MODELING DICTIONARY DATA

by

Robert F. Simmons

and

Robert A. Amsler

NATURAL LANGUAGE RESEARCH FOR CAI

Sponsored by

THE NATIONAL SCIENCE FOUNDATION

Grant GJ509E

## ABSTRACT

Forms and structures of definitions in Merriam-Webster's
dictionaries are presented to derive models of sense selection
contexts, sense meanings, and hierarchical relations among verbs.
The sense meaning model is presented as a case-role semantics
accompanied by time-ordered sets of assertions marked for truth
value. Systematic extraction of these types of models from
dictionary data is argued to be an encouraging line of research.

I think we are all agreed that since about 1968 we have had some very interesting natural language processing systems that have very limited but powerful capabilities within those limitations. The critical thing that has been going on in these is to choose a microworld as Ralph Grishman suggested [reference to preceding talk] where there is very little that in fact can be said, and then you are able to manage the semantics.

Now the situation of trying to understand text is quite different. You know, examples of those microworlds are The Woods Airline Guide [2], and the lunar rocks data base [3], the Winograd hand-and-blocks [4], and the Heidorn trucks at a gas depot [5], and Schank's John and Mary -- with John's attitude towards Mary. [laughter] And there are others. I think that in text processing we are nowhere near in such good shape. I have seen interesting experiments by Sager about 3 years ago [6] in terms of getting some semantic structure out of documents and I recall early work by Harris [7], to get kernel sentences out of documents. And most recently we have some work by Charniak on children's stories [8] in which the great contribution was, I think, to show how incredibly difficult it is to solve the problem of reference even there, in order to find the antecedent for a pronoun, 'it'. There are a couple of examples where you have to have the full strength of some kind of a problem solving/theorem prover system.

Most recently, a paper that hasn't gotten around very much, is a very ambitious effort by Roger Schank [9] on setting up discourse structure for paragraphs. That paper will probably be circulating pretty widely by Spring. In that, he very ambitiously takes the folk tale from Eskimo literature — several paragraph stories — about 5 I think — and he makes a causal chain of the relationships among the deep conceptual dependency structures that he derives. It's easily criticizable. There are many things where each of us would differ on whether that is the path or exactly what is going on. But I admire it very much because it is a very ambitious attempt to deal with a fairly large piece of text. The point is we don't have models for dealing with large amounts of text, or even with several paragraphs of text. We don't really know how to handle the semantics of text discourse.

The basic problems are: How do you use sentences to understand following sentences? How do you solve problems of reference? And what is an adequate lexical structure to do this? What is an adequate semantic structure to represent the resulting meanings? That whole family of questions is seeking answers.

We are very much concerned in my group at Texas to study all these things. We have one student (we haven't got very much budget for this) but one student is working on discourse analysis; a dissertation by Hendrix [10] is about finished on modelling techniques -- an approach toward representing meanings as set-theoretic expressions. It takes semantic nets right down to the abstract algebraic description of what the meaning of sentences is.

And what I will talk about today is Robert Amsler's studies of the Merriam-Webster dictionaries.

With that introduction, let me shift to the business of computational lexicology. I think 'lexicology' is a fine enough word. It is in contrast to lexicography -- studying how to put together meanings to make a dictionary is lexicography. Lexicology is, as best I understand it, the study of how the meanings are organized; what is the hierarchical structure; how much in the way of loops and things of this type occur; and how can you take and transform from this list, 150,000 entries, alphabetically organized, into some organization that is computationally more significant. That, essentially, is the study. The desired outcome is to transform the dictionary into data that is reasonably computable; and hopefully to get semantic models out with considerably less difficulty than is currently the case.

Now I want to talk about the kinds of semantic models. There are really three kinds of data that we are able to get, from the English lexicon. One is a model for the meaning of a word. And I will show examples of what I mean by model shortly. The next is a model that for each sense meaning of a word can detect the context patterns that will select that sense rather than some other. The third thing of course is we want to organize the whole set of words that use that word in its definition. Now these things will, I hope, become clearer in a little bit.

First, to the dictionary. I think we know from a linguistic point of view that the proper source of semantic information is a large corpus of ordinary usage of the language -- a very large corpus; the Brown-Kucera corpus [11] is one example; perhaps larger ones are needed. In about 1966, while I was still at System Development Corporation, with John Olney, John put in a proposal to NIH to keypunch Webster's Collegiate Dictionary. He managed to get the grant, and it cost about $35,000 and about two years to get the thing keypunched with great accuracy -- 1 error in a thousand strokes, or something of this level was the quality control on it. And he incidentally keypunched Webster's Pocket Dictionary [12]. Now the mass of data here is quite significant. The Webster's Collegiate currently resides on eight tapes. Now several years later and we could probably get that on 2 or 3 tapes of high density form. The pocket dictionary fits comfortably on one tape, but it pretty much fills it up. Whenever one does some computation with these data bases it is likely to be very expensive. So I have really spent a great deal of time over the past 5 years at Texas discouraging people from running these tapes through the computer unless they had a really strong hypothesis and a well-formulated plan.

Now, why this particular lexicon? I think somewhere we are aware of the fact, for example, that Random House also has a keypunched version of a rather large dictionary [13]. And some work has been done on that too. John Olney liked the Webster's

dictionary because he went up and visited the people in Springfield; Gove was editor at that time; and John was very much impressed by the way that they went about making their dictionary. What they have done, is to accumulate citations for perhaps fifty or one hundred years, and they have room after room of files filled with citations of usages of words. So there is the corpus on which that dictionary is based. Now many not-so-careful dictionary makers don't do that. You know their offices are typically -- a dozen filing cabinets and a large shelf of dictionaries [laughter] and they sort of rewrite the definitions from one dictionary on an intuitive basis or something into another form. Well Merriam-Webster's is a much more careful operation. People collect these citations and the process of making a dictionary from the citations is as follows:

Let us suppose that for the word 'move' you might have a thousand citations. And so you start putting them into heaps of similiar usages. You sort them out into pile and pile, pile, pile; and then you decide how to describe each pile. OK, the resulting description is a sense meaning for one usage of the word 'move'. Well, in Webster's Third International [14], in addition to presenting that description, a usage example is also included. And in the Collegiate it is very frequently the case that with the sense meaning there will be an example.

So what I am suggesting is that the Merriam-Webster's dictionaries are in fact based on a large corpus of English usages that are carefully collected. Now I suppose the weakness here is that it is all written usages, I imagine, I am not sure -- rather than spoken language. When a sense description is compiled with a typical usage example one is not quite sure how typical that example is and not quite sure how free it is from particular biases of a particular lexicographer who is liable to be an English teacher brought in every few years to go through the massive effort of producing a new edition of the dictionary. I think there is a new one due out in 1976, a Fourth International. Well, I am not too close to that, but linguistic quality is the reason that John Olney chose the Merriam-Webster's dictionary as a data base that was suitable for keypunching.

Now we have been poking around with dictionaries for quite some time, and they are really hard to work with. I might add by the way there are about 20 copies of these tapes around the country in various research groups, and not too much has been done in the seven years that that corpus has been available. The reason is because it is very very difficult and expensive to work with it, and very hard to know what you are doing while you work with it. And I am sad to say you may get that idea very clearly by the time this lecture is finished.

Now let me show you some of the things that emerge from Amsler's work on this dictionary. One of the first things that we have all noticed in the dictionary is that there appear to be hierarchical properties. Words are defined in terms of other words. And there is a tendency to go up the tree to an increasingly more abstract word. So 'march' might be defined in terms of 'walk' and 'walk' eventually gets up to 'move'. Figure 1 is an example of that. 'Retort' is defined "to answer sharply". And 'answer' is defined in its turn "to write or speak in response". 'write' is "to communicate in print" and to speak is "to communicate by voice." Finally, "to communicate" goes to "to make known" so we notice that the hierarchy is quite regular and we can now define "to retort" as "to make known by voice, in response, sharply."

Now what has been done? Well, we have taken advantage of the lexicographer's rather careful behavior to define a word in terms of a superclass word, a higher, more abstract term, while carefully setting out differentia as modifying phrases. So "to retort" is "to answer sharply." "to answer" is "to write or speak in response." So we are able in this hierarchy simply to lift up the differentia and to define everything in terms of the top-level verb. Well this is rather interesting. Many of us are aware of Schank's assertion that fourteen primitive verbs will account for all of the verb usages in English [15]. I doubt the number fourteen. On the other hand, it is fairly clear that a single primitive accounts for hundreds of verbs.

insert Figure 1 here

In these studies, we first looked at some 600 verbs of communication. In order to do that we inverted the dictionary so that for each word that occurred in a definition we had an entry; and associated with that entry were the definitions that it had occurred within. So by inverting the dictionary file we were able to sort things and discover just how words go up the tree — and sometimes around in circles. In the 600 verbs of communication there are 3 senses of meaning for communication and we haven't studied them in exhaustive enough detail to be absolutely sure, but it is my impression that 2 or 3 models of 'make known', i.e. 'communicate', will account for some 600 verb meanings by following this pattern of carrying up the differentia.

Most recently, Amsler has studied 200 verbs of motion; and the way he got these was a little bit simpler. The communication verbs are a deep hierarchy. On the verbs of motion, he simply took from the inverted file all of the verbs that use the word 'to move' as the defining kernel, — i.e. as the superclass; and those go up to what appear to be some 5 or 6 sense meanings of 'move.' I will show you these sense meanings shortly.

FIGURE 1.

USING A HIERARCHY OF VERBS

RETORT  -  TO ANSWER SHARPLY

ANSWER  -  TO WRITE OR SPEAK IN RESPONSE

WRITE  -  COMMUNICATE IN PRINT

SPEAK  -  COMMUNICATE BY VOICE

COMMUNICATE  -  MAKE KNOWN


TO  RETORT  $\Longrightarrow$ MAKE KNOWN, BY VOICE, IN RESPONSE, SHARPLY

$\Longrightarrow$ MAKE KNOWN, IN PRINT, IN RESPONSE, SHARPLY

7

The point of the hierarchy of verbs, particularly with reference to collecting semantic information, is: how can we sort the material that occurs in a good dictionary so that conceivably we can make models for 20, 50 or 100 units of meaning, and classify everything that is defined in terms of those words underneath them. There will be a great savings computationally if we are able to do this.

Let's see what I mean by a model. Let's consider a sentence, "John wired a greeting to Mary." We are still in the verbs of communication here. "to wire" is "to telegraph". "to telegraph" is to communicate, _instrument_ telegraph, _medium_ telegram. Now the representation of meaning that I found quite useful for answering questions is the semantic net structures that I have talked about many times but the notation is simpler now. In Figure 2 C2 is a token of 'communicate'; the verb used was 'wire'. The tense is past, Actant or Agent - John, Theme - a greeting. Source - John. Goal - Mary. Instrument - telegraph. Medium - telegram. That is what we can get out of a shallow level analysis of the sentence "John wired a greeting to Mary". Now that is just the first stage. I am not quite sure of terminology yet, but I think of that as a shallow semantic analysis or a shallow semantic structure. I used to call it a deep case structure but it is really not very deep at all. Those of you who are acquainted with the dependency structure can see that this head verb is simply in a flat tree dominating all of these arguments. Now the structure is really quite syntactic although it has semantic aspects. I also assume, by the way, that we have selected the particular sense meanings of these words.

insert Figure 2 here

Now we want more. First of all, associated with the definition of 'wire' are things that will take it up to 'communicate', and add arguments to the meaning of communicate. But now associated with communicate also in addition to the syntactic data that is needed to sort out its arguments will be a set of assertions. If somebody is to communicate something it is an event that occurs in time. So at the initial time, t1, the agent or actant knows the theme. John knows the greeting he sends. I think 'know' will be a primitive.

At t2 he sends the medium which in this case is the telegram. It might have been a letter. It might have been anything that carries information in regard to a communicate verb. At time t3 it is possible the goal person will receive that telegram; at time t4 it is possible that the goal person will know the message. So the final values of each assertion are: true, true, true, possible, possible. It is also the case that John, the agent, wanted the goal person to know the message and he wanted this at least from t1 to t3 when he sent it. Now that is a fair but not absolutely necessary type of inference. So what we are doing is to associate with 'communicate' a set of assertions defined over the argument

FIGURE 2.

A COMMUNICATION AND SOME OF ITS MEANING


JOHN WIRED A GREETING TO MARY


TO WIRE  $\Rightarrow$  TO TELEGRAPH

TO TELEGRAPH  $\Rightarrow$  TO COMMUNICATE, INST: TELEGRAPH, MED: TELEGRAM


(C1 TOK  COMMUNICATE, VB  WIRE, TENSE  PAST,
    A JOHN, TH  GREETING, S  JOHN, G MARY,
    INST  TELEGRAPH, MED  TELEGRAM)


COMMUNICATE: ASSERT (( KNOW  A  TH  T1  TN  T)
                     ( SEND  A  MED  T2  T3  T)
                     ( GET   G  MED  T3  T4  P)
                     ( KNOW  G  TH  T4  TU  P)
                     ( WANT  A  ( KNOW  G  TH) T1  T3  T))


INSTANTIATION:

    (( KNOW  JOHN  GREETING  T1  TN  T)
     ( WANT  JOHN  ( KNOW  MARY  GREETING) T1  T3  T)
     ( SEND  JOHN  TELEGRAM  T2  T3  T)
     ( GET   MARY  TELEGRAM  T3  T4  P)
     ( KNOW  MARY  GREETING  T4  TU  P))

9

variables.  And when these are applied to the particular  usage  we
discover that John knows the greeting at time t1 to tn, which means
that  he  knew the greeting all along, and John wanted Mary to know
the greeting from t1 at least until he sent it, and John  sent  the
telegram  from  time interval t2 - t3; it is possible that Mary got
the telegram during time t3 - t4; and  it  is  possible  that  Mary
knows the greeting at some time t4 to tu'.  It is also important to
point  out  that  t1 is before t2, t2 is before t3, etc., and tu is
before now, are all assertions that hold.  So  what  is  going  on
here?  This  is  really what we mean by a model.  We mean that the
assertions that are explicitly made in the sentence are one part of
the model of the meaning of that sentence and the  assertions  that
are  implied  in the setence are another part of the model.  So the
lexical  structure  that  we  seek  would  make  those  explicitly
available.  Now  we  know  from  the work of Charniak and Schank's
students, Riesbeck[16],  Rieger[17],  and  Goldman[18]  that  these
things  are  definitely  going  to be useful in discourse analysis,
although I confess I don't know how to do discourse  analysis  with
this  type  of  assertion yet.  That is another line of research we
are very much concerned with.

     When we analyze a mass of verbs one of the things that we  are
attempting  to  do  is to discover the case arguments that go along
with each.  For the verbs of communication (see  Figure  3),  there
are:  agent - which is an animate organism; theme - contents of the
communication;  source - the person or system that sends it; goal -
the receiver; manner - for example, sharply, harshly; instrument  -
telephone,  telegraph,  voice; depth and  length  -  long, short;
frequency - repeatedly; intensity - loudly; and medium -  the  form
of the communication.


                        insert Figure 3 here

     I think that is enough for communication verbs, as a notion to
get started with.  Figure 4 is an example of the data that has been
put  together — This is one page from the 200 verbs of motion, and
what has gone on here, is shown for example in 'paddle 2.1'; moves;
theme - the hands and feet; path - about; medium - presumably water
or something of this sort; So what is done is to take the verbs and
using an ordinary editor program go through with as  little  change
in  the  definition  as  possible,  and  mark  the category of each
distinguishing argument for the usage.  In fact, it takes a lot  of
sorting  in  heaps before one can get categories organized and then
one begins to  edit  the  definitions  to  get  them  into  uniform
structure.


                        insert Figure 4 here

# Figure 3.

## Arguments of Communication Verbs

AGENT   -   Person or Animate

THEME   -   Contents of Communication

SOURCE   -   Place or Person

GOAL   -   Person

MANNER   -   e.g. Sharply, Harshly

INSTRUMENT   -   e.g. Telephone

DEPTH/LENGTH   -   e.g. Long ... Descant

FREQUENCY   -   e.g. Repeatedly ... Nag

INTENSITY   -   e.g. Loudly ... Bellow

MEDIUM   -   e.g. Form of Communication - Letter, Telegram

Figure 4.

HURRY 3 - MOVE OR ACT (SPEED: WITH HASTE)
HURTLE 1 - MOVE WITH (SOUND: A RUSHING SOUND)(SPEED: RUSH)
INCH 2.1B - MOVE (SPEED: SLOWLY) <CARS INCHING ALONG THE SLIPPERY ROAD>
JERK 2 - MOVE IN (SUB-ACTS: (DISTANCE: SHORT)(ACCELERATION: ABRUPT) MOTIONS)
JIGGLE 1 - MOVE WITH (SPEED: QUICK)(DISTANCE: LITTLE)(STEADINESS: JERKY) MOTION
JOLT 1 - MOVE WITH A (ACCELERATION: SUDDEN)(STEADINESS: JERKY) MOTION
KEDGE 1 - MOVE (THEME: A SHIP) BY (SUB-ACT: HAULING ON (INSTRUMENT: A LINE ATTACHED TO A SMALL ANCHOR DROPPED AT THE DISTANCE ...
KEDGE 1 ... AND IN THE DIRECTION DESIRED)(NON-MOVERS: <2>>))
LASH 1.1 - MOVE (EMOTIVE/FORCE: VIGOROUSLY)
LABOR 2 - MOVE WITH (EMOTIVE/FORCE: GREAT EFFORT)
LUMBER 1 - MOVE (FORCE/EMOTIVE: HEAVILY) OR (EMOTIVE: CLUMSILY)
LURK 1 - MOVE (EMOTIVE: FURTIVELY) ; /SNEAK/
MARCH 1 - MOVE (PATH: ALONG) IN OR (EMOTIVE: AS IF IN MILITARY FORMATION)
MIGRATE 1 - MOVE FROM (SOURCE: ONE COUNTRY, PLACE, OR LOCALITY) TO (GOAL: ANOTHER COUNTRY, PLACE OR LOCALITY)
MILL 2 - MOVE (PATH: IN A CIRCLE) OR (MEDIUM: IN AN EDDYING MASS)
MOD 2 - MOVE (PATH: UP A DOWN) <THE TULIPS NODDED IN THE BREEZE>
NOSE 2 - TO PUSH OR MOVE WITH (INSTRUMENT: THE NOSE)
NOSE 6 - MOVE (PATH: AHEAD)(SPEED: SLOWLY) <THE SHIP NOSED INTO HER BERTH>
PADDLE 2.1 - MOVE (THEME: THE HANDS, A FEET)(PATH: ABOUT) IN (MEDIUM: SHALLOW WATER)
PADDLE 2.1 - MOVE (PATH: ON) OR (PATH: THROUGH)(MEDIUM: WATER) BY USING (INSTRUMENT: A PADDLE) OR (MANNER: AS IF BY USING A PADDLE)
PASS 1.3 - MOVE (PATH: PAST)(PATH: BEYOND), OR (PATH: OVER)(NON-MOVERS: <1>>)
PLAY 2.2 - MOVE (MOTIVE: AIMLESSLY)(PATH: ABOUT) ; /IMIFLE/ <PLAYS WITH A RING NERVOUSLY>
PLAY 2.6 - MOVE OR OPERATE IN A (SPEED: BRISK)(PATH)(PATH: STEADINESS/FREQUENCY: IRREGULAR), OR (STEADINESS/FREQUENCY: ALTERNATING) MANNER
POUND 4 - MOVE OR MOVE (PATH: ALONG) (FORCE/EMOTIVE: HEAVILY)
PROGRESS 1 - MOVE (PATH: FORWARD) ; /PROCEED/
PUTTER 1 - MOVE ON ACT (MOTIVE: AIMLESSLY) OR (MOTIVE: IDLY)
RATTLE 2 - MOVE WITH (SOUND: A CLATTERING SOUND) <RATTLE DOWN THE ROAD>
RECEDE 1 - MOVE (PATH: BACK) OR (PATH: AWAY) ; /WITHDRAW/
RECIPROCATE 1 - MOVE (PATH: BACKWARD A FORWARD) (STEADINESS/FREQUENCY: ALTERNATELY) <A RECIPROCATING MECHANICAL PART>
REMOVE 1.1 - MOVE FROM (SOURCE: ONE PLACE) TO (GOAL: ANOTHER PLACE) ; /TRANSFER/
REMOVE 1.2 - MOVE BY (SUB-ACT: LIFTING) OR (SUB-ACT: TAKING (PATH: OFF) OR (PATH: AWAY)(NON-MOVERS: <1>;)
REVOLVE 12 - MOVE OR CAUSE TO MOVE (PATH: IN AN ORBIT) ; <ALSO> ; /HUIATE/
RIDE 2 - TO FLOAT OR MOVE (PATH: ON (MEDIUM: WATER) <RIDE AT ANCHOR> ; <ALSO> ; TO MOVE (MANNER: LIKE A FLOATING OBJECT)
RISE 7 - MOVE (PATH: UPWARD) ; /ASCEND/
ROCK 1 - MOVE (PATH: BACK A FORTH) IN (INSTRUMENT: CRADLE) OR (PATH: AS IF IN A CRADLE)
ROLL 1 - MOVE BY (SUB-ACT: TURNING (PATH: OVER A OVER))
ROLL 3 - MOVE (PATH: BACK) OR (PATH: AWAY) ; /WITHDRAW/
RUN 15 - MOVE ON (INSTRUMENT: WHEELS)
RUN 7 - MOVE IN (QUANTITY: SCHOOLS) ESP. TO (GOAL: A SPAWNING GROUND) <SHAD ARE RUNNING>
RUSH 1 - MOVE OR (INSTRUMENT: WHEELS) OR (FORCE/RESISTANCE: AS IF ON WHEELS) ; PASS (RESISTANCE: FREELY)
SCOUR 1 - MOVE (PATH: FORWARD) OR ACT (EMOTIVE: WITH (EMOTIVE: GREAT HASTE OR EAGERNESS OR WITHOUT PREPARATION)
SCREW 3 - MOVE (SPEED: RAPIDLY)(PATH: THROUGH) ; /RUSH/
SCUD 1 - MOVE OR CAUSE TO MOVE (PATH: SPIRALLY)
SHAKE 2 - MOVE (SPEED: SPEEDILY)
SHAKE 2 - MOVE OR CAUSE TO MOVE (STEADINESS: JERKILY) OR (PATH/STEADINESS/FREQUENCY: IRREGULARLY); /QUIVER/
SHUFFLE 3 - MOVE WITH A (SPEED/RESISTANCE: SLIDING GAIT) OR (SPEED/RESISTANCE: DRAGGING GAIT)
SHUTTLE 1 - MOVE (PATH: BACK A FORTH)(SPEED: RAPIDLY) OR (FREQUENCY: FREQUENTLY)
SIDLE 1 - MOVE (PATH: SIDEWAYS) OR (ORIENTATION: SIDE FOREMOST)
SKIP 1 - MOVE WITH (SUB-ACTS: LEAPS A BOUNDS)
SKULK 1 - MOVE (EMOTIVE: FURTIVELY) ; /SNEAK/ ; /LURK/
SLIDE 1 - MOVE OR CAUSE TO MOVE (RESISTANCE: SMOOTHLY)(PATH: ALONG A SURFACE)
SLINK 1 - MOVE (EMOTIVE: STEALTHILY) OR (EMOTIVE: FURTIVELY)
SMACK 2 - MOVE (THEME: THE LIPS)(MOTIVE/SOUND: SO AS TO MAKE A SHARP NOISE)
SMASH 2 - MOVE (PATH: FORWARD) WITH (FORCE: FORCE) A (SUB-ACT: SHATTERING EFFECT)
SNEAK 1 - MOVE, ACT , OT TAKE IN A (EMOTIVE: FURTIVE) MANNER
SPIN 6 - MOVE (SPEED: RAPIDLY)(PATH: ALONG)
SPIRAL 1 - MOVE (PATH: IN A SPIRAL COURSE)
SPRING 1 - MOVE (ACCELERATION: SUDDENLY)(PATH: UPWARD) OR (PATH: FORWARD) ; /LEAP/ , /BOUND/
SPRING 2 - MOVE (SPEED: QUICKLY) BY (INSTRUMENT: ELASTIC (FORCE: FORCE))
SQUELCH 2 - MOVE IN (MEDIUM: SOFT MUD)
STEAM 3 - MOVE BY (INSTRUMENT: THE AGENCY OF STEAM) OR (MANNER: AS IF BY THE AGENCY OF STEAM)

We are developing methodology — so we are working with the pocket dictionary; and it is quite clear that the pocket dictionary has abbreviated many, many meanings to the point where they are not useful. I think "scud — to move speedily" is a good case in point. 'March' is defined without even using the notion that it is on foot. Abbreviated definitions are fine for some purposes, but the Collegiate is much better for model making and of course the Third International is, as far as we can detect without looking at it exhaustively, quite good.

We also use the dictionary to develop context patterns that will identify the sense meanings for each word. Figure 5, for example, is an analysis of 'move' from the Third International which includes sentence examples.


Has a physical object changed location?
Yes, Has a physical object changed ownership or status?
       Yes, e.g. The Christmas items were moving rapidly.
            We just moved to town.
Has the physical object changed ownership or status, No.
Move 1 — e.g. The cars moved down the road.
       The chess master moved the chess piece.


That one was puzzling for a while [laughter]


Has a part of the physical object changed location. Yes.
Move 2 — e.g. The trees moved gently in the breeze.
       He moved restlessly in his sleep.
       He pressed the button and the machine began moving.
       And so on.


Here is a continuation of that chart.


Has the rate at which something was happening changed? Yes.
       e.g. The plot moved quickly. The melody moves upward.
Has some action which is a part of the plan or procedure
been proposed or performed? So we get into Move 6.
       Move for a recess.
       Revolutionaries must make their moves carefully.
       Moves in social circles.
Has someone's emotional state changed?
       Move to tears.


insert Figures 5a,5b,5c here

Fig 5a

MOVE 4,5,6

②

has The rate
at which SomeThing
was happening changed?

MOVED
IN
SPEED

yes

no

MOVE4

The plot, melody
moves quickly

for a while There
was noThing to do,
but suddenly Things
really began to move

is some action
which is part of
a plan or procedure
been proposed
or performed?

MOVED
IN
THOUGHT

yes

no

MOVE6

moved for a recess

revolutionaries must
make Their moves
carefully

moves in different
social circles

Did someone
emotional
state
change?

yes

MOVE5

moved to
tears

Fig 56

15

FIGURE 5C.

MOVE CONTEXT PATTERNS

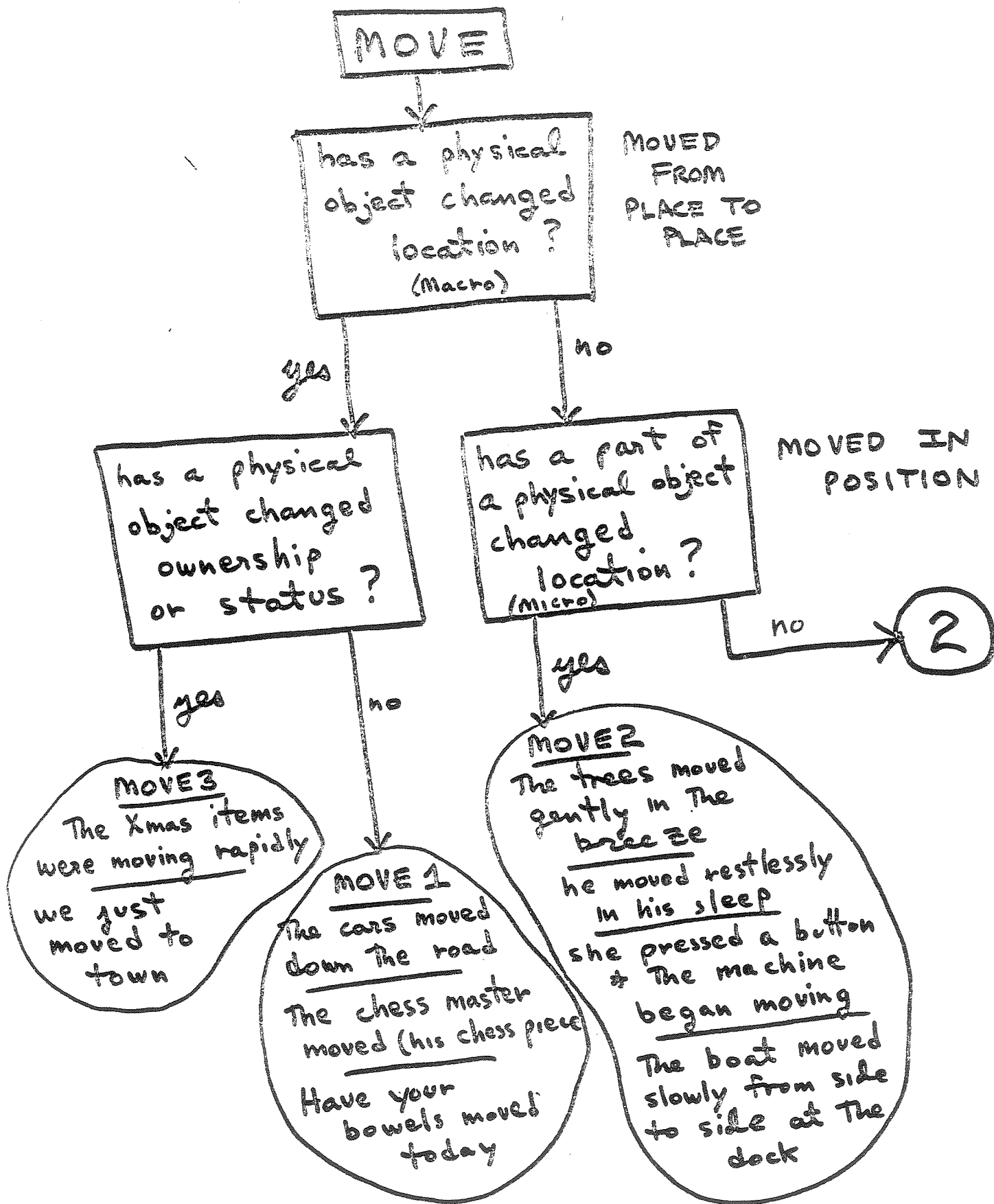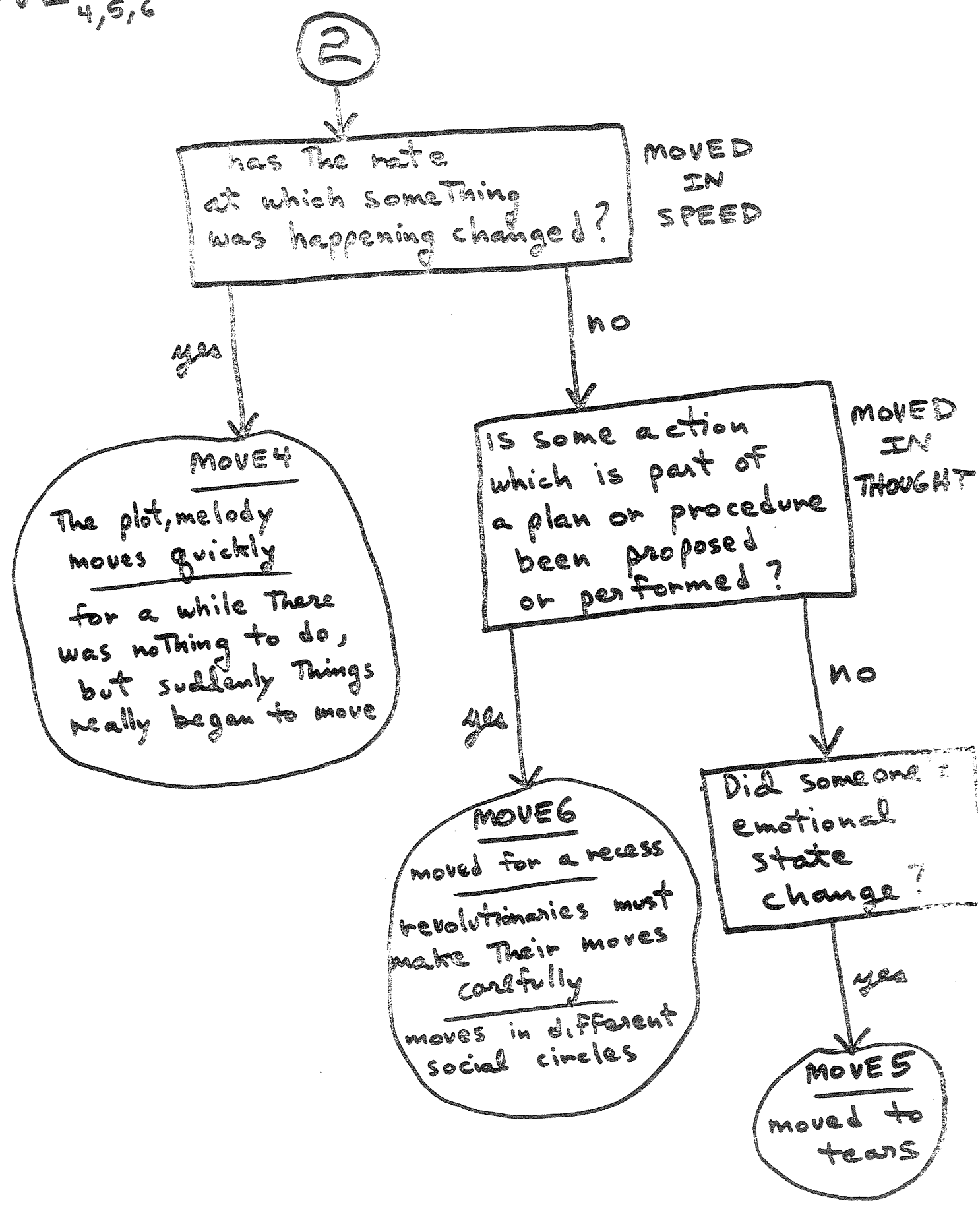| SENSE | AGENT | THEME | FROM SOURCE | TO GOAL | LOC | E.G. |
|---|---|---|---|---|---|---|
| MOVE1 | * | PHYS-OBJ | Loc1 | Loc2 | MEDIUM | TRAVEL, MARCH |
| MOVE2 | * | PHYS-OBJ PART | | | Loc1 | AGITATE, FIDGET |
| MOVE3 | ⟨ANIMATE⟩ | "CENTER OF ACTIVITY" | Loc1 | Loc2 | | MIGRATE1 |
| MOVE4 | | "PLOT, MELODY" | State1 | State2 | | |
| MOVE5 | | EMOTIONAL STATE | State1 | State2 | | TOUCH, PERSUADE, STIR2 |
| MOVE6 | ⟨HUMAN⟩ | "STATEMENT, RESOLVE" | | State1 | | "MAKE A MOTION" |

16

So there are six senses of 'move' that Amsler distinguished by studying the Third International. Now what does that mean computationally? In figure 5a we have move 1, 2, 3, 4, 5, 6. We have the basic, most important, arguments that usually occur. And here the agent must be human. Move 6, is to make a motion or resolution. So we must check the theme to make sure that it has some kind of a marker equivalent to 'statement' or 'resolve', a 'resolution' or something of this sort; and of course we will check the agent to determine that it is human. Emotional state, move 5; e.g. "the novel moved me to tears"; the criterion is that the Theme be a person and the Goal an emotional state. In terms of our experience with microworld modeling, it is usually sufficient to mark the nouns with particular semantic features in this kind of an area. But will that approach work all across the language? Semantic features are probably not sufficient, so after defining the criteria for selecting the particular senses of meaning; How does one mark the appropriate information on the nouns and distinguish the types of sentences that can be arguments for that sense? -- Well, we have learned a lot in terms of microworld models. And I think this will solve eventually.

In studying the verbs of motion Figure 6 is our current set of argument classifications: agent, theme, source, goal, instrument, path, medium, speed, acceleration, steadiness, continuity, force, resistance, and orientation. Now Amsler is making quite a detailed study of these arguments because he is still exploring the possibility that he will be able to take the definitions and push them onto a display scope and draw particular kinds of squiggles to represent in real time the meaning of the motion verb. Whether he will be able to do that or not remains to be seen. It is clear and easy for some things and not for others. Notice manner, of course, is liable to be an attitude, and 'thud' is "to move with a heavy sound". If one wants to make pictures of meanings, one needs more than a tempero-spacial frame of reference; one must also pick up the connotative, emotional frame of reference. In the 200 or so move verbs, there are at least a dozen that have sound and probably a couple dozen that have connotative things going on. In 'lurk', for example, 'slyly' is one of the connotations that goes along with it. So it is rather hard to draw pictures of all of the move verbs.


insert Figure 6 here


Figure 7 is an example of how we expect to use the move verb. "Arnold marched the army slowly through the countryside from New York to Montreal." First we give the shallow semantics of the sentence. C1 is a move. C1 represents this whole proposition, this whole idea. There is a moving going on. The verb that was used is "to march". Now I saved the verb because I just don't want to fuss around like Goldman does trying to find my way back to the

# FIGURE 6.

## ARGUMENTS OF MOTION VERBS

AGENT

THEME - (OBJECT MOVED)

SOURCE - FROM

GOAL - TO

INSTRUMENT - (USED FOR MOVING)

PATH - (COURSE OR PATH-OF-MOTION)

MEDIUM -

SPEED - (E.G. FAST, SLOW)

ACCELERATION - (E.G. SUDDEN)

STEADINESS - (E.G. STEADY, WITH JERKS)

CONTINUITY - (E.G. FREQUENTLY)

FORCE - (E.G. FORCEFULLY, SLIGHTLY)

RESISTANCE - (E.G. WITH FRICTION, AS IF ON WHEELS)

ORIENTATION - (SIDE FOREMOST - SIDLE)

surface representation. Some of you have read Goldman's excellent thesis[18] in terms of how you go from Schank's very deep conceptual dependency structure back into making sentences but in some cases one can save a lot of trouble by carrying the verb along, and not get into that. The tense - past; the agent - Arnold; the theme - C2 is a token of an army; now in fact, one needs tokens for all words but I am not showing the complete computational representation; I am trying to communicate it. I had to put a token for 'army' because there is a difference between what armies do in general and what this army did. So: 'move', by marching, in the past, done by Arnold, from New York, to Montreal, through the countryside, in military formation, speed - slowly. Now using the type of modelling that we discussed in the verbs of communication, the assertions associated with 'move' translate into this time-ordered series of assertions: - that it is probably true that Arnold was at New York from t1 to t-delta. It is true that the army was at New York t1 to t-delta. It is probably true that Arnold was in the countryside from t-delta to tn. It is true that the army was there. It is probably true that Arnold was in Montreal from tn to tu and it is true that the army was in Montreal from tn to tu. And t1 is before t-delta; t-delta is before tn; tn is before tu and tu is before now; are all true. So we have got time going on. Now 'slowly' is translated into "greater the move rate of armies than the move rate of C2"; C2 was the token of the army; it was a particular instantiation of the concept 'army' that occurred in this context. And "in military formation" during the whole period is a reasonable inference to make. Now I don't want to make this sound very cut-and-dried because there are a lot of other inferences that might be made. Notice, nobody said the army was on foot. And yet marched implied that it was. And nobody said that the army was tired when they got to Montreal. I don't know whether they were or weren't. But probably it is the case that you would want to make that inference to understand "Why did they sleep for 24 hours thereafter?" Well you need to make the possible inference that if the distance is from New York to Montreal and if one walks or if on foot, well then one will probably be tired. And if tired, one will want to sleep for a long time.


                    insert Figure 7 here



     I feel pretty good about how we are able to do things with verbs and I feel very hopeful that we can develop a methodology for sorting the dictionary to get this kind of information in fairly large quantities. With a generous sponsor some day we might just go at it and see how many primitives we in fact need for a particular purpose. The purpose is obviously important in terms of how you classify things.

     I don't know very much about nouns. We have struggled with nouns. They also occur in hierarchies in the lexicon. One

# Figure 7.

## A Move Model

Arnold Marched The Army Slowly Through The Countryside From
New York To Montreal.

(C1 TOK Move, VB March, Tense Past, A Arnold,
    TH (C2 TOK Army), MED Countryside, S New-York,
    G Montreal, MAN (in Military Formation), SPEED Slowly)

((ATP Arnold New-York T1 T-DELTA)
 (AT Army New-York T1 T-DELTA)
 (ATP Arnold Countryside T-DELTA TN)
 (AT Army Countryside T-DELTA TN)
 (ATP Arnold Montreal TN TU)
 (AT Army Montreal TN TU)
 (GR (Moverate Army)(Moverate C2))
 (IN Army Military-Formation T1 TN)
 (Before T1 T-DELTA)
 (Before T-DELTA TN)
 (Before TN TU)
 (Before TU Now) )

interesting fact is that the top of any noun hierarchy appears to
be an argument position in a verb. Once again, I haven't seen a
hundred nouns so I am not quite sure that this is always the case.
But in the few nouns that we have looked at and taken up the
hierarchy, we eventually reach an argument position. For example,
'message' is eventually defined at the top as 'that which is
communicated.' So it is an argument; it is defined as the theme of
a basic verb. Now a book (see figure 8) is all kinds of things.
It is print on pages of a given size. A handbook is a concise book
for memoranda or notes. A hardback is a book bound with hard
covers. A paperback is paper-covered. A primer is a book that is
used to teach children to read and so on. So summarized in a
network are all the things the dictionary describes about 'book'.
This doesn't focus on the hierarchical structure; this just shows
the main senses in which the word book can be used.


insert Figure 8 here


     To put it all together, as far as we know, we are going to end
up handling nouns the same way we handled verbs. That is, a noun
is going to be something that has a set of arguments. It'll be
defined in terms of the top level, a high level word, with
particular differentiating arguments from the meaning of that high
level word. And it too will be a predication like a verb. Except
of course it will fit into the argument position of the verb. So
the conclusion here is that I am talking about a good deal of work
that is in progress and I am not quite sure what the outcomes of it
will be. I think at this time it is really quite hopeful, because
it seems to be a fairly clear path to developing a methodology for
extracting from these large, large sources, carefully put together
of semantic information, that which we can use for question-
answering machines — machines that will eventually read text — I
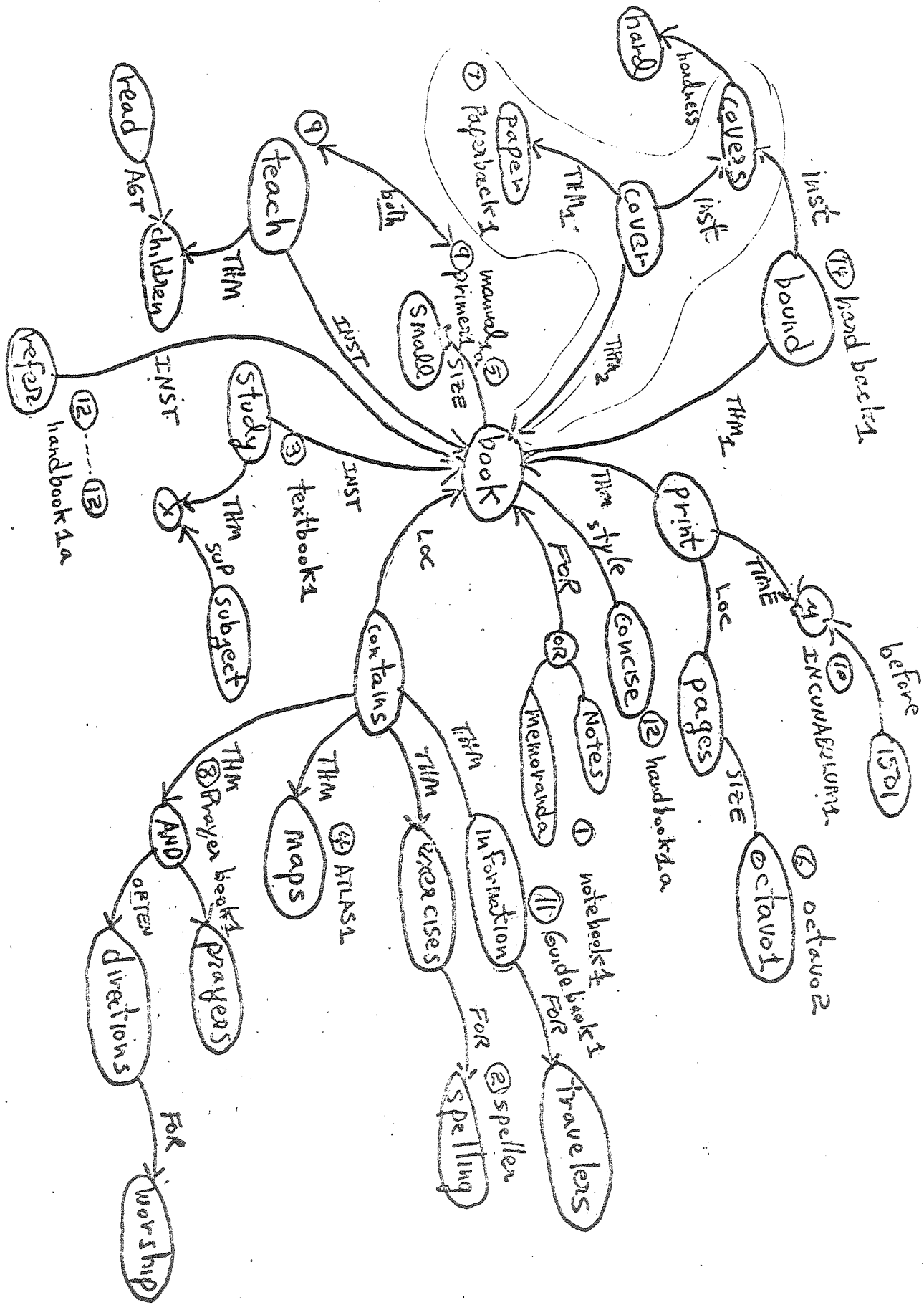still hope.


A P P L A U S E !  !  !

Fig. 8    "Book" Definitions

hard
hardness
covers
INST
cover
paper
Paperback1 ⑦
THM1
bound
THM2
INST
(T6) hand book1
(TP) INCUNABULUM
before
1501
INST
THM2
THM2
book
THM style
print
TIME
Loc
concise
Notes
OR
memoranda
FOR
pages
SIZE
octavo1 ⑥
octavo2 (K)
⑫ handbook1a
Loc
① notebook1
⑪ Guidebook1
FOR
② speller
FOR
spelling
travelers
Information
THM
exercises
THM
④ ATLAS1
maps
THM
contains
THM
⑧ Prayer book1
prayers
AND
OFTEN
directions
FOR
worship

teach
9
head
A6T
children
THM
both
manual1
primer1 size
INST
Small
INST
refer
⑫ --- ⑫
handbook1a
INST
Study
THM
sup
subject
③ textbook1
Loc

② octavo2

## Questions

**Hobbs**　　Concerning the decomposition of verbs, I guess the original motivation for Ogden's construction of Basic English [19] was reducing 1800 English verbs to 14 or so?

Is his work useful at all to you?

**Simmons**　　I've looked at that again and again and again; it has not yet been useful to me. I am not sure whether that is his fault or my fault. But I haven't yet been able to get use outof it.

**Hobbs**　　Concerning the nouns, is a reasonable way of looking at nouns to simply say, items in an index, or entries in an index, which point to a large number of facts? What you have to attempt is to organize these facts into various clusters of facts which are arranged in some sort of hierarchical kind of order according to the task that you are doing?

**Simmons**　　Sure, but isn't that the same thing that we are up against with verbs? We want to take the verb and transform it into a set of assertions, a set of facts. Similiarly we want to take the noun and transform it into a set of assertions, but sort of pick up the first level implications of its meaning. The critical question to me is how you use those sets of assertions to understand text. That is the problem of deep discourse analysis.

**Martin**　　And the first thing that is interesting to me is what a similiar path we have been down. I feel like I have done exactly what you have done except I had to do it by hand. I didn't have a dictionary. I pored over the dictionary at night, and my people have been poring over the dictionary. And just from hearing you talk I know that people here can't get an idea of how well this works. But this works actually quite well. It is very informative. And because I have done the same thing I can ask certain questions where I did things differently from you and I have wondered why. For one thing, in the case of 'book' - in the case of verbs too, I brought everything to a single level. Rather than saying there is a book and we know certain things about that; and then there is a handbook under that and there are certain things we know about that which are not true of books; it seems to be at the same level and the same thing with the

verbs. Rather than saying there are certain things about 'march' which have nothing to do with 'move', but yet they take some general properties from move; you have brought it all back to the same level; and that is one thing that I have done differently. I wondered why you did it that way. Do you see what I am saying? You tend to translate 'march' into 'move' plus a lot of arguments rather than put things directly on 'march' and then also know that it points up to 'move' and some of the properties come from there and some are directly on it.

Simmons      Yes. I am working on a hypothesis that we can describe the verbal meanings of English in something considerably less than the 20,000 verbs that English has. So I really want to be able to put the bulk of my effort on a few verbs and put the distinguishing features on the lower level verbs like 'march', and 'walk', and so on; and carry the distinguishing features up to the primitive model. I would like to have fifty verb models, if that is enough, and carry the distinguishing features up as added assertions. So that every word that is a descendent of one of these higher order verbs will have as true the assertions of the higher order verb plus its own distinguishing features.

Martin      Yes. That forces you to be able to do it in terms of models and features; that never allows you the out of saying that in fact I would be better off just to put a more or less description; but a much more ambitious way of trying to solve this whole problem.

Simmons      I think the difference, of course, is that you have a problem and an application directly under your fingers. You konw the pragmatics of what you are doing. We are in the more theoretical range of trying to see what we can do with the lexicon; and a general theme of how we are going to use models, rather than a particular application at this time. I think that would account for the differences.

Grishman      I would like to hold off on questions now, if I may. [due to schedule]

References

Grishman,R.,Sager,N.,Raze,C.  and Bookchin,B., "The Linguistic
String Parser." AFIPS. Conference Proceedings Vol. 42: AFIPS
Press, Muntvale,N.J.  1973, pp.  427-434.


[2]    Woods,William., "Procedural Semantics for a Question-Answering
       Machine," Fall Joint Computer  Conference, Proceedings 1968,
       pp.457-471.


[3]    Woods,W.A.;Kaplan,R.M.;Nash-Webber,B.  "The  Lunar  Sciences
       Natural Language Information System:  Final Report.  BBN
       Report No.  2378, BBN Cambridge, Mass.  June,1972.


[4]    Winograd,Terry.  Understanding Natural Language.   New York:
       Academic Press, 1972.


[5]    Heidorn,George E. "Natural Language Inputs  to  a  Simulation
       Programming  System."  NPS-55HD,  Naval Post Graduate School,
       Monterey, Calif.  1972.


[6]    Sager,Naomi, et.al., "An Application of Syntactic Analysis  to
       Information Retrieval," String Program Reports No.  6, N.Y.U.,
       Linguistic String Program, April,1970.


[7]    Harris,Zellig S., "Decomposition Lattices," T.D.A.P.  no.  70,
       U. of Pa., Linguistics Dept., 1967.


[8]    Charniak,Eugene  C.,  "Toward  a  Model  of  Children's  Story
       Comprehension."  AI TR-266, MIT, Cambridge,Mass., 1972.


[9]    Schank,Roger C., "Understanding Paragraphs," Instituto per gli
       Studi  Semantici  e  Cognitivi, Castagnola, Switzerland, 1974.
       (order from:  Centro di Documentazione della Fondazione Dalle
       Molle  per  gli  studi  linguistici  e  di  comunicazione
       internationale, Villa Barbariga, 30039 San  Pietro  di  Stra/
       Italy.) see also [20].


[10]   Hendrix,G., "Preliminary Constructs  for  the  Mathematical
       Modelling  of  English  Meanings."  University  of  Texas,
       Department of Computer Sciences, Working Draft, April 1975.
       (not for distribution).

[11]  Kucera,H.   and  Francis,N.,  Computational  Analysis  of
      Present-Day American English., Brown University Press., 1967.


[12]  Olney,J.;Revard,C.;Ziff,P.,   "Toward   the   Development   of
      Computational Aids for Obtaining a Formal Semantic Description
      of English," SDC document SP-2766/001/00, Oct., 1968.


[13]  Inquiries  regarding  the  availability  of  the  Random-House
      Dictionary of the English Language tapes should be directed to
      Mr.   Laurence Urdang, managing editor, Random House,Inc., 501
      Madison Avenue, New York 22, New York.


[14]  Gove,Philip  B.   (ed.)  Webster's  Third  New  International
      Dictionary (unabridged).   G.&C.   Merriam  Co.,  Publishers.,
      Springfield, Mass.  1971


[15]  Schank,Roger C., "The Fourteen  Primitive  Actions  and  their
      Inferences," Stanford Artificial Intelligence Laboratory, AIM-
      183, March,1973.


[16]  Riesbeck,  C.   "Computer  Analysis  of  Natural  Language  in
      Context," Ph.D.  Thesis, CS Department, Stanford, 1973.


[17]  Rieger,C.  "Conceptual Memory," Ph.D.  Thesis,  CS  Department,
      Stanford, 1973.


[18]  Goldman,N., "The Generation of English Sentences from  a  Deep
      Conceptual Base," Ph.D.  Thesis, Computer Science Departement,
      Stanford Univ.,Calif., 1973.


[19]  Ogden,C.K., The System of Basic English., Harcourt,Brace,  and
      Co., N.Y.  , 1934.


[20]  Schank,Roger C.  (ed.), Conceptual  Information  Processing,
      North-Holland Publishing Co., 1975 (inspress).