

# In-depth Exploitation of Noun and Verb Semantics to Identify Causation in Verb-Noun Pairs

Mehwish Riaz and Roxana Girju

Department of Computer Science and Beckman Institute

University of Illinois at Urbana-Champaign

Urbana, IL 61801, USA

{mriaz2, girju}@illinois.edu

## Abstract

Recognition of causality is important to achieve natural language discourse understanding. Previous approaches rely on shallow linguistic features. In this work, we propose to identify causality in verb-noun pairs by exploiting deeper semantics of nouns and verbs. Particularly, we acquire and employ three novel types of knowledge: (1) semantic classes of nouns with a high and low tendency to encode causality along with information regarding metonymies, (2) data-driven semantic classes of verbal events with the least tendency to encode causality, and (3) tendencies of verb frames to encode causality. Using these knowledge sources, we achieve around 15% improvement in F-score over a supervised classifier trained using linguistic features.

## 1 Introduction

The identification of cause-effect relations is critical to achieve natural language discourse understanding. Causal relations are encoded in text using various linguistic constructions e.g., between two verbs, a verb and a noun, two discourse segments, etc. In this research, we focus on identifying causality encoded between a verb and a noun (or noun\_phrase). For example, consider the following example:

1. At least 1,833 people **died** in **the hurricane**.

In example (1), the verb-noun\_phrase pair “died”-“the hurricane” encodes causality where event “died” is the effect of “hurricane” event.

Previously several approaches have been proposed to identify causality between two verbs (Bethard and Martin, 2008; Riaz and Girju, 2010; Do et al., 2011; Riaz and Girju, 2013) and discourse segments (Sporleder and Lascarides, 2008;

Pitler and Nenkova, 2009; Pitler et al., 2009). However, the problem of identifying causality in verb-noun pairs has not received a considerable attention. For example, Do et al. (2011) have studied this task but they worked only with a list of predefined nouns representing events. In this work, we focus on the linguistic construction of verb-noun (or noun\_phrase) pairs where noun can be of any semantic type.

Traditional approaches for identifying causality mainly employ linguistic features (e.g., lexical items, part-of-speech tags of words, etc.) in the framework of supervised learning (Girju, 2003; Sporleder and Lascarides, 2008; Bethard and Martin, 2008; Pitler and Nenkova, 2009; Pitler et al., 2009) and do not involve deeper semantics of language. Analysis of such approaches by Sporleder and Lascarides (2008) have revealed that the linguistic features are not always sufficient to achieve a good performance on the task of identifying semantic relations including causality. In this work, we propose a model that deeply processes and acquires the specific semantic information about the participants of a verb-noun\_phrase (v-np) pair (i.e., noun and verb semantics) to identify causality with a better performance over the baseline model depending merely on shallow linguistic features.

The work in this paper builds on our recent work reported in Riaz and Girju (2014). In that previous model, we identified the semantic classes of nouns and verbs with a high and low tendency to encode causation. For example, a named entity such as LOCATION may have the least tendency to encode causation. We leveraged such information about nouns to filter false positives. Similarly, we utilized the TimeBank’s (Pustejovsky et al., 2006) classification of verbal events (i.e., Occurrence, Perception, Aspectual, State, I\_State, I\_Action and Reporting) and their definitions to claim that the reporting events (e.g., say, tell, etc.)

just describe and narrate other events instead of encoding causality with them. We proposed an Integer Linear Programming (ILP) model (Roth and Yih, 2004; Do et al., 2011) to combine noun and verb semantics with the decisions of a supervised classifier which only relies on linguistic features.

In this paper, we extend our previous model by acquiring and exploiting the following three novel types of knowledge:

1. We learn the information about tendencies of various verb frames to encode causation. For example, our model identifies if the subject of verb “destroy” (“occur”) has a high (low) tendency to encode causation. Such information helps gain performance by exploiting causal semantics of each verb frame separately. We also learn and incorporate information about the verb frames in general e.g., how likely it is for the subject of any verb to encode causation with its verb.
2. In Riaz and Girju (2014), we utilized the TimeBank’s definition of reporting events to argue that such events have the least tendency to encode causation. Instead of relying on human judgment we now introduce a data intensive approach to identify the TimeBank’s classes of events with the least tendency to encode causation.
3. Although, information about the nouns with the least tendency to encode causation helps to filter false positives it can lead to false negatives when metonymic readings are associated with such nouns. Therefore, we introduce a metonymy resolver on top of our current model to avoid false negatives.

We provide details of our previously proposed model in section 3. We introduce new model and discuss its performance in sections 4 and 5. Section 6 concludes the current research.

## 2 Relevant Work

In Natural Language Processing (NLP), researchers are showing lots of interest in the task of identifying causality due to its various applications e.g., question answering (Girju, 2003), summarization (Chklovski and Pantel, 2004), future prediction (Radinsky and Horvitz, 2013), etc.

Several approaches have been proposed to identify causality in pairs of verbal events (Bethard and Martin, 2008; Riaz and Girju, 2010; Do et al., 2011; Riaz and Girju, 2013) and discourse

segments (Sporleder and Lascarides, 2008; Pitler and Nenkova, 2009; Pitler et al., 2009). However causality a pervasive relation of language can be encoded via various linguistic constructions. For example, verbs and nouns are the key components of language to represent events. Therefore in this work we focus on identifying causality in verb-noun pairs.

Previously researchers have followed the path of utilizing linguistic features in the framework of supervised learning (Girju, 2003; Bethard and Martin, 2008; Sporleder and Lascarides, 2008; Pitler and Nenkova, 2009; Pitler et al., 2009). Though linguistic features are important but other sources of knowledge are also critically required to achieve progress on the current task.

In recent years, researchers have proposed unsupervised metrics to identify causality between events (Riaz and Girju, 2010; Do et al., 2011). For example, Riaz and Girju (2010) and Do et al. (2011) introduced unsupervised metrics to learn causal dependencies between events. These metrics mainly depend on probabilities of co-occurrences of events and do not distinguish well causality from any other types of correlation (Riaz and Girju, 2013). In order to overcome this problem Riaz and Girju (2013) proposed some advanced metrics which combine probabilities of co-occurrences of events with the supervised estimates of cause and non-cause relations.

Considering the importance of employing rich sources of knowledge other than linguistic features for the current task, we have recently proposed a model that incorporates semantic classes of nouns and verbs with a high and low tendency to encode causation (Riaz and Girju, 2014). In this work, we exploit information about verb frames, data-driven verb semantics and metonymies to achieve more progress on our recent work.

## 3 Model for Recognizing Causality

In this section we provide an overview of our previous model (Riaz and Girju, 2014) for identifying causality in v-np pairs where v (np) stands for verb (noun\_phrase). This model works in the following two stages: (1) A supervised classifier is used to make binary predictions (i.e., the label cause (C) or non-cause ( $-C$ )) employing linguistic features, and (2) noun and verb semantics are then combined with the predictions of supervised classifier in the ILP framework to identify causality.

### 3.1 Supervised Classifier

To the best of our knowledge, there is no data set of v-np pairs with the labels C and  $\neg C$  available to us. For the current task we employ some heuristics to extract a training corpus of v-np pairs using FrameNet (Baker et al., 1998). FrameNet provides frame elements for the verbs and hand annotated examples (aka annotations) of these frame elements. Consider the following annotation from FrameNet “They **died** [*Cause* from **shotgun wounds**]” where the frame element “Cause” is given for the verb “died”. We remove the preposition “from” from the above annotation of frame element to acquire an instance of v-np (i.e., died-shotgun wounds) pair. We extract all annotations for verbs from FrameNet in which a frame element must contain at least one noun and no verb in it. We found such annotations for 729 distinct frame elements. We manually assigned the labels C and  $\neg C$  to these frame elements. Cause, Purpose, Reason, Result, Explanation are some examples of the frame elements to which we assigned the label C. Using the above mentioned assignments of labels C and  $\neg C$  to frame elements, we have acquired a training corpus of 4,141 (77,119) C ( $\neg C$ ) instances from FrameNet. In order to avoid class imbalance while training we employ an equal number of instances of both labels.

Due to space constraints, we refer the reader to Appendix A for the details of linguistic features to build the supervised classifier. We employ both Naive Bayes (NB) and Maximum Entropy (MaxEnt) algorithms to acquire predictions and probabilities of assignments of labels. We set up the following ILP using these probabilities:

$$Z_1 = \max \sum_{v-np \in I} \sum_{l \in L_1} x_1(v-np, l) P(v-np, l) \quad (1)$$

$$\sum_{l \in L_1} x_1(v-np, l) = 1 \quad \forall v-np \in I \quad (2)$$

$$x_1(v-np, l) \in \{0, 1\} \quad \forall v-np \in I \quad \forall l \in L_1 \quad (3)$$

Here,  $L_1 = \{C, \neg C\}$ ,  $I$  is the set of all v-np pairs.  $x_1(v-np, l)$  is a binary decision variable set to 1 only if the label  $l \in L_1$  is assigned to a v-np pair and only one label out of  $|L_1|$  choices can be assigned to a v-np pair (see constraints 2 and 3). In particular, we maximize the objective function  $Z_1$  (1) assigning the labels  $l \in \{L_1\}$  to v-np pairs depending on the probabilities of assignments (i.e.,  $P(v-np, l)$ ) obtained through the supervised classifier.

### 3.2 Noun and Verb Semantics

We automatically acquire and employ semantic classes of nouns and verbs with a high and low tendency to encode causation. Such information helps to reduce errors in predictions of the supervised classifier.

We derive two semantic classes of nouns for our purpose i.e.,  $C_{np}$  and  $\neg C_{np}$  where the class  $C_{np}$  ( $\neg C_{np}$ ) represents the noun\_phrases with a high (low) tendency to encode causation. For example, a noun\_phrase expression for a location has the least tendency to encode causation unless a metonymic reading is associated with it. In order to acquire these classes, we extract annotations of 936 distinct frame elements from FrameNet in which a frame element must contain at least one noun and no verb in it. These annotations of frame elements roughly represent instances of noun\_phrases (np). We manually assigned the labels  $C_{np}$  and  $\neg C_{np}$  to the frame elements. For example, we assign the label  $\neg C_{np}$  to the frame element “Place” which represents a location (see Appendix B for some examples of the frame elements with labels  $C_{np}$  and  $\neg C_{np}$ ). We also follow the approach similar to Girju and Moldovan (2002) to employ WordNet senses of nouns to acquire more instances of the classes  $C_{np}$  and  $\neg C_{np}$  (see Appendix B for the details). We have acquired a total of 280,212 instances of np (50% for each of the two classes i.e.,  $C_{np}$  and  $\neg C_{np}$ ) using both FrameNet and WordNet. Using these instances, we build a supervised classifier to identify the semantic class of np (see Appendix B for the details of features to build the classifier). We incorporate the knowledge of semantic classes of nouns by making the following additions to ILP:

$$Z_2 = Z_1 + \sum_{v-np \in I-M} \sum_{l \in L_2} x_2(f_{np}(v-np), l) P(f_{np}(v-np), l) \quad (4)$$

$$\sum_{l \in L_2} x_2(f_{np}(v-np), l) = 1 \quad \forall v-np \in I - M \quad (5)$$

$$x_2(f_{np}(v-np), l) \in \{0, 1\} \quad \forall v-np \in I - M \quad \forall l \in L_2 \quad (6)$$

$$x_1(v-np, \neg C) - x_2(f_{np}(v-np), \neg C_{np}) \geq 0 \quad \forall v-np \in I - M \quad (7)$$

Here  $L_2 = \{C_{np}, \neg C_{np}\}$ .  $f_{np}(v-np)$  is a function which returns np of a v-np pair.  $M$  is the set of v-np pairs with metonymic readings associated with np. Currently, this set is empty and in section 4.3 we introduce a metonymy resolver to pop-

ulate this set.  $x_2(f_{np}(v-np), l)$  is a binary decision variable set to 1 only if the label  $l \in L_2$  is assigned to np and only one label out of  $|L_2|$  choices can be assigned to np (see constraints 5 and 6). Constraint 7 enforces that if an np belongs to the class  $\neg C_{np}$  then its corresponding v-np pair is assigned the label  $\neg C$ . In particular, we maximize the objective function  $Z_2$  (4) subject to the constraints introduced till now. For each v-np pair, we predict the semantic class of np using our supervised classifier for the labels  $l \in L_2$  and set the probabilities – i.e.,  $P(f_{np}(v-np), l) = 1, P(f_{np}(v-np), \{L_2\} - \{l\}) = 0$  if the label  $l \in L_2$  is assigned to np. Also before running our supervised classifier, we run a named entity recognizer (Finkel et al., 2005) and assign the label  $\neg C_{np}$  to all noun\_phrases identified as named entities. We also determine association of metonymies with the noun\_phrases identified as named entities.

For the current task we also acquire two semantic classes of verbs i.e.,  $C_{ev}$  and  $\neg C_{ev}$  where the class  $C_{ev}$  ( $\neg C_{ev}$ ) contains the verbal events with a high (low) tendency to encode causation. In order to derive these two classes we exploit the TimeBank corpus (Pustejovsky et al., 2003) which provides seven semantic classes of verbal events – i.e., Occurrence, Perception, Aspectual, State, LState, LAction and Reporting. According to the definitions of these classes, we claim that the reporting events (e.g., say, tell, etc.) just describe and narrate other events instead of encoding causality with them. Using this claim, we consider that all instances of reporting verbal events of TimeBank belong to the class  $\neg C_{ev}$  and the rest of instances of verbal events lie in the class  $C_{ev}$ . After acquiring instances of the classes  $C_{ev}$  and  $\neg C_{ev}$ , we build a supervised classifier for these two classes. We use the features introduced by Bethard and Martin (2006) to build this classifier (see Bethard and Martin (2006) for the details). Employing predictions and probabilities of assignments of the labels  $C_{ev}$  and  $\neg C_{ev}$  we add the following two constraints to ILP: (1) if the event represented by v belongs to  $\neg C_{ev}$  then the corresponding v-np pair must be labeled with  $\neg C$  and (2) if a v-np pair is a causal pair then the event represented by v must be labeled with  $C_{ev}$ .

#### 4 Enriched Verb and Noun Semantics

This section describes the novel contributions of this work i.e., identification of semantics of verb

frames, semantic classes of verbal events via a data intensive approach and association of metonymic readings with noun\_phrases to identify causality with a better performance.

##### 4.1 Verb Frames

We introduce a method to acquire tendencies of various verb frames to encode causation. Consider the following two examples to understand the tendencies of verb frames of form  $\{v, gr\}$  to encode causation where v is the verb and gr is the grammatical relation of np with the verb v.

1. **The Great Storm of October 1987** almost totally **destroyed** the eighty year old pinetum at Nymans Garden in Sussex. (Cause (C))
2. **The explosion occurred** in the city’s main business area. (Non-Cause ( $\neg C$ ))

In above two examples the nps “The Great Storm of October 1987” and “The explosion” have the grammatical relations of subject with the verbs “destroyed” and “died”. In examples (1) and (2) the verb frames  $\{\text{destroy, subject}\}$  and  $\{\text{occur, subject}\}$  encode cause and non-cause relations. These examples reveal that each verb frame has its own tendency to encode causation. This type of knowledge helps gain performance by exploiting the semantics of each verb frame separately.

We leverage FrameNet annotations to acquire such type of knowledge. We collect all annotations of verbs from FrameNet and assign the labels C and  $\neg C$  to the frame elements as discussed in section 3.1. In FrameNet, example (1) is given as follows:

3. [*Cause* The Great Storm of October 1987] [*Degree* almost totally] **destroyed** [*Undergoer* the eighty year old pinetum at Nymans Garden in Sussex].

According to our assignments of labels C and  $\neg C$  to the frame elements, example (1) is given as “[C The Great Storm of October 1987] [ $\neg C$  almost totally] **destroyed** [ $\neg C$  the eighty year old pinetum at Nymans Garden in Sussex].”. After acquiring instances of the labels C and  $\neg C$  from example (1), we populate the fields of a knowledge base of verb frames (see Table 1). Fields of this knowledge base are  $\{v, gr\}$ ,  $\text{count}(\{v, gr\}, C)$  and  $\text{count}(\{v, gr\}, \neg C)$ . gr is the dependency relation of the frame element with the verb v. We use Stanford’s dependency parser (Marneffe et al., 2006) to collect dependency relations.  $\text{count}(\{v, gr\}, C)$  ( $\text{count}(\{v, gr\}, \neg C)$ ) is the count of the label C ( $\neg C$ ) of the frame  $\{v, gr\}$ . As shown in Table 1,

for the frame element ‘‘The Great Storm of October 1987’’, the word ‘‘Storm’’ has the dependency relation of ‘‘nsubj’’ with the verb ‘‘destroy’’. If there exists more than one dependency relations between the frame element and its verb then we choose the very first relation in the text order. According to the counts given in Table 1, {destroy, nsubj} has more tendency to encode a cause relation than the non-cause one. We have acquired 7,156 and 114,898 instances of the labels C and  $\neg C$  from FrameNet for populating the knowledge base of verb frames. We compute tendencies of verb frames to encode causality using the following scores:

$$\begin{aligned} S(\{v, gr\}, l) &= S_1(\{v, gr\}, l) \times S_2(\{*, gr\}, l) \quad (8) \\ S_1(\{v, gr\}, l) &= \frac{count(\{v, gr\}, l)}{count(\{v, gr\}, l) + count(\{v, gr\}, L_1 - \{l\})} \\ S_2(\{*, gr\}, l) &= \frac{count(\{*, gr\}, l)}{count(\{*, gr\}, l) + count(\{*, gr\}, L_1 - \{l\})} \end{aligned}$$

Counts of first component ( $S_1$ ) can be taken from the knowledge base of verb frames of form  $\{v, gr\}$ . The second component ( $S_2$ ) with counts  $count(\{*, gr\}, l)$  and  $count(\{*, gr\}, L_1 - \{l\})$  captures tendencies of verb frames in general. For example, what is the tendency of any subject to encode causality with its verb i.e., the score  $S_2(\{*, nsubj\}, C)$ . We populate the knowledge base of Table 1 with equal number of C and  $\neg C$  instances to calculate counts for  $S_2$ . We make the following additions to ILP to incorporate information about verb frames:

$$Z_3 = Z_2 + \sum_{\substack{v-np \in I \wedge \\ g(v-np) \in KB \wedge \\ f_{np}(v-np) \in C_{np}}} \sum_{l \in L_1} x_3(g(v-np), l) S(g(v-np), l) \quad (9)$$

$$\sum_{l \in L_1} x_3(g(v-np), l) = 1 \quad \forall \substack{v-np \in I \wedge \\ g(v-np) \in KB \wedge \\ f_{np}(v-np) \in C_{np}} \quad (10)$$

$$x_3(g(v-np), l) \in \{0, 1\} \quad \forall l \in L_1, \forall \substack{v-np \in I \wedge \\ g(v-np) \in KB \wedge \\ f_{np}(v-np) \in C_{np}} \quad (11)$$

$$x_3(g(v-np), l) \leq x_1(v-np, l) \quad \forall l \in L_1, \quad (12)$$

$$\forall \substack{v-np \in I \\ \wedge g(v-np) \in KB \\ \wedge f_{np}(v-np) \in C_{np}}$$

$$x_1(v-np, l) \leq x_3(g(v-np), l) \quad \forall l \in L_1, \quad (13)$$

$$\forall \substack{v-np \in I \wedge \\ g(v-np) \in KB \wedge \\ f_{np}(v-np) \in C_{np}}$$

Here,  $KB$  is the knowledge base of verb frames and  $g(v-np)$  is the function which returns the verb frame i.e.,  $\{v, gr\}$ . This function returns NULL value if there is no grammatical relation between  $v$  and  $np$  in an instance. The above changes in ILP are only applicable for the  $v-np$  pairs with

$\{v, gr\}$	$count(\{v, gr\}, C)$	$count(\{v, gr\}, \neg C)$
{destroy, nsubj}	1	0
{destroy, advmod}	0	1
{destroy, dobj}	0	1
[ $C$ The Great Storm of October 1987] [ $\neg C$ almost totally] <b>destroyed</b> [ $\neg C$ the eighty year old pinetum at Nymans Garden in Sussex].		

Table 1: A knowledge base of verb frames. This knowledge base is populated using the instances of C and  $\neg C$  labels given in this table.

$g(v-np) \in KB$  and  $np$  identified as of class  $C_{np}$  because we have already filtered the cases of  $np \in \neg C_{np}$  in section 3.2.  $x_3(g(v-np), l)$  is a binary decision variable set to 1 only if the label  $l \in L_1$  is assigned to  $g(v-np)$  and only one label out of  $|L_1|$  choices can be assigned to  $g(v-np)$  (see constraints 10 and 11). We add information about verb frames using constraints 12 and 13. These constraints enforce the predictions of the supervised classifier of causality (section 3.1) to be consistent with the predictions using tendencies of verb frames (i.e., score  $S(\{v, gr\}, l)$ ). We maximize objective function (9) subject to the above constraints. We remove those  $\{v, gr\}$  from KB which have  $count(\{v, gr\}, C) + count(\{v, gr\}, \neg C) < 5$  to avoid wrong predictions based on the small counts of verb frames.

## 4.2 Data-driven Verb Semantics

In section 3.2 we considered that reporting events belong to the class  $\neg C_{ev}$  with the least tendency to encode causation using the definition of these events in the TimeBank corpus. Instead of relying on definitions of events we now introduce a data intensive approach to automatically identify the class  $\neg C_{ev}$  of verbal events. In order to identify this class we extract training instances of verbal events encoding C and  $\neg C$  relations. Verbal events encode cause-effect relations using verb-verb (e.g., Five shoppers were **killed** when a car **blew up**.) and verb-noun linguistic constructions. Therefore for the current purpose we use the following two types of training instances: (A) a training corpus of 240K instances of verb-verb ( $v_i-v_j$ ) pairs encoding C and  $\neg C$  relations (named as Training $_{v_i-v_j}$ ) (we refer the reader to Riaz and Girju (2013) for the details of this training corpus) and (B) the training corpus  $v-np$  instances introduced in section 3.1 (named as Training $_{v-np}$ ).

Following is the procedure to derive  $V_{\neg C} \subseteq V$  where  $V = \{\text{Occurrence, Perception, Aspectual, State, I\_State, I\_Action, Reporting}\}$  and the set  $V_{\neg C}$  contains the TimeBank’s semantic classes

with the least tendency to encode a cause relation.

1. **Input:** Training corpus,  $V$
2. **Output:** Set  $V_{-C}$
3. For each training instance  $k$  employ the supervised classifier of Bethard and Martin (2006) to do the following:
  - (a) if  $k \in \text{Training}_{v_i-v_j}$  then identify the semantic class (sc) of both events represented by both verbs  $v_i$  and  $v_j$  and add this information to a set i.e.,  $T = T \cup (k_{v_i}, \text{sc}_{v_i}, l) \cup (k_{v_j}, \text{sc}_{v_j}, l)$  where  $\text{sc}_{v_i}$  is the semantic class of event of the verb  $v_i$  of instance  $k$  and  $l \in \{C, \neg C\}$ .
  - (b) Else if  $k \in \text{Training}_{v-np}$  then identify the semantic class (sc) of event represented by the verb  $v$  and set  $T = T \cup (k_v, \text{sc}_v, l)$ .
4. Using results of step 3, calculate tendency of each semantic class  $\text{sc} \in V$  to encode non-causality (i.e.,  $\text{score}(\text{sc}, \neg C)$ ) as follows:

$$\begin{aligned} \text{score}(\text{sc}, \neg C) &= \text{score}_1(\text{sc}, \neg C) \times \text{score}_2(\text{sc}, \neg C) \\ \text{score}_1(\text{sc}, \neg C) &= \left( \frac{\text{count}(\text{sc}, \neg C)}{\text{count}(\text{sc})} - \frac{\text{count}(\text{sc}, C)}{\text{count}(\text{sc})} \right) \\ \text{score}_2(\text{sc}, \neg C) &= \left( \frac{\text{count}(\text{sc}, \neg C)}{\text{count}(\neg C)} - \frac{\text{count}(\text{sc}, C)}{\text{count}(C)} \right) \end{aligned}$$

where  $\text{count}(m, n)$  is the number of instances of verbal events with the labels  $m$  and  $n$  and  $\text{count}(m)$  is the number of instances of verbal events with the label  $m$ .

5. Acquire a ranked list of semantic classes  $\text{list}_{\text{sc}} = [\text{sc}_1, \text{sc}_2, \dots, \text{sc}_m]$  s.t.  $\text{score}(\text{sc}_i, \neg C) \geq \text{score}(\text{sc}_{i+1}, \neg C)$ . From this list we remove the class  $\text{sc}_i$  if either  $\text{score}_1(\text{sc}_i, \neg C) < 0$  or  $\text{score}_2(\text{sc}_i, \neg C) < 0$ .
  - ▷ The following steps are used to determine the cutoff class  $\text{sc}_i \in \text{list}_{\text{sc}}$  s.t. the semantic classes  $\{\text{sc}_1, \text{sc}_2, \dots, \text{sc}_{i-1}\}$  have the least tendency to encode causation.
6.  $\text{result}_{\text{sc}_{-1}} = 0$  and  $\text{result}_{\text{sc}_0} = 0$ .
7. Remove  $\text{sc}_i$  from the front of  $\text{list}_{\text{sc}}$  and do the following:
  - (c) Predict the label ( $l$ )  $\neg C$  for all tuples of form  $(m, \text{sc}, l) \in T$  if  $\text{sc} \in \{\text{sc}_1, \text{sc}_2, \dots, \text{sc}_i\}$  and predict  $C$  for the rest of the tuples.
  - (d) Using predictions from step (c), calculate the  $\text{result}_{\text{sc}_i} = \text{F1-score} \times \text{accuracy}$  for the label  $l \in \{C, \neg C\}$ .
  - (e) If  $\text{result}_{\text{sc}_i} - \text{result}_{\text{sc}_{i-1}} < \text{result}_{\text{sc}_{i-1}} - \text{result}_{\text{sc}_{i-2}}$  then output  $\{\text{sc}_1, \text{sc}_2, \dots, \text{sc}_{i-1}\}$
  - (f) Else go to step 7.

Using the above procedure, we obtain the sets  $\{\text{Aspectual}\}$  and  $\{\text{Reporting, I\_State}\}$  with  $\text{Training}_{v_i-v_j}$  and  $\text{Training}_{v-np}$  corpora. We consider that the Aspectual, Reporting and I\_State events of the TimeBank corpus belong to the class  $\neg C_{e_v}$  and rest of the events lie in  $C_{e_v}$ . Using these semantic classes we apply the constraints introduced in section 3.2.

### 4.3 Metonymy Resolution:

Metonymy resolution is the task to determine if a literal or non-literal reading is associated with a

$\{v, \text{gr}\}$	$\text{count}(\{v, \text{gr}\}, C_{np})$	$\text{count}(\{v, \text{gr}\}, \neg C_{np})$
$\{\text{kill}, \text{nsubj}\}$	1	0
$\{\text{kill}, \text{dobj}\}$	0	1
$[C_{np} \text{Pissed off Angelus}] \text{ just kills } [\neg C_{np} \text{me}]$		

Table 2: A knowledge base of verb frames. This knowledge base is populated using the instances of  $C_{np}$  and  $\neg C_{np}$  labels given in this table.

natural language expression (Markert and Nissim, 2009). Consider the following example:

4. **The United States** has **killed** Osama bin Laden and has custody of his body. (Cause (C))

In example (4) “The United States” refers to a non-literal reading i.e., the event of “raid in Abbottabad on May 2, 2011 by the United States” rather than merely referring to a literal sense i.e., a country. The association of non-literal reading with “The United States” results in killing event. Previously, researchers have worked with hand-annotated selectional restrictions violation for this task (Markert and Nissim, 2009). In the example (4) a country cannot “kill” someone and thus a metonymic reading is associated with it. In this work we identify association of metonymies with noun\_phrases via verb frames and prepositions as explained below in this section.

In the first part of our approach we employ violations of tendencies of verb frames to identify if a non-literal reading is associated with a noun\_phrases. Particularly, we build a knowledge base of verb frames using  $C_{np}$  and  $\neg C_{np}$  classes as discussed in section 4.1. Consider the knowledge base given in Table 2 populated using the following FrameNet annotations “[*Stimulus*Pissed off Angelus] just kills [*Experiencer*me].” with assignments of labels  $C_{np}$  and  $\neg C_{np}$  to the frame elements. We populate the knowledge base using only those FrameNet annotations in which a frame element does not contain a verb.

Now we introduce our method to identify the association of non-literal reading with the “The United States” in example (4). The supervised classifier predicts the class  $\neg C_{np}$  for the np “The United States”. However, in the current state of knowledge base (Table 2)  $P(\{\text{destroy}, \text{nsubj}\}, C_{np}) > P(\{\text{destroy}, \text{nsubj}\}, \neg C_{np})$  where  $P$  is the probability. The prediction of  $\neg C_{np}$  for “The United States” violates the above probabilities. Considering this violation, we predict the association of metonymy with np.

In the second part of our approach we identify tendencies of prepositions to encode causation

and use violation of these tendencies to identify metonymies. For this purpose, we use the training corpus of v-np pairs with 4, 141 C and 77, 119  $\neg$ C training instances (see section 3.1). We employ only those training instances in which a preposition appears between v and np and there appears no verb between them. From these instances, we acquire a set of prepositions that appear between v and np. Using this set of prepositions (PR) as input to the following procedure, we acquire a set of prepositions (PR<sub>C</sub>) with the highest tendency to encode causation:

1. **Input:** Training Corpus of v-np pairs, PR
2. **Output:** PR<sub>C</sub>
3. Calculate tendency of each preposition  $pr \in PR$  to encode causality (i.e.,  $score(pr, C)$ ) as follows:

$$score(pr, C) = score_1(pr, C) \times score_2(pr, C)$$

$$score_1(pr, C) = \left( \frac{count(pr, C)}{count(pr)} - \frac{count(pr, \neg C)}{count(pr)} \right)$$

$$score_2(pr, C) = \left( \frac{count(pr, C)}{count(C)} - \frac{count(pr, \neg C)}{count(\neg C)} \right)$$

4. Acquire a ranked list of prepositions  $list_{pr} = [pr_1, pr_2, \dots, pr_m]$  s.t.  $score(pr_i, C) \geq score(pr_{i+1}, C)$ . From this list we remove  $pr_i$  if either  $score_1(pr_i, C)$  or  $score_2(pr_i, C) < 0$ .
5.  $result_{pr_{-1}} = 0, result_{pr_0} = 0$
6. Remove  $pr_i$  from the front of the  $list_{pr}$  and do the following:
  - (a) Predict the label C for all v-np training instances with  $pr \in \{pr_1, pr_2, \dots, pr_i\}$  and assign the label  $\neg C$  to the rest of the instances.
  - (b) Using predictions from step (a) calculate the  $result_{pr_i} = F1\text{-score} \times accuracy$ .
  - (c) If  $result_{pr_i} - result_{pr_{i-1}} < result_{pr_{i-1}} - result_{pr_{i-2}}$  then output  $\{pr_1, pr_2, \dots, pr_{i-1}\}$ .
  - (d) Else go to step 6.

The above procedure outputs the set  $PR_C = \{\text{for, by}\}$ . Now we introduce method to identify association of non-literal reading for the example “All weapon sites in Iraq were **destroyed** by **the United States**” where “the United States”  $\in \neg C_{np}$  as identified by the supervised classifier. However, the preposition “by” has a high tendency to encode causation and thus “the United States” may encode causation. Therefore, there is a possibility that this noun\_phrase has a non-literal sense attached to it which results in encoding causality. Using this method, we predict metonymies only for the v-np instances where preposition appears between v and np and there appears no verb between them. If any of two methods of metonymy resolution predicts the association of metonymy with np then we add v-np to the set M used in ILP (see section 3.2).

## 5 Evaluation and Discussion

In this section we present experiments and discussion on the performance achieved for the current task. In order to evaluate our model, we generated a test set of instances of v-np pairs. For this purpose, we collected three wiki articles on the topics of Hurricane Katrina, Iraq War and Egyptian Revolution of 2011. We apply a part-of-speech tagger and a dependency parser on all sentences of these three articles (Toutanova et al., 2003; Marneffe et al., 2006). We extracted all v-np pairs from each sentence of these articles. For each of the these three articles, we selected first 500 instances of v-np pairs. Two annotators were asked to provide the labels C and  $\neg C$  to the instances of v-np pairs using the annotation guidelines from Riaz and Girju (2010). We have achieved a 0.64 kappa score for the human inter-annotator agreement on a total of 1,500 v-np instances. This results in a total of 1,365 instances of v-np pairs with 11.86% C pairs.

In this section, we present performance of the following models (see Table 3):

1. **Baseline:** NB and MaxEnt (McCallum, 2002) supervised classifiers using only the shallow linguistic features (see section 3.1).
2. **Basic noun and verb semantics:** ILP with the addition of semantic classes of nouns without metonymy (denoted by  $+N_{!M}$ ) and the addition of semantic classes of verbs where  $\neg C_{ev} = \{(R)eporing\ events\}$  (denoted by  $+N_{!M}+V_{\{R\}}$ ). These models represent the work proposed in Riaz and Girju (2014) (section 3).
3. **Noun semantics with metonymies:** ILP with the addition of noun semantics involving metonymies resolved via verb frames (denoted by  $+N_{M_1}$ ), metonymies resolved via verb frames  $\{v, gr\}$  where  $gr \in GR = \{csubj, csubpass, nsubj, nsubjpass, xsubj, dobj, iobj, pobj, agent\}$  a set of core dependency relations of subjects and objects (denoted by  $+N_{M_{!GR}}$ ) and metonymies resolved via both verb frames and prepositions (denoted by  $+N_{M_{!GR}+M_2}$ ).
4. **Verb frames and data-driven verb semantics:** ILP with the addition of information about verb frames (denoted by  $+N_M+VF$  where  $M = M_{!GR} + M_2$ ), data-driven verb semantics i.e.,  $\neg C_{ev} = \{(A)spectual, (R)eporing, (I)\_ (S)tate\ events\}$  (denoted by  $+N_M+V_{\{A,R,IS\}}$ ) and both verb frames and data-driven verb semantics (denoted by  $+N_M+VF+V_{\{A,R,IS\}}$ )

S	B	+N <sub>1M</sub>	+N <sub>1M</sub> +V <sub>{R}</sub>	+N <sub>M1</sub>	+N <sub>M1GR</sub>	+N <sub>M1GR</sub> +M <sub>2</sub>	+N <sub>M</sub> +VF	+N <sub>M</sub> +V <sub>{A,R,IS}</sub>	+N <sub>M</sub> +VF+V <sub>{A,R,IS}</sub>
A	28.86	71.86	73.40	71.35	71.42	71.64	72.96	75.16	76.19
P	13.52	26.18	27.21	26.29	26.34	27.54	28.39	29.93	30.82
R	92.59	75.30	74.07	78.39	78.39	85.18	83.95	81.48	80.86
F	23.60	38.85	39.80	39.37	39.44	41.62	42.43	43.78	44.63
A	61.46	80.73	81.17	80.65	80.73	81.02	81.39	81.75	82.05
P	19.46	32.02	32.72	32.41	32.52	34.09	34.66	35.25	35.64
R	71.60	55.55	55.55	58.02	58.24	64.19	64.19	64.19	63.58
F	30.60	40.63	41.18	41.59	41.68	44.53	45.02	45.51	45.67

Table 3: Performance of (B)aseline, +N<sub>1M</sub>, +N<sub>1M</sub>+V<sub>{R}</sub>, +N<sub>M1</sub>, +N<sub>M1GR</sub>, +N<sub>M1GR</sub>+M<sub>2</sub>, +N<sub>M</sub>+VF, +N<sub>M</sub>+V<sub>{A,R,IS}</sub> and +N<sub>M</sub>+VF+V<sub>{A,R,IS}</sub> (see text for details) in terms of (S)cores of (A)ccuracy, (P)recision, (R)ecall, (F)-score. The row 1 (2) of this table presents results over NB (MaxEnt) baseline supervised classifier, respectively.

Table 3 shows that MaxEnt gives a very high accuracy and F-score as compared with NB. Model +N<sub>1M</sub>+V<sub>{R}</sub> with basic noun and verb semantics introduced in section 3.2 results in more than 10% improvement in F-score over NB and MaxEnt classifiers relying only on shallow linguistic features. Model +N<sub>M</sub>+VF+V<sub>{A,R,IS}</sub> with enriched verb and noun semantics brings more than 4% improvement in F-score over +N<sub>1M</sub>+V<sub>{R}</sub> with MaxEnt as baseline. We perform statistical significance test using bootstrap sampling method given in Berg-Kirkpatrick et al. (2012) (see Berg-Kirkpatrick et al. (2012) for the details). +N<sub>M</sub>+VF+V<sub>{A,R,IS}</sub> brings significant improvement in F-score over +N<sub>1M</sub>+V<sub>{R}</sub> with p-value 0.0.

Though +N<sub>1M</sub> gives significantly better F-score over baseline, it drops recall by more than 16%. Metonymy resolution helps perform quite better by recovering more than 8% recall with +N<sub>M1GR</sub>+M<sub>2</sub> over +N<sub>1M</sub>. +N<sub>M1GR</sub>+M<sub>2</sub> also results in 3.9% improvement in F-score over +N<sub>1M</sub> with MaxEnt as baseline model (significant improvement with p-value 0.0). Metonymies resolved via verb frames with all and core grammatical relations (i.e., set GR) recover more than 2% recall and slightly improve F-score.

Model with the addition of information of verb frames (i.e., +N<sub>M</sub>+VF) brings 0.49% improvement in F-score over +N<sub>M1GR</sub>+M<sub>2</sub> using MaxEnt as baseline model (significant improvement with p-value 0.027). Model with the addition of data-driven verb semantics (i.e., +N<sub>M</sub>+V<sub>{A,R,IS}</sub>) results in 0.98% improvement in F-score over +N<sub>M1GR</sub>+M<sub>2</sub> using MaxEnt as baseline model (significant improvement with p-value 0.0021). Overall the model +N<sub>M</sub>+VF+V<sub>{A,R,IS}</sub> yields more than 16% (20%) F-score (accuracy) over the baseline models build via NB and MaxEnt.

## 5.1 Error Analysis

We performed error analysis for the model +N<sub>M</sub>+VF+V<sub>{A,R,IS}</sub> by randomly selecting 50 False Positives (FP) and 50 False Negatives (FN).

For 32% FP instances information about verb frames is not available in the knowledge base of verb frames. To avoid this problem researchers should exploit some abstractions e.g., {semantic sense of v, gr} frames. Our model fails to identify the class  $-C_{np}$  for noun\_phrases of 29% FP instances due to the lack of enough training data for the semantic classes of nouns. In 21% FP instances v and np are not even relevant to each other. Our model first needs to determine relevance between v and np before identifying causality. Remaining 18% instances have v and np in temporal only sense, comparison relation or both represent parts of same event. There is need to extract more knowledge sources to better distinguish causality from any other type of relation.

77% FN instances are classified as non-causal due to the lack of enough v-np training data and require more sources of knowledge e.g., background knowledge. On remaining 23% FN instances our model fails to identify  $C_{np}$  class due to the lack of enough training data for the semantic classes of nouns.

## 6 Conclusion

This work has revealed that enriched semantics of nouns and verbs help gain significant improvement in performance over a baseline relying only on shallow linguistic features. Through empirical evaluation and error analysis of our model we have highlighted strengths and weaknesses of our model for the current task. Our work has provided a novel direction to exploit semantics of participants of causal relations to solve the challenge of identifying causality.



## References

- Collin F. Baker, Charles J. Fillmore and John B. Lowe. 1998. The Berkeley FrameNet project. *In proceedings of COLING-ACL. Montreal, Canada.*
- Taylor Berg-Kirkpatrick, David Burkett, and Dan Klein. 2012. An empirical investigation of statistical significance in NLP. *In Proceedings of the Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL).*
- Steven Bethard and James H. Martin. 2006. Identification of Event Mentions and their Semantic Class. *In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP).*
- Steven Bethard and James H. Martin. 2008. Learning Semantic Links from a Corpus of Parallel Temporal and Causal Relations. *In proceedings of ACL-08: HLT.*
- Timothy Chklovski and Patrick Pantel. 2004. VerbOcean: Mining the Web for Fine-Grained Semantic Verb Relations. *In proceedings of Conference on Empirical Methods in Natural Language Processing (EMNLP-04). Barcelona, Spain.*
- Quang X. Do, Yee S. Chen and Dan Roth. 2011. Minimally Supervised Event Causality Identification. *In proceedings of EMNLP-2011.*
- Jenny R. Finkel, Trond Grenager, and Christopher Manning. 2005. Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling. *In Proceedings of the Association for Computational Linguistics (ACL).*
- Roxana Girju and Dan Moldovan. 2002. Mining Answers for Causation Questions. *In American Associations of Artificial Intelligence (AAAI), 2002 Symposium.*
- Roxana Girju. 2003. Automatic detection of causal relations for Question Answering. *Association for Computational Linguistics ACL, Workshop on Multilingual Summarization and Question Answering Machine Learning and Beyond 2003.*
- Katja Markert and Malvina Nissim. 2009. Data and models for metonymy resolution. *Language Resources and Evaluation Volume 43 Issue 2, Pages 123–138.*
- Marie-Catherine de Marneffe, Bill MacCartney and Christopher D. Manning. 2006. Generating typed dependency parses from phrase structure parses. *In Proceedings of the International Conference on Language Resources and Evaluation (LREC).*
- Andrew K. McCallum. 2002. MALLET: A Machine Learning for Language Toolkit. <http://mallet.cs.umass.edu>.
- Emily Pitler, Annie Louis and Ani Nenkova. 2009. Automatic Sense Prediction for Implicit Discourse Relations in Text. *In proceedings of ACL-IJCNLP, 2009.*
- Emily Pitler and Ani Nenkova. 2009. Using Syntax to Disambiguate Explicit Discourse Connectives in Text. *In proceedings of ACL-IJCNLP, 2009.*
- James Pustejovsky, Patrick Hanks, Roser Saur, Andrew See, Robert Gaizauskas, Andrea Setzer, Dragomir Radev, Beth Sundheim, David Day, Lisa Ferro and Marcia Lazo. 2003. The TIMEBANK Corpus. *In Proceedings of Corpus Linguistics.*
- Kira Radinsky and Eric Horvitz. 2013. Mining the Web to Predict Future Events. *In proceedings of sixth ACM international conference on Web search and data mining, WSDM '13.*
- Mehwish Riaz and Roxana Girju. 2010. Another Look at Causality: Discovering Scenario-Specific Contingency Relationships with No Supervision. *In proceedings of the IEEE 4th International Conference on Semantic Computing (ICSC).*
- Mehwish Riaz and Roxana Girju. 2013. Toward a Better Understanding of Causality between Verbal Events: Extraction and Analysis of the Causal Power of Verb-Verb Associations. *Proceedings of the annual SIGdial Meeting on Discourse and Dialogue (SIGDIAL).*
- Mehwish Riaz and Roxana Girju. 2014. Recognizing Causality in Verb-Noun Pairs via Noun and Verb Semantics. *Proceedings of the Workshop on Computational Approaches to Causality in Language EACL, 2014.*
- Dan Roth and Wen-tau Yih. 2004. A Linear Programming Formulation for Global Inference in Natural Language Tasks. *In Proceedings of the Annual Conference on Computational Natural Language Learning (CoNLL).*
- Caroline Sporleder and Alex Lascarides. 2008. Using automatically labelled examples to classify rhetorical relations: An assessment. *Journal of Natural Language Engineering Volume 14 Issue 3, July 2008 Pages 369–416.*
- Kristina Toutanova, Dan Klein, Christopher Manning, and Yoram Singer. 2003. Feature-Rich Part-of-Speech Tagging with a Cyclic Dependency Network. *In Proceedings of Human Language Technology and North American Chapter of the Association for Computational Linguistics (HLT-NAACL).*

## Appendix A. Supervised Classifier

In this appendix, we provide a set of linguistic features taken from Riaz and Girju (2014) to identify causality in v-np pairs employing a supervised classifier (see section 3.1 for the details).

- **Lexical Features:** verb, lemma of verb, noun\_phrase, lemmas of all words of noun\_phrase, head noun of noun\_phrase, lemmas of all words between verb and noun\_phrase.
- **Syntactic Features:** part-of-speech tags of verb and head noun of noun\_phrase.
- **Semantic Features:** We adopted this feature from Girju (2003) to capture semantics of nouns. The 9 noun hierarchies of WordNet i.e., entity, psychological feature, abstraction, state, event, act, group, possession, phenomenon are used as this feature. Each of these hierarchies is set to 1 if any sense of head noun of noun\_phrase lies in that hierarchy, otherwise set to 0.
- **Structural Features:** This feature is applied by considering both subject (i.e., sub\_in\_np) and object (i.e., obj\_in\_np) of verb (v). For example, for the pair v-np the variable sub\_in\_np is set to 1 if the subject  $\in$  np, set to 0 if the subject  $\notin$  np and set to -1 if the subject is not available in an instance. The subject and object of a verb are its core arguments and may sometime be part of an event represented by a verb. Therefore, these arguments may have high tendency to encode non-causation with their verb.
- **Pairs:** The following pairs (verb, head noun of noun\_phrase), (subject<sub>verb</sub>, head noun of noun\_phrase) and (object<sub>verb</sub>, head noun of noun\_phrase) are used to capture relations.

## Appendix B. Noun Semantics

In this appendix, some examples of the frame elements of FrameNet and the WordNet senses belonging to the classes  $C_{np}$  and  $\neg C_{np}$  are given in Tables 4 and 5 (see section 3.2 for the details). We employ training instances acquired using the FrameNet annotations and WordNet senses for building a supervised classifier for the classes  $C_{np}$  and  $\neg C_{np}$ . Following is the list of features we use for this supervised classifier:

- **Lexical Features:** All words of noun\_phrase, lemmas of all words of noun\_phrase, head noun of noun\_phrase, first two (three) (four) letters

SC	FrameNet Labels
$c_{np}$	Event, Goal, Purpose, Cause, Internal cause, External cause, Result, Means, Reason, Phenomena, Coordinated event, Action, Activity, Circumstances, Desired goal, Explanation
$\neg c_{np}$	Artist, Performer, Duration, Time, Place, Distributor, Area, Path, Direction, Sub-region Frequency, Body part, Area, Degree, Angle, Fixed location, Path shape, Addressee, Interval

Table 4: Some examples of the frame elements of FrameNet to which we assign the semantic classes  $C_{np}$  and  $\neg C_{np}$ .

SC	WordNet Senses
$c_{np}$	{act, deed, human action, human activity}, {phenomenon}, {state}, {psychological feature}, {event}, {causal agent, cause, causal agency}
$\neg c_{np}$	{time period, period of time, period}, {measure, quantity, amount}, {group, grouping}, {organization, organisation}, {time unit, unit of time}, {clock time, time}

Table 5: This table shows our selected WordNet senses of nouns belonging to classes  $C_{np}$  and  $\neg C_{np}$ . For example, using the information provided in this table we assume that any noun concept whose all senses of WordNet lie in the semantic hierarchy of the sense {time period, period of time, period} is of class  $\neg C_{np}$ . We use English Gigaword corpus to collect instances of noun (or noun\_phrases) and label them with  $C_{np}$  and  $\neg C_{np}$  according to their senses in WordNet.

of head noun of noun\_phrase, last two, (three) (four) letters of head noun of noun\_phrase.

- **Word Class Features:** part-of-speech tags of all words of noun\_phrase and head noun of noun\_phrase.
- **Semantic Features:** Frequent sense of head noun of noun\_phrase.