Session 9:    SEMANTIC RESOLUTION


THE NATURE OF MULTIPLE   MEANING[1]

Don R.   Swanson

Ramo-Wooldridge Laboratories


Now that the stage has been set in previous discussions by the picturing of polysemia as a "monster" or a "blank wall",   let me add that there isn't a great deal more to be  said about multiple meaning that isn't either obvious or else wrong,   but I shall in any case take this opportunity to prove the point--not that it hasn't been convincingly demonstrated by others  on many occasions in the past.

I propose,   first of all,   to cut the monster down to size by carefully defining a limited portion of the whole problem and ignoring the rest.   "Multiple meaning" here refers  specifically to a phenomenon in which several equivalents in one language are said to represent some particular word in another language; it is evident,   therefore, that the phrase can be defined only in terms of some bilingual dictionary.   It is further evident  that the concept "bilingual dictionary" itself implies an approximation,   since we are never guaranteed that an appropriate equivalent necessarily exists for all thinkable contexts of each source-language word.   Having made these two observations, since they are fundamental to the nature of "multiple meaning" problems,   we shall nonetheless accept the idea of a bilingual dictionary as a practical tool for language translation and raise no further questions on its inherent limitations.

Given a dictionary,   then,   it is presumed that most words are represented by multiple equivalents,   and the problem is that of developing mechanizable rules for selecting an equivalent appropriate to any given context.    The word "context" itself ordinarily must be taken in its broadest sense and is by no means necessarily limited to linguistic data.   However,   since most of the observable and tractable context resides in the linguistic environment of the word,  I shall not be further concerned here with factors such as the identity of the writer (or speaker) and the cultural or situational factors which may have influenced the meaning of the words that he chose.    Thus  I  am

---

arbitrarily ignoring part of the problem,    perhaps an encyclopedic part,    but cannot say in any useful way just how much,    unless at the same time I discuss extensively the purposes and goals of automatic translation.    You will note that I have avoided the term "semantics" and have carefully identified and defined a rather specific area within which to consider problems of "meaning".

Since much machine translation research has centered on texts and dictionaries limited to specific technical fields,    there is some tendency to lose perspective on the real magnitude of the multiple meaning problem.    It has been reasonably well proved by several groups that the problem is indeed rather minor if the subject matter is narrow enough,    and if not too much text is examined,    and provided one is disposed to be charitable in reading the output.    This phenomenon of limited multiple meaning in limited scientific subject matter has led many projects to what might be called the "Eureka" stage of operational translation.    Perhaps after laboriously programming word lookup (delivered many months or years late,    of course,    since we all have tended to underestimate this deceptive task),    cleaning it up by cleverly programming the machine to erase (somewhat arbitrarily) all but one of the dictionary equivalents for each multiple meaning word,    and then inserting "of" in front of genitive nouns,    and then finally,    when the program is checked out,    gazing in rapture at an output that indeed resembles English,    the project director leaps naked from the bathtub,    vaults the blank wall of polysemia,    and hurtles down the boulevard trailing like ticker tape in the wind some 300 yards of high-speed-printer run,    shouting "Bar-Hillel was wrong-- I've done it! "

Multiple meaning was not really intended to be ignored at the "Eureka" stage,    but it was taken more or less on faith that any kind of initial grip on the problem could be followed by a bit harder work on more and more text until some kind of vaguely asymptotic convergence was attained.    Now unless one's objectives and ambitions in machine translation research are narrow indeed,  I submit that the "Eureka stage" optimism on multiple meaning is unjustifiable and that initial success of a limited nature has tended to obscure the very great difficulties that must necessarily appear when the breadth of subject matter    encompassed as a translation objective and the desired quality are  significantly  increased.

Session 9:   SEMANTIC RESOLUTION

   Now to turn my attention from vague generalities to vague
specifics,   let me say a  bit about our research on multiple meaning
at Ramo-Wooldridge.

   We have until recently confined our experimental translation to
the field of physics,   but have designed and implemented a research
procedure appropriate to at least part of the far more complex prob-
lems of meaning which manifest themselves as broader subject
matter is translated.    Thus our ambitions and objectives in the field
of meaning extend well beyond the boundaries of the narrow subject
matter which has thus far largely formed the basis for our experi-
mental research and collection of data.

   I plan to discuss principally here the design of our experimental
research procedures and our formats for the systematic collection
and analysis of data on multiple meaning.   In the main,   our approach
is inductive and consists of examining data on numerous individual
multiple meaning problems and attempting to infer general patterns
therefrom.    This inductive procedure is supplemented by the investi-
gation of several hypothetical models and subsequent deductive analysis
to understand the consequences of these hypotheses.   I cannot present
at the present time any simple,   regular,   and understandable pattern
which forms the basis and framework for representing multiple mean-
ing problems; i. e. ,   we do not have the "answer".    But the patterns
that we expect to emerge,   and  which  we   make   adequate   provision
to recognize should they do so,   are those which are revealed by the
data format itself which we employ.    This format takes into account
resolution by purely syntactic phenomena,  and I shall not make a
point here of separating grammatical from non-grammatical multiple
meaning.

   First,   let us observe that with any given dictionary it may be
assumed that, for certain words,   the choice among the various e-
quivalents  listed doesn't really matter,   and these  ought to be  excluded
from multiple meaning analysis at the outset.    For other words the
choice matters for some contexts but not for all contexts,   and for
still other words a choice is always necessary.    Less evident and
more interesting is the fact that for certain words in some contexts
several equivalents are clearly better than any one alone,   since the
reader is given the opportunity to interpolate a meaning which may

not exist in the form of any single English word or phrase.    (To ascribe meaning to a collection of words above and beyond that attributable to any single word is a suggestion which, if taken seriously, could have as much influence on human translation as on machine translation.)   In our research procedure,   the foregoing types of problems are recognized and recorded for each  multiple  meaning occurrence during the postediting of machine-translated text; this in- formation  is thus made available in coded form for further machine analysis.

Since I am reporting here on work being done by a number of R-W people, I shall attempt to identify specific phases of this work with specific individuals wherever appropriate.

The work that I am describing now on the design of data- collection procedures and formats is being carried out by Mr. Steven Smith,  under Professor P. L. Garvin's direction.

For each particular occurrence of a given multiple-meaning word,  the following data are also observed and recorded.   First of all,   a single word or phrase determiner occurring elsewhere in the same sentence is sought,   and in most cases found; and essentially all of those  cases for which a single determiner cannot be found are capable of resolution through general knowledge of the subject matter, i. e. ,   the resolution is made on the basis of scientific usage.   I do not attempt to project this  observation beyond the bounds of the  physics text we have studied.   A class of problems occurs,  to be  sure,  for which a single determiner cannot be found and for which even within scientific usage a problem of multiple meaning exists.   We can say nothing further at the present time about this category of problem except to observe that it is separately  tabulated  for later analysis.

To return now to the collection of information on identifiable and specific determiners,   it is next ascertained whether or not the resolution is based purely on questions of syntax.   Such problems are relatively manageable but unfortunately in the minority.   If the determiners fall into a clearly limited class for which an exhaustive list can immediately be prepared,   we define such determiners to constitute a "closed" set.   If the determiner is a member of a potentially "infinite" class possessing certain easily recognizable attributes (for example,   the class of all inanimate nouns),   then we define  these determiners as  constituting an "open"  set.

### Session 9:   SEMANTIC RESOLUTION

A card is prepared for each Russian word showing:

a.  "Word class" (essentially part of speech) for both determiner and determinee

b.  Syntactic relationship within sentence between determiner and determinee

c.  The possibility of intervening material

d.  The possibility of determination by absence rather than presence of the determiner

e.  Translation information including instructions for suppression and rearrangement

f.   Hierarchical relations among multiple equivalents,   so that eventually the most general term, i. e. ,   one subsuming the others, can be selected if selection is not otherwise decidable

Prepositions are treated somewhat separately,   and the following data are recorded (this phase of our studies is being carried out by Prof. Garvin and Dr. Gerta Worth, assisted by Mr. Onischenko):

a.  The first important distinction made is whether or not the preposition is part of a governed structure.   Government by predicatives,  gerunds,  infinitives,  modifiers,  and nominals is common; but it is important to observe that not all occurrences of prepositions are associated with government structures.   Governed and non-governed occurrences of prepositions require significantly different treatment with respect to   both word order and selection of equivalent.    Most of our own emphasis in this study has been on government by predicatives.

b.  It is further recognized that some government relationships are mandatory and others weak or optional.

c.  Our research procedure is intended to test among other things the hypothesis that certain predicatives,  for example,   may govern whole classes of prepositions which,  with  respect  to   such government,   then behave similarly.

d.  It is recognized in the data format that prepositions may lose their character as such when bound into idiomatic forms; new word classes are assigned to these structures accordingly as being characteristic of the structure itself.

Now let me turn our attention from research procedures and formats for data collection to certain models and hypotheses on which we base an approach at least partly deductive in nature.

Session 9:    SEMANTIC RESOLUTION

It is almost self-evident that the longer the lexical unit stored in the dictionary,  the less is the problem of multiple meaning. Certainly,  if one could store whole sentences,  the problem would indeed be approaching a practical (though not theoretical) solution. The difficulty with storing all sentences (apart from practical considerations) is that one would like to be able to translate sentences not yet conceived.   At some point intermediate between words and sentences,  we might hope to minimize multiple meaning,  but still deal with a relatively small number of oft-repeated units,  by considering strings of words or word combinations.   This concept (словосочетание) is used by certain modern Russian grammarians, and is being extensively investigated within our project by Mr. Curtis Benster.

The basic assumption underlying this approach is that many words,  of both the Russian and English languages,  as of others,  do not in fact have a meaning apart from that which they derive from their association with other words having more definite meaning characteristics.   Alongside words that primarily convey meaning ("content words") are words that primarily function as  "structural words"--having syntactic or grammatical function but little or no independent meaning.

In Russian this distinction is seen in the tendency of form words (and inflectional forms) to cluster into structures around content words,  principally verbs,  although nouns and even adjectives also figure in this.    These structures have semantic properties that are not derivable from the sum of the component parts taken singly. Words,  for instance,  which on occasion do indeed have the meanings given in our textbooks,   have these meanings attenuated or quite transformed in such combinations.    In view of comparable processes in the English language,   they then cannot be rendered by the English equivalent that in other circumstances may translate them correctly. Numerous individual occurrences of word combinations are being studied with respect to both structure and meaning.

Quite independently from the foregoing research,   we are also investigating the possibility of discovering  semantic attributes   of words based upon certain automatically recognizable statistical features of the  context.    Our initial endeavor in this direction has

been to attempt to discover a classification system for nouns based upon their frequency spectrum of categories of modifying adjectives, these categories being automatically  recognizable.    Our objective is to determine whether the categories that emerge, compiled by a fully automatic process, then have significant impressionistic features in common and, more importantly, whether the categories so discovered are useful in resolving problems of multiple meaning.    In our initial attempt at this analysis, each of several hundred modified occurrences of a few dozen nouns (those which occurred most frequently in a body of 30, 000 words of text) were categorized according to the type of modifier.

The modifiers were divided into 10 machine-recognizable categories,   as follows:

<u>List of Modifier  Classes</u>

(1)  demonstrative pronouns

(2)  possessive pronouns

(3)  definite and indefinite pronouns

(4)  numerals

(5)  adjectives ending in    - НИЙ   (except   СИНИЙ)

(6)  adjectives ending in    - ШИЙ (except   ХОРОШИЙ)

(7)  adjectives ending in    - ОВЫЙ,   ЕВЫЙ,   ВСКЙЙ
    (except  ТАКОВОЙ,    НОВЫЙ,    ОДИНАКОВЫЙ)

(8) adjectives ending in  - ОННЫЙ  (except ИОННЫЙ, НУКЛОННЫЙ
    ЕКСИТОННЫЙ,), -РНЫЙ (except  ОДНОМЕРНЫЙ,  ТРЕХМЕРНЫЙ)
    -АЛЬНЫЙ (except  ОСТАЛЬНОЙ, НАЧАЛЬНЫЙ, ПЕРВОНАЧАЛЬНЫЙ),
    -ИВНЫЙ, -СТСКИЙ, -НТНЫИ, -ИЧЕСКИЙ .

(9)  remaining adjectives

(10) participles (counted only if immediately preceding the noun
    without intervening governed matter)

Impressionistic semantic common denominators of classes (5) - (9):

(5)  location in space or time, e.g. ,  НИЖНИЙ,  ПОСЛЕДНИЙ

(6)  comparative or superlative

(7)  having to do with nouns, i. e. , entities

(8)  technical terms

(9)  "wastebasket" class

Session 9:   SEMANTIC RESOLUTION

We are endeavoring to determine,   on a statistical basis,   which nouns have the greatest number of common attributes so far as the spectrum of modifiers is concerned.   This work is still in progress and the statistics at the present time are too meager to justify further comment.   It is planned to pursue this study further only when the tabulated list of modifier spectra can be produced as a relatively in-expensive by-product of our other more direct attacks on problems of multiple meaning.

A more detailed account of this work will appear in a compre-hensive progress report,  forthcoming within a few months.