ABSTRACT
        The development of a technique for a computer-based
content analysis of interview data is described. A preliminary
version of ANACONDA (ANAlysis of CONcepts by DAta-processing) is
presented, and empirical results are shown from the application of
the technique by independent coders to test material. Proposed
modifications and extensions of the system are also discussed.
(DGC)

# A COMPUTER-BASED CONTENT ANALYSIS OF INTERVIEW DATA: SOME PROBLEMS IN THE CONSTRUCTION AND APPLICATION OF CODING RULES

Bernhard Bierschenk

This report discusses the development of a technique for a computer-based content analysis. It presents a flow-chart of different stages in the de-signing of an Analysis of Concepts by Data-processing. The acronym ANACONDA is the name that has been given to this technique. A condensed preliminary version of ANACONDA is presented and empirical results are shown from the application of the technique by independent coders to test material. The entire test material has been checked, in order to obtain a faultlessly punched and coded text. Empirical data are presented from this check. In conclusion the next steps are discussed: (1) scaling of qualifiers and (2) construction of registers.

Keywords: Data-processing, concept analysis, intercoder agreement, interview texts, psycholinguistics

CONTENTS

## 1. SOME INTRODUCTORY NOTES

When one has a complex verbal material with a low degree of structuriza-
tion, a goal-oriented analysis becomes extremely onerous. As a result,
researchers often try to avoid this situation by constructing "assessment
scales" containing statements or questions with fixed alternative answers.
They are easier/to handle from the processing point of view, even though
the gathering of information in the form of e.g. an interview would have
been more suitable. Manual analyses of verbal data are in addition often
determined in far too great a degree by practical considerations rather
than by what is desirable from a scientific viewpoint. The result/of such
an approach is that the analysis is usually limited by simple frequency
comparisons, which means that much information that is relevant for the
investigation is lost. In order to avoid such rough analyses, the researcher
should create quantities of data with a high degree of structurization.

The computer can be used for the handling and processing of text. But
an automatization of text analyses presupposes a coding of text and com-
puter programs that direct the processing of natural language. Content
analysis techniques based on the use of computers have been developed
within widely separated fields (see Gerbner, Holsti, Krippendorff,
Paisley & Stone, 1969). If computer-based analysis techniques are de-
veloped, the content analysis process will become flexible and the technique
can thereby be used for processing large quantities of verbal data. In
addition it becomes possible to refine the analysis technique so that
better statistical models can be used than has so far been the case.
Computer-based analysis techniques permit a greater dispersion of in-
formation in complex material that is difficult to survey. This in its turn
produces much more detailed analyses than a manually conducted ana-
lysis would allow. But at the same time automatic storage and processing
of text presupposes that the storage takes place in accordance with a
given format. The format states how each individual element should be
stored so that different elements can be placed in relation to each other
by means of e.g. Boolean algebra.

## 2. COLLECTION AND PROCESSING OF EMPIRICAL DATA

The problem of gaining access to and a technique for computer-based
content analysis has featured largely in a research project on search and
steering strategies in educational and psychological research planning.
This project is financed by the Swedish Board of Education. The work
within this project was initiated with an interview study involving forty
randomly selected researchers working in departments of educational
and psychological research in Sweden (see B. Bierschenk, 1974). The
material was collected during May, 1973. The interviews were recorded
on audio-tape. The recorded interviews have then been written out by
secretaries, but without the use of any phonetical transcription. The
total material now comprises approximately 4,000 pages of text,
measuring 21x29 cm. These are first to be prepared manually by the
insertion of syntactical codes, so that the interview data can then be
processed by the computer UNIVAC CD 3600, without the original struc-
ture of the text being lost.

In order to be able to carry out an Analysis of Concepts by Data-
processing (ANACONDA), rules have been worked out for the manual
coding and a condensed preliminary version of ANACONDA will be pre-
sented in this report. Empirical results of application of ANACONDA to
test material by independent coders are presented. Furthermore, a
check of the punching and the punch cards has been carried out in order
to attain perfect material and the control is described in this report.

The check of the intercoder agreement in the application of ANA-
CONDA has been carried out by Berg and the results are reported in
detail in Berg (1974). The construction of rules for ANACONDA has
been worked out by the linguist I. Bierschenk. A preliminary version and
some evaluation data are presented in I. Bierschenk (1974).

## 3. PRINCIPLES FOR CODING NATURAL LANGUAGE

The text material of the interviews must be formalized, which means
breaking it down into smaller units such as clauses, phrases, groups of
words or single words. This does not mean, however, that this analysis
is based on syntax, but it is instead based on the conceptual context (mo-
dels, images) underlying an utterance. The conceptual basis is the meaning
of an utterance. We assume that there is a covered structure, namely the
speaker's thought (what he wishes to say) and an uncovered structure (that
which in reality is said and heard). If an utterance is to have "meaning",
it must be comprehensible, which need not be the same as being gramma-
tically correct.

Seen in this way, an utterance consists of a conceptualization. The
basic unit of the conceptualization is the concept. A single word can have
meaning (e.g. represent a conceptualization) when it is uttered in a par-
ticular situation, in a particular context to a particular person and it then
forms a concept.

In this context it can also be of interest to mention that in principle
the generative theories adopt the same point of view (cf. Müller, 1969,
p. 67). According to Schank (1972, p. 555), there are two elemental kinds
of concepts. A concept can be either (1) independent or (2) dependent. By
independent concepts are meant all that are interpretable in isolation. By
dependent concepts are meant attributes or modifiers. Components with
independent meaning are e.g. nominals (subject, object) and actions, often
expressed in verb form. Dependent concepts only become meaningful
through the concepts that are modified by means of the dependent concepts.
From the view point of psychology, an analysis of natural language means
that the elements of the analysis are regarded as concepts and not as words.
In his chapter on the "Identification of conceptualizations underlying natural
language" Schank (1973, p. 192) now speaks of "three elemental kinds of
concepts. A concept can be either nominal, an action or a modifier". The
independent concepts are those that alone produce conceptualizations.
Schank calls these "picture producers". In other words, they are different
kinds of nominals. While Schank (1972, p. 566) also considers verbs to be
independent concepts, he has changed his position in 1973 (see Chap. 5,
p. 192) where he writes: "An action is what a nominal can be said to be
doing". Verbs are classified now solely as expressions of action towards
an object or goal.

Language is an expression of process (actions, events, conditions and
relationships and associated persons, objects and abstractions). This

process takes place within a structure: the clause. The process itself is represented by the verb. Participators in the process are e.g. persons and objects. They take the role of agent and goal. This role-playing in relation to the verb is called transitivity.

This means that we cannot extract information from a text if we only work with individual words. When people utter a thought, this takes place as economically as the situation permits. In a dialogue between e.g. researchers with a common frame of reference, the researchers can easily communicate since they make use of such verbal representation that the same conceptualization is produced in both parties involved (cf. Miller, 1969, p. 67). Thus, the interview text concerned forms a manifest output of what is indicated by the speaker.

By conceptualization is meant the individual's use of certain rules for relating concepts, e.g. grammar. Conceptualizations may be simple or complex. In this way an utterance in an interview situation can be rich in simultaneously underlying conceptualizations and make it difficult to represent these in a sentence. Consequently a sentence in a text can contain many completely expressed ideas and idea relations. The condensed information, which is a result of the inherent economy in clause-linking thus only becomes available in a content analysis by means of a supplementation procedure. But in order that we may carry out our analysis, a starting point is needed from which we can build up the structure in an utterance. In this analysis we begin with the action or the verb. An action can be said to be something that an "agent" can achieve in relation to an "object". Agent is used in the sense "action centre" and object consists of the means or the goal of an action. In principle only two cases exist, namely (1) agent and object coincide and (2) agent and object consist of two separate units. Modifiers describe the attributes (qualities, relations) that characterize agent, object or action. For this reason adverbials will be grouped around the verb, while different attributes (qualifying and describing) are arranged around subject and object.

Early attempts at designing computer programs that used natural language were primarily concerned with syntactic analysis. Today researchers in the field of psycho-linguistics no longer maintain that syntactic analysis is essential for the development of computer programs. Some have said that it is perhaps not at all necessary to use syntactic analyses (see Schank, 1972, p. 555). But this is not to say that syntax cannot be a very useful aid in an analysis of complex interview material. Syntax implies sequence or the relation between the different parts of an utterance. There are fixed and mobile positions in this structure. If these positions are

made use of in an analysis of text by arranging the basic elements (concepts) in classes that have a certain defined relation to each other, this facilitates considerably the processes of producing sequences, irrespective of whether the sequences are produced manually or by computer, or whether syntactical relations or psychological relations are stipulated and used. Each concept category can thereby also be related to each of the others by means of the conditions that are specified for a certain defined analysis purpose. The purpose of our analysis is primarily to establish which actions (with or without explicitly stated objects) are carried out by researchers.

## 4. REPRESENTATION OF STATEMENTS

A formalization of spoken text naturally cannot take place independently of later stages in the planned analysis. Nor does it in any way replace categories. Categories are namely the links between the theoretical anchorage of the research problem and technical aspects of the content analysis. In principle a content analysis can be carried out in accordance with the three following basic models: (1) the association model, which presents information in the form of statistical correlations between observable and non-observable variables, (2) the discourse model, which studies information defined by means of linguistic relationships and presents these relationships in denotations and connotations and (3) the communication model, which describes information by means of process and control within a dynamic interaction system. The choice of model 3 includes models 1 and 2. But since there are no adequate mathematical models for the last (3) model, it is model 2 that is most appropriate, considering the interview strategy that has been used. The discourse model describes extra-linguistic phenomena. It reproduces (is representative of) events within the source of information (the researcher) and nominals occuring in the discussion ("discourse") that refer to, separate or connect non-grammatical objects or concepts (see Krippendorff, 1969, p. 102).

The model presupposes that rules are drawn up. It should be pointed out that this kind of structurization is not possible without a considerable investment of energy, labor and time. Only when a highly structured quantity of data exists can it be used in the practical research work and when we wish at accelerated tempo to extract different types of information. The creation of material characterized by a high degree of structurization should be of particular importance when it can be assumed that the material will be able to provide answers to future questions, an assumption which should be relevant in research contexts.

The researcher's perception (description) and evaluation of the process of formulating problems is of primary interest in this study, and consequently the theories and techniques presented by Osgood (1956, 1959), Stone (1966) and Holsti (1969) are well-suited for an analysis of the interview material, since it is on the individual researcher's subjective interpretation of a given situation and action that we wish to focus. The technique and program that are being developed are based on agent-action-goal-relationship and associated modifiers. In addition to this segmentation, "themes" are also extracted from the sentences, so that important information can be recovered or mediated in addition to the information that becomes

available by means of the basic paradigm. Figure 1 presents the different stages of the analysis.

## 4.1 Directions for writing out spoken text

The directions that were given for the writing down of the interview material have not included phonological transcription rules, but the importance of an authentic recording of the audio-tape material has been emphasized. All audible utterances have been written down, which means that the speaker's slips of the tongue, corrections and imcomplete sentences are included.

## 4.2 Marking of text

The main problem in an analysis of verbal data for the purpose of obtaining information is that it must be possible to select the relevant material (concepts and concept relationships) from a quantity of possibly relevant material. For this reason the parts of the text that are relevant for the analysis should be marked. The object of the analysis is the statements made by the persons interviewed and therefore in the marking phase all interview questions, arguments and counterarguments from the interviewer are denoted "non-relevant" text.

## 4.3 Segmentation of text

A collection of interview texts or any other texts can be extremely comprehensive. Information that is to be extracted from a quantity of text can in addition be very dispersed. In order to find relevant information, each individual interview must be gone through from beginning to end. The postulation that there should be a structure in a text perhaps seems somewhat trivial in this context, but it is not superfluous. It is the permanent structure (established among other things by the order of the interview questions) of the interview material that can be used for a division of the large amount of text into manageable sections. It was considered unrealistic (and proved to be so) to treat each individual interview as a unit when coding and for this reason the interview material was divided into seven question complexes.

## 4.4 The agent-action-object (AaO) paradigm

Even if there is no technique for an "objective" analysis that reflects "all" the dimensions of the language, it is nevertheless possible to make use of certain general paradigms for an analysis, treatment and structurization of verbal data so that information becomes available. Dimensionality is a
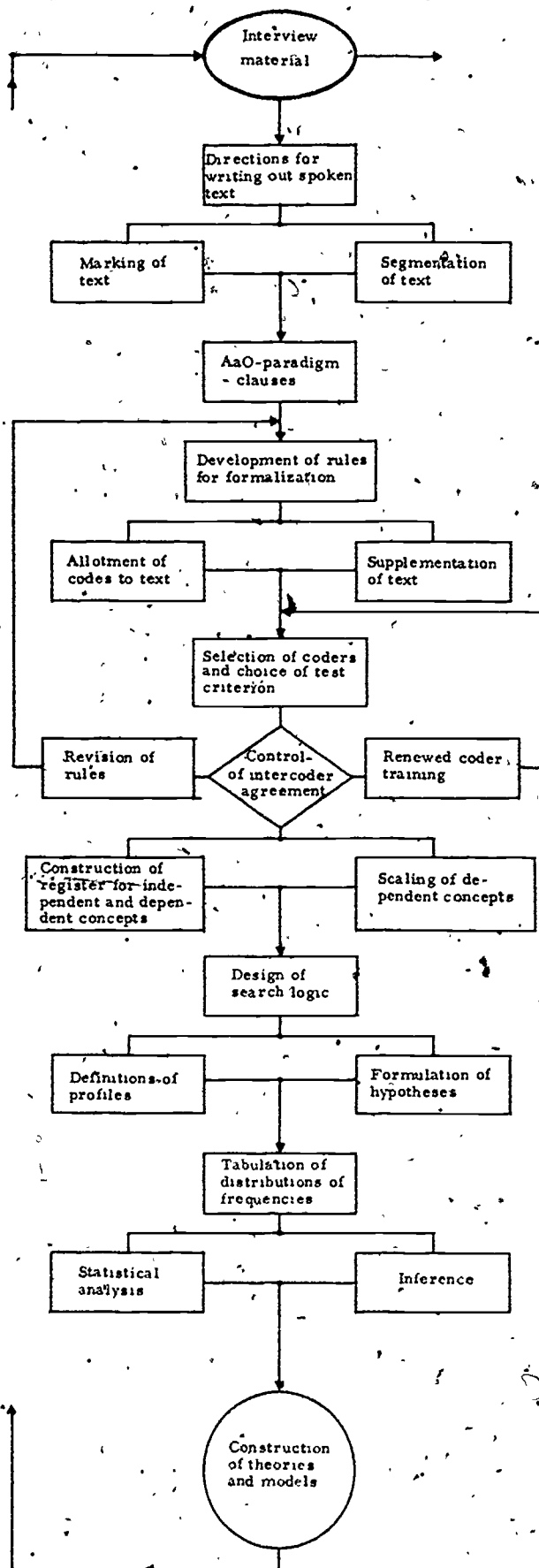
Figure 1. Flow-chart for designing a computer-based content analysis

central concept in behavioral science contexts. Horst (1968, p. 43) writes:
"The importance of this concept as a starting point for all psychological as
well as other scientific investigation is not generally recognized."

As has been mentioned earlier, every form of scientific investigation pre-
supposes that the researcher makes clear which aspects are to be mapped.
This in its turn leads to different models and data matrices, in which rows
usually represent the measurement objects of the investigation. Columns
refer to dimensions or attributes in which the objects of the investigation
("individuals") are to be measured. When using psychological tests, assess-
ment scales or questionnaires with fixed alte████████swers, one obtains
test values that can be used directly for sett█████data matrices. If on the
other hand one permits individuals to formulate their answers as they them-
selves wish, some technique is needed to help transform the answers into
scores. The development of such an analysis technique is planned. The
intended processing of the interview material is to lead to scores that can
be studied by means of statistical models for both "uni-dimensional" and
"multi-dimensional" analyses.

Each written or spoken text has a basic structure that consists of syn-
tactical units, which it should be possible to use in the construction of data
matrices. Each syntactical unit (phrases, clauses or sentences) consists
in its turn of "words" that are arranged in accordance with the structural
rules of the language. This relationship, together with the development of
computers, has led to a growing interest in recent years in the automatization
of content analyses. The computer's processing speed has resulted in
attempts also being made to develop algorithms with the help of which it
should be possible to identify "relevant information" as opposed to an iden-
tification of "words" as they occur in the text. (By algorithm is meant here
a mechanical method of approach for the transformation of utterances to
unambiguous analytical units.) For this purpose algorithmical codes must
be developed, i. e. codes based on rules for converting source material to
"equivalent terms".

The paradigm on which Osgood's "Evaluative Assertion Analysis"
(Osgood et al, 1956) is based is the AaO paradigm. Osgood presents a
method for separating attitude objects from "common meanings". Each
text is transformed into a succession of simple syntactical relationships.
The text is formalized according to the following model:

attitude object / connector / evaluating term
atitude object$_1$ / connector / attitude object$_2$

By the term attitude object, Osgood (1956, p. 47) denotes a message that can be limited in a general linguistic context ("common meaning materials"). Attitude objects are primarily nouns, which are placed in an evaluation scale dependent on predicative statements and attributes that are ascribed to these attitude objects in the text. In this way the method is considered to measure attitudes or evaluations of certain phenomena. The method is not efficient for our purposes, however, since it neither makes use of the entire text nor exploits the processing capacity of the computer.

A further development of Osgood's method into a computer-based analysis is presented by Holsti (1969). While the former method uses nouns and adjectives, the latter also takes verbs into consideration as being of importance for the evaluation of objects. Holsti's technique also permits coding of the theme of the sentence e. g. negation, tense and modality. The basis of this method is syntactical coding. The first step in Holsti's analysis is that one determines (1) agent and modifier, (2) action and modifier and (3) goal and modifier. By means of numerical codes, the agent (3) is linked with its qualifiers by e. g. a verb (4) with its qualifiers to the goal (7) with its qualifiers. This linking can be illustrated with the following concrete example from our material:

We (I) can say that / the resource situation / 3 has / 4 the whole time / 4 exercised pressure on / 4 problem definition / 7 problem limitation / 7. "I can say that" expresses the speaker's opinion on the continuation and need not be coded. Exercise pressure on cannot be separated, since the three words belong together within the expression (although pressure is not a verb). In addition the verb expresses direct action towards a goal. The 3-4-7 relationship states the direction of the action.

Holsti's analysis technique has been developed with the written text as starting-point. Such texts are carefully-prepared works, while our material is spoken text. Our starting point has been Holsti's agent-action-object-goal paradigm, but we wish to integrate Schank's psychological arguments into our continued work. When applying Holsti's method on the Swedish material, it soon became apparent however, that we needed to expand the codes. Codes for attributes to agent-action-goal and codes for different kinds of qualifiers were related to elemental concepts by means of a two-figure code system, in which the second figure states the respective modifier.

## 4.5 Development of rules for formalization of text

A system for the recovery of information should ideally be able to be used for every possible selection and not just for a selection of such material as

appears to be relevant on a certain given occasion. The demand on the
degree of structurization of the data quantity grows with increasing quanti-
ties of data, the frequency of use of the data quantity (search for information)
and increasing specification of problems. The concepts or statements
(conceptualizations) that exist in the data quantity must be predictable. A
computer-based presentation of relevant information is based on predictable
relationships between concepts and predictable statements. A preliminary
attempt at developing rules for an Analysis of Concepts by Data-processing
(ANACONDA) has been made (I. Bierschenk, 1974). Trial codings have been
carried out and intercoder agreement has been calculated (Berg, 1974,
I. Bierschenk, 1974). Using ANACONDA, about 10% of the interview mate-
rial has been prepared for manual processing. Some computer programs
for simpler processings are already in existence.

### 4.6 Allotment of codes to text

If we wish to guarantee that an analysis of texts leads to recovery of infor-
mation that is relevant to an investigation, we should choose for the ana-
lysis a technique that makes use of syntax. This means that the original
relationships that exist between concepts and statements are preserved.

In the analysis each concept holding information that has a different
function in the clause than the other parts has been picked out. Each "unit",
consisting of one or more words, is allotted a code. These codes are
divided into the two basic units of the analysis mentioned earlier, with the
following four main categories: 30, 40, 50, 70 (not 60, since 50 and 60 are
easily mistaken for each-other when data transcriptions are read). In
principle, these figures follow Holsti (1969). The first figure corresponds
to the following four categories: agent, action, object (means, goal). The
second figure states the independent or the dependent concept. The depen-
dent concepts that function as different kinds of qualifiers are allotted
figures that are immediately dominated by the independent concept: Agent,
for example, is given the figure 30 (0 = independent concept) and an adjec-
tive that describes the agent is given the figure 32. In addition to the indi-
vidual parts of the clause, the statement's tense, mood etc. (see pp 14-15)
are coded. In addition there are a number of codes for overall structures
that the coding of separate units cannot give. A theme has been defined as
a syntactic unit which consists of either (1) a main clause or (2) main
clause + one or more subordinate clauses. On the basis of ANACONDA a
test material has been prepared in accordance with the format presented
on pp 14-15.

| Identification codes | Theme codes | Text | | Functional codes |
|---|---|---|---|---|
| 123456789 10 11 12 13 14 | 15 - 25 | 26 | | 66 67 68 69 70 71 72 73 74 75 76 77 78 79 |
| 0101-0001 1 1 1 1 | | Den avgränsning | The limitation | 3 0 |
| 0 2 1 1 1 | | (inom projektet) | (within the project) | 3 3 |
| 0 3 1 2 1 | | som (avgränsning) | which (limitation) | 3 3 2.5 0 |
| 0 4 1 2 1 | | (inom projektet) | (within the project) | 5 3 |
| 0 5 1 2 1 | | vi (X-projektet) | we (project X), | 3 0 3 0 |
| 0 6 1 2 1 | | har gjort (avgränsat) | have made (limited) | 4 0 1 |
| 0 7 1 1 1 | | är | is | 4.1 |
| 0 8 1 1 1 | | i stort sett | largely | 4 2 |
| 0 9 1 1 1 | | att | that | 3 0.3 |
| 1 0 1 1 1 | | jag | I | 3 0 |
| 1 1 1 1 1 | | skriver om | write about | 4 0 |
| 1 2 1 1 1 | | intervjumaterialet | interview material | 5 0 |
| 1 3 1 1 1 | | när det gäller skolledar-biten | concerning the school leader section | 5 3 |
| 1 4 1 1 1 | | som (skolledarbiten) | which (the school leader section) | 5 3 4 3 0 |
| 1 5 1 1 1 | | leder fram till | leads to | 4 0 |
| 1 6 1 1 1 | | en beskrivning | a description | 5 0 |
| 1 7 1 1 1 | | utav vad typ av arbetsupp-gifter | of the type of tasks | 5 3 |
| 1 8 1 1 1 | | som (arbetsuppgifter) | which (the tasks) | 5 3 5 5 0 |
| 1 9 1 1 1 | | skolledarna | the school leaders | 7.0 |
| 2 0 1 1 1 | | ställs inför | are faced with | 4.0 |

Figur 2. Text on IBM-cards

16

The maximal unit is a sentence, which can be divided into clauses of different degrees. A sentence is complete as soon as it contains the two main constituents, subject and verb (phrase). This analysis works with the sequence of clauses. Therefore the main clause (the first column) can also be called the first clause. In the example (figure 2) the subject in the first clause has a postpositive qualification in the form of a whole clause, which must therefore be introduced in some other column. This is done by means of the figure 3 in column 68, which states that the whole sentence's first subordinate clause is to be found in columns 69 and 70. The first clause continues in columns 66 and 67, which is stated in column 71. The object in the clause is a that-clause and since it is "clause-worthy" it must be placed directly in the next empty column (3), which is stated in column 74. Thus, the object of the "main clause" consists of three clauses.

Since this analysis is to capture the function of each individual unit in a clause, it is sometimes necessary to have double coding in the form of clause subordinators. Here a "that" introduces a new clause. but has no other function. "Which", on the other hand, has a function. In the first and third cases, "which" changes function from qualifier to object, in the second case from qualifier to subject.

Figure 3 on p 17 presents the flow chart for a computer search for the text example shown in Figure 2. Only such structures as can be stated explicitly can be delegated to a computer-based system. Each desired facet cannot be stated in advance, nor be extracted from a text material. For this reason, we have, in addition to the segmentation discussed, also devised codes for the main theme, so that the fundamental information, which cannot be regained or mediated by means of this clause code, does not get lost.

## CODING SCHEDULE

| Col. | Variable Identification | Pos. | Comments |
|---|---|---|---|
| 1-2 | Interviewed person No. | 01-40 | |
| 3-4 | Question No. | 01-n | |
| 5-8 | Sentence No. | 0001-n | |
| 9-10 | Concept No. in sentence | 01-n | |

| Col. | Clause codes | Pos. | Comments |
|---|---|---|---|
| 11 | Source of statement | 1, 2 | 1 = the speaker himself<br>2 = someone other than the speaker has made the statement |

| 12 | Negation | 1, 2, 3 | 1 = not, no, none, nobody, nothing<br>2 = hardly or the like<br>3 = neither - nor |
|----|----------|---------|------------------------------------|
| 13 | Tense | 1, 2, 3 | 1 = present time, refers to when statement is made<br>2 = past time from the occasion when statement is made<br>3 = future time from the occasion when statement is made |
| 14 | Mood | 1, 2, 3, 4 | 1 = indicative<br>2 = imperative<br>3 = conjunctive<br>4 = modal auxiliaries |
| 15 | Condition | 1, 2 | 1 = the conditional clause<br>2 = the corollary |
| 16 | Cause | 1, 2 | 1 = the causal clause<br>2 = corollary |
| 17 | Concession | 1, 2 | 1 = the concession<br>2 = the restriction |
| 18 | Result or intention | 1, 2 | 1 = "the main clause"<br>2 = result/intention clause |
| 19 | Contrast | 1, 2 | 1 = either, admittedly, on the one hand<br>2 = --- or, --- but, on the other hand |
| 20 | Comparison | 1 | 1 = occurrence |
| 21 | Question | 1 | 1 = occurrence |
| 22 | Assumption | 1 | 1 = occurrence |
| 23 | Volition | 1 | 1 = occurrence |
| 24 | Number of cards (lines) with word units in chich there is not room for the unit on one line | 2-9 | |
| 25 | Interrelation between cards (lines) containing the same word unit | 1-9 | |
| 26-65 | <u>Text</u> | | |

<u>Summary: main categories</u>

| <u>Code</u> | <u>Content</u> |
|------|---------|
| 30 | Agent |
| 31 | Qualifier of agent (before) |
| 32 | Description of agent |

| 33 | Qualifier of agent (after) |
|----|----|
| 40 | Verb, verb phrase |
| 41 | Copula-verb |
| 42 | Clause adverbial (clause modifier) |
| 43 | Time adverbial (when?) |
| 44 | Place adverbial (where?) |
| 45 | Manner-degree adverbial (how?) |
| 46 | Adverb or other word stating comparison, contrast etc. |
| 50 | Direct object |
| 51 | Qualifier of direct object (before) |
| 52 | Description of direct object |
| 53 | Qualifier of direct object (after) |
| 70 | Object of goal |
| 71 | Qualifier of object of goal (before) |
| 72 | Description of object of goal |
| 73 | Qualifier of object of goal (after) |

| | Clause codes | Pos. | Comments |
|----|----|----|----|
| 66-67 | Main clause columns | | two-figure codes |
| 68 | Reference column for subordinate clause columns | 2, 3, 4 5 | |
| 69-70 | First subordinate clause columns | | two-figure codes |
| 71 | Reference column for other subordinate clause columns or a looping back to main clause | 1, 3, 4, 5 | |
| 72-73 | Second subordinate clause columns | | two-figure codes |
| 74 | Reference column for other subordinate clause columns or a looping back to main clause | 1, 2, 4, 5 | |
| 75-76 | Third subordinate clause columns | | two-figure codes |
| 77 | Reference column for other subordinate clause columns or a looping back to main clause | 1, 2, 3, 5 | |
| 78-79 | Fourth subordinate clause columns | | two-figure codes |
| 80 | Reference column for other subordinate clause columns or a looping back to main clause | 1, 2, 3, 4 | |

START

| 30 |
|----|
| 33 |

| 33 | 2 |
|----|---|

| 50 |
|----|
| 53 |
| 30 |

| 40 | 1 |
|----|---|

| 41 |
|----|
| 42 |

| 30 | 3 |
|----|---|

| 30 |
|----|
| 40 |
| 50 |
| 53 |

| 53 | 4 |
|----|---|

| 30 |
|----|
| 40 |
| 50 |
| 53 |

| 53 | 5 |
|----|---|

| 50 |
|----|
| 70 |
| 40 |

END

Figur 3. Search paradigm

## 4. 7 Supplementation of text

Some sentences can be fragments that cannot be supplemented into independent conceptualizations, i.e. agent-action-object (AaO). In these cases in which the coder does not understand an utterance, it is to be deleted. The utterance must be completely comprehensible, which means that different types of relation words (e.g. pronouns and adverbs) must be supplemented to their right meaning (referent) in the context. Supplements are placed in parenthesis, so that the analysis does not lose track of what the person interviewed in fact says. When choosing the words to be used in the supplements, those already used by the person interviewed are taken first, if the context does not make this impossible.

Some deletions are made when the material is segmented. When defining a sentence, one cannot always assume that each sentence in the text has been concluded with a fullstop. A unit between two fullstops can consist of several sentences, either separated by means of pauses that are marked in the transcription by a series of dots or fragments which can be supplemented and made into complete sentences. Another way of marking the beginning and end is by linking with "and" or other conjunctions, which in this analysis are taken as being the first unit in a sentence and coded as "having no meaning". (This does not apply to an "and" that links two objects in the same clause.) In the cases in which obvious corrections are made by the person interviewed, the utterance that is immediately corrected is not coded.

## 4. 8 Selection of coders and choice of criteria for intercoder agreement

The computer-based processing of text according to ANACONDA is based on pattern recognition and the treatment of manually inserted clause codes. The assessments which e.g. two independent assessors give for an "information unit" with respect to the same category can best be considered as parallel "tests", which at the same time assumes that both assessors have "identical" frames of reference or systems of relationships. An examination of the precision of the coding done by the assessors is one of the prerequisites if we are to be able to demonstrate the objectivity in content analytical processing of verbal material. The "reliability" of the assessors' coding is above all a problem of communication. i.e. the precision of the coding is dependent on the communicability of the criteria stated in ANACONDA. To summarize. it can be said that the reliability of the coding is a function of (1) the unequivocality of the information units, (2) the unequivocality of manual and category functions and (3) assessors' special

frame of reference, e. g. knowledge of linguistics and knowledge of the subject. The assessors form the measuring instrument in the analysis. In addition the unequivocality of the information contained in the basic units influences the reliability to a large degree. But since it is very difficult, if not impossible, to get the entire process under control, the possibility of increasing the reliability is usually limited to manipulation with the assessors and/or manual. For this reason it is more justifiable to use the term "intercoder agreement", at least as long as the allotment of codes cannot take place mechanically.

If we have estimated the intercoder agreement, irrespective of which method of estimation has been used, it is usually very difficult to judge whether the calculated index value can be considered satisfactory. It can be very difficult to determine a reasonable level of agreement, since there is no simple solution to this problem. Moreover, it is only possible to decide what can be considered a satisfactory "reliable" coding within the frame of a given problem.

Starting from a coded material of this kind, analyses can be carried out that are based on the discourse model (Krippendorff, 1969, p. 80). By means of this model, it becomes possible to represent events or ideas within the source of information (the researcher). The use of the model presupposes that independent coders can allot codes to the text with a satisfactory level of agreement. As the criterion we have stipulated 80% agreement.

Two methods of assessment were applied. The first method (Osgood et al, 1956, p. 57) states the proportional agreement. Segment markings; supplementations, coding of subject, object, verb and modifiers were assessed according to this method. Osgood's technique was applied primarily with the purpose of making it possible to compare our results with those presented by Osgood.

The second method is based on the binomial division hypothesis. By means of the binomial test (Siegel, 1956, p. 40), the extent to which the criterion (80% agreement) was fulfilled was studied. An observed value that is less than the test value of z-1.64 states that the intercoder agreement does not fulfill the criterion.

The persons who have coded the interview material have within the framework of the SÖH-project developed ANACONDA. In addition an attempt was made to train two additional coders. This attempt has for various reasons, however, been abandoned. The persons who were to code one part of the interview material were trained for about 15 hours in the use

of the first version of ANACONDA. They practised (1) separating relevant
from non-relevant text material, (2) segmenting text into meaningful units
and (3) identifying syntactical relationships. In addition they practised (4)
writing word units in sequence in accordance with a pattern, (5) differen-
tiating between types of clause according to definition and (6) interpreting
analysis units and arranging codes (see I. Bierschenk, 1974, p. 16).

By now the developmental work with ANACONDA has reached the point
at which, according to Figure 1, we have to estimate the intercoder agree-
ment. Some of the empirical results of the evaluation will be presented next.
But the steps below the "criterion of intercoder agreement" in Figure 1 will
not be pursued further here. Instead these steps will be discussed briefly
in Chapter 6.

## 5. INTERCODER AGREEMENT: SOME EMPIRICAL RESULTS

### 5.1 Intercoder agreement

If we are to be able to develop a technique for a computer-based, content
analysis, it will be necessary for us to create a system of rules that two
or more independent coders can use with a high degree of mutual agreement.
Berg (1974) calculated the agreement between two independent coders. This,
scrutiny concerned (1) supplementation and deletion, (2) segmentation and
(3) allotment of codes.

By means of a random table, four interview subjects (31, 2, 40 and 33)
were picked out from the interviewed sample of researchers. From the
respective interviews, four interview questions (5, 6, 7 and 8) concerning
information and documentation have been chosen. It can be assumed that
the information that will be extracted from the text will be relatively con-
crete and consequently easy to interpret. This should be an advantage in
the development of a new technique.

The interview question were to be coded in their entirety, so that the
context of the discussion could be used in supplementation. Spreading the
selection of text over the entire text or over all the subjects has been con-
sidered an unsuitable method of procedure.

The intercoder agreement was examined with regard to

1. segmentation of concepts. A check is made of whether both coders
   have supplemented and deleted identical words.

2. segmentation of clauses. A check is made of whether the coders have
   identical sentences.

3. allotment of codes to concepts. A check is made of whether both coders
   have allotted identical codes to one and the same concept.

4. allotment of codes to themes. A check is made of whether both coders
   have alloted identical codes to one and the same theme in a sentence.

All the comparisons are of the same type i. e. either there is agreement or
not. The number of common assessments has been noted. In addition the
total number of assessment  and the number of assessments each coder
has made separately have been calculated. A detailed scrutinization and
comprehensive documentation  may be found in Berg (1974). Here, how-
ever, only a summarizing table will be presented, with the values for
points 1 to 4 above. The values have been taken from Berg (1974, p. 30)
and will be presented in reorganized form.

Table 1. Summary of intercoder agreement in applying ANACONDA

| Steps in the analysis | | Interview person No. | | | |
|---|---|---|---|---|---|
| | | 31 | 2 | 40* | 33 |
| Segmentation of concepts (1) | $z$ | 3.92 | 2.20 | -.58 | 3.21 |
| | $i$ | .88 | .86 | .82 | .86 |
| | $N$ | 799 | 1098 | 237 | 1255 |
| Segmentation of clauses (2) | $z$ | 2.82 | 2.64 | .67 | -2.76 |
| | $p$ | | | .75 | |
| | $i$ | .94 | .93 | .92 | .84 |
| | $N$ | 165 | 227 | 47 | 246 |
| Allotment of codes 6 concepts (3) | $z$ | 7.64 | 9.42 | 1.16 | 8.51 |
| | $i$ | .91 | .92 | .83 | .90 |
| | $N$ | 841 | 1089 | 222 | 1190 |
| Allotment of codes to themes: source, time, mode (4) | $z$ | 7.33 | 5.51 | 1.40 | 4.37 |
| | $i$ | .98 | .93 | .93 | .93 |
| | $N$ | 320 | 397 | 83 | 422 |
| Segmentation of concepts before check on comparable text | $z$ | -9.89 | -13.08 | -4.71 | -17.60 |
| | $i$ | .77 | .76 | .76 | .74 |
| | $N$ | 1013 | 1377 | 272 | 1673 |
| Allotment of codes to concepts before check on comparable concepts | $z$ | -2.47 | -4.52 | -6.23 | -10.40 |
| | $i$ | .83 | .82 | .73 | .78 |
| | $N$ | 992 | 1328 | 283 | 1549 |

$z$    test value, binomial test
$p$    probability: $p < .05$ states
     that the criterion .80 has not
     been achieved
$i$    Osgood's index for agreement
$N$    total number of assessments

*IP 40 has given oral comments to question 5. Questions 6 and 7 were answered by filling in a questionnaire, while the IP did not comment on question 8.

The checks of the intercoder agreement in the steps of the analysis carried out so far show that segmentation can be done with a satisfactorily high level of agreement. As Table 1 shows, Osgood's index for agreement is between .74 and .98. Spiegelman, Terwilliger & Fearing (1953, p. 175) give as the minimum requirement an index value that is equal to or greater than .75, irrespective of the method by which the intercoder agreement has been estimated. Osgood (1956, p. 59) himself reports index values of between .64 and .88. Our result is by comparison very satisfactory, since the analysis this report is dealing with is much more detailed and comprehensive. In addition the interview material contains for natural reasons greater variations, while at the same time it is less complete than Osgood's printed material.

The binomial test shows, however, that the critical value .80 could not be established in every case. As is shown in Table 1, neither the "segmentation of concepts before a check on comparable text" nor the "allotment of codes to concepts before a check on comparable concepts" has resulted in satisfactory values. This is caused by the lack of unequivocal rules. If, for example, one coder uses the term "researcher" while another describes the same person as a "behaviorist", this leads to differences in supplementation. This difference can, however, be nullified by e.g. appropriate construction of registers and facetting. All the supplementations are marked in parenthesis, which makes it possible for us to analyse the material both with and without supplementations and thus investigate the extent to which this leads to different results.

The index values reported above the line are comparable with the results that we would have got by limiting concepts in written text. As can be seen from Table 1, the agreement is very good. though with the exception of "Segmentation of clauses" in interview No. 33. This is probably a result of there being a large number of unsupplemented clauses (see Berg, 1974, p. 23).

Attributes and adverbs have obviously caused most of the deviations in the coding. The agreement for attributes is admittedly over 80% but some of the deviations could be explained by the confusion that has occurred between the two categories. Thus, the coding of e.g. "Researcher A in Malmö" has partly been as an adverb of place "in Malmö" and partly as a postpositive attribute. In addition there has been confusion between adverbs of time and degree. In the clause "I read daily", the word 'daily' has been coded both as a statement of time (adverb of time) and as a statement of frequency (adverb of degree). Concerning the examples presented here, the rules will be improved.

## 5.2 Control of punch cards

It is important that the text material that is to form the basis for the further development of ANACONDA is faultless. Otherwise it would be very difficult if not impossible to determine whether a fault is caused by incorrect coding or by some deficiency in the test material. For this reason the text material transferred to punch cards has been checked both for faults despite correct coding and for faults resulting from incorrect coding. A detailed examination has been made and documented by I. Bierschenk (1974).

The test material comprises about 37,000 punch cards. The punching was carried out by the punching machine operator at the Department of Educational and Psychological Research, Malmö School of Education. A

selection of punch cards (10%) was handed to the Data Processing Centre for Research and Higher Education in Lund. The punchings were then examined for (1) identification faults (ip, no, question no, sentence, no, word no), (2) theme (source, negation, tense, case, other clause themes), (3) text (spelling, parenthesis, other text), (4) content (concepts, clause column).

The result of this examination is presented in condensed form in Table 2. (For more detailed information, see I. Bierschenk, 1974, p. 31).

Table 2. Punching and control-punching. Observed and relative frequency calculated on 70,260 punches

| Category | Specification | Punching | | Control-punching | |
|---|---|---|---|---|---|
| | | f | % | f | % |
| 1 | Identification | 7 | .01 | 4 | .01 |
| 2 | Theme | 3 | .00 | 8 | .01 |
| 3 | Text | 26 | .04 | 102 | .14 |
| 4 | Content | 22 | .03 | 32 | .05 |
| Σ | | 58 | .08 | 146 | .21 |

From Table 2 it emerges that the control-punchings have been carried out less well than the original punchings. The similarities are greatest within categories 1 and 2, while the differences are greatest within category 3. This can be explained by the fact that numerical codes (cat. 1 and 2) are more common for the machine punching operators and occur less frequently in this material. In addition there is a system in the theme codes. Source, tense and case are always punched, while negation and other clause themes are only punched when they occur. But category 4 also contains numerical punching. Incorrect punching has serious consequences if e.g. a verb is placed in some noun (object, subject) category. Moreover, it is an extremely time-consuming and difficult job to check all the concepts included in each respective code.

Mistakes in category 3 mean among other things that the parenthesis sign has been neglected. This sign is important, however, when we wish to keep apart actual statements and implied or imagined ones.

In order that we should be able to form an idea of the consequences of the content codes throughout the entire material, all the material was corrected. Thereby it became possible to draw up a protocol with all errors. The text of all forty interview persons was examined on questions 5, 6, 7 and 8, card by card, and every error was registered. The results of the examination are presented in condensed form in Table 3. (For more detailed information, see I. Bierschenk, 1974, p. 38.)

Table 3. Punching and coding errors in examination of the total punched
material: Observed and relative frequency calculated on
702,600 punches

| Category | Specification | Punching | | Coding | | Σ | |
|---|---|---|---|---|---|---|---|
| | | f | % | f | % | f | % |
| 1 | - | | | | | | |
| 2 | Theme | 6 | .00 | 53 | .00 | 59 | .00 |
| 3 | Text | 182 | .03 | 34 | .00 | 216 | .00 |
| 4 | Content | 88 | .00 | 99 | .00 | 187 | .03 |
| Σ | | 276 | .04 | 186 | .03 | 462 | .07 |

As can be seen in Table 3, the greater part of the errors depend on the
punching. Corrections within category 4 covariate with alterations in the
text. But since the examination made showed that we in future only need
calculate with approximately .04% punching errors and .03% coding errors,
they are with regard to the clause columns a negligible factor.

## 6. STRUCTURED REGISTERS

Before a computerized analysis of information can be realized, the researcher must state his theoretical standpoint, i.e. define his concepts. It is necessary to establish in advance which aspects of the material are to be paid attention. Categories form the link between the theoretical anchorage of the research problem and the technical aspects of the analysis of information. By means of registers, intended to be used for content analytic treatment of texts, natural language is transformed into formalized language. This transformation assumes a purpose or a theory. The questions that have guided the interviews with researchers at the departments of educational research are based on the assumption that the initial phase of the research is influenced by

1.  the qualities of individual researchers and the social system within which they have to act and react - constraints

2.  the interests, motivation and role-behaviors of individual researchers - intervening variables

3.  recommendations for changes or improvements of e.g. research planning - research policy

Our hope is that we shall be able to answer at least three questions:

1.  Which criteria guide the researcher's approach during the initial phase of the research, i.e. what values do the researchers have?

2.  What actions do the researchers take during the initial phase of the research and how are these evaluated?

3.  What steering mechanisms influence the development of the initial research phase?

But what situations and actions are to be extracted from the present material and how are they to be evaluated? These questions must be answered in connection with the construction of structured registers. In principle each individual concept in the text can form its own category.

Irrespective of how the registers are built up, it should be of great help if e.g. proper nouns and place references form separate facets. If it should prove to be desirable to have facetted registers, the analysis must begin with the recognition of individual concept patterns in the text according to a register. The analysis technique that is to be developed for the interview material demands at least four different registers: (1) Independent concepts (subject and object terms), (2) Dependent concepts (adjectives/ attributes), (3) Actions/kopula (verbs), (4) Dependent concepts (adverbs).

With the help of the computer, lists are produced of these parts of speech. Registers are then compiled on the bases of these lists. Criteria

for the selection of important concepts already exist to some extent in the
form of the question complexes that are dealt with in the interview. By
means of the KWIC program, the registers can be adapted very closely to
the verbal behavior of the researchers.

In constructing registers 2-4, Osgood's semantic differentials were
used. Each term is defined with regard to (1) evaluation, (2) activity and
(3) potency. The assessment is made according to seven-point and bipolar
scales with the respective pairs of adjectives (1) negative/positive,
(2) passive/active and (3) weak/strong. The advantage of this scaling tech-
nique is that it is simple to use and that we can study three independent
dimensions. By means of the evaluation dimension, the extent to which the
researcher assesses different aspects as good or bad can be studied. The
activity dimension measures the extent to which the researcher considers
that a particular aspect has influenced the development of project outlines
or behavior during the initial phase of the research process. The potency
dimension measures the researcher's sensitivity or responsiveness.
Dimensions two and three together express dynamics.

It is assumed that assessors can make reliable and valid assessments
of the tendency and intensity of a statement. An initial analysis for deter-
mining the reliability of the scaling of adjective and verb has been carried
out. The results are presented in Table 4:

Table 4.  Intraclass correlation in the scaling of modifiers and actions

| Class | Pairs of adjectives | | |
| | Negative/ positive | Passive/ active | Weak/ strong |
| --- | --- | --- | --- |
| Modifiers | .90 | .98 | 1.00 |
| Actions | .99 | .97 | .98 |

As is shown in Table 4, the words within each group are assessed very
similarly. An evaluation of the inter-assessor agreement has shown, how-
ever, that the assessors have differing ideas in the evaluation of both
modifiers and actions. One example of this dissimilarity in the assessment
can be given with the modifier "psychological". Depending on whether the
word stands before or after its main word, the assessment could be made
in different ways. If "psychological methods" are named in the text, the
word "psychological" can be treated as an adjective. If on the other hand
the text mentions "methods in psychology", the word is descriptive and in
the second case explanatory. The word "psychological" has been assessed
in the way shown in Table 5:

Table 5. Assessment scores of the modifier: psychological

| Assessor | Dimension 1 | 2 | 3 |
|----------|-------------|---|---|
| 1        | 5           | 5 | 5 |
| 2        | 6           | 6 | 6 |
| 3        | 4           | 4 | 6 |

The example given in Table 5 reflects the different references made by the assessors to the word "psychological". Even though the agreement in the assessment is low, the average nevertheless appears to be a good approximation of the sense in which the word is commonly used. The result shows that the word "psychological" as a descriptive term is not wholly neutral. The fact that the agreement between the assessors is low can be explained partly by their different backgrounds, partly by the circumstance that the assessment was carried out without any special instructings having been given. Thus, the way in which a person himself chooses an expression is our main concern. One could admittedly claim that (1) methods in psychology and (2) psychological methods are equivalent forms of expression for the same conceptualization. But the coding of the text must permit us to state content and not only position. The content is not determined until further steps have been carried out, i.e. the scaling.

## 6.1 Facets

A content analysis presupposes categories or facets. Concepts, idiomatic expressions or phrases represent thereby a variable according to a certain given theory. The fundamental means of approach in a content analysis is to identify these "signs", so that they can be coded according to the category to which they belong (Stone, 1966, pp. 170-186). In a content analysis, however, it is seldom so that the analysis concerns only one "category", but instead the interest concerns the relation between categories. The categories form the semantic and empirical anchorage of the analysis, since individual concepts are defined according to the researcher's conceptualization, e.g. category systems.

Registers for content analytical processing function as links between the natural language and a more formal, theory-oriented language. At the present stage of the analysis it is difficult to establish how the information made available with the help of the coding system presented and information such as specific theory-oriented concept groups should best be used. The building up of structured registers should, however, aim at

1.  the possibility of being able to test several theories
2.  clear divisions of relevant facets and
3.  if possible, the statement of the relations between individual facets, e.g. proper noun & institution, proper noun & role, etc.

Thus the construction of our registers involves both a classification of the material into evaluating categories and thematic classifications. If we are to be able to carry out the analysis successfully, it seems at present as if we must reach a decision on some form of facetting.

# 7. REFERENCES

Berg, M. Reliabilitetsprövning av en metod för innehållsanalys av intervju-text. /Reliability testing of a method of content analysis applied to interview texts./ Testkonstruktion och testdata, No. 26, 1974.

Bierschenk, B. Perception, strukturering och precisering av pedagogiska och psykologiska forskningsproblem på pedagogiska institutioner i Sverige. /Perception, structuring and definition of educational and psychological research problems at departments of education in Sweden./ Pedagogisk-psykologiska problem, No. 254, 1974.

Bierschenk, I. Konstruktion av ett regelsystem för en datorbaserad innehållsanalys av intervjutext: Preliminär manual och några utprövnings-resultat. /Construction of rules for a computer-based content analysis of interview texts: A preliminary manual and some evaluation data./ Testkonstruktion och testdata, No. 25, 1974.

Bobrow, D. G. Syntactic theories in computer implementation. In: Borko, H. Automated language processing. New York: Wiley, 1967. Pp. 215-251.

Gerbner, G., Holsti, O. R., Krippendorf, K., Paisley, W. J. & Stone, P. J. (Eds.) The analysis of communication content. Developments in scientific theories and computer techniques. New York: Wiley, 1969.

Holsti, O. R. Content analysis for the social sciences and humanities. Reading: Addison-Wesley, 1969.

Horst, P. Personality: The measurement of dimensions. San Francisco: Jossey-Bass, 1968.

Krippendorf, K. Models of messages: Three proto-types. In: Gerbner et al. The analysis of communication contents. Developments in scientific theories and computer techniques. New York: Wiley, 1969. Pp. 69-106.

Miller, G. A. Kommunikation och psykologi. /Communication and psychology./ Stockholm: Beckmans, 1969.

Osgood, Ch. E., Saporta, S. & Nunnally, J. C. Evaluative assertion analysis. Litera, 1956, 3, 47-102.

Schank, R. C. Conceptual dependency: A theory of natural language understanding. Cognitive Psychology, 1972, 3 (4), 552-631.

Schank, R. C. & Colby, K. M. (Eds.) Computer models of thought and language. San Francisco: Freeman, 1973.

Siegel, S. Nonparametric statistics for the behavioral sciences. New York: McGraw-Hill, 1956.

Spiegelman, M., Terwilliger, C. & Fearing, F. The reliability of agreement in content analysis. J. soc. Psychol., 1953, 37, 189-203.

Stone, P. J. The general inquirer: A computer approach to content analysis. Cambridge, Mass.: The MIT-Press, 1966.

**Department of
Educational and
Psychological Research**

**School of Education
Malmö, Sweden**

## Reference card

Bierschenk, B. A computer-based content analysis of interview
data: Some problems in the construction and application of
coding rules. Didakometry (Malmö, Sweden: School of Edu-
cation), No. 45, 1974.

## Abstract card

Bierschenk, B. A computer-based content analysis of interview
data: Some problems in the construction and application of
coding rules. Didakometry (Malmö, Sweden: School of Edu-
cation), No. 45, 1974.

This report presents a flow-chart of different steps in
the designing of an Analysis of Concepts by Data-pro-
cessing (ANACONDA). A condensed preliminary version of
ANACONDA is presented and empirical results of applica-
tion by independent coders are given.

Indexed:
1. Data-processing
2. Concept analysis
3. Intercoder agreement
4. Interview texts
5. Psycholinguistics