# FrameNet: A Knowledge Base for Natural Language Processing

**Collin F. Baker**
International Computer Science Institute
1947 Center St., Suite 600
Berkeley, California 94704 U.S.A.
collinb@icsi.berkeley.edu

## Abstract

Prof. Charles J. Fillmore had a life-long interest in lexical semantics, and this culminated in the latter part of his life in a major research project, the FrameNet Project at the International Computer Science Institute in Berkeley, California (http://framenet.icsi.berkeley.edu). This paper reports on the background of this ongoing project, its connections to Fillmore's other research interests, and briefly outlines applications and current directions of growth for FrameNet, including FrameNets in languages other than English.

## 1 Introduction

It was my honor to work closely with the late Charles Fillmore as part of the FrameNet project at the International Computer Science Institute in Berkeley, California (http://framenet.icsi.berkeley.edu) from 1997 until this year. It was a blessing to be in contact with that rare combination of a brilliant intellect, a compassionate heart, and genuine humility. This article will discuss where FrameNet fits in the development of Fillmore's major theoretical contributions (case grammar, frame semantics and construction grammar), how FrameNet can be used for NLP, and where the project is headed.

## 2 From Case Grammar to Frame Semantics to FrameNet

The beginnings of case grammar were contemporary with the development of what came to be called the "Standard Theory" of Generative Grammar (Chomsky, 1965), and related "through friendship" to the simultaneous development of Generative Semantics. Fillmore (1968) showed that a limited number of case roles could provide elegant explanations of quite varied linguistic phenomena, such as the differences in morphological case marking between nominative-accusative, nominative-ergative, and active-inactive languages, and anaphoric processes such as subject drop in Japanese. A year later (Fillmore, 1969), after explaining that verbs like *rob* and *steal* require three arguments, the culprit, the loser, and the loot, he continues in the next section to say

> It seems to me, however, that this sort of detail is unnecessary, and that what we need are abstractions from these specific role descriptions, abstractions which will allow us to recognize that certain elementary role notions recur in many situations,...Thus we can identify the culprit of rob and the critic of criticize with the more abstract role of Agent...in general...the roles that [predicates'] arguments play are taken from an inventory of role types fixed by grammatical theory.

But the search for the "correct" minimal set of case roles proved to be difficult and contentious, and it became apparent that some predicators, such as *replace* and *resemble*, required roles which did not fit into the usual categories. In fact, the original case roles (a.k.a. semantic roles, thematic roles, theta roles) were increasingly seen as generalizations over a much larger set of roles which provide more detailed information about the participants in a large variety of situations, described as **semantic frames** (Fillmore, 1976; Fillmore, 1977b).

Thus, the formulation of Frame Semantics should not be seen as a repudiation of the concept of case roles expounded in Fillmore 1968, but rather a recognition of the inadequacy of case roles as a characterization of all the different types of

interactions of participants that can be linguistically significant in using language to describe situations:

> ...[A]s I have conceived them, the repertory of cases is NOT identical to the full set of notions that would be needed to make an analysis of any state or event. ...[A] case frame need not comprise a complete description of all the relevant aspects of a situation, but only a particular piece or section of a situation. (Fillmore (1977a), emphasis in the original)

The concept of frames became part of the academic zeitgeist of the 1960s and 70s. Roger Shank was using the term **script** to talk about situations like eating in a restaurant (Schank and Abelson, 1977) and the term *frame* was being used in a more-or-less similar sense by Marvin Minsky (1974), and Eugene Charniak (1977).

**FrameNet as an Implementation of Frame Semantics**

During the late 1980s and early 1990s, much of Fillmore's effort went into joint work with Paul Kay, Catherine O'Connor, and others on the development of Construction Grammar, especially on linking constructions in which the semantic attributes of various constituents were represented by thematic roles such as Agent, Patient, Experiencer, Stimulus, etc., (cf. Levin (1993)). But semantic frames were always presupposed in Fillmore's discussion of Construction Grammar (e.g. Kay and Fillmore (1999)), just as Construction Grammar was always presupposed in discussions of Frame Semantics. In fact, some of the incidental references to semantic frames in the literature on construction grammar imply the existence of very sophisticated frame semantics. At the same time, Fillmore was becoming involved with the lexicographer Sue Atkins, and increasingly thinking about what the dictionary would look like, if freed from the limitations of publishing on paper (Fillmore and Atkins, 1994) and based on corpus data.

The FrameNet Project (Fillmore and Baker, 2010; Ruppenhofer et al., 2010a) at the International Computer Science Institute was launched in 1997, as an effort to produce a lexicon of English that is both human- and machine-readable, based on the theory of Frame Semantics and supported by annotating corpus examples of the lexical items. In part, FrameNet (FN) can be thought of as the implementation of a theory that was already well-developed, but, like other annotation projects, we have found that the process of annotating actual text has also pushed the development of the theory.

So what is a frame? Ruppenhofer et al. (2006) define a frame as "a script-like conceptual structure that describes a particular type of situation, object, or event along with its participants and props." Frames are generalizations over groups of words which describe similar states of affairs and which could be expected to share similar sets of roles, and (to some extent) similar syntactic patterns for them. In the terminology of Frame Semantics, the roles are called frame elements (FEs), and the words which evoke the frame are referred to as lexical units (LUs). A lexical unit is thus a Saussurian "sign", an association between a form and a meaning; the form is a lemma with a given part of speech, the meaning is represented as a semantic frame plus a short dictionary-style definition, which is intended to differentiate this lexical unit from others in the same frame. Each lexical unit is equivalent to a word sense; if a lemma has more than one sense, it will be linked to more than one LU in more than one frame; e.g. the lemma *run*.v (and all its word forms, *run, ran*, and *running*) is linked to several frames (**Self-motion, Operating a system**, etc.).

Some of this literature refers to two types of entities, **frames** and **scenes** (Fillmore, 1977c). However, early in the process of defining the FN data structure, it was recognized that more than two levels of generality might be needed, so it was decided to create only one type of data object, called a frame, and to define relations between frames at various levels of generality. Therefore, the term scene is not used in FrameNet today, although some frames which define complex events have the term **scenario** as part of their names, such as the **Employer's scenario**, with subframes **Hiring, Employing** and **Firing**.

In many cases, the framal distinctions proposed by Fillmore in early work are directly reflected in current FN frames, as in the pair of frames **Stinginess** and **Thriftiness**, discussed in Fillmore (1985). In other cases, the frame divisions in FN differ from those originally proposed, as in

the division of the original Commerce frame into three frames, **Commerce, Commerce_buy** and **Commerce_sell**, which are connected by frame-to-frame relations.

Because Frame Semantics began in the study of verbs and valences, there was emphasis initially on representing events, but the principle that a conceptual gestalt can be evoked by any member of a set of words also applies to relations, states, and entities, and the evoking words can be nouns, adjectives, adverbs, etc., as well as verbs. For example, the **Leadership** frame contains both nouns (*leader, headmaster*, *maharaja*), and verbs (*lead*, *command*); FEs in the **Leadership** frame include the LEADER and the GOVERNED, as in [LEADER Kurt Helborg] is the CAPTAIN [GOVERNED of the Reiksguard Knights].

## 3 Applications of FrameNet

Underlying other applications is the need for middle-ware to carry out automatic semantic role labeling (ASRL). Beginning with the work of Gildea and Jurafsky (2000; 2002), many researchers have built ASRL systems trained on the FrameNet data (Erk and Padó, 2006; Johansson and Nugues, 2007; Das et al., 2013), some of which are freely available. Other groups have built software to suggest new LUs for existing frames, or even new frames (Green, 2004)

Typical end-user applications for FrameNet include Question answering (Sinha, 2008) and information extraction (Mohit and Narayanan, 2003), and using FrameNet data has enabled some improvements on systems attempting the RTE task (Burchardt, 2008). The FrameNet website lists the intended uses for hundreds of users of the FrameNet data, including sentiment analysis, building dialog systems, improving machine translation, teaching English as a second language, etc. The FrameNet team have an active partnership with Decisive Analytics Corporation, which is using FN-based ASRL as for event recognition and tracking for their govenment and commercial clients.

## 4 Some Limitations and Extensions of the FrameNet Model

FrameNet works almost entirely on edited text, so directly applying the ASRL systems trained on current FN data will probably give poor results on, e.g. Twitter feeds or transcribed conversation.

FrameNet also works strictly within the sentence, so there is no direct way to deal with text coherence, although FrameNet annotation does indicate when certain core FEs are missing from a sentence, which typically indicates that that they are realized elsewhere in the text. This feature can be used to link arguments across sentences (Ruppenhofer et al., 2010b).

**Technical terms and Proper Nouns:**

FrameNet has taken as its mandate to cover the "core" lexicon of English, words in common use, whose definitions are established by their usage. The number of senses per word is known to increase with the frequency of occurrence Zipf (19491965), so the most frequent words are likely to be the most polysemous and therefore both the most important and the most challenging for NLP. In general, the FrameNet team have assumed that technical vocabulary, whose definitions are established by domain experts, will be handled in terminologies for each domain, such as the Medical Subject Headings of the U.S. National Library of Medicine (`https://www.nlm.nih.gov/mesh/meshhome.html`) and the Department of Defense Dictionary of Military Terms (`http://www.dtic.mil/doctrine/dod_dictionary/`). For similar reasons, FrameNet does not annotate proper nouns, also known in NLP as named entities. FrameNet cannot and has no reason to compete with the on-line resources for these domains, such as Wikipedia, lists of male and female personal names, and gazetteers. On the other hand, Frame Semantic resources have been produced in several specialized domains: Thomas Schmidt created a Frame-Semantic analysis of the language associated with soccer (in German, English, and French) (Schmidt, 2008), `http://www.kictionary.com`; and lexica in the legal domain have been produced for Italian (Venturi et al., 2009) and Brazilian Portuguese (Bertoldi and Oliveira Chishman, 2012).

**Negation and Conditionals:**

FrameNet does not have representations for negation and conditional sentences. The words *never*.adv and *seldom*.adv are LUs in the **Frequency** frame, but there is no recognition of their status as negatives. The general approach which the FrameNet team has proposed would be to treat negative expressions as parts of constructs li-

censed by constructions which have a "negation" frame as their meaning pole, and license negative polarity items over some scope in the sentence, but defining that scope is a notoriously difficult problem. We are just beginning to work a mental spaces approach to the related problem of conditional sentences, cf. Dancygier and Sweetser (2005) and Sweetser (2006). FrameNet does not include the word *if*, but does include both LUs and annotation for a number of modal verbs and other types of nouns and adjectives which can be used to express conditionality, incuding the following:

| Frame | : LUs |
|---|---|
| **Possibility** | : *can, could, might, may* |
| **Capability** | : *able*.a, *ability*.n, *can*.v, *potential*.n/a, . . . |
| **Likelihood** | : *likely*.a, *might*.v, *may*.v, *must*.v, *possible*.a, . . . |

## 5 Future directions: Expert curation vs. rapid growth

After almost two decades of work at varying levels of intensity, depending on funding, FrameNet contains almost 1200 Semantic Frames, covering almost 13,000 word senses (Lexical Units) , documented with almost 200,000 manual annotations. This is bigger than a toy lexicon, but far fewer LUs than WordNet or other lexicons derived automatically from the web. By virtue of expert curation, the FrameNet lexical database contains a wealth of semantic knowledge that is unique. The database is freely available from the FrameNet website.

One challenge we face now is finding a way to greatly expand FrameNet in a more cost-effective way while preserving the accuracy and richness of the annotation. We have recently done some small-scale experiments on crowd-sourcing various parts of the process in partnership with colleagues at Google, and the preliminary results are encouraging.

Another challenge comes as a result of the success of Frame Semantics as an interlingua (Boas, 2009). There are now projects building FrameNet-style lexical databases for many different languages; funded projects are creating FrameNets for German, Spanish, Japanese, Swedish, Chinese, French and Arabic; smaller efforts have created Frame Semantics-based resources for many other languages, including Italian, Korean, Polish, Bulgarian, Russian, Slovenian, Hebrew, and Hindi.

Some are produced almost entirely via manual annotation, while others are being created semi-automatically. The good news is that the general result seems to be that the frames devised for English can be used for the majority of LUs in each of these language. The challenge is finding a way to integrate the frame semantic work being done around the world, to create a truly multi-lingual FrameNet.

For more information on all these topics, please visit

`http://framenet.icsi.berkeley.edu`

## References

Anderson Bertoldi and Rove Luiza Oliveira Chishman. 2012. Developing a frame-based lexicon for the Brazilian legal language: The case of the criminal process frame. In Monica Palmirani, Ugo Pagallo, Pompeu Casanovas, and Giovanni Sartor, editors, *AI Approaches to the Complexity of Legal Systems*, volume 7639 of *Lecture Notes in Computer Science*, pages 256–270. Springer Berlin Heidelberg.

Hans C. Boas, editor. 2009. *Multilingual FrameNets in Computational Lexicography: Methods and Applications*. Mouton de Gruyter.

Aljoscha Burchardt. 2008. *Modeling Textual Entailment with Role-Semantic Information*. Ph.D. thesis, Universität des Saarlandes.

Eugene Charniak. 1977. Framed PAINTING: The representation of a common sense knowledge fragment. *Cognitive Science*, 1(4):235–264.

Noam Chomsky. 1965. *Aspects of the Theory of Syntax*. MIT Press, Cambridge, MA.

Barbara Dancygier and Eve Sweetser. 2005. *Mental spaces in grammar: conditional constructions*. Cambridge University Press, Cambridge, UK; New York.

Dipanjan Das, Desai Chen, André F. T. Martins, Nathan Schneider, and Noah A. Smith. 2013. Frame-Semantic parsing. *Computational Linguistics*, 40(1).

Katrin Erk and Sebastian Padó. 2006. Shalmaneser – a flexible toolbox for semantic role assignment. In *Proceedings of the fifth International Conference on Language Resources and Evaluation (LREC-2006)*, Genoa, Italy.

Charles J. Fillmore and B.T.S. Atkins. 1994. Starting where the dictionaries stop: The challenge for computational lexicography. In Antonio Zampolli and Sue Atkins, editors, *Computational Approaches to the Lexicon*. Oxford University Press.

Charles J. Fillmore and Collin F. Baker. 2010. A frames approach to semantic analysis. In Bernd Heine and Heiko Narrog, editors, *Oxford Handbook of Linguistic Analysis*, pages 313–341. OUP.

Charles J. Fillmore. 1968. The case for case. In E. Bach and R. Harms, editors, *Universals in Linguistic Theory*. Holt, Rinehart & Winston, New York.

Charles J. Fillmore. 1969. Toward a modern theory of case. In David A Reibel and Sanford A. Shane, editors, *Modern Studies in English: Readings in Transformational Grammar*, pages 361–375. Prentice-Hall, Englewood Cliffs, New Jersey.

Charles J. Fillmore. 1976. Frame semantics and the nature of language. *Annals of the New York Academy of Sciences: Conference on the Origin and Development of Language and Speech*, 280(1):20–32.

Charles J. Fillmore. 1977a. Frame semantics. pages 111–137.

Charles J. Fillmore. 1977b. The need for a frame semantics in linguistics. In Hans Karlgren, editor, *Statistical Methods in Linguistics*. Scriptor.

Charles J. Fillmore. 1977c. Scenes-and-frames semantics. In Antonio Zampolli, editor, *Linguistic Structures Processing*, number 59 in Fundamental Studies in Computer Science. North Holland Publishing.

Charles J. Fillmore. 1985. Frames and the semantics of understanding. *Quaderni di Semantica*, 6(2):222–254.

Daniel Gildea and Daniel Jurafsky. 2000. Automatic labeling of semantic roles. In *ACL 2000: Proceedings of ACL 2000, Hong Kong*.

Daniel Gildea and Daniel Jurafsky. 2002. Automatic labeling of semantic roles. *Computational Linguistics*, 28(3):245–288.

Rebecca Green. 2004. *Inducing Semantic Frames from Lexical Resources*. Ph.D. thesis, University of Maryland, College Park.

Richard Johansson and Pierre Nugues. 2007. LTH: Semantic structure extraction using nonprojective dependency trees. In *Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)*, pages 227–230, Prague, Czech Republic, June. Association for Computational Linguistics.

Paul Kay and Charles J. Fillmore. 1999. Grammatical constructions and linguistic generalizations: The what's x doing y? construction. *Language*, 75:1–33.

Beth Levin. 1993. *English Verb Classes and Alternations: A Preliminary Investigation*. University of Chicago Press, Chicago.

Marvin Minsky. 1974. A framework for representing knowledge. Memo 306, MIT-AI Laboratory, June.

Behrang Mohit and Srini Narayanan. 2003. Semantic extraction with wide-coverage lexical resources. In Marti Hearst and Mari Ostendorf, editors, *HLT-NAACL 2003: Short Papers*, pages 64–66, Edmonton, Alberta, Canada, May 27 - June 1. Association for Computational Linguistics.

Josef Ruppenhofer, Michael Ellsworth, Miriam R. L. Petruck, Christopher R. Johnson, and Jan Scheffczyk. 2006. *FrameNet II: Extended Theory and Practice*. International Computer Science Institute, Berkeley, California. Distributed with the FrameNet data.

Josef Ruppenhofer, Michael Ellsworth, Miriam R. L. Petruck, Christopher R. Johnson, and Jan Scheffczyk. 2010a. *FrameNet II: Extended Theory and Practice*. FrameNet Project, September.

Josef Ruppenhofer, Caroline Sporleder, Roser Morante, Collin Baker, and Martha Palmer. 2010b. Semeval-2010 task 10: Linking events and their participants in discourse. In *Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions (SEW-2009)*, pages 106–111, Boulder, Colorado, June. Association for Computational Linguistics.

Roger C. Schank and Robert P. Abelson. 1977. *Scripts, Plans, Goals and Understanding: an Inquiry into Human Knowledge Structures*. Lawrence Erlbaum, Hillsdale, NJ.

Thomas Schmidt. 2008. The Kicktionary: Combining corpus linguistics and lexical semantics for a multilingual football dictionary. In Eva Lavric, Gerhard Pisek, Andrew Skinner, and Wolfgang Stadler, editors, *The Linguistics of Football*, number 38 in Language in Performance, pages 11–23. Gunter Narr, Tübingen.

Steve Sinha. 2008. *Answering Questions about Complex Events*. Ph.D. thesis, EECS Department, University of California, Berkeley, Dec.

Eve Sweetser. 2006. Negative spaces: Levels of negation and kinds of spaces. In Stéphanie Bonnefille and Sébastien Salbayre, editors, *Proceedings of the conference "Negation: Form, figure of speech, conceptualization"*, Tours. Groupe de recherches anglo-américaines de l'Université de Tours, Publications universitaires Fran cois Rabelais.

Giulia Venturi, Alessandro Lenci, Simonetta Montemagn, Eva Maria Vecchi, Maria Teresa Sagri, and Daniela Tiscornia. 2009. Towards a FrameNet resource for the legal domain. In *Proceedings of the Third Workshop on Legal Ontologies and Artificial Intelligence Techniques*, Barcelona, Spain, June.

George Kingsley Zipf. 1949[1965]. *Human behavior and the principle of least effort: an introduction to human ecology*. Hafner Pub. Co., New York.