# 國 立 清 華 大 學

## 碩 士 論 文

題目　結合語料庫與知識庫之動詞語意框
解歧與習得

Combining corpus statistics and knowledge
base to disambiguate and acquire verb frames

系別　資訊工程學系　　組別　＿＿＿＿＿＿

學號姓名　100062654　高定慧（Ting-Hui Kao）

指導教授　張俊盛 博士（Dr. Jason S. Chang）

中華民國 一百零二 年 七 月

# 摘要

在句法中，動詞扮演著舉足輕重的角色。對語言學習者而言，學習相關的動詞架框（verb frame），更是一個重要的課題。 不幸的是，現今大部分的線上字典，例如朗文字典，對於動詞架框都只提供非常粗淺、概略性的表示法，也就是一般列舉的某事（something）與某人（somebody）。然而，這樣粗略的標籤並無法有效的幫助語言學習者。在本論文中，我們提出一個自動習得動詞語意框的方法，並提供更加全面、更容易理解的表示法。我們首先根據從語料庫中得到的句法關係，抽取出與動詞有關的參數並組合出動詞參數組（verb argument tuple），再藉由知識庫（knowledge base）決定各個動詞參數（verb argument）的語意類別。接著透過組合並估計動詞語意組（verb semantic tuple）出現的機率以消除歧義，最後便得到動詞架框（verb frame） 。我們開發了一個雛形系統 *FrameFinder*，根據上述方法自動產生動詞語意框，並採用一個人工編輯的動詞型態（verb pattern）字典作為標準答案，評估結果也顯示此方法對於常見動詞，可以得到令人滿意的準確度。


**關鍵字：動詞語意框習得、語意解歧、知識庫、最大期望演算法**

i

# ABSTRACT

Verb frames are very important for language learners, since they capture the semantics and word usages associated with verbs. Unfortunately, most online dictionaries such as Longman Dictionary show verb frames with broad semantic categories (i.e., *something* and *somebody*) which are not very informative. In this work, we introduce a method for automatically generating more comprehensive verb frames. The method involves extracting verb argument tuples based on grammatical relations acquired from a parsed corpus, obtaining intended semantic categories for each argument based on a knowledge base, estimating the probabilities of each semantically labeled tuples, and finally generating verb frames. We present a prototype system, *FrameFinder*, that applies the method to generate verb frames automatically. Evaluation on a set of verbs with manually compiled semantic patterns shows that the method is able to extract with high accuracy for the important high frequnecies verbs for language learning.

**Keyword: Verb Frame generation, Word Sense Disambiguation, Knowledge Base, Estimation-Maximization algorithm**

# 致 謝 辭

感謝我的指導教授，張俊盛老師，帶領我進入自然語言處理的領域。在清華的兩年，老師指引我做研究的方向，並鼓勵我探索開發自己有興趣的題材。每當遇到困難與瓶頸，老師一向樂觀面對，並與我一同積極、嚴謹的尋求解決之道，這份做研究的態度與熱忱，是我效法的目標，萬分感謝老師。而在口試期間，承蒙國立台灣大學陳信希教授與國立中央大學張嘉惠教授的諸多寶貴建議，讓本論文更趨於完備，僅此獻上誠摯感謝。

在自然語言處理實驗室這個溫暖的大家庭裡，特別感謝鑑城學長，從程式的規劃到實驗設計、分析評估結果，學長總是給我最實質的指引。感謝玫樺在百忙之中幫助我建立資料，使本論文得以採用自動化評估。感謝 Neo 分享相關研究領域的知識與文獻，讓我對前人的方法有更廣更深入的了解。謝謝國父不厭其煩的跟我說明程式的版本管理，並在電腦出狀況的時候即時伸出援手，讓本論文可以如期順利完成。感謝與我一同並肩努力的小草和小冷，幫忙張羅、準備口試的一切。謝謝 Joanne、小蘭、絢紋、祐維、文斌、正朋與 Nicolas，一路上你們的加油打氣，給我極大的動力。

感謝我的女朋友，在我忙於研究時，不斷給予鼓勵、包容與體諒。最後，感謝我的家人，一直以來的陪伴與支持，你們是我最堅實的後盾。這篇論文，獻給我的師長、同學、好友，以及所有關心我的各位。

<div align="right">

高定慧 敬書

自然語言處理實驗室

民國一百零二年

</div>

# Table of Contents

# List of Figures

# List of Tables

# CHAPTER 1

# INTRODUCTION

As researchers have begun using the Web as corpus to investigate language-related issues (Kilgarriff and Grefenstette, 2003), more and more users also turn to the Web as resource and tools for learning a language. An increasing number of Web services provide language related search. For example, Unified Verb Index[1] merges resources from several natural language processing projects. As for semantic information, most online dictionaries such as Longman Dictionary[2] show verb frames with broad semantic categories (i.e., *something* and *somebody*) which are not very informative.

In addition to online dictionaries, several corpus based reference systems focus on providing dictionary-like syntactical information for learners. For example, *StringNet* (Wible and Tsao, 2010) presents lexico-grammatical hybrid-ngrams, i.e., multiword patterns, and their relations to each other. Similarly, *GRASP* (Huang et al., 2011) provides general syntactic patterns as well as lexical bundles anchored at the query words. However, syntactic patterns even with broad semantics categories are not sufficient to helping the learner choose the right words (e.g., a verb) in a certain

---

[1] http://verbs.colorado.edu/verb-index/index.php

[2] www.ldoceonline.com

context. Consider the verbs *build* and *construct* sharing the same syntactic pattern "SOMEBODY *verb* SOMETHING." It is very easy for a learner to confuse the two verbs, leading to word choice errors (i.e., *construct relationship*). Learner can learn how to distinguish confusable words and make better word choices, if there is a reference tool that provides semantic patterns of word usage (e.g., "SOMEBODY construct ARTIFACT" and "SOMEBODY build RELATION"). Intuitively, these semantic categories (of arguments of a verb), such as SOMEBODY, ARTIFACT, and RELATION can be obtained by analyzing the grammatical relations and consulting a lexical semantics knowledge base (e.g., WordNet). Thus, we generate semantic verb frames by using a parser (e.g., Stanford Parser) to obtain verbal argument tuples composed of semantic roles, such as subjects, objects, and prepositional phrases, and disambiguate their semantic categories of these argument using a method similar to existing unsupervised WSD algorithms (e.g., Yarowsky 1992).

Consider the verb argument tuple *<workers, abandon, plant>* with two possible verb frame candidates, <**PERSON** *abandon* **ARTIFACT**> or <**PERSON** *abandon* **PLANT**>. The intended verb frame is <**PERSON** *abandon* **ARTIFACT**>. Intuitively, this verb frame can be obtained by comparing the numbers of ARTIFACT arguments (e.g., car, ship, and train) and PLANT arguments (e.g., flower) of the verb *abandon*.

We present a new system, *FrameFinder*, that automatically acquire semantic verb frames for a given verb. An example of *FrameFinder* retrieving patterns for the verb *abandon* is shown in Figure 1. *FrameFinder* has parsed the sentences containing the verb *abandon* in a corpus (e.g., "I abandoned the wild idea.") and obtained verb argument tuples from the grammatical relations in the parses (e.g., *<I, abandon, idea>*) and their likely semantic verb frames (e.g., <**PERSON** *abandon* **CONGNITION**> and <**PERSON** *abandon* **COMMUNICATION**>). *FrameFinder* uses a set of semantic categories such as PERSON, COGNITION obtained from lexical semantics knowledge bases such as WordNet. Finally, *FrameFinder* disambiguates each verb argument tuple by estimating the probability of semantic verb frame via iterative re-estimation. We describe the probability estimation model in more detail in Chapter 3.

*FrameFinder* iterates between model estimation and disambiguation and finally determines a set of semantic verb frames with probability. In our prototype, *FrameFinder* returns the frames to the user directly for language learning (see Figure 1); alternatively, the frames returned by *FrameFinder* can be used as knowledge source for a grammatical error correction system.

Input:　　　　　| abandon |

Knowledge base

　I, rescuers, governor, worker, Mary, friends, mother (PERSON)

　idea (COMMUNICATION, COGNITION)

　hope (FEELING, COGNITION, ATTRIBUTE)

　plan (COGNITION, ARTIFACT)

　plant (ARTIFACT, PLANT)

　ship (ARTIFACT)

　attempt (ACT)

Parsed sentences and candidate patterns

1. I finally abandoned the wild idea.

   nsubj(abandon, **I**) dobj(abandon, **idea**) → tuple <**I**, **abandon**, **idea**>

   → verb frame (PERSON, abandon, COMMUNICATION ), (PERSON, abandon, CONGNITION)

2. Polish rescuers abandon hope.

   nsubj(abandon, **rescuers**) dobj(abandon, **hope**) → tuple <**rescuers**, **abandon**, **hope**>

   → verb frame (PERSON, abandon, FEELING), (PERSON, abandon, CONGNITION)

3. The governor abandoned the plan to build a school there.

   nsubj(abandon, **governor**) dobj(abandon, **plan**) → tuple <**governor**, **abandon**, **plan**>

   → verb frame (PERSON, abandon, COGNITION), (PERSON, abandon, ARTIFACT)

4. Workers briefly abandoned the nuclear plant.

   nsubj(abandon, **workers**) dobj(abandon, **plant**) → tuple <**worker**, **abandon**, **plant**>

   → verb frame (PERSON, abandon, ARTIFACT), (PERSON, abandon, PLANT)

5. Mary had seemed to abandon her frivolous society friends.

   nsubj(abandon, **Mary**)　　dobj(abandon, **friends**) → tuple <**Mary**, **abandon**, **friends**>

   → verb frame (PERSON, abandon, PERSON)

6. An anxious mother would abandon the attempt until all the worry has died down.

   nsubj(abandon, mother)　　dobj(abandon, attempt) → tuple <**mother**, **abandon**, **attempt**>

   → verb frame (PERSON, abandon, ACT)

Part of semantic verb frames for *abandon*:

　PERSON abandon COGNITION, 13 %

　PERSON abandon ACT, 8.8%

　PERSON abandon PERSON, 5.0%

　PERSON abandon ARTIFACT, 2.3%

Figure 1. The process of acquiring semantic verb frames from corpus sentences

4

Consider the sentence, "I finally abandon the wild idea." The arguments *I* and *idea* of the verb *abandon* can be obtained from the dependency relations provided by the Stanford Parser. By introducing the knowledge base labels, the subject (nsubj) *I* and object (dobj) *idea* are then replaced with the concept labels. After that, by applying a disambiguation process, we determine the intended verb frame of *abandon* as "PERSON abandon COGNITION."

The rest of this thesis is organized as follows. We review the related works in the next chapter. Then in Chapter 3, we present our method for combining corpus statistics and knowledge base to disambiguate and acquiring verb frames. Chapter 4 describes the experimental setting and evaluation metholodgy. In Chapter 5, we compare the quality of the verb frames, which are acquired from corpus and knowledge base, with a manually compiled gold standard, Corpus Pattern Analysis (CPA), over a set of verbs listed in CPA. Finally, we discuss the future research direction and make a conclusion of this thesis in Chapter 6.

# CHAPTER 2

# RELATED WORK

Generating verb frame (also known as *verb valency frame* or *verb patterns*) has been

an important task in the fields of computational linguistics and language learning

(Abel et al. 2003.) The techniques of verb frame generation can be applied to many

natural language processing applications, such as lexicon construction (e.g., Baker et

al., 1998; Schuler, 2005; Hlaváčková et al., 2005; Hanks and Pustejovsky, 2005;

Herbst 2004), verb clustering (Sun and Korhonen, 2009; Walde et al., 2008),

grammatical error correction (Faure and Nedellec, 1998), synonym differentiation

(Mackay, 1980; Chen and Lin, 2009) and proposition inferring (Neverilová and Grác,

2012). In our work, we focus on automatically generating verb frames with correct

syntactic argument and coarse, informative semantic categories.

Previous research in the verb frame acquisition resorts to human efforts and

requires expert linguists to manually derive patterns from a corpus (e.g., Hanks 2004.)

Hornby (1942) observes that each verb is associated with different sets of syntactic

patterns and acquires 32 verb patterns by pick out structural threads and establishing

templates for pattern analyses. The *Oxford Advanced Learner's Dictionary* (Hornby

1989) uses these verb patterns to describe a usage for each verb sense. However, to derive such patterns are painstaking, time-consuming and might not achieve high coverage for many verbs a learner has to master. In addition, dictionaries do not provide quantitative information such as how often each verb frame is used and with what semantic categories (de Marken, 1992). At the same time, they typically overly simplified semantic roles (e.g., somebody or something) for a verb frame, which may confuse a learner and lead to word choice error.

To resolve this issue, Pustejovsky et al. (2004) use semi-automsatic bootstrapping process to compile a list of verb patterns based on the Theory of Norms and Exploitations (Hanks., 2004b). The verb patterns consist of argument or valency structure, as well as semantic values for each of the arguments. For example, for the verb *abandon*, the most frequent pattern is **[Human] abandon [plan or activity]**. In this pattern, *Human* and *Plan* are semantic categories, i.e., we expect the semantic type of the subject and object to fall into these ontological categories. Learners can use verbs more precisely throughing these semantic information.

An alternative approach to verb frame extraction is based on human defined rules. Bojar (2003) defines a series of strict constraints to filter out the sentences with complex structures, and generates the verb frames by applying the rules on these simple sentences from free text (e.g., untagged corpora.)

Furthermore, Boisson et al. (2013) proposes an unsupervised approach using Latent Dirichlet Alocation (LDA) to automatically classify the verb n-grams and to acquire verb frame clusters labeled with the most frequent words in the cluster. For example, by querying the verb and its object: **cultivate $N**[3], one would obtain a cluster containing the verb frames such as *cultivate crops*, *cultivate plants* and *cultivate vegetables*, where crops, plants, and vegetables contain the same semantic category respectively.

In a study more closely related to our work, Materna (2012) proposes a probabilistic approach based on LDA to identify semantic frames from semantically unlabeled text corpora. They compare probability distribution over semantic frames and group lexical units with similar meaning. For example, for the verb *consume*, the subject including *people*, *child* and the object *food*, *meal,* and Materna derives *people consume food*.

In contrast to the previous research, we automatically derive verb arguments by utilizing grammatical relations from a parsed corpus. At the same time, we introduce a knowledge base to obtain semantic categories and disambiguate the sense of intended category of each argument using the Expectation-Maximization (EM) algorithm (Dempster et al., 1977).

---

[3] "cultivate $N" is a *Linggle* (http://linggle.com/) query used for acquiring *verb-object* clusters.

# CHAPTER 3

# METHOD

Looking up verb frames, or verb uses (e.g., construct), from dictionaries often obtain insufficient information to help a learner employing the verb precisely. Traditional dictionaries typically present broad semantic frames correspond to the verb (e.g., somebody construct something.) Unfortunately, these verb frames with general labels may confuse a learner and lead to a word choice error. A proper verb frame such as <**PERSON** *construct* **ARTIFACT**> (instead of <**PERSON** *construct* **RELATION**>) could be acquired from plenty example sentences in the dictionary. To derive an appropriate verb frame, a promising approach is to utilize dependency relations and automatically generate and disambiguate verb frames with specific semantic roles based on a lexical-semantic knowledge database.

## 3.1 Problem Statement

We focus on the issue of generating semantic verb frames: deriving sequences of semantic arguments for typical usage of a given verb. These frames are then returned to assist a language learner in learning how to use a verb correctly. The frames can be

```
(1) Extract verb-argument tuples based on grammatical relations
    acquired from a corpus.
        (Section 3.2)
(2) Obtain semantic roles for each argument in verb-argument tuples
    from a knowledge base.
        (Section 3.3)
(3) Estimate the probabilities of each semantically labeled tuples
    and generate verb frames.
        (Section 3.4)
```

Figure 2. Outline of the verb frame acquisition process of *FrameFinder*.

also by a grammatical checker in detecting and correcting word choice errors (e.g.,

Faure and Nedellec, 1998). Thus, it is crucial that the semantic roles precisely reflect

the intended meaning of arguments. At the same time, the set of derived verb frames

must be large enough to cover most usage of the given verb. Therefore, our goal is to

return a comprehensive set of verb frames with very precise semantic roles. We now

formally state the problem that we are addressing.

**Problem Statement**

We are given a general purpose corpus $C$, a knowledge base $KB$, and a verb $v$. Our

goal is to derive a set of verb frames for $v$ with semantically labeled arguments to

account for the usage of $v$. For this, we generate a tuple of semantically labeld

arguments, $<role_0, role_1, …, role_n>$ for each instance of $v$ in $C$, $<arg_0, v, arg_1, …,$

$arg_n>$ such that $role_i$ is a possible sense of $arg_i$ in *KB*, and $role_i$ most likely reflects the intended meaning of $arg_i$.

In the rest of this section, we detail our approach to this problem. First, we describe our method to extract argument tuples for each instance (Section 3.2.) This extracting process relies on dependency relations, provided by a lexicalized dependency parser (which will be described in Section 4.3). In the next section, we show how to obtain semantic roles for each verb argument (Section 3.3.) Then we disambiguate the semantic roles and obtain verb frames by iteratively disambiguating and estimating the probabilities of these argument tuples by adopting the *Expectation-Maximization* (EM) algorithm. (Section 3.4)

## 3.2 Extracting tuples of verb arguments

In the first stage of the verb frame generating process (Step (1) in Figure 2), we collect a set of grammatical relations that contain the syntactical structures information of sentences. Consider the sentences "I finally abandon the wild idea". The subject *I* and the object *idea* for the verb *abandon* can be extracted and combined to form a tuple of verb arguments, *<I, abandon, idea>*. Table 1 shows the corresponding grammatical relations provided by a dependency parser for Example 1 and 2.

Example 1. *I finally abandon the wild idea.*

Example 2. *My friends shouldn't accuse me of hiding anything.*

Example 3. *The deer would eat the plants.*

|  | Grammatical relations | Verb argument tuples |
|---|---|---|
| Example 1 | nsubj(abandon, I)[4]<br>dobj(abandon, idea) | <I, abandon, idea> |
| Example 2 | nsubj(accuse, friends)<br>dobj(accuse, me)<br>prep_of(accuse, hiding) | <friends, accuse, me, of, hiding> |
| Example 3 | nsubj(eat, deer)<br>dobj(eat, plants) | <deer, eat, plants> |

Table 1. The grammatical relations and the verb arguments for Example 1, 2 and 3.

The input to this stage is a set of sentences with their corresponding grammatical relations. As we will describe in Section 4.3, we use a lexicalized dependency parser to parse the input sentences in a corpus to obtain these grammatical relations.

The output of this stage is a set of verbal argument tuples that can be used to generate candidate semantic verb frames for the given verb.

The arguments of a verb can be extracted based on several relations including *subject*, *object* and *modifier*. Consider the sentence "My friends shouldn't accuse me of hiding anything". We can obtain the subject *friends*[5], the object *me*, the prepositional modifier *of* as well as the prepositional complement *hiding* of the verb *accuse*, and

---

[4] The annotations of grammatical relations we used are described more detail in *Appendix A*.
[5] In this paper, we only consider the head of a noun phrase as the subject or the object arguments.

yield the verbal argument tuple <*friends*, *accuse*, *me*, *of*, *hiding*> ordered by the lexical position in the original sentence.

## 3.3 Obtaining semantic roles

In the second stage of the generating process (Step (2) in Figure 2), we assign semantic roles to each verb argument extracted in Section 3.2, and generate a set of semantically labeled tuples that are verb frame candidates as input to the next stage.

Many approaches have been proposed to automatically label semantic roles. In our work, we mainly focus on obtaining semantic roles for a lexical term rather than constructing a model to generate such labels. For this, we exploit a lexical-semantic knowledge base, which contains organized lexical items into semantic categories (e.g., PERSON, ARTIFACT or COGNITION), to facilitate the process of labeling verbal arguments.

With such knowledge base, we can efficiently obtain semantic roles of an argument. For example, consider the verbal argument tuple <*friends, accuse, me, of, hiding*> extracted from Example 2. By looking up the knowledge base, we can easily recognize the argument *friends* as representing a PERSON and *hiding* as representing an ACT. At the same time, we exploit the semantic information of pronouns such as *you*, *she* and *me*. In this example, the argument *me* is unambiguously recognized as

PERSON, and finally yield a semantically labeled tuple <PERSON, *accuse,* PERSON, *of*, ACT>.

In some cases, a lexical term is ambiguous and has many senses, (i.e., belonging to more than one category in a knowledge base.) Consider another verbal argument tuple <*deer, eat, plants*> extracted from Example 3. We obtain unambiguously ANIMAL for *deer*, and PLANT (plant life) and ARTIFACT (buildings) for *plants*. Thus, we obtain two different semantic tuples, <ANIMAL, *eat*, PLANT> and <ANIMAL, *eat*, ARTIFACT>. Figure 3 below illustrates the process for two example sentences.

To a human, it is obvious that the intended meaning of *plants* is PLANT rather than ARTIFACT in the context of *deer* and *eat*. We need to resolve this issue to determine the proper semantic role of arguments, i.e., <ANIMAL, *eat*, PLANT> for <*deer, eat, plants*> to acquire the precise semantic verb frame.

Input*:*    <deer, eat, plants>

Knowledge Base:
    deer: ANIMAL
    plants: PLANT, ARTIFACT

Output:
    <ANIMAL, *eat*, PLANT>
    <ANIMAL, *eat*, ARTIFACT>

Figure 3. The process of obtaining semantic categories for each argument in <*deer, eat, plants*>

## 3.4 Generating verb frames

Once all possible semantic labels for each argument are obtained, we then

disambiguate and label arguments with semantic roles. The disambiguation procedure

is shown in Figure 4.

In Step (1a), (1b) and (1c), since our approach is based on the EM algorithm, we

initialize the probability of an assignment, $p(frame|v)$, given the semantic tuples and

the verb uniformly. Table 2 shows an example containing semantic tuples and

corresponding initial probability for the verb *eat*. Step (2) starts with reseting

$count(frame|v)$, the occurrence of each semantic tuples (Step (2a)). Then in Step 2(b),

the system get all semantic tuples according to the argument tuple, and we collect the

normalized probabilities $prob(frame)/nf$ and add to the count for each semantic tuple

(Step 2(c) and 2(d)).

15

```
Procedure GenerateVerbFrames(ArgumentTuples)

(1a)   FrameCandis = {}

(1b)   allFrames = {getVerbFrames(t) for each tuple t in ArgumentTuples}

(1c)   p(frame|v) = 1 / |allFrames|

       while not convergent

(2a)     count(frame|v) = 0 for frame in allFrames

           for each tuple t in ArgumentTuples

(2b)         tupleFrame = getVerbFrames(t)

(2c)         nf = sum(prob(frame) for all frame in tupleFrames)

               for each frame in tupleFrames

(2d)             count(frame) += prob(frame)/nf

           for all frame in allFrames

(3)          p(frame|v) = count(frame|v)/|ArgumentTuples|

       for each tuple t in ArgumentTuples

(4)        append argmax(frame) to FrameCandis

(5)    RankedFrames = Sort(frames in FrameCandis in decreasing order
                           of frequency)

       return top K ranking RankedFrames
```

Figure 4. the procedure of generating verb frames from argument tuples

| Argument tuple | Semantic tuple | Probablity |
|---|---|---|
| deer, eat, plants | < ANIMAL, *eat*, PLANT> | 0.25 |
| deer, eat, plants | < ANIMAL, *eat*, ARTIFACT> | 0.25 |
| birds, eat, seeds | < ANIMAL, *eat*, PLANT> | 0.25 |
| Birds, eat, seeds | < ANIMAL, *eat*, PERSON> | 0.25 |
| ant, eat, sugar | < ANIMAL, *eat*, FOOD> | 0.25 |

Table 2. Samples of semantic tuple with uniformly initial probabilities.

| Semantic tuple | Probablity |
|---|---|
| < ANIMAL, *eat*, PLANT> | 0.89 |
| < ANIMAL, *eat*, FOOD> | 1.0 |
| < ANIMAL, *eat*, ARTIFACT> | 0.11 |
| < ANIMAL, *eat*, PERSON> | 0.11 |

Table 3. The samples of semantic tuples with corresponding probability at the converge state of EM.

| Semantic tuple | Count |
|---|---|
| < ANIMAL, *eat*, PLANT> | 2 |
| < ANIMAL, *eat*, FOOD> | 1 |
| < ANIMAL, *eat*, ARTIFACT> | 0 |
| < ANIMAL, *eat*, PERSON> | 0 |

Table 4. The samples of semantic tuples with corresponing count in the final step.

Next, in Step (3), by normalizing these counts, we can estimate a new probability distribution $p(frame|v)$ iteratively until the output converged. For each argument tuple, we assign the corresponging semantic tuple with highest probability while the EM process is converged. (Step (4)) Finally, in Step (5), the system counts the semantic tuples and outputs the verb frames ordered by their frequencies. For example, in Table 3, **ANIMAL** *eat* **PLANTS** is the verb frame with highest probability for the argument tuple <deer, *eat,* plant>. Table 4 shows the count of each semantic tuple by aggregating the argument tuples with their corresponding semantic tuple which has the highest probability, and the system outputs **ANIMAL** eat **PLANTS** and **ANIMAL** eat **FOOD**.

# CHAPTER 4

# EXPERIMENTAL SETTING

*FrameFinder* was designed to acquire for a given verb semantic verb frames with arguments and semantic labels from a corpus. As such, *FrameFinder* is developed using a parser and external semantic knowledge base. Furthermore, since the goal of *FrameFinder* is to derive a comprehensive set of verb frames with informative semantic labels for language learning, we evaluate *FrameFinder* on the frame level rather than the argument level. Finally, as it is very time-consuming and error-prone to prepare answer keys for evaluation, we use an existing source of verbs and verb frames developed under the *Corpus Pattern Analysis Project* to evaluate *FrameFinder*. In this section, we first compare two verb frame generation methods (Section 4.1). Then, Section 4.2 introduces the evaluation metrics for the performance of *FrameFinder*, and details of the corpus statistic, tools we used and relevance judgement are reported in Section 4.3.

## 4.1 Verb frame Generation Methods Compared

*FrameFinder* starts with a given verb, the parsed sentences contain the verb with grammatical relations. The output of *FrameFinder* is a set of semantic verb frames with arguments and semantic labels, which can be shown to language learners to teach how to use a verb, or used by a grammar checker to correct grammatical errors.

Recall that we use a knowledge base to provide probable semantic categories for each verb argument. Some arguments might be ambiguous and have more than one semantic category to choose from. There are many ways to disambiguate and identify the intended meaning of an argument. We focus on evaluating the proposed method with two alternative disambiguation modules:

— Most Frequent Sense (**MFS**): This module always assigns the semantic category related to the most frequent sense to an argument.

— Expectation-Maximization (**EM**): This module considers probabilities of semantic categories for all arguments of the given verb and disambiguates by iteratively re-estimating the probability.

Additionally, we assess the performance of the models by the percentage with which the generated frames cover the most frequent usage of the given verb, or more specifically, the percentage of the most frequent verb usage instances covered.

## 4.2 Evaluation metric

Information extraction systems are usually evaluated based on the quality and completeness of the extracted information, traditionally using two metrics: *precision* and *recall* (Salton, 1989). However, the common version of the recall metrics is not well suited to evaluate a system like *FrameFinder*, where we aim to acquire the commonly used verb frames for the average language learners. To evaluate *FrameFinder*, we take into account that the extracted frames for a given verb should be valid and cover common usages of the verb. Therefore, the performance of our model is evaluated using the three metrics: precision, weighted recall, and F-measure. Each verb frame extracted by the model is judged to be either correct or incorrect against a gold standard. Precision is calculated as the fraction of correct verb frames among the frames extracted, and weighted recall measures the sum of coverage percentages of the gold-standard verb frames extracted by the system. The definitions of the metrics are described below:

*Definition 4.1*. The precision for a given verb $v$ is the percentage of verb frames that are correct in arguments as well in semantic labels among the verb frames returned by the model for the verb $v$.

*Example* 4.1. Consider a given verb $v$ and the 5 verb frames extracted by the model for this verb. If we 3 of these 5 verb frames are correct (i.e., in the gold standard for $v$) then the precision for $v$ is $3/5 = 0.6$.

*Definition 4.2.* The weighted recall for a given verb *v* is the accumulated percentage of gold-standard verb frames which the model extracts for the verb *v*.

*Example* 4.2. Consider a given verb *v* with 6 verb frames in the gold standard with the coverage decreasing percentage of .30, .25, .20, .15, .08, and .02 respectively. If the first, third, and fourth gold-standard frames match the 5 verb frames extracted by the model. Then, the weighted recall for *v* is .30 + .20 + .15 = 0.65.

*Definition 4.3.* The F-measure of the model for a verb *v* is defined as 2 PR / (P+R), where P = precision and R = weighted recall.

*Example* 4.3. Consider the result of a verb *v* in Example 4.1 and 4.2. Then, the F-measure = $2 \times 0.6 \times 0.65 / (0.6+0.65) = 0.62$

## 4.3 Evaluation verb frame and Relevance judgments

We evaluate the performance of the methods described in Section 4.1 using some evaluation metrics (Section 4.2). We draw from the Princeton *WordNet* (Miller et al., 1990) to provide semantic categories for argument. *WordNet* is the most widely used lexical semantic knowledge database and has been used as a word sense inventory for word sense disambiguation research. We used the 25 unique noun beginners (sometimes called *supersense* or lexical file names) (see Table 5), as the semantic

| ACT | FOOD | PROCESS |
|---|---|---|
| ANIMAL | GROUP | QUANTITY |
| ARTIFACT | LOCATION | RELATION |
| ATTRIBUTE | MOTIVATION | SHAPE |
| BODY | OBJECT | STATE |
| COGNITION | NATURAL | SUBSTANCE |
| COMMUNICATION | PERSON | TIME |
| EVENT | PLANT | |
| FEELING | POSSESSION | |

Table 5. The lexical file classes in WordNet[6]

| Generic relations | Annotation | Generic relations |
|---|---|---|
| subject | nsubj | normal subject |
| | xsubj | controlling subject |
| object | dobj | direct object |
| | iobj | indirect object |
| | pobj | object of a preposition |
| modifiers | prt | phrasal verb particle |
| | prep | prepositional modifier |
| | prepc | prepositional clausal modifier |

Table 6 The types of dependencies we focused in Stanford Dependency

labels. The portion of the British National Corpus (BNC) from which we extract verb

argument tuples are filtered using the following criteria:

*Criterion* 1:  The sentence must contain a subject.

*Criterion* 2:  The verb in question has no highly ambiguous arguments.

*Criterion* 3:  The sentence contains verbs, which is not the verb, *to be*

---

6  More details about lexical files are shown in *Appendix* B

After filtering, we obtained 4,693,767 sentences which were then parsed using *Stanford Dependency Parser* (Marneffe et al., 2006) to produce a list of grammatical relations (see Table 6) from which we form an argument tuple for each verb.

To evaluate the extracted verb frames, we tested the 812 verbs and 3,100 verb argument patterns developed under *Corpus Pattern Analysis Project*[7]. We converted the concepts in CPA semantic patterns to WordNet supersenses and derived gold standard for evaluating semantic verb frames. Table 7 shows a sample of the verbs and semantic verb frames. To convert CPA concepts (such as HUMAN, INSTITUTION, PLAN, ACTIVITY), we asked an English lecturer to map each concept to one or more WordNet synset[8], and expand the alternative concepts in a CPA patterns (e.g., HUMAN|INSTITUTION abandon PLAN|ACTIVITY) into multiple verb frames without alternative semantic categories (e.g., <**PERSON** *abandon* **COGNITION**> and <**PERSON** *abandon* **ACT**>) After the mapping and expansion, we obtained 13,540 verb frames. Table 8 shows sample verb frames after the conversion. From the parsed sentences, we extracted 1,007,510 verb argument tuples for the 812 verbs listed in CPA. For each of these verbs, verb frames from argument tuples were derived and evaluated against gold standard. An extracted verb frame is judged as "correct" if and only if it is present in the gold standard based on CPA patterns.

---

[7] http://deb.fi.muni.cz/pdev/
[8] The CPA-WordNet mapping table is listed in *Appendix* C

| Verb | Sample of Lexical Patterns |
|---|---|
| abandon | [[Human | Institution]] abandon [[Activity | Plan]] |
| | [[Human | Institution]] abandon [[Attitude]] |
| | [[Human]] abandon [[Artifact]] |
| acquire | [[Human | Institution]] acquire [[Asset | {Physical Object = Valuable}]] |
| | [[Human | Institution | Animal | Entity]] acquire [[Property]] |
| | [[Human | Animal]] acquire [[Information = Knowledge | Skill]] |

Table 7. Sample verbs and lexical patterns from the CPA

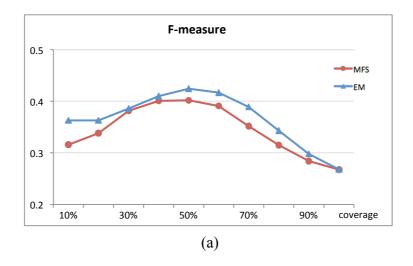| Original Pattern | Transformed Patterns |
|---|---|
| [[Human | Institution]] abandon [[Activity | Plan]] | [PERSON] abandon [ACT] |
| | [GROUP] abandon [ACT] |
| | [PERSON] abandon [COGNITION] |
| | [GROUP] abandon [COGNITION] |

Table 8. Sample of transforming CPA patterns into WordNet format

# CHAPTER 5

## EVALUATION RESULTS

In this chapter, we report the results of the experimental evaluation on two systems, **MFS** (most frequent sense) and **EM** (Expectation and Maximization), using the methodology described in the previous section.

During the evaluation, 812 verbs were automatically evaluated based on the manually compiled CPA-WordNet mapping rules. Recall that we will evaluate precision at different coverage levels. Figure 5(a) shows the average F-measure for the **MFS** and **EM** over the 812 verbs for varying level of coverage on most frequent usage (MFU). The **EM** system achieved higher F-score overall. This shows that it is somewhat effective to perform sense disambiguation in very limited context (some WSD disambiguatation algorithms consider wider context of 100-word windows). The average values of F-measure are around 0.42 when covering 50% of most frequent usage.
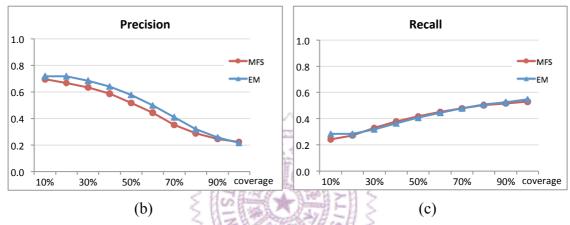
(a)



(b)



(c)

Figure 5. The average f-score, precision and recall at each coverage range for **EM** and **MFS**
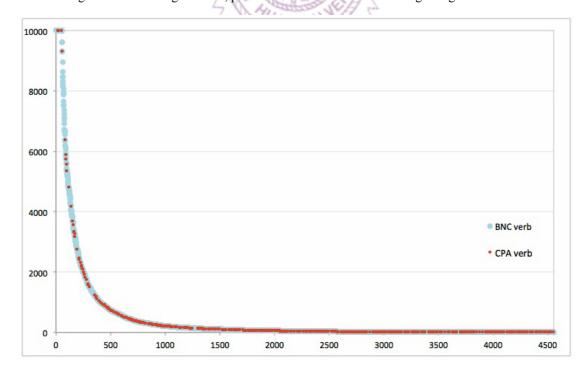


Figure 6. Number of occurrences of verbs in BNC and CPA

26

Figures 5(b) and 5(c) shows the average precision and weighted recall rates for varying level of coverage on MFU. Again, the **EM** system had the higher precision and weighted recall at all levels of coverage than the **MFS** system. The highest average precision is around 60% when covering top 50% of MFU for all verbs evaluated.

However, recall that our goal is to provide verb frames to facilitate language learning. With that in mind, some verbs might be more important for our purpose, and the precision rate for specific verbs (e.g., high-frequency verbs in General Service List (West, 1953)) is more relevant. In order to investigate how the system performs for different verbs (especially high frequency verbs not listed in the CPA). Hence, we investigated the verb and patterns in our gold standard, i.e., CPA, especially how the frequency of a given verb affects the precision rate of the systems. Thus, we compared the frequency distributions of 812 verbs in CPA and top 4,500 verbs obtained from BNC. Figure 6 shows that CPA verbs tends to concentrated in the mid to low frequency range. Note that we only plotted the verbs that occur more than 10 times, and group verbs that occur more than 10,000 to fit every thing in the figure.

The low frequency values of verbs and small numbers of instances used in *FrameFinder* which might negatively effect the precision.

| Group | Verb count criteria | # of verb in CPA |
|-------|---------------------|------------------|
| VH | > 6000 | 26 |
| H | 3000 – 6000 | 50 |
| M | 500 – 3000 | 85 |
| L | 150 – 500 | 77 |
| VL | < 150 | 575 |

Table 9. Five grouping criteria and the count of verb in CPA.

As we can see in Figure 6, very few common verbs that are important to language learning are present in the *CPA* list of verbs. To investigate how the systems perform for high and low-frequency verbs, we divided 812 verbs into five groups: **VH**, **H**, **M**, **L** and **VL**. Table 9 shows the range of counts and the number of CPA verb for each frequency group. We then calculate the precision rate for each of these group. Figure 7 shows the precision for all groups at 10% coverage of MFU. Figure 7(a) and 7(b) show the proposed systems achieved near 100% precision for almost all verbs in the **VH** and **H** groups, which contain more important verbs for intermediate to advanced language learners. On the other hand, Figure 7(e) shows that *FrameFinder* have much lower precision rate for verbs in the **VL** group.
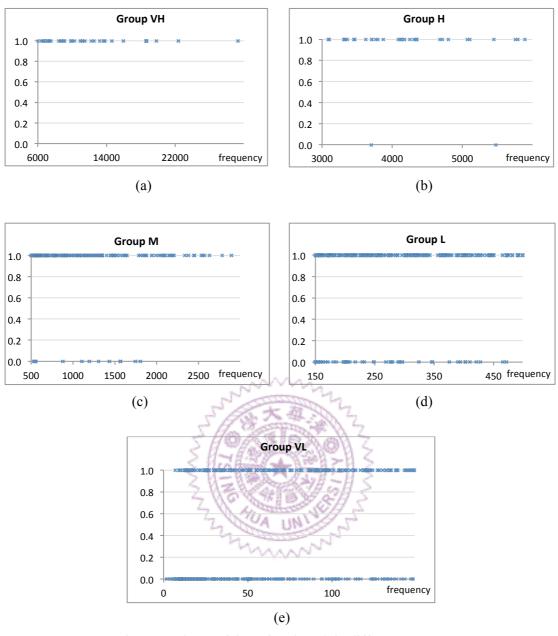
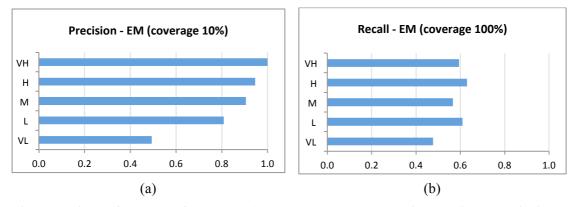Figure 7. The precision of each verb in different groups.



Figure 8. The performance of *FrameFinder* at coverage 10% (a) and 100% (b) respectively.

Figure 8 summarized the performance for each group at two coverage level of 10% and 100% respectively. From the figure, we can conclude that the systems perform much better for very high and high-frequency verbs covering the top 10% of the most frequent usage of verbs. Note that very infrequent usage of verbs is seldom useful for for intermediate to advanced language learners.

We also conducted error analysis of this automatic evaluation using a gold standard not constructed specifically for this research. After examine closely some of failed cases, we found that in many cases the failure is not due to the effectiveness of the systems, but might probably caused by **semantic category inconsistency** and **optional adverbial**.

Recall that we automatically transform CPA semantic categories to WordNet annotations. Consider the argument tuple *<Iraq, abandon, program>* and its relevant verb frame **<LOCATION** *abandon* **COMMUNICATION>**. In this case, we obtained the unique semantic category for *Iraq* from WordNet is **LOCATION**, which stands for a country or a particular place. However, CPA uses different semantic categories to distinguish a country and a place, i.e., **GROUP** and **LOCATION**. In CPA, the best pattern to describe this argument tuple is **<GROUP** *abandon* **COMMUNICATION>**. As a consequence, the extracted verb frame was labeled as *incorrect* in evaluation. As a result of semantic categories inconsistency, the system

performance could to some extent underestimated based on the automatically generated transformation. In the future, a human cross validation is necessary for assessing the actual performance.

Additionally, we observed that some verb frames only partially matched the gold standard. That is, in some cases, *FrameFinder* generated verb frames with an optional preposition. For example, we obtained the verb frame <**PERSON** *abominate* **COMMUNICATION** *as* **ARTIFACT**>, whereas the gold standard is <**PERSON** *abominate* **COMMUNICATION**> not containing the "as **ARTIFACT**" slot. In the literature of linguistic student, such prepositional phrase or adverbial are optional and is sometimes called *outer case* to distinguish it from obligatory inner cases (e.g., the prepositional phrase, "on the table" for the verb "put"). The outer-case argument show not be part of a verb frame, leading to lower precision for the proposed systems that do not exclude outer case in the first stage. We could solve this issue considering the strength of association between the prepositional phrase (or adverbial) and the verb in question.

These two issues seriously affect the performance of *FrameFinder* especially for the low-freqency verbs in the **VL** group, but have negligible effect on the high-frequency verbs in the **VH** and **H** groups. As indicated in Figure 8(a),

*FrameFinder* could achieve near perfect precision for frequent verbs (100 for **VH** verbs and 95% **H** verbs).

Thus, we have reason to believe that the performance would be greatly improved by extension of the systems in terms of training data and algorithm. More specifically, we envision that with a much larger corpus than the BNC which we have used in this study, the precision and weighted recall rates to increase scientifically. Further improvement is very likely with outer-case arguments excluded using collocation extraction methods (e.g., Smadja 1993), mutual information (Hank and Church, 1990), and other statistical tests such as log likelihood ratio (Dunning, 1993).

# CHAPTER 6

## FUTURE WORK AND SUMMARY

Many avenues exist for future research and improvement of our system. For example, existing web corpora, e.g., UKWaC with more than 2 billion words and ClueWeb09 with 503,903,810 English web pages, more than 70 billion words, can be exploited. More sophisticated word sense disambiguation methods could be employed. More grammatical relations (especially the passive constructs) should also be used. Additionally, an interesting direction to explore is using the extracted verb frames as selectional preference for verb clustering, or to interpret noun-noun sequences (Vanderwende 1994). Yet another direction of application would be to use the verb frame information for learning or teaching the difference of verb usage. For example, the derived verb frames could help identify prefered contexts of near-synonym verbs, which has been a challenging task for language learners.

In summary, we have introduced a method for generating verb frames automatically. The method involves extracting verb-argument tuples based on grammatical relations acquired from a parsed corpus, obtaining probable semantic categories for each argument in verb-argument tuples from a lexical semantic

knowledge base, estimating the probabilities of each semantically labeled tuples, and

generating verb frames. Evaluation on a set of verbs listed in a manually compile verb

patterns shows that the method is able to extract verb frames with reasonable

precision and recall of the most frequent usages.

# REFERENCES

Abel, A., Gamper, J., Knapp, J., & Weber, V. (2003). Describing Verb Valency in an Electronic Learner's Dictionary: Linguistic and Technical Implications. In World Conference on Educational Multimedia, Hypermedia and Telecommunications (Vol. 2003, No. 1, pp. 1202-1209).

z

Baker, C. F., Fillmore, C. J., & Lowe, J. B. (1998, August). The berkeley framenet project. In Proceedings of the 17th international conference on Computational linguistics-Volume 1 (pp. 86-90). Association for Computational Linguistics.

Boisson, J., Kao, T.H., Wu, J.C. Yen, T.H., and Chang, J. S. (2010). Linggle: a Web-scale Linguistic Search Engine for Words in Context. Association for Computational Linguistics System Demo 2013.

Bojar, O., Semecký, J., & Benesová, V. (2005). VALEVAL: Testing Vallex Consistency and Experimenting with Word-Frame Disambiguation. Prague Bull. Math. Linguistics, 83, 5-18.

Chen, M. C., & Lin H. (2009). Self-efficacy, foreign language anxiety as predictors of academic performance among professional program students in a general English proficiency writing test. Perceptual and Motor Skills, 2009(109), 420-430.

Chen, M. H., Huang, C. C., Huang, S. T., & Chang, J. S. (2010). GRASP: Grammar-and Syntax-based Pattern-Finder for Collocation and Phrase Learning. In PACLIC (pp. 357-364).

Church, K. W., & Hanks, P. (1990). Word association norms, mutual information, and lexicography. Computational linguistics, 16(1), 22-29.

De Mareken, C.G. "Parsing the LOB Corpus". In Proceedings of the 28th Annual Meeting of the ACL, 1990, pp. 243-251.

De Marneffe, M. C., MacCartney, B., & Manning, C. D. (2006, May). Generating typed dependency parses from phrase structure parses. In Proceedings of LREC (Vol. 6, pp. 449-454).

Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. Journal of the Royal Statistical Society. Series B (Methodological), 1-38.

Dunning, T. (1993). Accurate methods for the statistics of surprise and coincidence. Computational linguistics, 19(1), 61-74.

Faure, D., & Nédellec, C. (1998, May). A corpus-based conceptual clustering method for verb frames and ontology acquisition. In LREC workshop on adapting lexical and corpus resources to sublanguages and applications (Vol. 707, No. 728, p. 30).

Hanks, P. (2004a). Corpus pattern analysis. In Euralex Proceedings (Vol. 1, pp. 87-98).

Hanks, P. (2004b). The syntagmatics of metaphor and idiom. International Journal of Lexicography, 17(3), 245-274.

Hanks, P., & Pustejovsky, J. (2005). A pattern dictionary for natural language processing. Revue Française de linguistique appliquée, 10(2), 63-82.

Herbst, T. (Ed.). (2004). A Valency Dictionary of English: A Corpus Based Analysis of the Complementation Patterns of English Verbs, Nouns and Adjectives (Vol. 40). Walter de Gruyter.

Hlaváčková, D., & Horák, A. (2005). Verbalex–new comprehensive lexicon of verb valencies for czech. In Proceedings of the Slovko Conference, Bratislava, Slovakia.

Hornby, A. S., Gatenby, E. V., & Wakefield, H. (1942). Idiomatic and Syntactic English Dictionary.

Hornby, A. S. (ed.). Oxford Advanced Learner's Dictionary of Current English. Oxford, UK: Oxford University Press, 1989.

Im Walde, S. S., Hying, C., Scheible, C., & Schmid, H. (2008). Combining EM Training and the MDL Principle for an Automatic Verb Classification Incorporating Selectional Preferences. In ACL (pp. 496-504).

Kilgarriff, A., & Grefenstette, G. (2003). Introduction to the special issue on the web as corpus. Computational linguistics, 29(3), 333-347.

Materna, J. (2012). LDA-Frames: an unsupervised approach to generating semantic frames. In Computational Linguistics and Intelligent Text Processing (pp. 376-387). Springer Berlin Heidelberg.

McKay, S. (1980). Teaching the syntactic, semantic and pragmatic dimensions of verbs. TESOL Quarterly, 17-26.

Miller, G. A., Beckwith, R., Fellbaum, C., Gross, D., & Miller, K. J. (1990). Introduction to wordnet: An on-line lexical database. International journal of lexicography, 3(4), 235-244.

Nevěřilová, Z., & Grác, M. (2012, January). Common Sense Inference Using Verb Valency Frames. In Text, Speech and Dialogue (pp. 328-335).

Pustejovsky, J., Hanks, P., & Rumshisky, A. (2004, August). Automated induction of sense in context. In Proceedings of the 20th international conference on Computational Linguistics (p. 924).

Salton, G. (1989). Automatic Text Processing: The Transformation, Analysis, and Retrieval of. Addison-Wesley.

Schuler, K. K. (2005). VerbNet: A broad-coverage, comprehensive verb lexicon.

Smadja, F. (1993). Retrieving collocations from text: Xtract. Computational linguistics, 19(1), 143-177.

Sun, L., & Korhonen, A. (2009, August). Improving verb clustering with automatically acquired selectional preferences. In Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 2-Volume 2 (pp. 638-647). Association for Computational Linguistics.

Vanderwende, L. (1994, August). Algorithm for automatic interpretation of noun sequences. In Proceedings of the 15th conference on Computational linguistics-Volume 2 (pp. 782-788). Association for Computational Linguistics.

West, M. P. (1953). A general service list of English words: with semantic frequencies and a supplementary word-list for the writing of popular science and technology. Longmans, Green.

Wible, D., & Tsao, N. L. (2010, June). StringNet as a computational resource for discovering and investigating linguistic constructions. In Proceedings of the NAACL HLT Workshop on Extracting and Using Constructions in Computational Linguistics (pp. 25-31). Association for Computational Linguistics.

Yarowsky, D. (1992, August). Word-sense disambiguation using statistical models of Roget's categories trained on large corpora. In Proceedings of the 14th conference on Computational linguistics-Volume 2 (pp. 454-460). Association for Computational Linguistics.

# APPENDIX A - Stanford Dependency Annotation

*dobj : direct object*

The direct object of a VP is the noun phrase which is the (accusative) object of the verb.

> "She gave me a raise"  *dobj*(gave, raise)

> "They win the lottery"  *dobj*(win, lottery)

*iobj : indirect object*

The indirect object of a VP is the noun phrase which is the (dative) object of the verb.

> "She gave me a raise"  *iobj*(gave, me)

*nsubj : nominal subject*

A nominal subject is a noun phrase which is the syntactic subject of a clause. The governor of this relation might not always be a verb: when the verb is a copular verb, the root of the clause is the complement of the copular verb, which can be an adjective or noun.

> "Clinton defeated Dole"  *nsubj*(defeated, Clinton)

> "The baby is cute"  *nsubj*(cute, baby)

*pobj: object of a preposition*

The object of a preposition is the head of a noun phrase following the preposition, or the adverbs "here" and "there". (The preposition in turn may be modifying a noun, verb, etc.) Unlike the Penn Treebank, we here define cases of VBG quasi-prepositions like "including", "concerning", etc. as instances of pobj. (The preposition can be called a FW for "pace", "versus", etc. It can also be called a CC – but we don't currently handle that and would need to distinguish from conjoined prepositions.) In the case of preposition stranding, the object can precede the preposition (e.g., "What does CPR stand for?").

"I sat on the chair"          *pobj*(on, chair)


*prep: prepositional modifier*

A prepositional modifier of a verb, adjective, or noun is any prepositional phrase that serves to modify the meaning of the verb, adjective, noun, or even another prepositon. In the collapsed representation, this is used only for prepositions with NP complements.

"I saw a cat in a hat"          *prep*(cat, in)
"I saw a cat with a telescope"  *prep*(saw, with)
"He is responsible for meals"   *prep*(responsible, for)


*prepc: prepositional clausal modifier*

In the collapsed representation, a prepositional clausal modifier of a verb, adjective, or noun is a clause introduced by a preposition which serves to modify the meaning of the verb, adjective, or noun.

"He purchased it without paying a premium" *prepc_without*(purchased, paying)


*prt: phrasal verb particle*

The phrasal verb particle relation identifies a phrasal verb, and holds between the verb and its particle.

"They shut down the station"    *prt*(shut, down)


*xsubj : controlling subject*

A controlling subject is the relation between the head of an open clausal complement (xcomp) and the external subject of that clause.

"Tom likes to eat fish"         *xsubj*(eat, Tom)

# APPENDIX B - WordNet lexicographer classes with descriptions.

| lexicographer file name | brief description |
| --- | --- |
| ACT | nouns denoting acts or actions |
| ANIMAL | nouns denoting animals |
| ARTIFACT | nouns denoting man-made objects |
| ATTRIBUTE | nouns denoting attributes of people and objects |
| BODY | nouns denoting body parts |
| COGNITION | nouns denoting cognitive processes and contents |
| COMMUNICATION | nouns denoting communicative processes and contents |
| EVENT | nouns denoting natural events |
| FEELING | nouns denoting feelings and emotions |
| FOOD | nouns denoting foods and drinks |
| GROUP | nouns denoting groupings of people or objects |
| LOCATION | nouns denoting spatial position |
| MOTIVE | nouns denoting goals |
| OBJECT | nouns denoting natural objects (not man-made) |
| PERSON | nouns denoting people |
| PHENOMENON | nouns denoting natural phenomena |
| PLANT | nouns denoting plants |
| POSSESSION | nouns denoting possession and transfer of possession |
| PROCESS | nouns denoting natural processes |
| QUANTITY | nouns denoting quantities and units of measure |
| RELATION | nouns denoting relations between people or things or ideas |
| SHAPE | nouns denoting two and three dimensional shapes |
| STATE | nouns denoting stable states of affairs |
| SUBSTANCE | nouns denoting substances |
| TIME | nouns denoting time and temporal relations |

# APPENDIX C - CPA to WordNet Transformation.

| Object | Match (Yes: 1, No: 0) | WordNet Sense | |
|---|---|---|---|
| Speech Act | 1 | (1) | \<communication\> |
| Speech Act | 1 | (2) | \<communication\> |
| Speech Act | 1 | (3) | \<communication\> |
| Horse | 1 | (1) | \<animal\> |
| Horse | 0 | (2) | \<artifact\> |
| Horse | 0 | (3) | \<group\> |
| Goal | 1 | (1) | \<cognition\> |
| Goal | 0 | (2) | \<location\> |
| Goal | 0 | (3) | \<artifact\> |
| Self | 1 | (1) | \<cognition\> |
| Self | 1 | (2) | \<person\> |
| Projectile | 1 | (1) | \<artifact\> |
| Projectile | 1 | (2) | \<artifact\> |
| Building Part | 1 | (1) | \<artifact\> |
| Building Part | 0 | (2) | \<act\> |
| Building Part | 1 | (3) | \<act\> |
| Crime | 1 | (1) | \<act\> |
| Crime | 1 | (2) | \<act\> |
| Asset | 1 | (1) | \<attribute\> |
| Vehicle | 1 | (1) | \<artifact\> |
| Vehicle | 0 | (2) | \<communication\> |
| Vehicle | 0 | (3) | \<substance\> |
| Document | 1 | (1) | \<communication\> |
| Document | 1 | (2) | \<artifact\> |
| Document | 1 | (3) | \<possession\> |
| Bird | 1 | (1) | \<animal\> |
| Bird | - | (2) | \<food\> |
| Bird | - | (3) | \<person\> |
| Document Part | 1 | (1) | \<communication\> |
| Document Part | 1 | (2) | \<artifact\> |
| Document Part | 1 | (3) | \<possession\> |
| Body | 1 | (1) | \<body\> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense |
|---|---|---|
| Body | 0 | (2) <group> |
| Body | 0 | (3) <body> |
| Offer | 1 | (1) <communication> |
| Offer | 1 | (2) <communication> |
| Offer | - | (3) <act> |
| Gas | 0 | (1) <state> |
| Gas | 1 | (2) <substance> |
| Gas | 1 | (3) <substance> |
| Broadcast | 1 | (1) <communication> |
| Broadcast | 1 | (2) <communication> |
| Glass | 1 | (1) <substance> |
| Glass | 1 | (2) <artifact> |
| Glass | 0 | (3) <quantity> |
| Plane | 1 | (1) <artifact> |
| Plane | 0 | (2) <shape> |
| Plane | 0 | (3) <state> |
| Obligation | 1 | (1) <act> |
| Obligation | 1 | (2) <state> |
| Obligation | 1 | (3) <relation> |
| Software | 1 | (1) <communication> |
| Institution | 1 | (1) <group> |
| Institution | 0 | (2) <artifact> |
| Institution | 1 | (3) <cognition> |
| Sound | 1 | (1) <attribute> |
| Sound | 1 | (2) <cognition> |
| Sound | 1 | (3) <phenomenon> |
| Blemish | 1 | (1) <attribute> |
| Drug | 1 | (1) <artifact> |
| Stuff | - | (1) <substance> |
| Stuff | 1 | (2) <artifact> |
| Stuff | 1 | (3) <possession> |
| Visible Feature | 1 | (1) <cognition> |
| Visible Feature | 1 | (2) <body> |
| Visible Feature | 1 | (3) <communication> |
| Character Trait | 1 | (1) <attribute> |
| Surface | 1 | (1) <artifact> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense |
|---|---|---|
| Surface | 0 | (2) &lt;location&gt; |
| Surface | 1 | (3) &lt;object&gt; |
| Entity | - | (1) &lt;Tops&gt; |
| Event | 1 | (1) &lt;Tops&gt; |
| Event | 1 | (2) &lt;state&gt; |
| Event | 0 | (3) &lt;phenomenon&gt; |
| Storm | 1 | (1) &lt;phenomenon&gt; |
| Storm | 1 | (2) &lt;state&gt; |
| Storm | 1 | (3) &lt;act&gt; |
| System | 1 | (1) &lt;artifact&gt; |
| System | 1 | (2) &lt;group&gt; |
| System | 1 | (3) &lt;substance&gt; |
| State | 0 | (1) &lt;location&gt; |
| State | 1 | (2) &lt;Tops&gt; |
| State | 0 | (3) &lt;group&gt; |
| Location | 1 | (1) &lt;Tops&gt; |
| Location | 0 | (2) &lt;act&gt; |
| Location | 1 | (3) &lt;act&gt; |
| Aperture | 0 | (1) &lt;artifact&gt; |
| Aperture | 0 | (2) &lt;object&gt; |
| Aperture | 1 | (3) &lt;artifact&gt; |
| TV Program | 1 | (1) &lt;communication&gt; |
| TV Program | 0 | (2) &lt;artifact&gt; |
| Dog | 1 | (1) &lt;animal&gt; |
| Dog | 0 | (2) &lt;person&gt; |
| Dog | 0 | (3) &lt;person&gt; |
| Physical Object | 1 | (1) &lt;Tops&gt; |
| Physical Object | 0 | (2) &lt;cognition&gt; |
| Physical Object | 0 | (3) &lt;communication&gt; |
| Plant Part | 1 | (1) &lt;plant&gt; |
| Plant Part | 0 | (2) &lt;shape&gt; |
| Plant Part | 0 | (3) &lt;person&gt; |
| Animal Group | - | (1) &lt;Tops&gt; |
| Meat | 1 | (1) &lt;food&gt; |
| Meat | 1 | (2) &lt;plant&gt; |
| Meat | 0 | (3) &lt;cognition&gt; |

| Object | Match (Yes: 1, No: 0) | WordNet Sense | |
|---|---|---|---|
| Picture | 1 | (1) | <artifact> |
| Picture | 1 | (2) | <artifact> |
| Picture | 0 | (3) | <cognition> |
| Medium | 1 | (1) | <communication> |
| Medium | 0 | (2) | <location> |
| Medium | 0 | (3) | <communication> |
| Relationship | 1 | (1) | <relation> |
| Relationship | 1 | (2) | <state> |
| Relationship | 1 | (3) | <state> |
| Permission | 1 | (1) | <communication> |
| Permission | 1 | (2) | <act> |
| Fire | 1 | (1) | <event> |
| Fire | 1 | (2) | <act> |
| Fire | 1 | (3) | <process> |
| Material | 1 | (1) | <substance> |
| Material | 0 | (2) | <communication> |
| Material | 1 | (3) | <artifact> |
| Agreement | 1 | (1) | <communication> |
| Agreement | 0 | (2) | <attribute> |
| Agreement | 1 | (3) | <state> |
| Eventuality | Entity | 1 | (1) | <Tops> |
| Eventuality | Entity | 1 | (2) | <state> |
| Eventuality | Entity | 1 | (3) | <phenomenon> |
| Flag | 1 | (1) | <artifact> |
| Flag | 0 | (2) | <communication> |
| Flag | 0 | (3) | <plant> |
| Part | 1 | (1) | <relation> |
| Part | 0 | (2) | <artifact> |
| Part | 1 | (3) | <object> |
| Numerical Value | 1 | (1) | <attribute> |
| Numerical Value | 1 | (2) | <quantity> |
| Numerical Value | 0 | (3) | <communication> |
| Route | 1 | (1) | <location> |
| Route | 1 | (2) | <artifact> |
| Word | 1 | (1) | <communication> |
| Word | 1 | (2) | <communication> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense |
|---|---|---|
| Word | 1 | (3) <communication> |
| Eventuality | 1 | (1) <event> |
| Particle | 1 | (1) <substance> |
| Particle | 1 | (2) <object> |
| Particle | 0 | (3) <communication> |
| Light | 1 | (1) <phenomenon> |
| Light | 1 | (2) <artifact> |
| Light | 0 | (3) <cognition> |
| Metal | 1 | (1) <substance> |
| Metal | 1 | (2) <substance> |
| Speech Sound | 1 | (1) <attribute> |
| Speech Sound | 1 | (2) <cognition> |
| Speech Sound | 1 | (3) <phenomenon> |
| Pace | 1 | (1) <time> |
| Pace | 1 | (2) <quantity> |
| Pace | 1 | (3) <attribute> |
| Command | 1 | (1) <communication> |
| Command | 0 | (2) <group> |
| Command | 1 | (3) <attribute> |
| Mathematical Value | 1 | (1) <cognition> |
| Mathematical Value | 1 | (2) <attribute> |
| Mathematical Value | 1 | (3) <possession> |
| Narrative | 1 | (1) <communication> |
| Disease | 1 | (1) <state> |
| Property | 1 | (1) <possession> |
| Property | 1 | (2) <attribute> |
| Property | 1 | (3) <location> |
| Emotion | 1 | (1) <feeling> |
| Concept | 1 | (1) <cognition> |
| Fetus | 1 | (1) <animal> |
| Weight | 1 | (1) <attribute> |
| Weight | 0 | (2) <artifact> |
| Weight | 0 | (3) <attribute> |
| Money | 0 | (1) <possession> |
| Money | 1 | (2) <possession> |
| Money | 1 | (3) <possession> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense |
|---|---|---|
| Energy | 0 | (1) &lt;phenomenon&gt; |
| Energy | 1 | (2) &lt;attribute&gt; |
| Energy | 1 | (3) &lt;attribute&gt; |
| Ceramic | 1 | (1) &lt;artifact&gt; |
| Road Vehicle | 1 | (1) &lt;artifact&gt; |
| Road Vehicle | 0 | (2) &lt;communication&gt; |
| Road Vehicle | 0 | (3) &lt;substance&gt; |
| Water Vehicle | 1 | (1) &lt;artifact&gt; |
| Water Vehicle | 0 | (2) &lt;communication&gt; |
| Water Vehicle | 0 | (3) &lt;substance&gt; |
| Body Part | 1 | (1) &lt;body&gt; |
| Body Part | 1 | (2) &lt;group&gt; |
| Body Part | 0 | (3) &lt;body&gt; |
| Computer | 1 | (1) &lt;artifact&gt; |
| Computer | 1 | (2) &lt;person&gt; |
| Money Value | 0 | (1) &lt;possession&gt; |
| Money Value | 1 | (2) &lt;possession&gt; |
| Money Value | 1 | (3) &lt;possession&gt; |
| Human | 1 | (1) &lt;animal&gt; |
| Deity | 1 | (1) &lt;person&gt; |
| Firearm | 1 | (1) &lt;artifact&gt; |
| Sense Organ | 1 | (1) &lt;body&gt; |
| Sense Organ | 0 | (2) &lt;group&gt; |
| Sense Organ | 0 | (3) &lt;artifact&gt; |
| Physical Object Part | 1 | (1) &lt;Tops&gt; |
| Physical Object Part | 0 | (2) &lt;cognition&gt; |
| Physical Object Part | 0 | (3) &lt;communication&gt; |
| Motorbike | 1 | (1) &lt;artifact&gt; |
| Animal | 1 | (1) &lt;Tops&gt; |
| Privilege | 1 | (1) &lt;attribute&gt; |
| Privilege | 1 | (2) &lt;attribute&gt; |
| Privilege | 1 | (3) &lt;attribute&gt; |
| Signal | 1 | (1) &lt;communication&gt; |
| Signal | 1 | (2) &lt;motive&gt; |
| Signal | 1 | (3) &lt;phenomenon&gt; |
| Thread | 1 | (1) &lt;artifact&gt; |

| Object | Match (Yes: 1, No: 0) | WordNet Sense |
|---|---|---|
| Thread | 1 | (2) <object> |
| Thread | 1 | (3) <cognition> |
| Cloth | 1 | (1) <artifact> |
| Plan | 1 | (1) <cognition> |
| Plan | 1 | (2) <cognition> |
| Plan | 0 | (3) <artifact> |
| Artwork | 1 | (1) <communication> |
| Device | 1 | (1) <artifact> |
| Device | 1 | (2) <communication> |
| Device | 1 | (3) <act> |
| Information Source | 1 | (1) <communication> |
| Information Source | 1 | (2) <cognition> |
| Information Source | 0 | (3) <communication> |
| Musical Instrument | 1 | (1) <artifact> |
| Musical Instrument | 0 | (2) <act> |
| Musical Instrument | 0 | (3) <person> |
| Language | 1 | (1) <communication> |
| Language | 1 | (2) <communication> |
| Language | 1 | (3) <communication> |
| Weapon | 1 | (1) <artifact> |
| Weapon | 0 | (2) <communication> |
| Human | Animal | 1 | (1) <Tops> |
| Rule | 1 | (1) <cognition> |
| Rule | 1 | (2) <cognition> |
| Rule | 1 | (3) <communication> |
| Beverage | 1 | (1) <food> |
| Flame | 0 | (1) <process> |
| Watercourse | 1 | (1) <object> |
| Watercourse | 1 | (2) <object> |
| Watercourse | 1 | (3) <artifact> |
| Action | 1 | (1) <act> |
| Action | 0 | (2) <state> |
| Action | 1 | (3) <act> |
| Quantity | 1 | (1) <Tops> |
| Quantity | 1 | (2) <attribute> |
| Quantity | 0 | (3) <cognition> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense | |
|---|---|---|---|
| Illness | 1 | (1) | <state> |
| Bicycle | 1 | (1) | <artifact> |
| Process | 1 | (1) | <act> |
| Process | 1 | (2) | <cognition> |
| Process | 1 | (3) | <communication> |
| Garment | 1 | (1) | <artifact> |
| Attitude | 1 | (1) | <cognition> |
| Attitude | 1 | (2) | <attribute> |
| Attitude | 0 | (3) | <act> |
| Wood | 1 | (1) | <substance> |
| Wood | 1 | (2) | <group> |
| Wood | 1 | (3) | <person> |
| Abstract Entity | 1 | (1) | <Tops> |
| Soil | 0 | (1) | <state> |
| Soil | 1 | (2) | <substance> |
| Soil | 1 | (3) | <object> |
| Limb | 1 | (1) | <body> |
| Limb | 0 | (2) | <plant> |
| Limb | 0 | (3) | <location> |
| Fish | 1 | (1) | <animal> |
| Fish | 0 | (2) | <food> |
| Fish | 0 | (3) | <person> |
| Information | 1 | (1) | <communication> |
| Information | 1 | (2) | <cognition> |
| Information | 1 | (3) | <communication> |
| Container | 1 | (1) | <artifact> |
| Movie | 1 | (1) | <communication> |
| Machine | 1 | (1) | <artifact> |
| Machine | 0 | (2) | <person> |
| Machine | 0 | (3) | <group> |
| Land | 1 | (1) | <possession> |
| Land | 1 | (2) | <object> |
| Land | 0 | (3) | <location> |
| Injury | 1 | (1) | <state> |
| Injury | 1 | (2) | <event> |
| Injury | 1 | (3) | <event> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense | |
|---|---|---|---|
| Vehicle Group | 1 | (1) | \<artifact\> |
| Vehicle Group | 0 | (2) | \<communication\> |
| Vehicle Group | 0 | (3) | \<substance\> |
| Hair | 1 | (1) | \<body\> |
| Hair | 0 | (2) | \<quantity\> |
| Hair | 0 | (3) | \<plant\> |
| Plant | 0 | (1) | \<artifact\> |
| Plant | 1 | (2) | \<Tops\> |
| Plant | 0 | (3) | \<person\> |
| Human | Institution | 1 | (1) | \<Tops\> |
| Human | Institution | 0 | (2) | \<substance\> |
| Human | Institution | - | (3) | \<cognition\> |
| Liquid | 1 | (1) | \<substance\> |
| Liquid | 0 | (2) | \<state\> |
| Liquid | 1 | (3) | \<substance\> |
| Power | 1 | (1) | \<attribute\> |
| Power | 0 | (2) | \<phenomenon\> |
| Power | 1 | (3) | \<cognition\> |
| Musical Performance | 0 | (1) | \<communication\> |
| Musical Performance | 1 | (2) | \<act\> |
| Musical Performance | 1 | (3) | \<act\> |
| Food | 1 | (1) | \<Tops\> |
| Food | 1 | (2) | \<food\> |
| Food | 0 | (3) | \<cognition\> |
| Abstract | 1 | (1) | \<cognition\> |
| Abstract | 1 | (2) | \<communication\> |
| Artifact | 1 | (1) | \<Tops\> |
| Proposition | 1 | (1) | \<communication\> |
| Proposition | 1 | (2) | \<communication\> |
| Proposition | 0 | (3) | \<communication\> |
| Building | 1 | (1) | \<artifact\> |
| Building | 1 | (2) | \<act\> |
| Building | 1 | (3) | \<act\> |
| Opportunity | 1 | (1) | \<state\> |
| Time Period | 0 | (1) | \<event\> |
| Time Period | 1 | (2) | \<time\> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense | |
|---|---|---|---|
| Time Period | 1 | (3) | <time> |
| Structure | 0 | (1) | <artifact> |
| Structure | 1 | (2) | <attribute> |
| Structure | 1 | (3) | <cognition> |
| Vapour | 1 | (1) | <substance> |
| Vapour | 0 | (2) | <process> |
| Language Part | 1 | (1) | <communication> |
| Language Part | 1 | (2) | <communication> |
| Language Part | 1 | (3) | <communication> |
| Resource | 1 | (1) | <possession> |
| Resource | 1 | (2) | <attribute> |
| Resource | 1 | (3) | <cognition> |
| Room | 1 | (1) | <artifact> |
| Room | 0 | (2) | <quantity> |
| Room | 0 | (3) | <state> |
| Music Part | 1 | (1) | <communication> |
| Music Part | 1 | (2) | <cognition> |
| Music Part | 1 | (3) | <act> |
| Request | 1 | (1) | <communication> |
| Request | 1 | (2) | <communication> |
| Human Group | 1 | (1) | <group> |
| Human Group | 1 | (2) | <group> |
| Human Group | 0 | (3) | <group> |
| Human Role | 1 | (1) | <group> |
| Human Role | 1 | (2) | <group> |
| Human Role | 0 | (3) | <group> |
| Activity | 1 | (1) | <act> |
| Activity | 0 | (2) | <state> |
| Activity | 0 | (3) | <process> |
| State Of Affairs | 1 | (1) | <cognition> |
| State Of Affairs | 0 | (2) | <state> |
| State Of Affairs | 1 | (3) | <event> |
| Chord | 0 | (1) | <shape> |
| Chord | 1 | (2) | <communication> |
| Alcoholic Drink | 1 | (1) | <food> |
| Alcoholic Drink | 1 | (2) | <substance> |

| Object | Match (Yes: 1, No: 0) | WordNet Sense |
| --- | --- | --- |
| Reputation | 1 | (1) &lt;state&gt; |
| Reputation | 1 | (2) &lt;state&gt; |
| Reputation | 1 | (3) &lt;cognition&gt; |
| Phrase | 1 | (1) &lt;communication&gt; |
| Phrase | 0 | (2) &lt;communication&gt; |
| Phrase | 0 | (3) &lt;communication&gt; |
| Waterway | 0 | (1) &lt;object&gt; |
| Waterway | 1 | (2) &lt;artifact&gt; |
| Light Source | 1 | (1) &lt;phenomenon&gt; |
| Light Source | 1 | (2) &lt;artifact&gt; |
| Light Source | 1 | (3) &lt;cognition&gt; |

γ