

## A TAXONOMY FOR ENGLISH NOUNS AND VERBS

Robert A. Amsler  
Computer Sciences Department  
University of Texas, Austin, TX 78712

**ABSTRACT:** The definition texts of a machine-readable pocket dictionary were analyzed to determine the disambiguated word sense of the kernel terms of each word sense being defined. The resultant sets of word pairs of defined and defining words were then computationally connected into two taxonomic semi-lattices ("tangled hierarchies") representing some 24,000 noun nodes and 11,000 verb nodes. The study of the nature of the "topmost" nodes in these hierarchies, and the structure of the trees reveal information about the nature of the dictionary's organization of the language, the concept of semantic primitives and other aspects of lexical semantics. The data proves that the dictionary offers a fundamentally consistent description of word meaning and may provide the basis for future research and applications in computational linguistic systems.

### 1. INTRODUCTION

In the late 1960's, John Olney et al. at System Development Corporation produced machine-readable copies of the Merriam-Webster New Pocket Dictionary and the Merriam-Webster Seventh Collegiate Dictionary. These massive data files have been widely distributed within the computational linguistic community, yet research upon the basic structure of the dictionary has been exceedingly slow and difficult due to the significant computer resources required to process tens of thousands of definitions.

The dictionary is a fascinating computational resource. It contains spelling, pronunciation, hyphenation, capitalization, usage notes for semantic domains, geographic regions, and propriety; etymological, syntactic and semantic information about the most basic units of the language. Accompanying definitions are example sentences which often use words in prototypical contexts. Thus the dictionary should be able to serve as a resource for a variety of computational linguistic needs. My primary concern within the dictionary has been the development of dictionary data for use in understanding systems. Thus I am concerned with what dictionary definitions tell us about the semantic and pragmatic structure of meaning. The hypothesis I am proposing is that definitions in the lexicon can be studied in the same manner as other large collections of objects such as plants, animals, and minerals are studied. Thus I am concerned with enumerating the classificational organization of the lexicon as it has been implicitly used by the dictionary's lexicographers.

Each textual definition in the dictionary is syntactically a noun or verb phrase with one or more kernel terms. If one identifies these kernel terms of definitions, and then proceeds to disambiguate them relative to the senses offered in the same dictionary under their respective definitions, then one can arrive at a large collection of pairs of disambiguated words which can be assembled into a taxonomic semi-lattice.

This task has been accomplished for all the definition texts of nouns and verbs in a common pocket dictionary. This paper is an effort to reveal the results of a preliminary examination of the structure of these databases.

The applications of this data are still in the future. What might these applications be?

First, the data should provide information on the contents of semantic domains. One should be able to determine from a lexical taxonomy what domains one might be in given one has encountered the word "periscope", or "petiole", or "petroleum".

Second, dictionary data should be of use in resolving semantic ambiguity in text. Words in definitions appear in the company of their prototypical associates.

Third, dictionary data can provide the basis for creating case grammar descriptions of verbs, and noun argument descriptions of nouns. Semantic templates of meaning are far richer when one considers the taxonomic inheritance of elements of the lexicon.

Fourth, the dictionary should offer a classification which anthropological linguists and psycholinguists can use as an objective reference in comparison with other cultures or human memory observations. This isn't to say that the dictionary's classification is the same as the culture's or the human mind's, only that it is an objective datum from which comparisons can be made.

Fifth, knowledge of how the dictionary is structured can be used by lexicographers to build better dictionaries.

And finally, the dictionary if converted into a computer tool can become more readily accessible to all the disciplines seeking to use the current paper-based versions. Education, historical linguistics, sociology, English composition, etc. can all make steps forward given that they can assume access to a dictionary is immediately available via computer. I do not know what all these applications will be and the task at hand is simply an elucidation of the dictionary's structure as it currently exists.

### 2. "TANGLED" HIERARCHIES OF NOUNS AND VERBS

The grant, MCS77-01315, "Development of a Computational Methodology for Deriving Natural Language Semantic Structures via Analysis of Machine-Readable Dictionaries", created a taxonomy for the nouns and verbs of the Merriam-Webster Pocket Dictionary (MPD), based upon the hand-disambiguated kernel words in their definitions. This taxonomy confirmed the anticipated structure of the lexicon to be that of a "tangled hierarchy" [8,9] of unprecedented size (24,000 noun senses, 11,000 verb senses). This data base is believed to be the first to be assembled which is representative of the structure of the entire English lexicon. (A somewhat similar study of the Italian lexicon has been done [2.11]). The content categories agree substantially with the semantic structure of the lexicon proposed by Nida [15], and the verb taxonomy confirms the primitives proposed by the San Diego LNR group [16].

This "tangled hierarchy" may be described as a formal data structure whose bottom is a set of terminal disambiguated words that are not used as kernel defining terms; these are the most specific elements in the structure. The tops of the structure are senses of words such as "cause", "thing", "class", "being", etc. These are the most general elements in the tangled hierarchy. If all the top terms are considered to be





This essential property, the inability to write a definition explaining a word's meaning without using another member of some small set of near synonymous words, is the basis for describing such a set as a PRIMITIVE. It is based upon the notion of definition given by Wilder [21], which in turn was based upon a presentation of the ideas of Padoa, a turn-of-the-century logician.

The definitions are given, the disambiguation of their kernel's senses leads to a cyclic structure which cannot be resolved by attributing erroneous judgements to either the lexicographer or the disambiguator; therefore the structure is taken as representative of an undefinable primitive concept, and the words whose definitions participate in this complex structure are found to be undefinable without reference to the other members of the set of undefined terms.

The question of what to do with such primitives is not really a problem, as Winograd notes [22], once one realizes that they must exist at some level, just as mathematical primitives must exist. In tree construction the solution is to form a single node whose English surface representation may be selected from any of the words in the primitive set. There probably are connotative differences between the members of the set, but the ordinary pocket dictionary does not treat these in its definitions with any detail. The Merriam-Webster Collegiate Dictionary does include so-called "synonym paragraphs" which seem to discuss the connotative differences between words sharing a "ring".

While numerous studies of lexical domains such as the verbs of motion [1,12,13] and possession [10] have been carried out by other researchers, it is worth noting that recourse to using ordinary dictionary definitions as a source of material has received little attention. Yet the "primitives" selected by Donald A. Norman, David E. Rumelhart, and the LNR Research Group for knowledge representation in their system bear a remarkable similarity to those verbs used most often as kernels in The Merriam-Webster Pocket Dictionary and Donald Sherman has shown (Table 4) these topmost verbs to be among the most common verbs in the Collegiate Dictionary as well [19]. The most frequent verbs of the MPD are, in descending order, MAKE, BE, BECOME, CAUSE, GIVE, MOVE, TAKE, PUT, FORM, BRING, HAVE, and GO. The similarity of these verbs to those selected by the LNR group for their semantic representations, i.e., BECOME, CAUSE, CHANGE, DO, MOVE, POSS ("have"), TRANSF ("give", "take"), etc., [10.14.18] is striking. This similarity is indicative of an underlying "rightness" of dictionary definitions and supports the proposition that the lexical information extractable from study of the dictionary will prove to be the same knowledge needed for computational linguistics.

The enumeration of the primitives for nouns and verbs by analysis of the tangled hierarchies of the noun and verb forests grown from the MPD definitions is a considerable undertaking and one which goes beyond the scope of this paper. To see an example of how this technique works in practice, consider the discovery of the primitive group starting from PLACE-1.3A.

place-1.3a - a building or locality used for a special purpose

The kernels of this definition are "building" and "locality". Looking these up in turn we have:

building-1a - a usu. roofed and walled structure (as a house) for permanent use

locality-0a - a particular spot, situation, or location

Table 4. 50 Most Frequent Verb Infinitive Forms of W7 Verb Definitions (from [19]).

1878	MAKE	157	FURNISH
908	CAUSE	154	TURN
815	BECOME	150	GET
599	GIVE	150	TREAT
569	BE	147	SUBJECT
496	MOVE	141	HOLD
485	TAKE	137	UNDERGO
444	PUT	132	CHANGE
366	BRING	132	USE
311	HAVE	129	KEEP
281	FORM	127	ENGAGE
259	GO	127	PERFORM
240	SET	118	BREAK
224	COME	118	REDUCE
221	REMOVE	112	EXPRESS
210	ACT	107	ARRANGE
204	UTTER	107	MARK
190	PASS	106	SEPARATE
188	PLACE	105	DRIVE
178	COVER	104	CARRY
173	CUT	101	THROW
169	PROVIDE	100	SERVE
166	DRAW	100	SPEAK
163	STRIKE	100	WORK

This gives us four new terms, "structure", "spot", "situation", and "location". Looking these up we find the circularity forming the primitive group.

structure-.2a - something built (as a house or a dam)

spot-1.3a - LOCATION, SITE

location-.2a - SITUATION, PLACE

situation-.1a - location, site

And finally, the only new term we encounter is "site" which yields,

site-.0a - location <\* of a building> <battle \*>

The primitive cluster thus appears as in Figure 5.

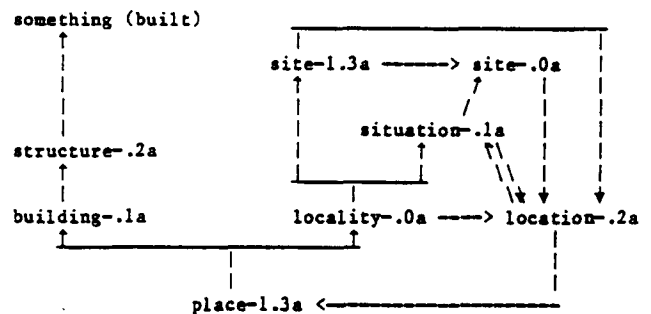


Figure 5. Diagram of Primitive Set Containing PLACE, LOCALITY, SPOT, SITE, SITUATION, and LOCATION

### 2.3 NOUNS TERMINATING IN RELATIONS TO OTHER NOUNS OR VERBS

In addition to terminating in "dictionary circles" or "loops", nouns also terminate in definitions which are actually text descriptions of case arguments of verbs or relationships to other nouns. "Vehicle" is a fine

example of the former, being as it were the canonical instrumental case argument of one sense of the verb "carry" or "transport".

vehicle - a means of carrying or transporting something

"Leaf" is an example of the latter, being defined as a part of a plant.

leaf - a usu. flat and green outgrowth of a plant stem that is a unit of foliage and functions esp. in photosynthesis.

Thus "leaf" isn't a type of anything. Even though under a strictly genus/differentia interpretation one would analyze "leaf" as being in an ISA relationship with "outgrowth", "outgrowth" hasn't a suitable homogeneous set of members and a better interpretation for modeling this definition would be to consider the "outgrowth of" phrase to signify a part/whole relationship between "leaf" and "plant".

Hence we may consider the dictionary to have at least two taxonomic relationships (i.e. ISA and ISPART) as well as additional relations explaining noun terminals as verb arguments. One can also readily see that there will be taxonomic interactions among nodes connected across these relationship "bridges".

While the parts of a plant will include the "leaves", "stem", "roots", etc., the corresponding parts of any TYPE of plant may have further specifications added to their descriptions. Thus "plant" specifies a functional form which can be further elaborated by descent down its ISA chain. For example, a "frond" is a type of "leaf".

frond - a usu. large divided leaf (as of a fern)

We knew from "leaf" that it was a normal outgrowth of a "plant", but now we see that "leaf" can be specialized, provided we get confirmation from the dictionary that a "fern" is a "plant". (Such confirmation is only needed if we grant "leaf" more than one sense meaning, but words in the Pocket Dictionary do typically average 2-3 sense meanings). The definition of "fern" gives us the needed linkage, offering,

fern - any of a group of flowerless seedless vascular green plants

Thus we have a specialized name for the "leaf" appendage of a "plant" if that plant is a "fern". This can be represented as in Figure 6.

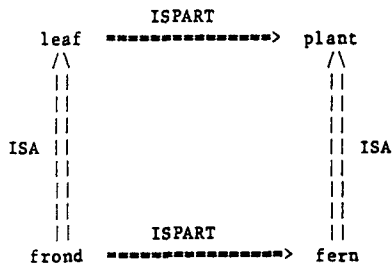


Figure 6. LEAF:PLANT::FROND:FERN

This conclusion that there are two major transitive taxonomies and that they are related is not of course new. Evens et al. [6,7] have dealt with the PART-OF relationship as second only to the ISA relationship in importance, and Fahlman [8,9] has also discussed the

interaction of the PART-OF and ISA hierarchies. Historically even Raphael [17] used a PART-OF relationship together with the ISA hierarchy of SIR's deduction system. What however is new is that I am not stating "leaf" is a part of a plant because of some need use this fact within a particular system's operation, but "discovering" this in a published reference source and noting that such information results naturally from an effort to assemble the complete lexical structure of the dictionary.

## 2.4 PARTITIVES AND COLLECTIVES

As mentioned in Section 2.3, the use of "outgrowth" in the definition of "leaf" causes problems in the taxonomy if we treat "outgrowth" as the true genus term of that definition. This word is but one example of a broad range of noun terminals which may be described as "partitives". A "partitive" may be defined as a noun which serves as a general term for a PART of another large and often very non-homogeneous set of concepts. Additionally, at the opposite end of the partitive scale, there is the class of "collectives". Collectives are words which serve as a general term for a COLLECTION of other concepts.

The disambiguators often faced decisions as to whether some words were indeed the true semantic kernels of definitions, and often found additional words in the definitions which were more semantically appropriate to serve as the kernel -- albeit they did not appear syntactically in the correct position. Many of these terms were partitives and collectives. Figure 7 shows a set of partitives and collectives which were extracted and classified by Gretchen Hazard and John White during the dictionary project. The terms under "group names", "whole units", and "system units" are collectives. Those under "individuator", "piece units", "space shapes", "existential units", "locus units", and "event units" are partitives. These terms usually appeared in the syntactic frame "An \_\_\_\_\_ of" and this additionally served to indicate their functional role.

1	QUANTIFIERS	3	EXISTENTIAL UNITS
1.1	GROUP NAMES pair.collection.group cluster.bunch. band (of people)	3.1	VARIANT version.form.sense
1.2	INDIVIDUATORS member.unit.item. article.strand, branch (of science, etc.)	3.2	STATE state.condition
2	SHAPE UNITS	4	REFERENCE UNITS
2.1	PIECE UNITS sample.bit.piece, tinge,tint	4.1	LOCUS UNITS place.end,ground, point
2.2	WHOLE UNITS mass.stock,body. quantity.wad	4.2	PROCESS UNITS cause.source.means. way.manner
2.3	SPACE SHAPES bed,layer.strip,belt, crest,fringe,knot. knob,tuft	5	SYSTEM UNITS system.course.chain. succession.period
		6	EVENT UNITS act,discharge, instance
		7	EXCEPTIONS growth.study

Figure 7. Examples of Partitives and Collectives [3]

#### ACKNOWLEDGEMENTS

This research on the machine-readable dictionary could not have been accomplished without the permission of the G. & C. Merriam Co., the publishers of the Merriam-Webster New Pocket Dictionary and the Merriam-Webster Seventh Collegiate Dictionary as well as the funding support of the National Science Foundation. Thanks should also go to Dr. John S. White, currently of Siemens Corp., Boca Raton, Florida; Gretchen Hazard; and Drs. Robert F. Simmons and Winfred P. Lehmann of the University of Texas at Austin.

#### REFERENCES

1. Abrahamson, Adele A., "Experimental Analysis of the Semantics of Movement." in Explorations in Cognition, Donald A. Norman and David E. Rumelhart, ed., W. H. Freeman, San Francisco, 1975, pp. 248-276.
2. Alinei, Mario, La struttura del lessico. Il Mulino, Bologna, 1974.
3. Amsler, Robert A. and John S. White. "Final Report for NSF Project MCS77-01315, Development of a Computational Methodology for Deriving Natural Language Semantic Structures via Analysis of Machine-Readable Dictionaries," Tech. report. Linguistics Research Center, University of Texas at Austin, 1979.
4. Amsler, Robert A., The Structure of the Merriam-Webster Pocket Dictionary. PhD dissertation, The University of Texas at Austin, December 1980.
5. Calzolari, N., "An Empirical Approach to Circularity in Dictionary Definitions," Cahiers de Lexicologie, Vol. 31, No. 2, 1977, pp. 118-128.
6. Evens, Martha and Raoul Smith. "A Lexicon for a Computer Question-Answering System." Tech. report 77-14, Illinois Inst. of Technology. Dept. of Computer Science, 1977.
7. Evens, Martha, Bonnie Litowitz, Judith Markowitz, Raoul Smith and Oswald Werner. Lexical-Semantic Relations: A Comparative Survey. Linguistic Research, Carbondale, 1980.
8. Fahlman, Scott E., "Thesis progress report: A system for representing and using real-world knowledge," AI-Memo 331, M.I.T. Artificial Intelligence Lab., 1975.
9. Fahlman, Scott E., A System for Representing and Using Real-World Knowledge. PhD dissertation, M.I.T., 1977.
10. Gentner, Dedre, "Evidence for the Psychological Reality of Semantic Components: The Verbs of Possession," in Explorations in Cognition, Donald A. Norman and David E. Rumelhart, ed., W. H. Freeman, San Francisco, 1975, pp. 211-246.
11. Lee, Charmaine, "Review of La struttura del lessico by Mario Alinei." Language, Vol. 53, No. 2, 1977, pp. 474-477.
12. Levelt, W. J. M., R. Schreuder, and E. Hoenkamp, "Structure and Use of Verbs of Motion," in Recent Advances in the Psychology of Language, Robin Campbell and Philip T. Smith, ed., Plenum Press, New York, 1976, pp. 137-161.
13. Miller, G., "English verbs of motion: A case study in semantic and lexical memory." in Coding Processes in Human Memory, A.W. Melton and E. Martins, ed., Winston, Washington, D.C., 1972.
14. Munro, Allen. "Linguistic Theory and the LNR Structural Representation." in Explorations in Cognition, Donald A. Norman and David E. Rumelhart, ed., W. H. Freeman, San Francisco, 1975, pp. 88-113.
15. Nida, Eugene A., Exploring Semantic Structures. Wilhelm Fink Verlag, Munich, 1975.
16. Norman, Donald A., and David E. Rumelhart. Explorations in Cognition. W.H. Freeman, San Francisco, 1975.
17. Raphael, Bertram. SIR: A Computer Program for Semantic Information Retrieval, PhD dissertation. M.I.T., 1968.
18. Rumelhart, David E. and James A. Levin. "A Language Comprehension System." in Explorations in Cognition, Donald A. Norman and David E. Rumelhart, ed., W. H. Freeman, San Francisco, 1975, pp. 179-208.
19. Sherman, Donald, "A Semantic Index to Verb Definitions in Webster's Seventh New Collegiate Dictionary." Research Report. Computer Archive of Language Materials, Linguistics Dept., Stanford University, 1979.
20. Sparck Jones, Karen. "Dictionary Circles," SDC document TM-3304, System Development Corp., January 1967.
21. Wilder, Raymond L., Introduction to the Foundations of Mathematics, John Wiley & Sons, Inc., New York, 1965.
22. Winograd, Terry, "On Primitives, prototypes, and other semantic anomalies." Proceedings of the Workshop on Theoretical Issues in Natural Language Processing, June 10-13, 1975, Cambridge, Mass., Schank, Roger C., and B.L. Nash-Webber, ed., Assoc. for Comp. Ling., Arlington, 1978, pp. 25-32.