11p

# Computer Representation of Semantic Information

by Bertram Raphael

April 3, 1963

# Computer Representation of Semantic Information*

## I.  Introduction

A major obstacle in the development of learning machines, mechanical translation, advanced information retrieval systems, and other areas of artificial intelligence, has been the problem of defining, encoding, and representing within a computer the "meaning" of the text data being processed.  Various devices have been used to avoid this problem, but very little work has been done toward solving it.  The purpose of this memo (and the thesis research with which it is associated) is to describe one possible solution, and report on a computer program which demonstrates its feasability.

## II.  Semantics:  meaning and models

"Semantics" has variously been defined in the literature as "the study of meanings" or "the study of models".  In this section we shall try to explain, though not necessarily resolve, this apparent ambiguity.

Under "meaning" we include all those vague notions which are sometimes distinguished as the "meanings" of words, the "significance" of phrases or sentences, the "standard interpretation" of symbols in an artificial language.  Our basic definitions follows:  Definition 1: The meaning of a message is the set of ideas about objects and relations among objects which the message evokes in a typical human receiver. (We shall leave open, for the present, several questions which this definition suggests, such as how to define "typical human", whether one needs a concept of "meaning" for animals (and computers), and why "set of ideas" is a better intuitive notion than "meaning").

---

A <u>message</u> is generally expressed in some <u>language</u>, which has its own <u>grammar</u> (or <u>syntactic structure</u>). The meaning of a message, however, is a set of mental images; while it clearly must depend on the syntax of the message, the dependence may be quite indirect and obscure. Therefore, in many specialized areas of information (message) processing people have found it useful to employ <u>models</u>. <u>Definition 2</u>: A <u>model</u> for a class of messages is an artificial representation of at least some of the information content of the messages.

A <u>useful</u> model is one which has the following properties:

a. Changes in the model reflect changes in the original messages in a simple, well-defined manner.

b. Changes in the model reflect changes in the <u>meanings</u> of the original messages in a simple, well-defined manner.

Thus one may better understand the changes in meanings of certain messages, under certain transformations, by studying the effects of corresponding changes upon an appropriate useful model.

<u>Examples</u>:

1. The meaning of a verbal statement of a plane-geometry problem includes the ideas of line -segments, connections, shapes, etc. The usual model is a (pencil or chalk) diagram. This is a useful model because it satisfies properties a) (every student learns how to translate between diagrams and good verbal descriptions) and b) (for most people verbal discussions of geometric relations evoke visual images which correspond quite closely to drawn diagrams).

2. The nature of problem-solving ability in human beings at present is not known (the answer to the question, "Why can people solve problems?" does not have any general accepted meaning). However, the GPS computer program of Newell, Shaw, and Simon reflects the behavior of human subjects in certain problem-solving tasks in a direct way. If we accept GPS as a useful model of human behavior, we must also accept it as a theory of

the meaning of problem-solving ability (i.e. a close representation of
the information-processing procedures involved in human problem-solving).

3. Logicians develop and study formal systems which have no meanings
other than their syntactic structures. However, occasionally systems are
developed in order to study the properties of external (usually mathematical)
relationships. On these occasions one says that statements in the formal
systems have meaning "under standard interpretation", in terms of the
ideas of corresponding realtions between objects (usually numbers).
Models are invented which satisfy our definition of useful models (they
correspond in a well-defined manner to statements in the formal system,
and their properties match our intuitive notions of the meanings, under
standard interpretation of those statements). Semantics, in Mathematical
Logic, is the study of these models. (A description of the structure of
these models is beyond the scope of the present memo.)

4. The purpose of communication in natural language is to convey
meaning, usually about objects in the real world. If one can find a use-
ful model applicable to a significant part of the information conveyed
in normal conversation, then one will have a valuable tool for dealing
with the problem discussed in section I of this paper. The remainder of
this paper describes one such model.

III  **A Model for Natural Language**

Words may be considered the basic symbols in most natural languages.
Certain words (usually nouns) are thought of as denoting objects in the
real world, or classes of such objects. Other words (usually verbs and
prepositions) denote relationships between real objects. We see no need
for defining "meaning" of such words in a more fundamental manner than
as the thought about real objects (or object-relations). (We neglect for
the time being the problem of abstractions.)

In our model the basic objects are words themselves; in particular,
those words which usually denote objects or classes of objects. The
structure of our model for relationships is similar to that used in

Mathematical Logic, where an n-ary relation (predicate) is represented as the set of those ordered n-tuples of objects for which the relation holds. However, in our model each object (word) shall be labelled with the names of those relations to which it belongs and those objects to which it is thus related. Therefore, the model consists of a set of words, and, associated with each word, a formal description of how it is related to other words in the model. Precise representations for these descriptions are discussed in section IV below.

To demonstrate the usefulness of this model, we must show

a) how statements expressed in natural language may be translated into the model, and

b) how the model may be related to the meanings of the statements.

For (a) above, it is clear that classical grammatical parsing of English sentences is not very useful. The part-of-speech classes and relations usually used are syntactic elements only remotely related to the relations associated with meaning. However, the fact that many forms of English sentences unambiguously determine various relations is intuitively clear. Formal systems of linguistic analysis based on ideas similar to the above have been suggested, e.g., in books by Fries and Reichenbach. Much work remains to be done in this field of "semantic parsing", the details of which we are happy to leave to ambitious, enlightened linguists. As a first step, we shall content ourselves with the following, trivial form of analysis: A list of sentence formats will be developed, each of which determines a particular relation (or, for ambiguous sentences, choice of relations) among the unbound variables in the format (e.g., the format "x is a y" indicates that the ordered pair $<x,y>$ is in the "subset" relation). A new text sentence may then be analyzed by searching the list of formats for one which matches it in a specified manner.

The answer to (b) above is apparent. The meaning of a word and its description, in the model, is just the meanings of the various elements of the description, associated with the meaning of the word. For example,

"the blue chair next to the lamp" could be represented in the model as the word "chair" and an associated description including the property-word "blue", the relation-words "next-to" and their object "lamp", and possibly other words whose meanings are ideas of the size, shape, etc., of chairs in general.

## IV    Computer Representations

Since the model consists of words and word-descriptions, computer representations are most easily developed with the aid of a symbol-manipulating programming language. For various reasons LISP was chosen as the principle language for coding this system. Since the original input data is to be sentences in natural English, the COMIT language was used, in an earlier version of this system, as a pre-processor. However, in view of the limited grammatical processing required by the "format" method described above, LISP functions have now been written which also handle this original processing.

The use of a model for meanings in natural language involves two tasks: translating from sentences into the model, and interpreting meanings from the model. The model should be chosen so as these tasks are of comparable difficulty, so that hopefully neither one will be prohibitive. Schemes which store complete input text (as the computer model of the text) must solve the complete problem of extracting meaning from text whenever interpretation of the model is required. At the other extreme, models which very closely mirror meanings are suitable in certain specialized contexts (e.g., the geometry diagram), but for arbitrary general text would place a tremendous burden on the translation process from the text into the model. The property (description) list structure seems to be a suitable intermediate model. In this case the input text processor establishes, by some sort of semantic parsing, what relations should be stored in the model. These relations are then placed on property-lists associated with object-words in the form of attribute-value pairs; the attribute is the name of the relation or property, and its value indicates the other object-words associated with the given one

in that relation (if any). Interpretation of the model can be then carried out by functions which operate on these property-lists.

A simple information retrieval (automatic question-answering) system based upon the above ideas is being programmed, and operates as follows: as each sentence is read, the input routines perform the semantic analysis (by matching formats) and invoke appropriate routines which act on the model. If the sentence is not recognized, it is rejected with an appropriate comment. If a declarative sentence is recognized, the annexing routine for that particular relation adds information to the model (after first checking to avoid contradiction or excess redundancy). If an interrogative sentence is recognized, the appropriate routines are called which search the model (and indicate the required response, or reason for failure). Imperative sentences may also be used if routines are provided which will modify the model appropriately.

## V   Present state of the system

The present system will recognize about twenty sentence formats for unambiguous declarative and interrogative statements about set relations, part-whole relations including numbers, and left-to-right spacial relations. Associated annexing and searching routines are provided, including some which consider interactions between the set and the part relations. The input routine always indicates which other programs are invoked, and those programs provide informative responses describing the actions they perform. One of the spacial-relations search routines constructs a diagram of the ordering, in which case inner parentheses indicated objects which are known to be adjacent. The ability of the system is illustrated by the following sample dialogue between a human (H) and the computer program (C). One should keep in mind that the bulk of C's effort is devoted to con-structing and searching property-lists.

H:  Every boy is a person.

C:  SETR (BOY PERSON).

I UNDERSTAND THE SUPERSET RELATION BETWEEN PERSON AND BOY.

H:  John is an example of a boy.

C:  SETR (JOHN BOY).

I UNDERSTAND THE SUPERSET RELATION BETWEEN BOY AND JOHN.

H:  Is John a person?

C:  SETRQ (JOHN PERSON).

YES.

H:  Is a person a boy?

C:  SETRQ (PERSON BOY).

SOMETIMES.

H:  Any person has two hands.

C:  PARTRN (HAND PERSON 2).

I UNDERSTAND THE SUPERPART RELATION BETWEEN PERSON AND HAND.

I REALIZE THE NUMBER RELATION BETWEEN 2 AND PERSON.

H:  How many fingers does John have?

C:  PARTRNQ (FINGER JOHN).

I DON'T KNOW WHETHER FINGER IS PART OF JOHN.

H:  A finger is part of a hand.

C:  I UNDERSTAND THE SUPERPART RELATION BETWEEN HAND AND FINGER.

H:  How many fingers are on John?

C:  PARTRNQ (FINGER JOHN).

HOW MANY FINGER PER HAND?

H:  Every hand has 5 fingers.

C:  I KNOW THE SUPERPART RELATION BETWEEN HAND AND FINGER.

I REALIZE THE NUMBER RELATION BETWEEN 5 AND HAND.

H:  How many fingers does John have?

C:  PARTRNQ (FINGER JOHN).

THE ANSWER IS 10.

H:  Who is president of the United States?

C:  SENTENCE FORM NOT RECOGNIZED.

H:  The card is just to the right of the pad.

```
C:    JRIGHT (CARD PAD).

      I REALIZE THE JRIGHT RELATION BETWEEN CARD AND PAD.

      I REALIZE THE JLEFT RELATION BETWEEN PAD AND CARD.

H:    The card is to the right of the pad.

C:    RIGHT (CARD PAD).

      THE ABOVE RELATION IS KNOWN.

H:    The card is to the left of the pad.

C:    RIGHT (PAD CARD).

      THE ABOVE RELATION IS IMPOSSIBLE.

H:    The pen is to the left of the pad.

C:    RIGHT (PAD PEN).

      I UNDERSTAND THE RIGHT RELATION BETWEEN PAD AND PEN.
      I UNDERSTAND THE LEFT RELATION BETWEEN PEN AND PAD.
H:    Is the pen just to the left of the pad?
C:    JRIGHTQ (PAD PEN).

      INSUFFICIENT INFORMATION.

H:    Is the card to the right of the pen?

C:    RIGHTQ (CARD PEN).

      YES.

H:    The pencil is to the right of the pen.

C:    RIGHT (PENCIL PEN).

      I UNDERSTAND THE RIGHT RELATION BETWEEN PENCIL AND PEN.

      I UNDERSTAND THE LEFT RELATION BETWEEN PEN AND PENCIL.

H:    Where is the pen?

C:    WHERE (PEN).

      SOMEWHERE TO THE LEFT OF THE FOLLOWING:  (PAD, PENCIL).

H:    What is the relative position of the pen?

C:    LOCATE (PEN).

      THE LEFT-TO-RIGHT ORDER IS:  (PEN (PAD CARD)).

      TO FURTHER SPECIFY POSITION YOU MUST INDICATE WHERE THE PENCIL
          IS WITH RESPECT TO THE PAD.

H:    The pencil is to the left of the card.

C:    RIGHT (CARD PENCIL).
```

C:    I UNDERSTAND THE RIGHT RELATION BETWEEN CARD AND PENCIL.

        I UNDERSTAND THE LEFT RELATION BETWEEN PENCIL AND CARD.

H:    What is the relative location of the pen?

C:    LOCATE (PEN).

        THE LEFT-TO-RIGHT ORDER IS:  (PEN PENCIL (PAD CARD)).


## VI   Future Goals

The objects of the research being reported here are to define _meaning_ and _useful model_, to derive an example of a useful model for processing information presented in the form of natural language, and to discover the feasibility, desirability, and limitations of that model by means of a computer program. The object is _not_ to produce a complete, working computer program for any particular application; therefore, we ask forgiveness for the trivial and fragmentary appearance of the examples included above and proposed below. We hope, however, that the ideas for memory organization and information-processing procedures discussed here (and in future reports on this work) will find wide application in many sticky areas of artificial intelligence research.

In the near future we plan to study the applicability of the property-list model for natural language processing with respect to the following problems:

1. Interaction: Each relation added to the repertoire of the system may interact, as far as the search routines are concerned, with many other relations. For example we would like the present system to respond correctly to: "Newsweek is a magazine. A magazine has pages. Newsweek is to the right of the telephone. Where are the pages?" We hope to get some feeling for the magnitude of the difficulties involved in coordinating relations.

2. Specific _vs._ general information: It is important to distinguish between facts which are generally true ("A boy has ten fingers.") and those which refer to specific instances ("A boy is running."). The cases may even be contradictory ("That boy has eleven fingers."). We believe

this can easily be handled by labelling attributes <u>specific</u> or <u>general</u> when necessary, but more study in this direction is necessary.

3. Imperatives: At present the model can be grown and searched, but not modified. Inclusion of a set of imperative statement formats, and associated routines which change existing model structure, is certainly called for. These are necessary to permit erasure of specific information as well as correction of false general information.

4. Wider applicability: The present system deals only with certain binary relations. We should also describe how to encode in the model trinary relations, unary modifiers, and other forms. (Note that we will not consider the analysis of tenses, compound clauses, etc., which are part of the semantic parser's problem. We shall pick up his presumed results and decide how to treat them in our model.)

5. Ambiguities: Certain syntactic ambiguities may frequently be resolved on the basis of semantic information which has previously been stored in the model through the use of unambiguous sentences. We expect to illustrate this principle in the computer program.

# Scanning Agent Identification Target