In

<sup>6</sup>Note that we also slightly overload *states* to include all kind of properties such as mental states

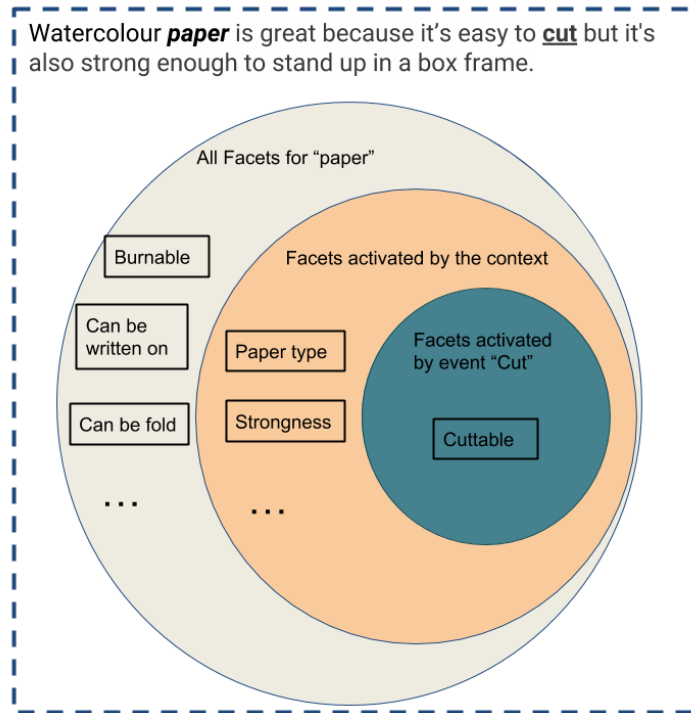


Figure 7.3: A Venn Graph showing the relations between different sets of facets.

example 7.6, we can identify two active facets from the event “pay”: the first one shows that the “payer” will give up the money, and the second one shows that the “goal” to resolve claims is to be achieve. Similarly, there are facets that are not emphasized in this sentence, such as who are the “receiver”.

(7.6) Google agreed to pay 125 million dollars to resolve outstanding claims and establish an independent ”Book Rights Registry”.

Readers who are familiar with frame semantics will realize that “payer”, “goal” and the “receiver” are the core frame elements. We consider the core frame elements to be closely related to the event facets. To be specific, we consider the facets of events are the state changes caused (or to be caused) by the event. Since each core argument may have been changed by the event, each one of them could be corresponding to one facet.

## 7.4 Relations to Prior Research

The identity problems have been studied directly or indirectly from many different perspectives. Here we introduce some related research to the best of our knowledge.

### 7.4.1 Quasi-Identity

Recasens et al. [128] first formally propose the notion of quasi-identity (a.k.a. near-identity), which fills the gap between full coreference and non-coreference. Recasens et al. [127] provides a list of typology for quasi-entity coreference, which list 15 possible subtypes of quasi-identity under 4 high-level types (i.e. Name metonymy, meronymy, class and spatio-temporal function). Relevant corpus study have been conducted guided by the theory and the typology. Our theoretical analyses are largely inspired and based on their theoretical foundation.

Similar analyses have later been applied to events by Hovy et al. [74]. Besides full coreference (exact coreference), they propose quasi-identity should be establish between event mentions when “most aspects are the same, but some additional information is provided for one or the other that is not shared”. They further identify 2 types of quasi-identity between event mentions, namely *membership* and *subevent*. A corpus study is also performed, and later computationally modelled by Araki et al. [3]. The event sequencing work in Chapter 4 also follows this line. In these prior work, the interaction between event and entity mentions are not adequately addressed, which is instead the focus of the proposed work.

### 7.4.2 Script Modeling

One research line relevant to the quasi-identity of events is on script modeling, which similar to the *subevent* relation defined in [74]. Recent work on script modeling is pioneered by [16, 19, 23]. Rich computational work has then been devoted to this area [8, 17, 20, 22, 30, 76, 77, 85, 100, 101, 111, 113, 117, 118, 119, 131, 132, 133, 138, 139, 151, 154]. The common nature of these work uses the event context to predict other events, which is quite similar to language modelling. As discussed in the previous chapters, the facet view can help fill some important gaps in script modeling, by changing the basic share frame assumption in these approaches.

### 7.4.3 Area Similar to Quasi-Identity

There are several specific types of relation between discourse elements being studied. Some of them are very similar to quasi-identity, and can be roughly classified into the quasi-identity typology. It is important to note the connections and differences of these relations and our problem.

**Bridging** is a phenomenon related to entity coreference. Example 7.7 is an excerpt taken from the OntoNotes [73] corpus. More general phenomenon are also linguistically studied in Prince [123].

(7.7) ... as much as possible of the Polish center will be made from aluminum, steel and glass recycled from Warsaw’s abundant rubble. ... The windows will open. The carpets won’t be glued down and walls will be coated with non-toxic finishes.

**Metonymy Resolution** is another problem that shares many common characteristics with quasi-identity. The problems studies xxx. A notable computation area in metonymy resolution is location metonymy[66].



#### **7.4.4 Context Dependent Modeling**

More generally, modeling the facets of entity mentions is related to some research on context dependent models. Gillick et al. [63] introduces a task named context dependent entity type detection. Dasigi et al. [44] determine the granularity of meaning of a word in context, in order to improve prepositional phrase attachment. Peters et al. [116] introduce a new type of deep contextualized word representation to incorporate context information, such as language usage and polysemy.

In the next section, we briefly propose some experiments for both aspects.

September 25, 2018  
DRAFT

# Chapter 8

## Proposed Experiments

In this chapter, we propose the experiments to validate the theoretical hypotheses discussed in Chapter 7, and to show their potential impact to some existing NLP problems.

One challenge to model complex semantic phenomenon is data sparsity. Though this is a general problem for NLP, it becomes severe for complex semantic tasks such as coreference and event sequencing: large amount of training data is required for them, but the annotation tasks are also very difficult for human at the same time. For example, the facets of entities and events are not easy to be defined from the first place, which makes manual annotation almost infeasible. In this section, we sketch our proposed solution on how to solve these problems. We try to seek indirect supervision signals, which often come in abundance and is closer to the application end. We also attempt to incorporate existing knowledge sources to simplify the problem.

### 8.1 Facets Understanding from a Static View

In this section, we propose several computational experiments to validate the first hypothesis: *only partial facets of entities and events are being activated given a discourse context*. To achieve this, we first need to identify the facets of the entity and events, and then show that by focusing on some facets, we can improve the performance of some downstream tasks.

#### 8.1.1 Facet Representation

To the best of our knowledge, there is few direct annotation for facets. Instead, we propose to use some indirect methods to represent the facets. At this step, our main goal is to learn a space for facets, where similar facets are close by.

#### Entity Event Dual Representation

In our previous chapter, we have shown that events are useful for identifying the active facets for the entities. We can collect facets about the entity via events. These facets can be also represented using the events. We call these as “event dual representation” for the facets. For example, 8.1 shows that “*paper*” is burnable (flammable) as indicated by the event “burning”. 8.2 shows another facet that “*paper*” can be written on.

(8.1) As others have pointed out, *burning paper* inside can be dangerous.

(8.2) I found that *writing on paper* helped me to quit analyzing and create faster.

Such knowledge is in the form of *predicate*-able (e.g. burnable, writable), which can be obtained via raw text documents. We propose to map the *predicate*-able facets of the entities to a low dimensional embedding space, mimicking word embeddings. This is close to the line of work on learning frame embeddings. But one major difference is that we are connecting the frame embedding knowledge to the entities, instead of tying them to the predicates. Further, we are plan to further connect these frame embeddings with some existing background knowledge.

## Incorporating Background Knowledge

Mining patterns from raw text arguments is often limited. The distribution learnt about the entities are often skewed. For example, common sense knowledge, which is quite important for reasoning, is rarely mentioned in text. Further, the proposed dual representation only covers sparse predicate related information, we further propose to incorporate other sources of knowledge to this work, with a special focus to common sense. The background knowledge can help form the facets for entities, while the event centric approach will help identify active ones. Learning both knowledge into the same space should help showing the connection between the predicate knowledge and the entity facets.

Specifically, we have investigated and found the following resources to be useful for our purpose.

**Wikipedia Sections:** Wikipedia articles contain rich information for many real world entities. For example, the entity *paper* has an entry in Wikipedia (<https://en.wikipedia.org/wiki/Paper>). Each section in the Wikipedia page describes different facets of *paper*, including: *History*, *Early sources of fibre*, *Etymology*, *Papermaking*, *Applications*, *Types*, *thickness and weight*, *Paper stability*, *Environmental impact*. The content in each section contains rich information around the particular facet. We can utilize the additional entities and events in the sections to get more information of the entity of interest.

**ConceptNet:** While Wikipedia provides detailed descriptions to entities, there are some knowledge that are more subtle to human. For example, the *paper* entry in Wikipedia does not mention that paper is flammable. Alternatively, we propose to augment the data with ConceptNet, which is a knowledge graph of common sense entities. In ConceptNet, the entry for *paper* (<http://www.conceptnet.io/c/en/paper>) contains interesting facts. For example, it shows that *paper* can be burnt, cut, written on and more. It also shows that *paper* is related to *writing*, *academic*, *book*, which covers a variety of aspects of *paper*.

**Procedural Knowledge Base:** There are also rich procedural knowledge datasets available online, such as WikiHow. Such datasets provide rich information on common sense events. This will be another useful source of information to augment event facets.

### 8.1.2 Facet Identification

This step is closely connected to the facet representation step. In fact, if we can successfully learn some embedding representations for the facets, it will be easier to identify the active facets given

the context. However, since the active facets are mostly related to the context and purpose, the validation tasks will affect the active ones. We now describe several candidate validation tasks to prove utility of the facet representations. Note that we are not trying to perform all three tasks in our thesis, but will focus on whether a task can prove our hypotheses.

### **Event Coreference**

The core reason for active facet hypothesis is actually to solve the “Frame Inconsistence” problem. Event coreference is clearly a great candidate to validate our approach. We plan to improve our event coreference system by incorporating the argument information in the facet form.

### **Entity and Quasi-Entity Coreference**

Since the facet based representation is inspired by quasi-entity coreference. It would also be useful to show that the facet representation can help establish an estimation for quasi-entity coreference. The NIDENT [130] corpus contains a set of 60 documents with quasi-coreference links. These links are assigned with an “identity degree” value from 1 to 3 (with 3 being exact coreference). We plan to investigate how well can our facet based representation approximate these annotations.

### **Implicit Semantic Role Labeling (SRL) Selection**

Directly modeling quasi-entity coreference is informative. However, the NIDENT dataset is small and may be difficult to use even as a test set. We find implicit SRL to be another useful validation task. The implicit SRL task studies the following problem: the argument for a predicate may not appear in the sentence, but it can be found in the same document. The task thus requires a system to find such arguments. Unlike traditional SRL, a system has no access to the syntax information under the setting of implicit SRL. The task can be considered as a “Null coreference problem”, where the missing argument in a sentence can be considered as a “Null argument filler”, and the task is to find in the document which other argument phrase can be coreferent with it.

There are several dataset available for this task, including the implicit NomBank [62] and the SemEval 2010 SRL task [134]. The identity between a Null argument and the filler argument is closer to our definition of quasi-identity: they are only required to be functionally identical.

## **8.2 Understanding States and Facets from a Dynamic View**

We now present our proposed experiments to study the state change perspective for events.

### **8.2.1 State Modeling**

The core proposed work is to capture the state change of entities after a series of events described by the text document.

To show the effectiveness of state modeling, we attempt to work on one of the following validation task:

## **Script Induction**

Script Induction (a.k.a. Event Schema Induction) is a task to mine prototypical event chain patterns from raw text. We consider

## **State Aware Selectional Preference**

Unfortunately, there is no dataset annotated in this form. However, it is convenient to use a cloze style task to validate such argument filling problem [29].

## **Chapter 9**

### **Timeline**

September 25, 2018  
DRAFT



# Bibliography

- [1] David Ahn. 2006. The stages of event extraction. *ARTE '06 Proceedings of the Workshop on Annotating and Reasoning about Time and Events*, pages 1–8.
- [2] Jun Araki. 2018. *Extraction of Event Structures from Text*. Ph.D. thesis, Carnegie Mellon University.
- [3] Jun Araki, Zhengzhong Liu, Eduard Hovy, and Teruko Mitamura. 2014. Detecting Subevent Structure for Event Coreference Resolution. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 4553–4558, Reykjavik, Iceland. European Language Resources Association (ELRA).
- [4] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. DBpedia: A Nucleus for a Web of Open Data. In *ISWC'07/ASWC'07 Proceedings of the 6th international The semantic web and 2nd Asian conference on Asian semantic web conference*, pages 722–735.
- [5] A. Bagga and B. Baldwin. 1998. Entity-based cross-document coreferencing using the vector space model. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics-Volume 1*, pages 79–85. Association for Computational Linguistics.
- [6] Amit Bagga and Bagga Baldwin. 1999. Cross - Document Event Coreference: Annotations, Experiments, and Observations. In *ACL-99 Workshop on Coreference and Its Applications*.
- [7] CF Baker, CJ Fillmore, and JB Lowe. 1998. The berkeley framenet project. *Proceeding ACL '98 Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics*, pages 86–90.
- [8] Niranjan Balasubramanian, Stephen Soderland, and OE Mausam. 2013. Generating Coherent Event Schemas at Scale. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*.
- [9] Breck Baldwin and Thomas S Morton. 1998. Dynamic Coreference-Based Summarization. *Proceedings of the Third Conference on Empirical Methods in Natural Language Processing (EMNLP-3)*, pages 1–6.
- [10] Regina Barzilay and Mirella Lapata. 2008. Modeling Local Coherence: An Entity-Based Approach. In *Proceedings of the 43rd Annual Meeting of the ACL*, pages 1–34.
- [11] Cosmin Adrian Bejan and Sanda Harabagiu. 2010. Unsupervised event coreference resolution with rich linguistic features. In *Proceedings of the 48th Annual Meeting of the*

- Association for the Association for Computational Linguistics*, July, pages 1412–1422.
- [12] Anders Björkelund and Jonas Kuhn. 2014. Learning Structured Perceptrons for Coreference Resolution with Latent Antecedents and Non-local Features. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, pages 47–57.
  - [13] Johan Bos and Jennifer Spenader. 2011. An annotated corpus for the analysis of VP ellipsis. *Language Resources and Evaluation*, 45(4):463–494.
  - [14] L Breiman. 2001. Random forests. *Machine learning*, pages 5–32.
  - [15] Susan Windisch Brown, Claire Bonial, Leo Obrst, and Martha Palmer. 2017. The Rich Event Ontology. In *Proceedings of the Events and Stories in the News Workshop*, pages 87–97.
  - [16] Nathanael Chambers. 2011. *Inducing Event Schemas and Their Participants from Unlabeled Text*. Doctoral, Stanford University.
  - [17] Nathanael Chambers. 2013. Event Schema Induction with a Probabilistic Entity-Driven Model. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807.
  - [18] 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.
  - [19] Nathanael Chambers and Dan Jurafsky. 2008. Unsupervised Learning of Narrative Event Chains. In *ACL ’08 Meeting of the Association for Computational Linguistics*, pages 789–797.
  - [20] Nathanael Chambers and Dan Jurafsky. 2009. Unsupervised learning of narrative schemas and their participants. In *ACL ’09 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*, pages 602–610. Association for Computational Linguistics.
  - [21] Nathanael Chambers and Dan Jurafsky. 2010. A Database of Narrative Schemas. In *LREC 2010, Seventh International Conference on Language Resources and Evaluation*.
  - [22] Nathanael Chambers and Dan Jurafsky. 2011. Template-Based Information Extraction without the Templates. In *HLT ’11 Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*, pages 976–986.
  - [23] Nathanael Chambers and Dan Jurafsky. 2012. Learning the Central Events and Participants in Unlabeled Text. In *International Conference on Machine Learning (ICML 2012)*.
  - [24] Jonathan Chang, Sean Gerrish, Chong Wang, and David M Blei. 2009. Reading Tea Leaves: How Humans Interpret Topic Models. In *Advances in Neural Information Processing Systems 22*, pages 288—296.
  - [25] Chen Chen and Vincent Ng. 2013. Chinese Event Coreference Resolution: Understanding the State of the Art. In *Proceedings of the 6th International Joint Conference on Natural Language Processing*, October, pages 822–828.

- [26] Chen Chen and Vincent Ng. 2015. Chinese Event Coreference Resolution : An Un-supervised Probabilistic Model Rivaling Supervised Resolvers. In *Human Language Technologies: The 2015 Annual Conference of the North American Chapter of the ACL*, pages 1097–1107.
- [27] Zheng Chen and H Ji. 2009. Graph-based event coreference resolution. *Proceedings of the 2009 Workshop on Graph-based Methods for Natural Language Processing*, (August):54–57.
- [28] Zheng Chen, H Ji, and Robert Haralick. 2009. A pairwise event coreference model, feature impact and evaluation for event coreference resolution. In *Proceedings of the Workshop on Events in Emerging Text Types*, 3, pages 17–22.
- [29] Pengxiang Cheng and Katrin Erk. 2018. Implicit Argument Prediction with Event Knowledge. In *NAACL 2018*, 2012.
- [30] JC Kit Cheung, H Poon, and Lucy Vanderwende. 2013. Probabilistic Frame Induction. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL/HLT 2013)*.
- [31] Nancy Chinchor. 1992. MUC-5 EVALUATION METRIC. In *Proceedings of the 5th Conference on Message Understanding*, pages 69–78.
- [32] Timothy Chklovski and Patrick Pantel. 2004. VERB OCEAN : Mining the Web for Fine-Grained Semantic Verb Relations. In *In Proceedings of Conference on Empirical Methods in Natural Language Processing (EMNLP-04)*, pages 33–40.
- [33] Prafulla Kumar Choubey and Ruihong Huang. 2017. Event Coreference Resolution by Iteratively Unfolding Inter-dependencies among Events. In *EMNLP 2017*.
- [34] Prafulla Kumar Choubey, Kaushik Raju, and Ruihong Huang. 2018. Identifying the Most Dominant Event in a News Article by Mining Event Coreference Relations. In *NAACL 2018*.
- [35] Michael Collins. 2002. Discriminative Training Methods for Hidden Markov Models: Theory and Experiments with Perceptron Algorithms. In *Proceedings of the Conference on Empirical Methods in NLP (EMNLP 2002)*, July, pages 1–8.
- [36] Ronan Collobert, Jason Weston, Leon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural Language Processing (almost) from Scratch. *Journal of Machine Learning Research*.
- [37] Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. 2006. Online Passive-Aggressive Algorithms. *Journal of Machine Learning Research*, 7:551–585.
- [38] Agata Cybulska and Piek Vossen. 2012. Using Semantic Relations to Solve Event Coreference in Text. In *SemRel2012 in conjunction with LREC2012*, pages 60–67.
- [39] Agata Cybulska and Piek Vossen. 2014. Using a sledgehammer to crack a nut? Lexical diversity and event coreference resolution. *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 4545–4552.
- [40] Mary Dalrymple, Stuart M. Shieber, and Fernando C N Pereira. 1991. Ellipsis and higher-

order unification. *Linguistics and Philosophy*, 14(4):399–452.

- [41] Laurence Danlos. 2003. Event coreference between two sentences. *Computing Meaning*, 77:271–288.
- [42] D Das, Nathan Schneider, Desai Chen, and NA Smith. 2010. Probabilistic frame-semantic parsing. In *In Proceedings of the North American Chapter of the Association for Computational Linguistics Human Language Technologies Conference (NAACL 2010)*, volume 3, Los Angeles.
- [43] Dipanjan Das and NA Smith. 2011. Semi-Supervised Frame-Semantic Parsing for Unknown Predicates. In *HLT ’11 Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*, volume 1, pages 1435–1444.
- [44] Pradeep Dasigi, Waleed Ammar, Chris Dyer, and Eduard Hovy. 2017. Ontology-Aware Token Embeddings for Prepositional Phrase Attachment. In *ACL 2017*.
- [45] Ofer Dekel, C Manning, and Yoram Singer. 2004. Log-linear models for label ranking. *Advances in neural information processing systems*, 16:497–504.
- [46] Quang Xuan Do, Yee Seng Chan, and Dan Roth. 2011. Minimally supervised event causality identification. In *EMNLP ’11 Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 294–303.
- [47] Quang Xuan Do, Wei Lu, and Dan Roth. 2012. Joint Inference for Event Timeline Construction. *EMNLP-CoNLL ’12 Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, (July):677–687.
- [48] M. Dojchinovski, D. Reddy, T. Kliegr, T. Vitvar, and H. Sack. 2016. Crowdsourced corpus with entity salience annotations. In *Proceedings of the 10th International Conference on Language Resources and Evaluation, LREC*, pages 3307–3311.
- [49] Jesse Dunietz and Daniel Gillick. 2014. A New Entity Salience Task with Millions of Training Examples. In *Proceedings of the European Association for Computational Linguistics*, pages 205–209.
- [50] Greg Durrett and Dan Klein. 2013. Easy victories and uphill battles in coreference resolution. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1971–1982.
- [51] Greg Durrett and Dan Klein. 2014. A Joint Model for Entity Analysis: Coreference, Typing, and Linking. *Proceedings of the Transactions of the Association for Computational Linguistics*.
- [52] Aymen Elkhilfi and Rim Faiz. 2009. Automatic annotation approach of events in news articles. *International Journal of Computing & Information Sciences*, 7(1):40–50.
- [53] Günes Erkan and Dragomir R Radev. 2004. LexRank : Graph-based Lexical Centrality as Salience in Text Summarization. *Journal of Artificial Intelligence Research*, 22:457–479.
- [54] Gilles Fauconnier. 1997. *Mappings In Thought and Language*. Cambridge University Press.

- [55] Christiane Fellbaum. 1998. WordNet.
- [56] Eraldo Rezende Fernandes, Cícero Nogueira dos Santos, and Ruy Luiz Milidiú. 2012. Latent structure perceptron with feature induction for unrestricted coreference resolution. *Joint Conference on {EMNLP} and {CoNLL-Shared} Task*, pages 41–48.
- [57] Paolo Ferragina and Ugo Scaiella. 2010. Tagme: on-the-fly annotation of short text fragments (by wikipedia entities). In *CIKM 2010*.
- [58] Francis Ferraro and Benjamin Van Durme. 2016. A Unified Bayesian Model of Scripts , Frames and Language. *Proceedings of the 30th Conference on Artificial Intelligence (AAAI 2016)*, pages 2601–2607.
- [59] JR Finkel, Trond Grenager, and Christopher Manning. 2005. Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL 2005)*.
- [60] Joseph Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76.
- [61] Radu Florian, John F Pitrelli, Salim Roukos, and Imed Zitouni. 2010. Improving mention detection robustness to noisy input. *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, (October):335–345.
- [62] M Gerber and J Y Chai. 2010. Beyond NomBank: A Study of Implicit Arguments for Nominal Predicates. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, July, pages 1583–1592.
- [63] Dan Gillick, Nevena Lazic, Kuzman Ganchev, Jesse Kirchner, and David Huynh. 2014. Context-Dependent Fine-Grained Entity Type Tagging. Technical report.
- [64] Kartik Goyal, Sujay Kumar Jauhar, Mrinmaya Sachan, Shashank Srivastava, Huiying Li, and Eduard Hovy. 2013. A Structured Distributional Semantic Model for Event Co-reference. In *Proceedings of The 51st Annual Meeting of the Association for Computational Linguistics*, Bulgaria.
- [65] Joseph Evans Grimes. 1975. *The Thread of Discourse*. New York.
- [66] Milan Gritta, Mohammad Taher Pilehvar, Nut Limsopatham, and Nigel Collier. 2017. Vancouver Welcomes You! Minimalist Location Metonymy Resolution. *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1248–1259.
- [67] Mark Hall, Geoffrey Holmes, Bernhard Pfahringer, Ian Witten, Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H Witten. 2009. The WEKA data mining software : An update The WEKA Data Mining Software : An Update. *ACM SIGKDD Explorations Newsletter*, 11(1):10–18.
- [68] Michael Halliday and Ruqaiya Hasan. 1976. *Cohesion in English*.
- [69] Daniel Hardt. 1992. An algorithm for VP ellipsis. In *Proceedings of the 30th annual meeting on Association for Computational Linguistics*, January, pages 9–14.
- [70] Daniel Hardt. 1997. An empirical approach to VP ellipsis. *Computational Linguistics*,

23(4):525–541.

- [71] Daniel Hardt. 1998. Improving Ellipsis Resolution with Transformation-Based Learning. *AAAI Fall Symposium*, pages 41–43.
- [72] Laura Hasler and Constantin Orasan. 2009. Do coreferential arguments make event mentions coreferential? In *Proceedings of DAARC*, pages 151–163.
- [73] Eduard Hovy, Mitchell Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. 2006. OntoNotes: the 90% solution. In *Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, June, pages 57–60.
- [74] Eduard Hovy, T Mitamura, F Verdejo, J Araki, and Andrew Philpot. 2013. Events are Not Simple: Identity, Non-Identity, and Quasi-Identity. In *The 1st Workshop on EVENTS: Definition, Detection, Coreference and Representation, NAACL-HLT 2013 Workshop*, pages 21–28, Atlanta.
- [75] Kevin Humphreys, Robert Gaizauskas, and Saliha Azzam. 1997. Event coreference for information extraction. In *Proceedings of a Workshop on Operational Factors in Practical, Robust Anaphora Resolution for Unrestricted Texts, 35 th Annual Meeting of Assoc. for Computational Linguistics*, pages 75–81, Madrid.
- [76] Bram Jans, Steven Bethard, I Vulić, and MF Moens. 2012. Skip n-grams and ranking functions for predicting script events. *EACL ’12 Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 336–344.
- [77] Niels Kasch and Tim Oates. 2010. Mining script-like structures from the web. In *Proceedings of the NAACL HLT 2010 First International Workshop on Formalisms and Methodology for Learning by Reading*, pages 34–42.
- [78] Diederik P. Kingma and Jimmy Lei Ba. 2015. Adam: A Method for Stochastic Optimization. In *ICLR 2015*.
- [79] Heeyoung Lee, Yves Peirsman, and Angel Chang. 2011. Stanford’s Multi-Pass Sieve Coreference Resolution System at the CoNLL-2011 Shared Task. In *Proc. of the 15th Conference on Computational Natural Language Learning: Shared Task*, June, pages 28–34.
- [80] Heeyoung Lee, Marta Recasens, Angel Chang, Mihai Surdeanu, and Dan Jurafsky. 2012. Joint Entity and Event Coreference Resolution across Documents. In *Proceedings of the 2012 Conference on Empirical Methods in Natural Language Processing and Natural Language Learning*.
- [81] Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. End-to-end Neural Coreference Resolution. In *EMNLP 2017*.
- [82] Peifeng Li, Guodong Zhou, Qiaoming Zhu, and Libin Hou. 2012. Employing compositional semantics and discourse consistency in Chinese event extraction. In *EMNLP-CoNLL ’12 Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, July, pages 1006–1016.
- [83] Peifeng Li, Qiaoming Zhu, and Xiaoxu Zhu. 2011. A Clustering and Ranking Based Approach for Multi-document Event Fusion. *12th ACIS International Conference on Software*

*Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing*, pages 159–165.

- [84] Qi Li, Heng Ji, and Liang Huang. 2013. Joint Event Extraction via Structured Prediction with Global Features. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (ACL2013)*.
- [85] Zhongyang Li, Xiao Ding, and Ting Liu. 2018. Constructing Narrative Event Evolutionary Graph for Script Event Prediction. In *IJCAI 2018*.
- [86] Linguistic Data Consortium. 2008. ACE ( Automatic Content Extraction ) English Annotation Guidelines for Entities, Version 6.6 2008.06.13. Technical report.
- [87] Linguistic Data Consortium. 2015. DEFT Rich ERE Annotation Guidelines: Events v.2.6. Technical report.
- [88] Tie-Yan Liu. 2009. Learning to Rank for Information Retrieval. *Foundations and Trends in Information Retrieval*, 3(3):225–331.
- [89] Zhengzhong Liu, Jun Araki, Eduard Hovy, and Teruko Mitamura. 2014. Supervised Within-Document Event Coreference using Information Propagation. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 4539–4544, Reykjavik, Iceland. European Language Resources Association (ELRA).
- [90] Jing Lu and Vincent Ng. 2017. Joint Learning for Event Coreference Resolution. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 90–101.
- [91] Jing Lu, Deepak Venugopal, Vibhav Gogate, and Vincent Ng. 2016. Joint Inference for Event Reference Resolution. In *Proceedings of the 26th International Conference on Computational Linguistics*.
- [92] Xiaoqiang Luo. 2005. On coreference resolution performance metrics. *HLT ’05: Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing*, (October):25–32.
- [93] Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J Bethard, and David McClosky. 2014. The Stanford CoreNLP Natural Language Processing Toolkit. *Proceedings of 52nd Annual Meeting of the ACL: System Demonstrations*, pages 55–60.
- [94] Daniel Marcu. 1999. Discourse Trees are Good Indicators of Importance in Text. *Advances in Automatic Text Summarization*, pages 123–136.
- [95] Katie McConky, Rakesh Nagi, Moises Sudit, and William Hughes. 2012. Improving event co-reference by context extraction and dynamic feature weighting. *2012 IEEE International Multi-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support*, pages 38–43.
- [96] PN Mendes, Max Jakob, Andres Garcia-silva, and Christian Bizer. 2011. DBpedia spotlight: shedding light on the web of documents. *Proceedings of the 7th Conference on Semantic Systems*, pages 1–8.
- [97] Tomas Mikolov, I Sutskever, and Kai Chen. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*

26, pages 3111–3119.

- [98] Teruko Mitamura, Zhengzhong Liu, and Eduard Hovy. 2015. Overview of TAC KBP 2015 Event Nugget Track. In *TAC KBP 2015*, pages 1–31.
- [99] Teruko Mitamura, Zhengzhong Liu, and Eduard Hovy. 2018. Events Detection, Coreference and Sequencing: What’s next? Overview of the TAC KBP 2017 Event Track. In *TAC 2017*, pages 1–42.
- [100] A Modi and Ivan Titov. 2014. Inducing Neural Models of Script Knowledge. In *Eighteenth Conference on Computational Natural Language Learning (CoNLL 2014)*.
- [101] Ashutosh Modi and Ivan Titov. 2013. Learning Semantic Script Knowledge with Event Embeddings.
- [102] Jane Morris and Graeme Hirst. 1991. Lexical cohesion computed by thesaural relations as an indicator of the structure of text. *Computational Linguistics*, 17:21–48.
- [103] Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A Corpus and Evaluation Framework for Deeper Understanding of Commonsense Stories. In *Proceedings of NAACL-HLT 2016*, pages 839–849.
- [104] Martina Naughton. 2009. *Sentence Level Event Detection and Coreference Resolution*. Ph.D. thesis, National University of Ireland, Dublin.
- [105] Vincent Ng and Claire Cardie. 2002. Improving machine learning approaches to coreference resolution. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, July, pages 104–111, Philadelphia.
- [106] Thien Huu Nguyen and Ralph Grishman. 2015. Event Detection and Domain Adaptation with Convolutional Neural Networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 365–371.
- [107] Leif Arda Nielsen. 2003. Using Machine Learning techniques for VPE detection. In *Proceedings of Recent Advances in Natural Language Processing*, pages 339–346.
- [108] Leif Arda Nielsen. 2004. Robust VPE detection using automatically parsed text. In *Proceedings of the Student Workshop, ACL*, pages 31–36.
- [109] Leif Arda Nielsen. 2004. Using Automatically Parsed Text for Robust VPE detection. In *Proceedings of the fifth Discourse Anaphor and Anaphora Resolution conference (DAARC)*.
- [110] Leif Arda Nielsen. 2005. *A corpus-based study of Verb Phrase Ellipsis Identification and Resolution*. Doctor of philosophy, King’s College London.
- [111] Simon Ostermann, Michael Roth, Stefan Thater, and Manfred Pinkal. 2017. Aligning Script Events with Narrative Texts. In *Conference: Proceedings of the 6th Joint Conference on Lexical and Computational Semantics (\*SEM 2017)*.
- [112] Ted Pedersen and Jason Michelizzi. 2004. WordNet :: Similarity - Measuring the Relatedness of Concepts. *HLT-NAACL-Demonstrations ’04 Demonstration Papers at HLT-NAACL 2004*, (July):38–41.



- [113] Haoruo Peng, Snigdha Chaturvedi, and Dan Roth. 2017. A Joint Model for Semantic Sequences : Frames, Entities, Sentiments. *CoNLL 2017*, (CoNLL):173–181.
- [114] Haoruo Peng, Daniel Khashabi, and Dan Roth. 2015. Solving Hard Coreference Problems. In *Human Language Technologies: The 2015 Annual Conference of the North American Chapter of the ACL*, pages 809–819, Denver, Colorado.
- [115] Haoruo Peng, Yangqi Song, and Dan Roth. 2016. Event Detection and Co-reference with Minimal Supervision. In *EMNLP 2016*.
- [116] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *NAACL 2018*.
- [117] Karl Pichotta. 2016. Learning Statistical Scripts With LSTM Recurrent Neural Networks. In *In Proceedings of the 30th AAAI Conference on Artificial Intelligence*, February.
- [118] Karl Pichotta and Raymond J. Mooney. 2016. Using Sentence-Level LSTM Language Models for Script Inference. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 279–289.
- [119] Karl Pichotta and RJ Mooney. 2013. Statistical Script Learning with Multi-Argument Events. In *In Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2014)*, 2012, pages 220–229, Gothenburg, Sweden.
- [120] Marco Ponza, Paolo Ferragina, and Francesco Piccinno. 2018. SWAT: A System for Detecting Salient Wikipedia Entities in Texts.
- [121] Sameer Pradhan, Alessandro Moschitti, and Olga Uryupina. 2012. CoNLL-2012 Shared Task : Modeling Multilingual Unrestricted Coreference in OntoNotes. In *Conll*, June, pages 1–40.
- [122] Sameer S. Pradhan, Lance Ramshaw, Ralph Weischedel, Jessica MacBride, and Linnea Micciulla. 2007. Unrestricted Coreference: Identifying Entities and Events in OntoNotes. *International Conference on Semantic Computing (ICSC 2007)*, pages 446–453.
- [123] Ellen F Prince. 1981. Toward a taxonomy of given-new information. *Radical pragmatics*, 3(1979):223–255.
- [124] 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.
- [125] Altaf Rahman and Vincent Ng. 2009. Supervised models for coreference resolution. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, August, pages 968–977.
- [126] Marta Recasens and Eduard Hovy. 2011. Blanc: Implementing the rand index for coreference evaluation. *Natural Language Engineering*, 1(1).
- [127] Marta Recasens, Eduard Hovy, and M. Antònia Martí. 2010. A Typology of Near-Identity Relations for Coreference (NIDENT). *7th International Conference on Language Resources and Evaluation (LREC-2010)*, (i):149–156.

- [128] Marta Recasens, Eduard Hovy, and M. Antonia Marti. 2011. Identity, non-identity, and near-identity: Addressing the complexity of coreference. *Lingua*, 121(6):1138–1152.
- [129] Marta Recasens, Mc De Marneffe, and Christopher Potts. 2013. The Life and Death of Discourse Entities: Identifying Singleton Mentions. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, June, pages 627–633.
- [130] Marta Recasens, M Antònia Martí, and Constantin Orasan. 2012. Annotating Near-Identity from Coreference Disagreements. In *Proceedings of The Eight International Conference on Language Resources and Evaluation (LREC 2012)*, pages 165–172.
- [131] Michaela Regneri, A Koller, and M Pinkal. 2010. Learning script knowledge with web experiments. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 979–988.
- [132] Michaela Regneri, Alexander Koller, J Ruppenhofer, and Manfred Pinkal. 2011. Learning Script Participants from Unlabeled Data. In *Proceedings of Recent Advances in Natural Language Processing*, pages 463–470.
- [133] Rachel Rudinger, Pushpendre Rastogi, Francis Ferraro, and Benjamin Van Durme. 2015. Script Induction as Language Modeling. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1681–1686.
- [134] Joseph Ruppenhofer, Caroline Sporleder, Roser Morante, Collin F Baker, and Martha Palmer. 2010. SemEval-2010 Task 10: Linking Events and Their Participants in Discourse. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, July, pages 45–50.
- [135] Evan Sandhaus. 2008. The New York Times Annotated Corpus.
- [136] S Sangeetha and Michael Arock. 2012. Event Coreference Resolution using Mincut based Graph Clustering. *International Journal of Computing & Information Sciences*, pages 253–260.
- [137] Roger C Schank and Robert P Abelson. 1977. *Scripts, Plans, Goals and Understanding*. Lawrence Erlbaum Associates.
- [138] Lei Sha, Sujian Li, Baobao Chang, Zhifang Sui, and Computer Science. 2008. Jointly Learning Templates and Slots for Event Schema Induction.
- [139] Tomohide Shibata and Sadao Kurohashi. 2011. Acquiring Strongly-related Events using Predicate-argument Co-occurring Statistics and Case Frames. In *Proceedings of the 5th International Joint Conference on Natural Language Processing*, pages 1028–1036.
- [140] Stuart M. Shieber, Fernando C. N. Pereira, and Mary Dalrymple. 1996. Interactions of scope and ellipsis. *Linguistics and Philosophy*, 19(5):527–552.
- [141] E Skorochood’ko. 1971. Adaptive Method of Automatic Abstracting and Indexing. In *Proceedings of the IFIP Congress 71*.
- [142] Zhiyi Song, Ann Bies, Stephanie Strassel, Tom Riese, Justin Mott, Joe Ellis, Jonathan Wright, Seth Kulick, Neville Ryant, and Xiaoyi Ma. 2015. From light to rich ERE: annotation of entities, relations, and events. In *Proceedings of the 3rd Workshop on*

*EVENTS at the NAACL-HLT*, pages 89–98.

- [143] A Sun, R Grishman, and S Sekine. 2011. Semi-supervised relation extraction with large-scale word clustering. *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics*, pages 521–529.
- [144] Stephen Tratz and Eduard Hovy. 2011. A fast, accurate, non-projective, semantically-enriched parser. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2010, pages 1257–1268.
- [145] Naushad UzZaman. 2012. *Interpreting the Temporal Aspects of Language*. Ph.D. thesis, University of Rochester.
- [146] Naushad Uzzaman and James F Allen. 2010. TRIPS and TRIOS System for TempEval-2: Extracting Temporal Information from Text. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, July, pages 276–283.
- [147] Naushad UzZaman, Hector Llorens, Leon Derczynski, Marc Verhagen, James Allen, and James Pustejovsky. 2013. Semeval-2013 task 1: {T}empeval-3: Evaluating time expressions, events, and temporal relations. In *Second joint conference on lexical and computational semantics (\* SEM)*, volume 2, pages 1–9.
- [148] Zeno Vendler. 1957. Verbs and times. *The Philosophical Review*, 66(2):143–160.
- [149] Piek Vossen and Tommaso Caselli. 2015. Storylines for structuring massive streams of news. In *Proceedings of the First Workshop on Computing News Story Lines*, pages 40–49.
- [150] Bonnie Lynn Webber. 1978. *A Formal Approach to Discourse Anaphora*. Garland Press, New York.
- [151] Noah Weber, Leena Shekhar, Niranjan Balasubramanian, and Nathanael Chambers. 2018. Hierarchical Quantized Representations for Script Generation. In *EMNLP 2018*.
- [152] Chenyan Xiong, Zhengzhong Liu, Jamie Callan, and Tie-Yan Liu. 2018. Towards Better Text Understanding and Retrieval through Kernel Entity Saliency Modeling. In *SIGIR 2018*.
- [153] Chenyang Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan Liu, and Russell Power. 2017. End-to-End Neural Ad-hoc Ranking with Kernel Pooling. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 55–64.
- [154] Wenlin Yao and Ruihong Huang. 2018. Temporal Event Knowledge Acquisition via Identifying Narratives. In *ACL 2018*.
- [155] Congle Zhang, Stephen Soderland, and Daniel S Weld. 2015. Exploiting Parallel News Streams for Unsupervised Event Extraction. volume 3, pages 117–129.