

Carnegie Mellon HUNT LIBRARY

ILLiad TN: 68577



Lending String: *PMC,VA@,VA@

Patron: Rapaport, William

Journal Title: Dictionaries in the electronic age ;
proceedings of the Conference, September 18-19, 1989,
St. Catherine's College, Oxford, England.

Volume: Issue:
Month/Year: 1989
Pages: ???

Article Author: University of Waterloo. Centre for the
New Oxford English
Dictionary. Conference (5th ; 1989 ; Oxf

Article Title: Martin Kay; The Concrete Lexicon and
the Abstract Dictionary

Imprint: Waterloo, Ont. ; UW Centre for the New O

ILL Number: 3333868

Call No.: 413.028 U58P 1989 1

Location: OVRSQ-2

Item No.:

ARIEL

7pgs

Borrower: BUF

Shipping Address:

234 Lockwood Memorial Library, ILL
SUNY at Buffalo
Buffalo, NY 14260

Fax:

Ariel: 128.205.111.1

The Concrete Lexicon and the Abstract Dictionary

Martin Kay

Xerox Palo Alto Research Center and Stanford University

The way I understand it, there is a very big difference between a dictionary and a lexicon. A dictionary is a book that white-collar workers keep on their desks and that is second in popularity only to the Bible among domestic bibliographic knickknacks. A lexicon is an abstraction; something that we all have in our heads by virtue of having learnt a language. *The* lexicon — as opposed to *a* lexicon — is even more abstract; to the extent that a language, like English, is a well defined notion, and to the extent that we can decide who speaks it and who does not, there is a lexicon — *the lexicon of English* — that is in the heads of all of them. It is an idealization, like a frictionless pulley or a perfect gas which the lexicons in our heads presumably approximate.

If the lexicons in our heads more nearly approximated the ideal lexicon, then we might very well not be here — here at this meeting in Oxford, I mean. I assume we have dictionaries to thank for bringing us here. Because, if our lexica were more similar, we should not need dictionaries on our desks and in our houses. The reason I need a dictionary is to help me understand something that another person said or wrote that used something from his lexicon, or hers, and which I do not have in mine. I turn to the dictionary in the hope that the lexicographer, who is professionally committed to having a large lexicon and to writing it down, will be able to supply my lack. The obverse of this, of course, is the situation in which I am the one trying to communicate and suspect that there is lexical material known to other speakers of my language that would either clarify my message or impress my audience.

I hasten to point out that some dictionaries, and most notably the *Oxford English Dictionary*, do not fit my picture. The *OED* is not a desk dictionary intended mainly to stop gaps in people's mental lexica. Most of the information it contains is not part of the lexicon of the language, or the lexicon of any individual. Information about when a word was first used in a particular sense, and by whom, or about its etymology, is not the kind of information of which lexica are made. These are not things a person knows by virtue of being a speaker, or writer, of English — even a very good one. They are things a person would know by virtue of being a scholar, and that is something quite different.

There is more than casual interest these days in the question of whether a lexicon could exist elsewhere than in the head of a person. To put it another way, can the abstract lexicon be made concrete? The scientific interest of the question is clear; it is always easier to study something *in vitro* than to grope for it in the dark recesses of a person's brain. The question is especially interesting to the extent that its structure and content turned out to be different from the structure and content of a dictionary. A concrete lexicon would also be of great interest to linguistic technologists, that new breed of engineers that promise to build us machines that can understand us when we speak, translate from one language into another, and generally masquerade as people. The lexicon that would have to be lodged in the chips of such a device would be no abstraction.

You may say that the kind of device I am imagining, and to the design of whose forerunners I devote much of my time, needs only what people need. If such a device required a concrete lexicon to function, then so do people, and the lexicon is not the abstraction I am claiming it is at all. Of course this is true. If we could build an effective translation machine, or cause computers to understand our speech, without building all the other properties of human kind into these devices, we should have evidence that the lexicon can be made concrete and could at least speculate that it was concrete somewhere in our own heads. The trouble is that it may not be possible, either in principle or in practice to disengage the lexicon from everything that makes up our culture and our cognitive system as a whole. The lexicon is, after all, the locus of de Saussure's "arbitraire du signe" — the place at which language makes contact with the world, and most of the world it makes contact with is not the domain of the physicist, it is in our heads. However, though I see some cause to suppose that the lexicon may suffuse our whole mental being, I also see reason to suppose that some of what we know is very specifically lexical, and that is the line I want to concentrate on here.

If I am right, and dictionaries exist to supply what is lacking in a person's lexicon, it is natural to ask how similar the concrete lexicon would be to a dictionary or, at least, to the structure that would presumably be revealed if one could abstract away from the particular notational conventions used in it. This is what I mean by the *abstract* dictionary. For example, Hornby¹ shows a verb as taking a "to+infinitive" by writing "VP2" against it. This sequence of characters has nothing to do with the abstract dictionary that is made concrete in Hornby's book. So, I have two questions that I want to explore a little. One concerns the extent to which we should expect to profit from abstract dictionaries in the enterprise of creating the concrete lexicon. And the other is the obverse of that: to what extent could abstract, and hence concrete, dictionaries be improved by causing them to be the nearest things we can make to the concrete lexicon?

In recent years, a noticeable investment has been made by computational linguists in the attempt to turn standard published dictionaries to their own purposes. There has been some success, for example, in conducting grammatical analysis of sentences based on no more information about the words in them than can be obtained from *Longman's Dictionary of Contemporary English*² Of course, this success is measured only against analyses using specially prepared files of lexical information and this success is, as everyone knows, very limited.

Linguists of the most diverse persuasions, who might agree about little else in their field, have over the last few years been coming to the view that a great deal of the information that was previously placed elsewhere properly belongs in the lexicon. This is what is behind grammatical proposals such as Lexical Functional Grammar, Head Driven Phrase Structure Grammar, the rebirth of categorial grammar, and much of what has been done by the followers of Chomsky. It is not hard to see how this comes about. When a linguist examines the following pair of sentences, he notices, in effect, that the verbs in them belong to Hornby's VP3 class — they take two objects, the second of which is a to+infinitive.

¹ *The Advanced Learner's Dictionary of English*, Oxford University Press, 1963.

² Longmans, 1978.

- (1) I promised him to be there
 (2) I expected him to be there

Then he notices a difference between them, namely that the person whose presence is predicted in (1) is me, whereas in (2), it is him. He concludes that the VP3 class of verbs needs to be split up into the VP3A verbs, like "promise", and the VP3B verbs, like "expect".

Each time a new distinction is made in the lexicon, like the one between VP3A and VP3B, a parallel distinction must presumably be made in the grammar, for it is the grammar that must make sense of categories like VP3A and VP3B. When we notice the distinction between sentences (1) and (2), we are not noticing two separate phenomena, one grammatical and one lexical. What speakers know about their language that underlies the distinction is apparently one kind of fact which nevertheless seems to reside partly in their lexicon and partly in their grammar. We can get around the dilemma that this seems to put us in by making the grammar more of an interpreter of lexical entries than the repository of very particular grammatical facts. Specifically, whereas we might previously have been inclined to put annotations like VP3A and VP3B into lexical entries, and to write general rules that made explicit how those annotations were to be interpreted, we would now put the interpretation in the lexicon and write much more general rules which know nothing of fine distinctions but which require constraints in lexical entries to be maintained in sentences. So, for example, the lexical entries for "promise" and "expect" might contain something like the following as part of the grammatical description of the words:

promise	np, (np), inf(1)
expect	np, (np), inf

I make no claims for the particular notation. What it is intended to mean is that both verbs take three arguments, a subject, and object, and an infinitive phrase. The object is optional in both cases and accordingly, the second "np" is placed in parentheses in both cases. By default, the noun phrase most recently preceding an infinitive functions as the subject of the infinitive phrase, so that, in "I expected to go" it is I whose departure is foreseen, whereas in "I expected you to go", it is yours. The verb "promise" does exhibit the default behavior and the infinitive is marked to show that it is the first argument of the verb that will be its subject.

There are two obvious, and closely related reasons why early attempts to use abstract dictionaries as models for the construction of a concrete lexicon have started with syntax. They both have to do with the fact that syntax is, *par excellence*, the aspect of language for which we have the best abstractions. We have formalized systems of rules and parsing algorithms with well known mathematical and computational properties and, as a consequence, syntax is a central area of concern for linguists in general and computational linguists in particular. The second reason is that, precisely because of this status that syntax has, it is an area in which lexicographers feel most at home using formal devices like VP2 and VP3.

Suppose that English had another verb — say “cromise” with exactly the same meaning as promise, but with the grammatical description of “expect”. In particular, suppose that its subject was the assumed subject of the infinitive object. Notice that this would have some interesting purely syntactic consequences. For example, sentences (3) and (4) would be grammatical, whereas (5) and (6) would not.

- (3) The children promised their mother to behave themselves.
- (4) The children cromised their mother to behave herself.
- (5) *The children promised their mother to behave herself.
- (6) *The children cromised their mother to behave themselves.

Now, consider example (4) more carefully. Just what is it that the children are doing? Since “cromise” has the same meaning as “promise”, they are promising to do something. What is it that they are promising to do? Only two possibilities seem to be open, neither particularly plausible, and neither according in any way with my intuitions as a native speaker. The first possibility is that they are promising that they will behave themselves. The occurrence of the word “herself” in the sentence can only be accounted a unique grammatical anomaly. The other is that the children are giving an undertaking that their mother will behave herself. But, in this case, we can hardly claim that “cromise” really has the same meaning as “promise”, because to promise is to give an undertaking on ones own behalf. It is not transferable to anyone else, even in the same family. The upshot, I claim, is that a verb with the meaning of “promise” and the grammatical description of “expect” could not be part of the English language. If this is true, then it would be redundant for a dictionary or a lexicon to make the grammatical distinction explicit. The justification for distinguishing VP3A from VP3B would therefore have to be that of moving into the more formal part of the entry what would otherwise be available only to an agent that was able to interpret the English of the definition.

The big question, of course, concerns the extent to which we expect to be able to formalize that part of the dictionary that is now written in English. There is more to this than finding a more obscure language in which to capture this information. Many researchers, especially members of the artificial intelligentsia, will claim that the sheer amount of information in a typical dictionary entry is simply too little to support the kinds of inference that people base on them.

Consider the word “give”, in what I take to be its primary meaning. This is sense 3 in Hornby:

allow (sb.) to have; allow or cause (sb. or sth.) to pass into the care or keeping of.

Needless to say, it requires considerable skill to interpret a definition like this. Suppose I have before me the sentence “The man gave his son some money”. I know that give, in sense 3, fits patterns 18A and 19A, and 18A fits this sentence. But it is a long way from here to lining up the parts of the sentence with the parts of the definition. If I get it right, it will turn out that the subject of “give” is also the intended subject of the verb phrases in the definition. So it is the man who is allowing and causing. But, who is the “sb.” (somebody) in the first clause, and the “sb. or sth.” (somebody or something) in the second. In fact, the man’s son is the sb. in the first clause, but the money is the sth. in the second. The money shows up in the first clause as the implied object of the verb “have” and, in the son is the implied object of the preposition “of” in the second.

Certainly there is work to do in making the definitions more perspicuous³ Of more immediate concern is the question of the amount of information that the definitions contain.

A style of definition that has appealed as more useful to some workers in artificial intelligence runs more along the following lines:

A person (p_1) gives an object (q) to another person (p_2) if there is some time (t_1) at which (p_1) has (q) and (p_2) does not have (q), and a time (t_2) at which (p_1) no longer has the object, but (p_2) does. Furthermore, this change in the state of affairs was occasioned by some action on the part of (p_1).

Actually, this would probably be considered hopelessly inadequate because it appeals to the word "have", one of the least precise words in the language. But the flavor is there and that is all I need. For all its imprecision, this parody of an artificial intelligence definition is more restrictive than Hornby's sense 3. Hornby's permissive rubric leaves room for such examples as "The machine gave me my money back" in which (p_1) is not a person, "I gave the machine my last dime", where (p_2) isn't, "He gave me all the information I needed", in which (p_2) still has (q) at (t_2) and "His words gave me courage" in which (p_1) not only lacks the properties of a person, but almost certainly did not have (q) at (t_i), for any (i).

So, doubtless, as many have long suspected, the lexicon will have to be a much larger compendium of information than the dictionary. Perhaps we need a degree of precision in the lexicon that would be out of place in the dictionary. But if we do, we will almost certainly find it more difficult to account for one of the most remarkable properties of language, namely its great adaptability, especially in matters touching the lexicon. The more we insist on lexical entries that match minute details of the situations they are used to describe, the less easily we shall be able to turn them to unfamiliar uses and to account for the fact that they are still understood.

I think it possible that Hornby's entry, at least for a word like "give", suffers from overspecification, rather than the opposite. Many of his senses would be subsumed by a definition that said " (p_1) gives (p_2) (q) if (p_1) causes (p_2) to have (q)", where the word "have" is intended to cover most, if not all, of the meanings it can have in English. Certainly it covers sentences as diverse as the following:

We just gave the car some new tires.

That noise gives me a head ache.

I would like to give you a crack at it.

I can give you ten minutes to finish it.

How much will you give me for my old car? (*Hornby, sense 2*)

³ Some interesting moves in this direction were made by the designers of the Cobuild dictionary.

You should give the other boys a good example. (*Hornby, sense 4*)

You've given me your cold. (*Hornby, sense 5*)

If a lexicon of this sort could be constructed, it might well go further towards capturing the kinds of generalizations on which people presumably depend to make what would otherwise be a bewildering array of essentially unrelated facts memorable. At the same time, it would give the beginnings of an account of the adaptability of language to new situations. Notice that, on the account we have given, a change in the meaning of "have" has the side effect of changing one of the meanings of "give".

The verb "give" does, of course, have other senses, for example "He gave his life for his friends". He did not cause anyone to have his life; he simply ceased to have it himself. This is a much more restricted sense. Notice how it differs from "He gave his life to his friends" or, even more strikingly, from "He gave his fortune to the movement". Hornby points out that VP19 does not apply to these last two examples — we do not say "He gave his friends his life", and "He gave the movement his fortune", though possible, is not preferred.

I think Hornby would include "He gave his life for his friends" under his sense 6 (*devote, dedicate*), though it also partakes somewhat of sense 1 (*hand over as a present*). The parodied artificial intelligence definition was a combination of this sense — to relinquish — and my original sense — to cause to have. Now, the sentence "He gave his fortune to the movement" can be interpreted under both of these senses and I would claim that another component of the flexibility of the lexicon is that it encourages just this sort of usage, different aspects of the same word being invoked in parallel.

Much of what needs to be done in order to construct the concrete lexicon is essentially anthropological in nature because it consists in mapping the ontology of the culture and articulating the relationships between that map and the vocabulary of the language. The notion of giving is not given by the physical world; it is enshrined in our culture and, presumably many others. The notion of an exchange, or a transaction, builds on the notion of giving: one thing is given and another is given in exchange. The commercial transaction is still more specific, and limited possibly to a smaller range of cultures. Here that which is given in exchange is constrained to be money. Money, in its turn, is, by definition, the secondary object of transfer in the commercial transaction. Once it has been established that the ontology contains such a thing as the commercial transaction, two essentially independent enterprises can be undertaken in relation to it. One is to seek criteria according to which instances of it can be recognized in the real world. The other is to look for its reflections in the vocabulary of the language.

Verbs like "buy", "sell", "pay", "cost" and "charge", and nouns like "cost", "price", "merchandise", "buyer", and "seller" all have lexical entries that declare them to be names of the commercial transaction. The difference between "buy" and "sell" is not in the kinds of thing they can be used to refer to — they both refer to commercial transactions — but in the syntax that they make available for referring to the participants in the transaction and, consequently, the perspective in which they place those participants. There are four participants that must be present in any commercial transaction, a buyer, a seller, some merchandise, and some cash. The presence of all of them is inferred whenever any of these words is used, though they do not have to be named. If "buy" is used in the active voice, the buyer must be mentioned, and that person is cast in the role

of agent, or instigator, of the transaction. The merchandise must always be mentioned when "buy" is used, as object of the active verb or subject of the passive. The point is that it is the participants one wishes to name and the light in which one wishes to cast them that determines the word one chooses rather than its meaning, or the kind of thing to which it can refer.

I have several reasons for suspecting that the concrete lexicon will not be a dictionary in the usual sense, and will not serve the makers of dictionaries very much except as an abstraction. One is that it will be impossible to achieve the level of specificity that it will require without large amounts of technical terminology and specialized notational devices. It is therefore likely to be accessible only to people with special knowledge and training. The second reason is similar, though it may seem to contradict the first. It is that the concrete lexicon will in some ways be deceptively — often quite misleadingly simple. The sense (give₁) (sense 1 of "give") is defined to be "(cause₁ to have₁) where "(cause₁)" is sense 1 of "cause", and "(to have₁)" is sense 1 of "have". More accurately stated, "(give₁)", "(cause₁)", and "(to have₁)" are types of objects in the underlying ontology. Ignoring the subscripts, the idea that comes across is disarmingly simple, but the subscripts are crucial to the intention. Finally, I suspect that the number of richly meaningful lexical items for which nothing can be said, except in the part of the lexicon I have called the "sensory map" will be fairly large. The items, about which we shall have nothing to say will be the ones to which the most space is usually devoted, despite the fact that nobody reads the articles. Several of the senses of words like "be", "have", "set", and the like fall in this category. The reason why the definitions are so long is, in short, that they are undefinable, a fact which the makers of the concrete lexicon will have to face. The reason that nobody reads these entries that have been agonized over at such length, is that these are items with which speakers of the language need no assistance. They know more about them than the lexicographer could ever say.

I believe, as a growing number of others seem to, that the time is ripe to start on at least a preliminary version of the concrete lexicon because of what it could contribute to our understanding of language and because it would facilitate technological steps that we are ready to take. But, for reasons some of which I have tried to outline, I do not believe that dictionaries will contribute much to the enterprise or that it can be accomplished in any important measure by processing dictionaries automatically. I believe that the greatest impediment will be literal mindedness — the almost irresistible temptation to make of the layman's lexicon the physicist's encyclopedia. In the lexicon, giving is simple, once you know about having, and, unicorns are just as real as horses.