

391

PRR

VOLUME I | NR. 3 | OCTOBER 1968

RECEIVED

JUL 3 '69

U. OF R. LIBRARY

COMPUTER  
STUDIES

*in the Humanities  
and Verbal Behavior*

MOUTON

12  
05  
107  
1#3

# COMPUTER STUDIES

## *in the Humanities and Verbal Behavior*

### CO-EDITORS

FLOYD R. HOROWITZ, University of Kansas  
LEWIS SAWIN, University of Colorado  
SALLY Y. SEDELOW, University of North Carolina

### BOARD OF EDITORS

<i>Anthropology</i>	GEORGE L. COWGILL, Brandeis University	<i>Mathematics</i>	PAUL R. HALMOS, University of Michigan
<i>Archaeology</i>	PAUL S. MARTIN, Field Museum of Natural History	<i>Music</i>	LEJAREN HILLER, University of Illinois
<i>Art</i>	LESLIE MEZEI, University of Toronto	<i>Philosophy</i>	LARRY TRAVIS, University of Wisconsin
<i>Bibliography</i>	ERIC BOEHM, American Bibliographical Center, Santa Barbara	<i>Political Science</i>	OLE R. HOLSTI, University of British Columbia
<i>Classics</i>	T. M. ROBINSON, University of Calgary	<i>Programming</i>	DAVID BRIDGER, Washington University
<i>Education</i>	ROBERT SCHMIEL, Goucher College	<i>Psycholinguistics</i>	DANIEL E. BAILEY, University of Colorado
<i>Folklore</i>	ELLIS B. PAGE, University of Connecticut	<i>Social Psychology</i>	JOHN B. CARROLL, Educational Testing Service, Princeton
<i>History</i>	JOHN Q. ANDERSON, University of Houston	<i>Sociology</i>	HANS E. LEE, Michigan State University
<i>Library Science</i>	THEODORE K. RABB, Princeton University	<i>Speech</i>	WALTER A. SEDELOW, University of North Carolina
<i>Linguistics</i>	RALPH H. PARKER, University of Missouri		EDWARD E. DAVID, Bell Telephone Laboratories
<i>Literature and Modern Languages</i>	SHELDON KLEIN, University of Wisconsin	<i>Statistics</i>	GERALD M. SIEGEL, University of Minnesota
	ROBERT S. WACHAL, University of Iowa	<i>Theater Translation</i>	JULIET SHAFFER, University of Kansas
	BERTRAND AUGST, University of California		HAROLD P. EDMUNDSON, University of Maryland
	JESS BESSINGER, JR., New York University		GARY GAISER, Indiana University
	WILLIAM INGRAM, University of Michigan		SILVIO CECCATO, Università di Milano
	HENRY KUČERA, Brown University		DAVID A. DINEEN, University of Kansas
	JAMES W. MARCHAND, Cornell University		
	STEPHEN PARRISH, Cornell University		
	ALICE POLLIN, New York University		
<i>Mass Communication</i>	WILLIAM J. PAISLEY, Stanford University		

### EDITORIAL ADDRESS

*Computer Studies in the Humanities and Verbal Behavior* is jointly sponsored by the universities of Colorado, Kansas and North Carolina. Manuscripts: LEWIS SAWIN, Department of English, University of Colorado, Boulder, Colorado, U.S.A. General correspondence: FLOYD R. HOROWITZ, Department of English, University of Kansas, Lawrence, Kansas, U.S.A.

*Computer Studies in the Humanities and Verbal Behavior* appears quarterly in issues of approximately sixty pages. Four issues constitute a volume. The subscription price is \$ 10.00/f 36.00 per year per volume; single issues cost \$ 3.00/f 10.50. Back issues as well as subscriptions and single issues can be ordered from every bookseller or subscription agency, or directly from Mouton & Co., P.O. Box 1132, The Hague, The Netherlands.

## Contents

JUDITH E. SELVIDGE and THEODORE K. RABB, DATA-TEXT: A Simple and Flexible Programming System for Historians, Linguists, and Other Social Scientists . . . . .	107
GEORGE PSATHAS, Computer Analysis of Dyadic Interaction . . . . .	115
J. ZVI NAMENWIRTH, Some Long and Short Term Trends in One American Political Value: A Computer Analysis of Concern with Wealth in Sixty-two Party Platforms . . . . .	126
HOWARD P. IKER and NORMAN I. HARWAY, A Computer Systems Approach Towards the Recognition and Analysis of Content . . . . .	134
REVIEW ARTICLES . . . . .	155

---

*Some articles for forthcoming issues:*

- EVERETT ALLDREDGE, Preservation of Documentation for Conventional and Automated Systems
- GARY BERLIND and GEORGE W. LOGEMANN, An Algorithm for Musical Transposition
- ERWIN DANZIGER, Tutorial on Computer Hardware
- ALLEN FORTE, Analytic Methods of Machine Storage
- DOUGLAS HINTZMAN, Learning and Memory in a Discrimination Net
- OLE R. HOLSTI, Computer Content Analysis for Measuring Attitudes: The Assessment of Qualities and Performance
- CHARLES H. KELLOGG, Data Management in Ordinary English: Examples
- ALASTAIR MCKINNON and ROGER WEBSTER, A Method of 'Author' Identification
- R. NARASIMHAN, Intelligence and Artificial Intelligence
- ELLIS B. PAGE, A Program for the Automatic Evaluation of Student Essays
- JAMES B. RHOADS, The Role of the National Archives in Facilitating Statistical Inquiry
- WALTER A. SEDELOW, JR., History as Language (Part I); Verbal Structure of Hume's *History of England* (Part II)
- HUNG-CHING TAO, A Chinese Computer Alphabet for Automatic Machine Processing
- WAYNE TOSH, Machine Translation, 1969

# A Computer Systems Approach Towards the Recognition and Analysis of Content\*

HOWARD P. IKER and NORMAN I. HARWAY

*Department of Psychiatry, School of Medicine and Dentistry  
University of Rochester, Rochester, N.Y.*

*This paper explores a method of content analysis which allows the user to discover what his data are about without having to furnish a priori categorizations within which to classify these data. Based originally on problems encountered in psychotherapy — the question of the relationship between the substance and structure of oral communications between psychotherapist and patient — the paper also explores the possibility of adapting this method of content analysis to other kinds of data. The initial task of this analysis is to generate informational categories with which many content analytic systems begin. Based on an associational approach, the basic unit of information is the word itself. A matrix of intercorrelations among words is eventually factor analyzed to determine systematically the common factors which may account for the matrix in a meaningful way. To reduce the number of words to a manageable size, a system of programs (WORDS) has been developed. This paper then considers the current use of the WORDS System, its goals and structure, the results and implications of some current research, and plans for the future.*

The key question which we have been exploring for almost six years (2, 3, 4) is whether there exists a method for content analysis which will allow the user to discover what his data are about without having to furnish *a priori* categorizations within which to classify these data. Any other currently used content-analytic system of which we are aware requires, at the least, that the user furnish a set of categories to which the various portions of the text are to be allocated. It is our purpose to demonstrate that there is an alternative to this approach.

Our original interest in the area was in the process of psychotherapy, in the changes in cognitive organization which occur as a result of treatment, and in the process of change. Psychotherapy, if we exclude for the moment certain of the behavioral therapies, involves the user of oral communication to modify, among other things, a person's perspective of himself and of his world. To the extent that this occurs, this change should be reflected in the individual's verbal behavior. Since the assumption is that the communications of the psychotherapist are influential in accomplishing this change, we are led to the

question of the relationship between the substance and the structure of the oral communications among psychotherapist and patient, their change with time and with progress in treatment.

Over the years during which we have been developing the system for making such analyses, it has become clear to us that the system can be applied to kinds of data other than those embodied in psychotherapy materials.

The basic assumption upon which this system rests is that there exists sufficient meaning within the word and within the temporal associations among and between words to allow the elicitation of major content materials and categories. In short, it is the initial task of our approach to *generate* the kinds of informational categories with which many content analytic systems begin.

Utilizing an associational approach as the foundation for our method, we take, as our unit of information, the word itself. Dividing an input document into segments of time, or segments of equal length, or equal numbers of sentences, etc., it is possible to count the frequency with which each word occurs in each such segment. Using these data, intercorrelations among all words may be obtained; operationally, these intercorrelations represent the degree of co-occurrence, i.e., association, between words as they are observed across successive units of the data-base. This matrix of intercorrelations may then be factor analyzed to determine, in a systematic fashion, if

\* This is a draft of a paper commissioned by and prepared for the National Conference on Content Analysis, November 16-18, 1967 at the Annenberg School of Communications, University of Pennsylvania, Philadelphia, Pa. The research was supported in part by U.S. Public Health Service Grant MH-10444, National Institute of Mental Health. We wish to acknowledge also the aid of Miss Janet Barber, Mr. Gerald Leibowitz, and Dr. Edward Ware.

there are common factors which can account for the obtained associational matrix in an efficient and meaningful way.

In a factor-analytic approach towards this kind of data, it is necessary to reduce the number of different words that are examined. This criterion, a function of machine-size, program running time, and factor interpretability demands that the number of variables (words) analyzed be held to some reasonable minimum. It is the job of the system of programs which we have developed, WORDS, to allow this reduction and the subsequent statistical analyses that are necessary.

This paper will discuss the current implementation of the WORDS System, its goals and structure, the results and implications of some current research, and will then discuss our plans for the future concerning changes in the system and directions for future research.

## IMPLEMENTATION OF THE WORDS SYSTEM

### *Computing Facility*

WORDS has been developed at the University of Rochester Computing Center. Until recently, all large scale work in this facility has been run on an IBM 7074 computer. The 7074, a second generation, medium-speed machine, has characteristics including extremely powerful input/output (I/O) hardware and logic, and highly flexible scatter read/write commands. It is a fixed word length machine (five characters or ten digits with sign per machine word) with decimal arithmetic and hardware supported floating point operations. The configuration at the University of Rochester operates with a 10K core and is supported by eight 729-IV tape drives and a 1301 disk file. Final output from the system must be to tape with all printing and punching done offline on an IBM 1401 configuration.

This configuration runs under control of the Computing Center's resident monitor which operates the system on a batched queue basis. As such, all daytime runs operate under closed-shop conditions with all operations handled by a computing-center machine operator.

### *System Goals*

Before beginning actual programming on WORDS, a set of system goals was developed. The considerations dictating choice of these goals were based on the several uses to which WORDS could be put. WORDS had to be capable of large scale, repetitive, methodological investi-

gations; as such program running times had to be as efficient as possible. WORDS had to be useful as a production device; as such, ease of system use and good turnaround time were needed. Finally, WORDS would probably be used by other members of the University and therefore should not be difficult to learn.

With these considerations, as well as those dictated by the method itself, a set of generalized criteria were developed to help in directing the systems work and programming applications to follow.<sup>1</sup>

### *User-Orientation*

WORDS was designed to be as easy to use and to learn as possible. Despite this desire, the final results are far short of the mark as can be attested to by a 110 page user's manual. While the system is not inordinately difficult to use and does not require any extensive experience with computers or programming, its use clearly requires some amount of training. Considering the complexity of the total system itself, this requirement is not unreasonable; nevertheless, both we and the users would be happier were the method easier to understand. Future implementations will result in a much more easily used system. This decrease in difficulty will come about for two reasons; the first is a direct result of our experience with the current system and knowledge of how we could make things easier even within the current system itself; the second reason stems from the greatly increased sophistication and flexibility of third-generation software which substitutes the operating system of the computer for much programming effort on the part of the developer.

Because WORDS does not demand any computer experience of the user it must be able to run under control of the target computer's resident monitor in a closed-shop environment. Since WORDS contains a large number of programs any, or all, of which can be called in any order appropriate to the purpose of the run, WORDS itself must be monitored. For a second generation resident monitor such a consideration demands a systems monitor capable of operating with and communicating with the resident.

### *Flexibility*

WORDS is designed so that the user should have no great

<sup>1</sup> Throughout the development and programming of the WORDS System, the advice and assistance of Mrs. Barbara Rothe has been invaluable.

difficulty in manipulating his data as is necessary for production of appropriate output. Thus, WORDS has been written to allow the user to configure a run with as many or as few programs as necessary, to utilize, wherever possible, mnemonics rather than numeric information for control purposes, to afford as effortless a handling of I/O as possible, and to allow a wide range of output formats.

### *Efficiency*

Typical data-bases in our current research involve files of approximately 25,000 words. Typical run configurations often run more than 20 separate program calls for successive manipulations of this data. A total run, then, can easily involve the computer in the manipulation of well over one million records. On a machine with the speed of the 7074, then, a high degree of program efficiency is a clear necessity.

### *Protectiveness*

Use of WORDS almost invariably demands repetitive runs on the computer with the input for any given run being based, at least in part, on output from prior runs; in such a case, some devices are needed to protect prior data from possible destruction during a run. Additionally, a complete run configuration for any extensive processing by WORDS can generate a complicated calling sequence with massive I/O operations so that a software or hardware failure is always a possibility. WORDS is designed to fail-safe and to produce as much diagnostic information for the user as possible.

Appendix I may be consulted for further information on the system. It details the structure of the WORDS System: its systems organization, data organization, program organization and a list and description of each of the major programs in the system.

## CURRENT RESEARCH

### *Methodologic Issues*

#### *The UHH Approach*

Over the past several years, apart from the complete re-programming of WORDS, our major research emphasis has been to investigate some of the methodological problems confronting this technique. One of the most

important of these issues, both on a practical and theoretical level, derives from the technique used for the reduction of the number of different words found in the raw data base. The question which we have investigated was whether alternative methods for reduction of these words could be found that did not involve the extensive use of synonymization.

In our analysis of a typical data base — five psychotherapy interviews — we encounter about 1300 different words making up the total 25,000 word protocol. Since the maximum matrix with which we can work is 215 variables, we must make reductions comprising 83% of the total number of different words.

The method we initially employed was based on a four phase process. Utilizing the PARSE programs, all articles, prepositions, conjunctions, etc., were removed. The next phase used STRIP to change all words into root form. The third, heavily employing synonymization, combined words having the same basic meaning and being used in the same fashion. Lastly, a list of remaining different words, sorted by frequency of occurrence, was generated; beginning with the word of highest frequency, this list was downcounted to reach 215 different words which were then subjected to matrix analysis. It is the synonymization phase which we have found most difficult in implementation and most dangerous in terms of objectivity. Synonymization requires two basic commitments from the user: the first being extensive amounts of time and the second being the use of subjective and potentially unreliable judgments. In many respects, synonymization places the same demands on the user as does the typical content-analytic method of *a priori* categorization. Except for the fact that the data itself, rather than the initial interests of the investigator, generates the potential categories (generic word), we are faced in synonymization with the same degree of subjectivity and inefficiency as we would be were we to have begun word reduction with a set of categories into which the data was to be allocated.

With a synonymizing procedure, the time demands placed upon us in reduction of the data in a typical set of five interviews was considerable. At least a month was required, from both investigators, with a heavy investment of twenty-seven pre-factor runs on WORDS being necessary before producing the final rotated factor structure representing the content-recognition portion of the analysis.

More important than time, however, was the constant requirement that we exercise our own judgment as to when two words were being used in a fashion and with a meaning such that they could be combined into one. While WORDS is implemented to make this task easier and more reliable

than before (cf. Appendix 1: HSTRY, IDIOM, TEXT 1 and TEXT 2, PARSE, etc.), there is little question but that our ability to maintain a high degree of freedom from *a priori* needs and ideas in the selection of combinations decreased as the time wore on and the combinations became more and more difficult to locate.

Our experience, then, in the analysis of data using *ad hoc* synonymization procedures to carry a large bulk of the reduction process demanded that we investigate other techniques that were more efficient and less subjective. Accordingly, we began work with a different approach in which synonymization played a minimal role and in which major reductions were accomplished by deletion procedures.

This new technique, which we have euphemistically labeled UHH (Untouched by Human Hands) as opposed to SYN (Synonymization), operates according to a generalized set of rules. The application of these rules, rather than *ad hoc* decisions deriving from our own inspection of the data, clearly reduces both the time required and the subjectivity involved in the reduction process.

UHH rules currently fall into two phases: pre- and post-factoring. In the pre-factor phase, we first parse and then subject the data to a common STRIP run for the purposes of de-inflection and a change to root form of all comparatives. Following this, EDIT is applied. EDIT is used to make four kinds of change: (1) deletion by part of speech so that all articles, prepositions, conjunctions, etc., are removed; (2) deletion of certain words which carry very little meaning outside of context, e.g., *sort, still, be, thing, ago*, etc.; (3) deletion by the combination of word/speech categories so that, for example, *kind* (adjective) is retained while *kind* (noun) is deleted, or *like* (verb) is retained while all other forms of the word are dropped; (4) lastly, a low level of *pre-determined* synonymization is applied in which generic words are created to subsume a set of other highly related high frequency words, e.g., NO is held as a generic word which will contain all occurrences of *neither, never, nobody, none, nor, not, nothing, and nowhere*.

It is of importance to the UHH approach to note that this kind of synonymization rule is applied to data prior to any analysis of the data and, where feasible, is applied prior to any inspection of the data itself.

After subjecting the data to this set of pre-factor rules, a downcount is taken on a frequency ordered list and the 215 words with the highest frequencies are then subjected to factoring. Post-factoring rules serve a common goal: the improvement of the obtained factor structure. Such improvement can come about in two major ways within the constraints of our methods and techniques; the first kind of improvement obtains factors which have a better

statistical structure such that loadings are improved, amount of variance extracted is increased, and factor independence is better. The second kind of improvement, obviously not independent of the first, is to improve the "meaningfulness" of the factors.

In our attempt to improve both structure and content, we have begun to investigate post-factor rules for the deletion and/or synonymization of words. There are basically two such rules. The first derives from clusters within the factors themselves. Thus, one factoring run yielded a factor with the days of the week heavily loaded within it; we utilized this information to produce a generic "TIME" term subsuming the days and thus opened the matrix size for the inclusion of seven additional words should this be indicated.

The second rule which we make use of in post-factoring derives from the fact that common usage holds loadings less than .30 as being fundamentally uninterpretable. Thus, we are also investigating the results from re-factoring after having dropped all words which never obtain a loading greater than .29 anywhere in the obtained factor structure. Obviously, both of these rules result in a reduction of total matrix size; one of the questions, within a UHH approach, with which we are concerned is whether we better obtain our goal of improved structure by re-factoring with a smaller size matrix or whether it is more fruitful to include additional words (formerly uncludable) now available because of the open slots in the matrix. Our initial results suggest that replacements, rather than a size reduction per se, is the more appropriate technique.

Using a UHH approach, the time for complete analysis of the same set of data mentioned earlier, has changed to approximately four hours of investigators' time as opposed to almost a month using SYN and, as would be expected, the number of computer runs has sharply reduced. Typically, four runs carry us through initial screening, parsing, deinflection, editing, factoring and re-factoring. Depending upon turnaround times and daily load requirements at the university computing-center, we can reasonably expect to complete analyses in something less than a normal work week. With SYN, turnaround time and repetitive runs, made six weeks a minimum.

The results, with UHH, have been quite encouraging. Before beginning the UHH factoring of interview 23-27 of subject PI. we had available the SYN results on that same data.<sup>2</sup> Accordingly, we used, as a partial criterion, the

<sup>2</sup> This data consists of a set of 462 consecutive psychoanalytic treatment sessions which were tape recorded and made available to us by Dr. F. Gordon Pleune. We are grateful for his help and his cooperation in the analysis of this data.

comparability of the two factor structures to make some judgment about the viability of a UHH approach. While the factor structures were not identical, there was sufficient similarity between the two to make us believe that the method should be refined further.

Table 1 illustrates the kind of structure and the basic similarity between the two approaches.<sup>3</sup>

TABLE 1

Part 1. *Comparison of SYN and UHH Factors\**

(Varimax rotated loadings truncated at  $< .30$ )

SYN FACTOR 3		UHH FACTOR 12	
<i>Old</i>	95	<i>Old</i>	94
<i>Change</i>	94	<i>Clothes</i>	91
<i>Dress</i>	85	<i>High</i>	91
<i>Look</i>	83	<i>School</i>	83
<i>Friend</i>	79	<i>Friend</i>	82
<i>Okay</i>	77	<i>Dress</i>	62
<i>School</i>	77	<i>Look</i>	53
<i>Clothes</i>	75	<i>Enjoy</i>	50
<i>Sloven</i>	72	<i>Definite</i>	47
<i>Attract</i>	70	<i>See</i>	45
<i>Relax</i>	63	<i>Long</i>	43
<i>Keep</i>	56	<i>Apologize</i>	41
<i>Apologize</i>	41	<i>Always</i>	34
<i>Good</i>	38	<i>Keep</i>	34
<i>Shave</i>	38	<i>Lot</i>	34
<i>Meet</i>	36	<i>Attention</i>	32
<i>Differ</i>	34	<i>Normal</i>	30
<i>Peculiar</i>	33		
<i>Attention</i>	32		
<i>See</i>	32		

\* SYN refers to word reductions via synonymization as well as deletion; UHH data is reduced without synonymization.

Part 2. *SYN Factor 3 High Scoring Segment*

Friend of mine and I and every time I see this old high school friend, I am always dressed in old clothes. And one day, I sort of apologized for always seeing him this way and he said he never saw anybody look as good or as well in old clothes as I do. I don't know if I do it for, I know that I have done this, that I have gone to parties not shaven and dressed this way and would go weeks in high school, not weeks but a whole week, with dressing this way or slovenly. But I, I guess it, it would attract attention. I know if I see somebody dressed this way and unshaven for a whole week I would look at them myself.

<sup>3</sup> In this, and in all other analyses reported in this paper, a combination factor of 5 has been used in preparing the data for correlation. Thus, if psychotherapy data is being analyzed each successive set of five minutes is combined and analyzed as one observation; analogously, in analysis of the book data to be reported later, each successive set of five pages has been combined and analyzed as one observation.

*Factor Scoring*

One of the more important uses for WORDS, in a post-factor environment, is for content analysis. It is of importance to us that we be able to investigate content changes over time, across speakers, under different circumstances, etc. Our first step for implementing this goal, involves locating those portions of the data base from which specific content and/or thematic areas were being elicited. The SCORE program in WORDS is designed to accomplish this task.

With the final factor structure loadings as the prime data, SCORE will scan the data base, from which these loadings were obtained, and will assign to each of the observations (segments or combinations thereof) a numeric factor score. Using only factor loadings greater than .49, this score is obtained, for each factor, by using the loading of all words on that factor, as a multiplier for the frequency of occurrence of that word in the observation being scored; these separate word-scores are then summed over all the selected words on that factor yielding a factor score. In order to afford comparability between factor scores, SCORE computes standard- as well as raw-scores for each factor on each observation.

Table 1 also illustrates the use of SCORE. The segment presented in Table 1 is the highest scoring segment for the SYN factor seen in that table. This is a typical result and there seems little question that SCORE can quite well locate that portion of the data base which is heavily saturated with the material that is helping to elicit the factor which is being scored and that the material located is consonant with the content of the factor.

*Illustrative Results*

As part of our continuing research program, we have recently begun application of the system to data other than the series of over four-hundred continuous psychotherapeutic interviews which comprised the initial data base from which the system was developed. As a suitable vehicle, we chose a set of two psychotherapy interviews recorded over fifteen years ago (1). Our choice of these two interviews was dictated by the fact that they form the nucleus of a book, *Comparative Psycholinguistic Analysis of Two Psychotherapeutic Interviews*, which was edited by Gottschalk and published in 1961: the purpose of the symposium, from which this book was generated, was to bring together several workers in the area of content analysis and to bring their different approaches and skills to bear on the same set of two interviews. The results of



these different analyses form the major part of the book; our hope was that analysis of these same two interviews by WORDS would yield data which would illustrate the utility of the system by allowing comparison of our results with those of some of the members of the symposium. Our purpose in presenting these results is not to offer or interpret further information as regards the particular case under analysis; rather, by showing some portion of our results we aim only towards demonstrating the utility of WORDS as a method for content analysis.

The data incorporated in the two interviews analyzed in the symposium is based upon two separate meetings, interviews number 8 and 18 of the patient under treatment. Following a description of the patient and the comments of the psychotherapist as to the content of the interviews, the book then details the interviews themselves, a set of physiological observations as to skin temperature and heart rate for both therapist and patient (on a minute-by-minute basis) and then presents papers by Strupp, Jaffe, Mahl, Gottschalk, et al., and DiMascio. The papers by the last four authors present materials which are reported in tabular or graphic form in the text in a quantifiable fashion; thus, Jaffe uses the verbal diversification index (type-token ratio) and a percent present-tense index, Mahl defines a speech disturbance ratio and silence quotient, Gottschalk presents and scores categories for anxiety, hostility, and schizophrenic disorganization, and DiMascio makes use of the physiological indices listed above.

In analyzing the data of these interviews by WORDS, we used the UHH approach. Following the usual pre-factor rules, we submitted a list of a hundred-and-ninety different words for factoring. With these results as a baseline we then included, as an additional set of twenty-five variables, the various indices derived from the papers presented by the members of the symposium; inclusion was done by representing the variables as though they were words with a frequency equivalent to their "score". Thus, for example, if the patient had had a heart rate of 75 in segment four of interview eight, a variable PHR was included seventy-five times within that same segment; following usual WORDS reduction and summarization procedures PHR would yield a combined frequency of 75 in segment four of interview eight. After analysis of the basic plus auxilliary data matrix, we compared this with the matrix of words only and found little basic difference between the two sets of factors. The substantive data factors were almost the same; that auxilliary data which had been submitted simply loaded *within* various of these factors.

Although we extracted twenty factors, we found to our surprise (perhaps because we were dealing with a data base

only 40% of our usual size) that almost 100% of the variance had been extracted by the first nineteen factors.

The results of the analysis have been very provocative. The significant and frequent loadings of so much of the auxilliary data clearly indicates the extent to which content themes extracted by WORDS relate to indices extracted according to other widely different theoretical constructs. It is important, however, to note that some of the relationships are probably artifacts. In Factor 13 we have a loading of .70 for "OUTWARD HOSTILITY" and a dominant loading of .96 for the word *kill*; this simply suggests that the word *kill* plays an important role in scoring the category. On the other hand, such indices as the type-token ratio, physiologic measures such as heart rate and skin temperature, percent of present tense verbs, etc. cannot be reasonably explained in terms of content artifact.

The results of this analysis are too extensive to present in entirety. Rather, we shall show four factors of the second analysis (in which the auxilliary data was included) in order to illustrate the kind of materials produced.

In the beginning of the Gottschalk book, in a chapter by Kanter and DiMascio, a summarization of the content of the two interviews is furnished by the psychotherapist. Table 2 contains a portion of the therapist's description, the significant loadings on Factor 15 of the WORDS analysis, and the segments chosen by the SCORE program as the highest loaded in the data on that particular factor.

TABLE 2

*Gottschalk Data Example 1*

Part 1. *Therapist statement.* He continued to work at understanding his feelings. In the course of this, he told of blocking himself and hurting himself and handling the humiliation by clowning and playing the buffoon as his father had before him. With a feeling of horror, he told of his identification with being publicly humiliated before those who matter. (1, p. 20)

Part 2. *FACTOR 15.* (Varimax rotated loadings truncated at <.30)

ANXIETY	92	*THERAPIST HEART RATE	-46
HUMILIATE	89	OTHER	44
*EMOTIONAL DISCOMFORT	81	*INTRAPRNL SCHIZ WITHDRWL	-38
*FREE ANXIETY	73	*BLOCKED RELATIONS	-36
SPECIAL	72	*PATIENT SKIN TEMPERATURE	-36
FAIL	64	*TOTAL SCHIZ WITHDRWL	36
UPSET	54	WIFE	35
CERTAIN	50	ALL	-34
WANT	50	EXPECT	-32
*% PRESENT TENSE VERBS	50	MANY	-31
INTELLECTUALIZE	-47	*GRATIFYING RELATIONSHIPS	31
*TYPE TOKEN RATIO	-46		

\* Non-verbal categories.

40. being the therapist to this group and then this woman on a program which I identified with. Of course I was really upset for this woman in front of her child. I remember a similar feeling of the greatness of the movie "Bicycle Thief", the humiliation of the father in front of his son when he was caught stealing the bicycle. I thought that was the greatest part of the movie, very much aroused. I am sure I cried at that part of the movie. Uh that really affected me but ah I don't know, and then yesterday I started talking about my father and his humiliation and then his handling the humiliation by clowning

41. and being buffoon and so do I, I clown and play the buffoon. Humiliation is a real issue to me, humiliation in front of people, I think of tarring and feathering someone. I think of it with horror. The humiliation with which they treated collaborators who were stripped and had their heads shaved, especially the woman. It really deeply affects me. I must identify with that — being humiliated ... Gee, I always talk about my anxiety, ha ha, I can talk freely to people, I'm very anxious, I'm very anxious over such and such, almost ready to have them say no, you weren't, or to show them that I wasn't really ...

In a later paper, Mahl presents information as to his analysis of the data using the speech disturbance ratio and silence quotient. Mahl raises a question as to what causes the variations in the speech disturbance ratio; noting that his objective measures are not designed to answer this question, Mahl nevertheless attempts to pinpoint some of the variation by noting that "the therapist's increasingly prodding, insistent questions and comments on the patient's lateness ... are associated with the progressive rise in speech disturbance level ..." Mahl identifies the 29th-32nd minute of interview 18 as being these points. Table 3 presents Factors 11 and 18 which contain the two highest loadings of SDR across all of the extracted factors. Factor scoring for Factor 11 selects segments (minutes) 26-30; scoring on Factor 18 analogously retrieves segments 31-35.

As a final illustration of the results, it is worth considering some of the physiological materials presented and analyzed by DiMascio in another chapter of the book. Because each of the participants in the symposium had utilized his own methods and skills for construction of the various indices to be applied to the data, DiMascio's ability to relate the various physiologic indices to this other data was confined to a series of correlations. What WORDS allows, on the other hand, can be easily seen by inspection of Factor 4 presented in Table 4.

Again, we should note that this re-analysis of the data from the content analysis symposium is not an attempt to offer new information or conclusions about the data analyzed by that group, although it could indeed serve such a purpose. Rather, we cite this information and show these results to begin our attempts in making an assessment

TABLE 3

*Gottschalk Data Example 2*  
(Varimax rotated loadings truncated at < .30)

FACTOR 11		FACTOR 18	
RETALIATE	97	FEAR	78
IDENTIFY	81	*STRUCTURAL SCHIZ	
CONFIDENT	80	WITHDRWL	-75
REASON	72	MUST	69
*INWARD HOSTILITY	71	YES	58
CHANGE	70	*SPEECH DIST. RATIO	53
YOU (Therapist)	68	NEW	50
SPEAK	60	RESIST	49
*SPEECH DIST. RATIO	54	COURSE	42
*TYPE TOKEN RATIO	-54	MANY	-42
*GRATIFYING RLTSNHPs	53	CONCERN	-36
SEE	47	EXPECT	-36
REACT	46	*TOTAL SCHIZ WITHDRWL	-33
SAY	44	MAYBE	-31
CHANCE	38	LESS	30
FEEL	33		
TAKE	32		
ACTUAL	31		
TALK	-31		
*FREE ANXIETY	30		
*EMOTIONAL WELL BEING	-30		
*TOTAL SCHIZ WITHDRWL	-30		
High Scoring			
Segments or	26-30		31-35
Minutes			

\* Non-verbal categories

TABLE 4

*Gottschalk Data Example 3*

Part 1. *FACTOR 4.* (Varimax rotated loadings truncated at < .30)

FLY	98	HOPE	54
BOY	95	REMEMBER	51
LET	94	*OUTWARD HOSTILITY	47
MUMPS	93	GOOD	44
AGGRESS	91	*INTRAPRSNL SCHIZ	
BACK	88	WITHDRWL	42
*THERAPIST HEART RATE	-77	*SELF ESTEEM	40
GI	68	MAYBE	-39
*OUTWARD HOSTILITY		PRETTY	39
THEME	65	SHOULD	39
FEW	64	*SPEECH DIST. RATIO	-39
*PATIENT HEART RATE	-61	FAMILY	-32
*PATIENT SKIN		CAN	31
TEMPERATURE	-58	MANY	-30
*INTERPRSNL SCHIZ			
WITHDRWL	-56		

\* Non-verbal categories

Part 2. *High Scoring Segment 52.*

(Prior to this segment, patient describes a private beach used by his family and his discovering a group of soldiers there one day who refused to leave. He speaks of his maneuvers to get them to go, his hopes that they will contract his son's mumps, his wife's fear he will get into a fight; he talks of getting ready to flycast the beach area and remarks that if one is not skillful people behind the caster can sometimes get hurt. He states his inability to tell the soldiers:)

... would you mind moving I'd like to cast this area, I don't let my kids sit there, which I don't. Uh, I just didn't say anything. I let the fly fall a few times thinking I hope one of them gets the fly, and then being afraid though to really hook someone with a fly or whip them with it. Actually, it's a whipping, a good sharp slap. Being afraid I realized well here I'm going through all this indirect aggression and suffering it through and I can't really express it. Even if it were a physical fight I would have been proud of myself to have been able to do it and feeling that I could have carried it off. I was physically in wonderful shape. I'd been swimming and I'm usually in good shape anyway, thinking if it would just help establish a relationship between me and my son. I'd just seen the picture — in which ...

of the validational properties of WORDS and of the factors and factor structures which it produces.

We have been concerned about the validational properties of this method since its very inception. While there are many validational criteria that might be applied, two major aspects of validity have seemed of primary import to us. The first of these two criteria concerns the extent to which the factor "fits" the data. That is, considering each factor independently how well does it identify its portion of the data (as defined by high factor score) and to what degree is that portion of the data selected consonant with the factor. The second validational criterion concerns the factor configuration derived; how well can the data base, or portions of the data base, be described in terms of the factors.

Our attempts to answer these questions with the kinds of data for which the system was originally developed — psychotherapeutic protocols — is difficult. We used the Gottschalk data in our hopes that it would offer enough ancillary information to help us assess these questions. While the data did indeed make it easier for us to assess our factors — by utilizing descriptions of the protocols from the book — it demonstrated that, once again, we would be forced to buttress these factors by judgmental statements made by others as to what was indeed the content of the data base. Since it is precisely in order to avoid such dependence upon judgmental techniques that we developed the system, we felt dissatisfied with the results of the Gottschalk analysis. It seemed clear that what we needed, for analysis, was a data base that had clearly defined content that was known to large numbers of people. In short, we wanted a data base whose content was clear enough and public enough to make a "face validity" approach towards factor assessment a reasonable tenable procedure.

Accordingly, we turned our attention to famous children's books. Such books are relatively short, utilize somewhat restricted language, tend to have clearly defined content, and are usually known, in broad outline, to many people. Having examined a number of possibilities, we selected Frank Baum's *The Wizard of Oz*. This story

satisfies all the criteria mentioned above and is certainly one of the most well-known children's books of all time.

In the analysis of *Oz*, we felt that we were posing a fairly stringent but appropriate test for the WORDS System. The major themes of *Oz* are well known both from the book and the movie.<sup>4</sup> Certainly, if the method is viable, one must expect to see content themes and/or materials clustering around the Tin Woodman, Dorothy's flight in the cyclone to the land of Oz from her home in Kansas, the Cowardly Lion, the Scarecrow, the Wicked Witch, the Wizard himself, etc.

*Oz* was the largest single data base we have ever analyzed, numbering almost 42,000 words. The data was analyzed according to the UHH rules mentioned earlier. After assignment of parts of speech, and deinflection to root form, we were left with a data base of approximately 1400 different words. This list, in frequency order, was down-counted in order to obtain the 215 highest frequency words; these words constituted the UHH analysis. Twenty factors were extracted, rotated by varimax, and factor-scored for each chapter. The factors, presented in Appendix 2, account for 80% of the variance represented by the initial  $215 \times 215$  matrix submitted.

The two questions posed earlier were examined in the light of these results. The first of the questions, the validity of each factor for its eliciting content, may be examined by referring to Table 5. This table describes each of the factors both in terms of its four highest loading words and in terms of a brief description of the content of the total factor. With the presentation of each factor will be found that chapter yielding the highest factor score for the factor and the title of the chapter. Because of the nature of the chapter headings, it can be seen that twelve of the twenty factors can be immediately verified by inspection of the key loadings on the one hand and the chapter title (in which they occur most heavily) on the other; thus, Factor 2 is most heavily loaded into Chapter 5, "The Rescue of the Tin Woodman", Factor 8 in Chapter 14, "The Winged Monkeys", etc. The remaining eight factors do not automatically link with the chapter titles;

<sup>4</sup> *Oz*, now in the public domain, is available in numerous editions. We checked several to confirm the fact that ours was standard in content. Two changes were made to the book in our analysis: (1) Chapter 20, "The Dainty China Country", was omitted since it added many different words for a very small increment in total data; (2) Chapters 23 and 24 were combined since Chapter 24 is but one-half page long. It is worth noting that several differences exist between the movie and the book; we mention the major ones in order not to confuse the reader: The movie has ruby slippers, the book silver shoes; the movie pays much attention to Kansas and reproduces the Kansas characters in Oz while the book does neither; the movie has Oz as a "dream", in the book Oz is "real"; finally, the movie omits any reference to the Kalidahs, the Golden Cap, the Hammer Heads, and the Field Mice.

TABLE 5  
*Identification of Factors in Oz*  
 (For complete factors cf. Appendix 2)

Factor and four high loaded words	Factor content	Highest loading chapter and title	Factor and four high loaded words	Factor content	Highest loading Chapter and title
1. <i>ax</i> <i>oil</i> <i>Tin Woodman</i> <i>tin</i>	The Woodman, his body, one time romance with Munchkin girl.	5. "The Rescue of the Tin Woodman"	11. <i>spectacles</i> <i>Guardian/</i> <i>Gates</i> <i>want</i> <i>Emerald City</i>	The entry to the Emerald City via the Guardian of the Gates who requires all to wear green spectacles on entrance.	10. "The Guardian of the Gates"
2. <i>Oz</i> (HEAD) <i>Oz</i> (LADY) <i>kill</i> <i>send</i>	Meetings with Oz in his various disguises. Oz's demand that they kill Wicked Witch of West in order to receive their requests.	11. "The Emerald City of Oz"	12. <i>flower</i> <i>stork</i> <i>sleep</i> <i>bright</i>	The poppy field with its "deadly" fragrance which puts anyone going through to sleep forever.	8. "The Deadly Poppy Field"
3. <i>farmer</i> <i>brick</i> <i>scarecrow</i> <i>road</i>	The Scarecrow, his creation, stupidity, clumsiness, need for a brain.	3. "How Dorothy Saved the Scarecrow"	13. <i>terrible</i> <i>man</i> <i>home</i> <i>promise</i>	The wizard; his trickery of the group; his broken promise to each of them.	15. "The Discovery of Oz the Terrible"
4. <i>Munchkins</i> <i>Witch</i> <i>East</i> <i>woman</i>	People of the East; freed by Dorothy whose house killed the Witch.	2. "The Council with the Munchkins"	14. <i>silver shoes</i> <i>end</i> <i>water</i> <i>Wicked Witch</i>	The magic shoes gained by Dorothy from Witch of East and coveted by Witch of the West leading to her downfall.	12. "The Search for the Wicked Witch"
5. <i>Uncle Henry</i> <i>house</i> <i>Aunt Em</i> <i>bed</i>	Dorothy's aunt and uncle, her home in Kansas, her trip to the land of Oz.	1. "The Cyclone"	15. <i>room</i> <i>green</i> <i>soldier</i> <i>dress</i>	The palace of Oz in Emerald City where the Group await their audiences with Oz.	11. "The Emerald City of Oz"
6. <i>wolf</i> <i>lie</i> <i>crow</i> <i>die</i>	Attacks by animals on one or more members of group during book; chief such attack is instigated by Wicked Witch.	12. "The Search for the Wicked Witch"	16. <i>balloon</i> <i>air</i> <i>silk</i> <i>make</i>	The device which first brought Oz to the magic land and which he and Dorothy hope to use to return home again.	17. "How the Balloon was Launched"
7. <i>coward</i> <i>near</i> <i>Cowardly Lion</i> <i>heart</i>	The Cowardly Lion, his attack on Toto, Woodman, Scarecrow.	6. "The Cowardly Lion"	17. <i>pretty</i> <i>Hammer Heads</i> <i>kind</i> <i>dress</i>	The Hammer Heads, people with projectile heads who bar the group in their journey south to see Good Witch.	22. "The Country of the Quadlings"
8. <i>Gaylette</i> <i>Quelala</i> <i>time</i> <i>Winged</i> <i>Monkeys</i>	The Winged Monkeys, creation of their controlling agent — Golden Cap — by the sorceress Gaylette.	14. "The Winged Monkeys"	18. <i>pole</i> <i>river</i> <i>middle</i> <i>let</i>	The ubiquitous pole with which Scarecrow has much trouble. Chief problem: stuck on pole in middle of river.	8. "The Deadly Poppy Field"
9. <i>Winkies</i> <i>tinsmith</i> <i>set</i> <i>careful</i>	Winkies, people of West, slaves to Wicked Witch, freed by Dorothy; their special friendship to the Woodman.	13. "The Rescue"	19. <i>tree</i> <i>side</i> <i>branch</i> <i>seem</i>	The various forests and the scenery encountered by the group in their travels.	19. "Attacked by the Fighting Trees"
10. <i>mouse</i> <i>Queen Mouse</i> <i>safe</i> <i>turn</i>	The field mice; Queen Mouse, saved by Woodman directs mice to save Lion from poppy field.	9. "The Queen of the Field Mice"	20. <i>courage</i> <i>real</i> <i>brain</i> <i>very</i>	A generalized factor about the "needs" of the group: courage for the lion, a heart for the woodman, brains for scarecrow.	15. "The Discovery of Oz the Terrible"

however, all relate quite appropriately as can be seen in Table 5 under the "content" which relates the factor content to the chapter content. Thus, Factor 2 describing the various meetings of the group with the disguised Wizard occurs in Chapter 11, "The Emerald City of Oz"; it is, however, in that chapter that all of these particular

meetings take place. Likewise, Factor 17 achieves its main loadings in Chapter 22, "The Country of the Quadlings", a chapter almost completely devoted to the group's attempt to reach the country of the Quadlings and their difficulty in following their route since it is barred by the Hammer Heads whom they must circumvent.

Our reaction to these results, then, is that the factors are indeed relevant to the content areas which they identify. This approach has dealt only with the *highest* factor-loading score for each factor; factors, however, have a factor score for *each* of the chapters under analysis and we therefore turned our attention to those areas in which factors were occurring at some significantly high level. The most effective way to present this kind of analysis is on a chapter-by-chapter basis; thus, we are raising the question as to how well the *chapters* are described by the factors as opposed to how well a particular factor dovetails with its highest scored segment. Table 6 presents this chapter-by-chapter analysis. The chapter, its title, and a brief description of its content are furnished; for each chapter, the factor number, standard-score for that factor, and a mnemonic based on the basic factor structure are listed. All factors with a standard score of +2.00 or greater are presented. Factor scores with asterisks indicate the highest score ever obtained by the factor.

Again, the results are very encouraging. In Chapter 1 which is devoted to Dorothy's home in Kansas and her trip, the only factor scoring at 2.00 or more is Factor 5, "Home". In Chapter 2 which details her arrival in Oz, the meeting of the Munchkins, receipt of the Silver Shoes, and her plans to see the Wizard, Factors 4, 14, and 13 score at or above 2.00 thus identifying the "Munchkins", the "Silver Shoes", and the "Wizard".

Close inspection of Table 6 will reveal that certain themes, noted in the "contents" column are not being supported by appropriate factors and that one chapter, Chapter 21, has no factor scored at +2.00 or greater. This chapter is concerned with the cowardly lion's "election" to become king of beasts after he has killed a monster spider terrorizing the other animals in the forest. *Spider* and *monster*, the two words most closely associated with the chapter's theme, did not obtain frequencies high enough to be included in the 215 word list for factoring. Whether enough other words, themselves associated with these two key words are included in the 215 list and potentially available in factors beyond number 20, is presently something we cannot ascertain. The issue of what words are included, and the consequences of missing words, is a topic to which we shall later turn our attention. It is perhaps worth noting that the highest scoring factor in Chapter 21 with a standard-score of +1.29 is Factor 7, "Cowardly Lion".

After concluding this initial analysis of *Oz*, we decided to try another approach which would yield information as to the effects of word choice on factor extraction. In order to do this, we again started with the list of 1400 different words (used in the last analysis to yield the

TABLE 6

*Relationship between Factors and Chapters in Oz*  
(For complete factors cf. Appendix 2)

Chapter and major themes presented	Factor and mnemonic	Z-score
1. "The Cyclone". Aunt Em, Uncle Henry, Dorothy, Toto; trip in cyclone begins.	5. Home	7.81*
2. "The Council with the Munchkins". Arrival in Oz where Dorothy meets Munchkins and Good Witch of the East; learns she has killed Wicked Witch of East and freed Munchkins; gets silver shoes; decides to see Wizard to go back to Kansas.	4. Munchkins 14. Silver Shoes 13. Wizard	5.17* 2.50 2.40
3. "How Dorothy saved the Scarecrow". Begins trip, spends night Munchkin farm, meets Scarecrow, gets him off pole; he joins trip to seek brains.	3. Scarecrow 4. Munchkins 18. Pole	5.61* 3.19 2.33
4. "The Road through the Forest". While walking, Scarecrow tells story of his creation his failure as a scarecrow, his stupidity, etc. They find a cottage to spend the night.	3. Scarecrow 6. Animal Attack 13. Wizard 19. Trees	4.40 3.04 2.34 2.32
5. "The Rescue of the Tin Woodman". Woodman is found, oiled, freed from rust; tells how Wicked Witch of East enchanted his ax making him cut off arms, legs, etc. to prevent marriage to Munchkin girl; joins trip to seek a heart.	1. Tin Woodman 4. Munchkins 13. Wizard 3. Scarecrow 19. Trees	7.33* 4.78 3.31 2.15 2.10
6. "The Cowardly Lion". Lion attacks Woodman and Scarecrow, to no avail, turns to Toto, Dorothy slaps him, cowardice revealed; learns purpose of trip and joins to seek courage from Wizard.	7. Cowardly Lion 20. Needs	6.04* 2.91
7. "The Journey to the Great Oz". Must cross large ditches on back of Lion; attacked by Kalidahs; Lion delays them while Woodman chops down log bridge; escape; come to river and prepare to build a raft to cross.	19. Trees 18. Pole 7. Cowardly Lion	5.73 2.29 2.00
8. "The Deadly Poppy Field". Build raft, begin to cross, Scarecrow stuck on pole middle of river; ashore, others solicit stork to carry Scarecrow back; she does; come to poppies; Lion, Dorothy, Toto succumb; Woodman and Scarecrow carry Dorothy and Toto but cannot move heavy lion.	18. Pole 12. Poppies	8.99* 7.66*

TABLE 6 (Cont.)

*Relationship between Factors and Chapters in Oz*  
(For complete factors cf. Appendix 2)

Chapter and major themes Presented	Factor and mnemonic	Z-score	Chapter and major themes presented	Factor and mnemonic	Z-score
9. "The Queen of the Field Mice". Woodman saves Queen Mouse from wildcat, she offers help, Woodman makes "truck"; mice drag Lion from poppies.	10. Mice	6.62*	16. "The Magic Art of the Great Humbug". Oz puts pin and needles in Scarecrow to make him "sharp", gives the Woodman a heart inside his chest, feeds the Lion a dose of courage.	20. Needs	2.46
10. "The Guardian of the Gates". Lion awake, trip continues; spend night at cottage discussing hope for success with Oz; reach Emerald City, admitted by Guardian of Gates; all must wear spectacles to prevent blindness from brilliance of city.	11. Emerald City 13. Wizard 20. Needs	8.95* 2.52 2.20	17. "How the Balloon was Launched". Oz and Dorothy try to leave via another balloon they have built. Balloon, with Oz, leaves without Dorothy who is chasing Toto.	16. Balloon	8.10*
11. "The Emerald City of Oz". Enter city, reach palace, each assigned separate room pending their separate audiences with Oz as a different figure; Oz refuses each pending death of Wicked Witch of West; they agree to kill her.	2. Audience 15. Palace 7. Cowardly Lion 20. Needs 3. Scarecrow	6.95* 6.62* 3.17 2.69 2.40	18. "Away to the South". They ask Winged Monkeys if they can take Dorothy home; monkeys explain they are powerless outside of Oz; decide to go South to seek aid from Good Witch.	8. Winged Monkeys	3.45
12. "The Search for The Wicked Witch". Leave city, Guardian of Gates removes spectacles; in land of West seen by Wicked Witch who sends wolves, crows, etc. to kill; all fail; sends Winkies — slaves — also fail; uses Golden Cap and sends Winged Monkeys who destroy Woodman and Scarecrow; rest prisoners; Witch's greed for Silver Shoes angers Dorothy who tosses water on her and melts her.	6. Animal Attack 14. Silver Shoes 9. Winkies 8. Winged Monkeys 7. Cowardly Lion 11. Emerald City	8.42* 7.57* 4.57 3.00 2.40 2.39	19. "Attacked by the Fighting Trees". Leave city again and start south; in forest are attacked by trees; they escape and continue on.	19. Trees 11. Emerald City	7.92* 2.72
13. "The Rescue". Ask Winkies, freed from Witch, to rescue Woodman and Scarecrow; both saved; Winkies carefully repair Woodman; start for Emerald City; Dorothy takes pretty Golden Cap.	9. Winkies	9.97*	20. "The Dainty China Country". This chapter was not analyzed. Cf. fn. 4.	No factor score value > + 1.99	
14. "The Winged Monkeys". Lost; call field mice who suggest Golden Cap; call Winged Monkeys and begin flight to Emerald City; on way, Monkeys tell their story and explain charm behind Cap.	8. Winged Monkeys 10. Mice 13. Wizard 18. Pole	7.43* 3.22 2.49 2.29	21. "The Lion becomes the King of Beasts". Lion is asked by forest animals to destroy monster spider; does so; is elected "king" of beasts.	17. Hammer Heads	9.32*
15. "The Discovery of Oz the Terrible". Back at Emerald City Oz delays seeing them; agrees at threat of Monkeys; discovery of Oz as humbug; he tells of trip via balloon from Omaha; they still believe he can grant their wishes.	20. Needs 13. Wizard 16. Balloon 11. Emerald City	8.38* 4.11* 2.95 2.07	22. "The Country of the Quadlings." Prevented from crossing hill by Hammer Heads who knock them down with their projectile heads; they call Winged Monkeys who carry them over hill to castle of Good Witch of the South.	5. Home	2.01
			23. "Glinda grants Dorothy's Wish".		
			24. "Home Again". (Analyzed together; cf. fn. 4.) Witch arranges for Scarecrow to rule Emerald City, Woodman the Winkies, Lion the animals, returns Golden Cap to Monkeys, shows Dorothy use of Silver Shoes; Dorothy returns home to Kansas.		

\* Highest score achieved by the factor anywhere in the book.

downcounted 215 list) but on this occasion we approached the list in a completely different fashion. Disregarding any attempt at objectivity we carefully selected a set of 215 words corresponding to our knowledge of the book and our hopes as to what factors would be extracted. In short, we loaded the matrix for analysis with the *best* set of 215 words we could possibly choose. From that point on, analysis was conducted according to standard procedures and twenty factors were extracted.

Space does not permit us to present this complete factoring run here. Inspection of the two factor structures, that from the "downcounted" data and that from the "chosen" data was very encouraging; from our knowledge of the book we found no difficulty in matching sixteen of the twenty downcounted factors with those in the chosen word results. Table 7 presents these matchings and indicates those factors that could not be matched from one analysis to the other. While it was clear to us that, for example, downcounted Factor 15 was the "same" factor as chosen Factor 16, it was also clear that this congruence might not be apparent to someone who was not intimately familiar with the book. We therefore searched for a way in which to demonstrate this similarity. Our technique was to take the factor scores for Factor 15 across the entire book and correlate them with the factor scores for the (ostensibly) matched Factor 16. We did this for each of the fifteen unilateral matches in the data and the results were startling. The last column of Table 7 presents these correlations; there is one correlation of .49, one of .79, one more at .86 and the remaining twelve correlations all have values equal to or greater than .97.

Because of the size of these coefficients, we were concerned that some artifact might be operating. We could think of only one: the number of words in common between each two potentially matched factors. Thus, despite our different methods of choice, were a pair of factors yielding high correlations based on basically the same set of highly loaded words (loadings greater than .49), we would be building correlations based mainly on a common set of frequencies.

Table 7 also examines this possibility and demonstrates that it is untenable. Using the downcounted factors as the base, we established the number of downcounted factor-scoring words in common with the possibly matched chosen factor-scoring words. This data is presented both as a fraction and a percentage; inspection of the correlations clearly demonstrates the complete lack of effect of the "commonality" on the correlation.

The only other possibility for an artifact of which we were aware also hangs on the question of common words. Thus, if the total 215 matrix of downcounted data is

TABLE 7

*A Comparison Between Downcounted Word Factors and Chosen Word Factors*

Factor mnemonic	Downcount Factor	Chosen factor	Overlap	Correlation
Tin Woodman	1	2	7/14 50%	.99
Audience	2	3	5/11 45%	.97
Scarecrow	3	7,19		
Munchkins	4	5	5/9 56%	.97
Home	5	4	4/10 40%	.98
Animal Attack	6	1	4/8 50%	.86
Cowardly Lion	7	6	4/5 80%	.98
Winged Monkeys	8	20	3/8 38%	.79
Winkies	9	18	2/7 29%	.97
Mice	10	9	3/7 47%	.99
Emerald City	11	12	4/5 80%	.99
Poppies	12	13	4/9 44%	.98
Wizard	13	10	4/8 50%	.69
Silver Shoes	14	14	3/7 43%	.99
Palace	15	16	1/9 11%	.97
Balloon	16	8	1/7 14%	.49
Hammer Heads	17			
Pole	18			
Trees	19			
Needs	20			
King of Beasts		11		
Leaving Oz		15		
Kalidahs		17		

composed of a high proportion of the 215 chosen word data, its similarity in total words might allow extraction of basically identical factors with the few different words holding onto the high loading areas. It is sufficient to note that only 98 words in the downcounted data and chosen word data are in common, i.e., 46%.

DISCUSSION

The results presented in this paper have led us to four major conclusions. We shall discuss each in turn.

*Factor Validity*

We raised two questions concerning factor validity in our analysis of the *Oz* data. The first tested validity by asking the extent to which each factor "fit" the data; were the segments for which a factor was most heavily scored good representations of that factor and vice-versa. Secondly, we asked the extent to which the configuration of factors was able to describe the data base itself. Both of these questions were answered affirmatively in the *Oz* data and

while the second has never been tested on any other database, the question of "fit" has been examined now over some widely differing sources of material. We take it then that the approach we have pursued is indeed a viable one and that the factors which it extracts are valid representations of the data which elicits them.

#### *Data Base Description*

While raised mainly as a validation criterion, the ability of the factors to describe the *Oz* data on a chapter-by-chapter basis has suggested the possibility of a configurational approach towards the description of major content. In such an approach, the information utilized would be the set of factors and their standardized scores operating configurationally over the various units of the data base under analysis; it might thus be feasible to describe changes in a data base by noting the configurational changes taking place. We have done no research into this area as yet but there are several statistical techniques available for configurational analysis and we shall begin to investigate them soon.

#### *Word Selection*

The high degree of correlation between the *Oz* downcounted factors as contrasted with the chosen-word factors implies some degree of insensitivity, in the factor *structure*, to the words available for building these factors. We have no information, at present, as to the degree of that insensitivity other than that furnished by the analyses on *Oz*. Thus, we do not know whether the 46% overlap figure between the two sets of words analyzed represents a figure that is adequate because of the nature of the data; it is possible that a much greater overlap may be required when dealing with data whose interrelationships are more subtle and complex than is the case with *Oz*. Further, our UHH rules, operating on the downcount analysis, fairly well restrict the data analyzed to nouns, adjectives, verbs, and adverbs. Whether the kind of concordance obtained with *Oz* would hold when data is analyzed with a heavy preponderance of articles, prepositions, conjunctions, etc., is something we cannot presently answer. Nevertheless, *some* degree of insensitivity is a clearly demonstrable finding and the fact that this robustness does exist has considerable implications for our future lines of research.

We have noted earlier in the paper that the major critical

problem we face in the operational procedures of WORDS is in the selection of what words to retain for analysis. We have tacitly assumed that the deletion of words would cause changes in the factors obtained and in the configurational relationships among these factors. Our first indication that this assumed sensitivity might be overstated came when we began our investigation into a UHH rather than an SYN approach. As noted in the paper, both approaches yielded factors many of which could be related one to the other. Several of the factors, however, could not. The *Oz* data confirms this finding. Of the twenty factors extracted in both the downcounted and chosen-word analyses, four of the downcounted factors could not be matched to those in the chosen-word set while three of the chosen-word factors were not matchable; (the discrepancy exists because one of the downcounted factors — the Scarecrow — separated into two factors in the chosen-word data). Both sets of unmatched factors, the three from the chosen-word data and the four from the downcounted data, were equally "good". Table 7 shows that the downcounted data extracted unmatched factors for the Hammer Heads (Factor 17), The Pole (Factor 18), the Fighting Trees (Factor 19), and the Needs (Factor 20). Likewise, the chosen-word data allowed extraction of unmatched factors representing the Kalidahs (Factor 17), Dorothy's flight from the Land of Oz (Factor 15) and King of Beasts (Factor 11). This last factor is worth noting because it will be remembered that no factor in the downcounted analysis obtained a standard score of over +2.00 in Chapter 21, "The Lion Becomes the King of Beasts"; scoring Factor 11 from the chosen-word analysis, however, yields a factor score of +8.00 for Chapter 21. Clearly, then, as we have analyzed the data, a change in the set of words submitted for analysis *does* cause some change in the factors extracted. We are not sure but that some of these "missing" factors might not have been extracted had we continued factor extraction past our limit of 20 but this is problematic. Nevertheless, the fact remains that a considerable change in submitted words still yielded a match on sixteen of the twenty factors, a matched percentage of 80%.

The implication of this result seems clear: We have some margin of safety in our selection of what words will be submitted for analysis. This implication, coupled with the fact that the downcounted data furnished as good a description of the data as did the chosen-word analysis, strongly supports our feeling that a UHH approach, with its advantages of extreme speed and objectivity, is the correct way to pursue our future developments in the system. We shall speak more, later, about some of our plans for increasing the efficiency of the UHH approach.



## *Inclusion of Non-Word Materials*

It will be remembered that we did two analyses on the Gottschalk data presented earlier. In both analyses we submitted the same list of 190 different downcounted words for factor extraction but in the first analysis these words were examined by themselves whereas the second analysis added twenty-five nonverbal variables for analysis; these twenty-five variables were composed of various categories of content analysis, physiological data, etc. The results indicated that the factor structure of both analyses was almost identical insofar as the words and their loadings were concerned. What happened was that the non-verbal data tended to appear within factors already established, in its absence, during the first analysis.

We feel that this finding offers a line of development for the use of WORDS that is quite interesting. We can envision at least two major uses for the inclusion of non-verbal data along with the submitted list of words. On the one hand, we believe it would be interesting to see what descriptive value could be furnished by such non-verbal data in order to add to the utility of the extracted factors and to clarify further the segments of the data base chosen for examination because of high scores on that factor. On the other hand, we can see a use of WORDS in the developmental phases of a categorization system which would allow the developer or investigator of that content-analytic method to investigate the degree to which his content categories are intrinsically related to the various content materials uncovered by the factor structure itself.

## PLANS FOR THE FUTURE

### *Methodologic Issues*

There are three major research areas which we intend to pursue and which we should now like briefly to detail.

### *Statistical Word Selection*

The results presented in the paper have clearly suggested that we have some flexibility in the choice of words to be submitted for analysis. We have long been interested in objective methods for such word selection but were troubled because objectivity seemed to demand a price in factor interpretability and meaningfulness. The UHH approach has, however, given us sufficient encouragement to look into this issue further.

In the standard UHH approach as we had formulated

it, the final selection of words for analysis was done by downcounting a frequency ordered list of remaining words; this technique, of course, guaranteed that the highest frequency words would be included for analysis with the low frequency words being deleted. There is nothing compelling, however, about such an approach. Rather than depending on frequency selection, we plan to pay extensive attention to that correlation matrix which precedes the factoring run as a method for making word choices. Our reasoning is as follows: An intercorrelation matrix computation is very much faster than a factoring run given the same size of matrix input. With correlations on a  $215 \times 215$  matrix being computed across, say, a hundred observations, we can reasonably expect the correlation program to run at least eight to ten times faster than the factoring run that will follow it. Further, a factoring run demands more of the machine's available core capacity than does a correlational approach and it is therefore feasible to run larger matrices through a correlation program than through a factor analysis program on a machine of a given size. Putting these facts together, we intend to allow correlational runs on word matrices of orders running to about 800, a size which is usually capable of holding *all* different words left in a data base after deinflection to root form. We shall then utilize another program to inspect this matrix of intercorrelations and to choose from it the 215 words best meeting a set of criteria that will ensure the development of "good" factor structure if, indeed, such factor structure is inherent in the data. There are at least two criteria that make for "good" factors; one, which has been discussed extensively during the paper, is the validity of each of the factors and of the factor configuration. Another, stemming from common factor analytic usage, is simply the loadings on each extracted factor — how much variance do they extract — and, as a result of summation across the factors, how much variance does the total extracted factor set remove from the input matrix. We do not believe that these criteria are independent; we have found, in prior research, that good statistical factors tend to be the more valid factors for our use.

The statistical criterion for factorial "goodness" then is one approach that we can utilize. Since it is a fact that correlational matrices with very high overall correlations will yield better statistical factor structure than those with very low overall correlations, we should like to investigate the possibility that good factor structure can be obtained by eliminating words whose overall correlations are low in favor of those with high mean correlations. Another, non-independent approach, may lie with the variance of the correlations obtained between a given word and all

other words in the matrix. Other things being equal, high variance is better than low for factor analytic operations; such variance is obviously not independent of the mean correlational level associated with a word but may allow selection of words for retention from among other words with equal mean correlations.

Should such an approach prove productive, we would have a completely objective and very fast method for the selection of the "right" words to be included in a factoring run.

#### *Measurement of Specific Words*

We are interested in exploring a somewhat different approach towards the "measurement" of specific and key words in a data base. As it is currently, all words are potentially admissible for analysis in a WORDS run. If a word, for example *mother*, is in the data base then it stands a chance of admission for analysis that is independent of its meaning. We believe, however, that if the word *mother* is an important one for the user of the system, it might be fruitful to analyze the data base deliberately leaving out the word in any such analysis. We should then be interested to see what happens to factor scoring techniques as they are applied, for all factors, on those observational segments where *mother* does not appear vs. those where it does. We have no evidence on what the effect of such an operation will be but think the possibility of success interesting enough to give it some priority in our future research efforts.

#### *Reprogramming of WORDS*

The University of Rochester has recently acquired an IBM 360 model 50 computer and will, within another year, update that machine to a model 65. WORDS will be reprogrammed to run on the 360. Programming on the IBM 7074 had, of necessity, to be in assembler language because no higher level language existed capable of doing the job. PL/I has met that need and reprogramming for WORDS will be in that language. Because the 360 is a very popular machine, we shall, for the first time, have the ability and opportunity to make the WORDS System available to others outside the University proper.

While PL/I cannot come even close to matching the efficiency of assembler language coding, it allows us a high degree of programming efficiency and offers the distinct possibility, within the next two years, of being implemented on a number of other manufacturer's machines; this, of course, would allow even further dissemination

of the WORDS System. Further, with the increasing speed of third generation machines, the overhead generated by PL/I should be more than compensated for by the increased operating efficiency of the target computer.

In this reprogramming, we shall begin investigation of a data flow logic which we hope to implement. As originally constructed, WORDS was based on the concept of repetitive runs on the computer for the purpose of data reduction with each run taking its input from the prior run's output. With the marked success we have obtained in a UHH approach, with much faster and larger machines available, and with the possibility of word selection being accomplished by a statistical criterion embodied in the correlation matrix, we believe it will be possible to reduce the complete analysis of a data base into two runs on the computer. The first, and more trivial, of these two runs would be for purposes of correcting spelling and any other errors that have crept into the data base during punching and initial entry into the system. The second run would then take place in a fashion somewhat similar to the following: The data would have parts of speech assigned, would go to an analytically oriented deinflection routine, would have words changed and/or deleted according to pre-set rules, e.g., delete all non-verb forms of *like*, and would have all words whose frequency is equal to or less than some pre-set criterion deleted; the remaining words would then be readied for a complete intercorrelation matrix whose results, as mentioned earlier, would be used to select the *N* highest correlating words for submission to factoring. Results of the factor procedure would be automatically submitted for rotation and the rotated factors would be channeled through for factor scoring on the original data with results of the scoring being made available graphically (for plotting offline) as well as in their usual printed form.

While an automatic procedure of this kind must await the results of our research into the effect of using the correlation matrix as the statistical criterion for selection of the factoring matrix, the programming and systems logic embodied in the preceding description are well within the state of the art of both hardware and software of present third generation machines. Indeed, the entire second run we have described, assuming a thousand different words for initial screening, with final factoring on approximately a  $200 \times 200$  matrix should run in somewhat less than two hours on the model 65 IBM 360 that will soon be available at the University of Rochester. Thus, the possibility of utilizing WORDS on an almost completely automatic, and therefore almost completely objective, approach towards the analysis of major content cluster is potentially quite feasible.

## REFERENCES

- Gottschalk, L. A., Ed. *Comparative psycholinguistic analysis of two psychotherapeutic interviews* (New York, International Universities Press, 1961).
- Iker, H. P., and N. I. Harway, "Computer analysis of content in psychotherapy", *Psychological Reports*, 14, (1964), 720-722.
- , "Objective content analysis of psychotherapy by computer", in K. Enslein, Ed., *Data acquisition and processing in biology and medicine*, Vol. 4 (New York, Pergamon Press, 1966).
- , "A computer approach towards the analysis of content", *Behavioral Science*, 10 (1965), 173-183.

## APPENDIX I

### *Current Structure of Words*

WORDS currently consists of forty programs. Ignoring those programs which can only be called internally (by another program) and those called by a "package" call (causing internal manufactures of a set calling sequence), there are thirty-two programs available to the user. Each of these programs will be briefly described later. Initially, however, we will detail the basic system-, data-, and program-organization of WORDS.

### *Systems Organizations*

To use WORDS, a series of control cards must be prepared by the user indicating what programs are to be called, in what order, what each program is to do, where each program is to locate and leave its input and output data.

Since preparation of control cards can be complex, it is important that all such cards be extensively screened before allowing a run to begin. This function is satisfied by requiring the first program in any WORDS run to be CHECK. CHECK will subject every control card to an extensive series of generalized validity checks and then further check each card for any idiosyncratic forming peculiar to the particular program being called. CHECK allows the run to proceed only if all cards are error free. After receiving the last control card, CHECK then issues an internal call for the administrative WORDS monitor, MNTRA.

MNTRA uses the control cards, passed by CHECK, and (1) sets up a general calling configuration for the entire run, (2) replaces sort parameter mnemonics with actual sort fields for later use by the sort programs, (3) imposes extensive configurational checks on the run to rule out any logical impossibilities (any of which would then result in a dump of the entire run), (4) forms the entire calling sequence, package generated sub-sequences, sort mnemonics, etc., into a continuous block of data which

is written to the online 1301 disk file to form a common communications block and then (5) issues a call for the executive WORDS monitor, MNTRB.

MNTRB is normally called on completion of each program in the calling sequence. Upon obtaining control of the machine, MNTRB (1) retrieves the communications block, (2) maintains a record of elapsed time for the just completed program, (3) furnishes the next program in line with necessary data I/O information, (4) indicates where the program may obtain message space if it is needed and (5) where the program's specification card, carrying auxiliary data, may be found; (6) if the program is a sort, transfers necessary sort control information parameters to the core of the machine, (7) updates the communications block and re-writes it back to the disk and (8) issues a call for the next program in the run configuration.

Normally, each program called returns control to MNTRB after completion thus again beginning the series of operations noted. The program PRINT, however, the last called program in every run, does not release to MNTRB but rather to the resident monitor which then terminates the job, tallies total time, and moves on to the next job on the queue.

Throughout the entire run, extensive error trapping procedures are activated. Immediately following CHECK, MNTRA makes certain changes to the core-resident linkages into the resident monitor (later restored by PRINT) in order to prevent normal error recovery under control of the resident. Whenever an irremediable error occurs that is not trapped by a WORDS program itself (typically, I/O or arithmetic), a default branch is taken by the machine operator to keep the queue moving. The changes made to core-linkage by MNTRA overrides this default branch and forces control to another linkage (also provided by MNTRA) which automatically issues a call for PRINT and notes the name of the offending program. On entry, PRINT determines if it has been called normally (by MNTRB) or not. If not, appropriate warning messages of an impending dump are issued and PRINT then goes into normal end-of-job procedures. Inspection of the material furnished by PRINT usually allows the user to diagnose the cause and location of the malfunction. If an error is trapped by a WORDS program, a diagnostic message is issued and the series of events just noted take place by forcing a branch to the PRINT core-linkage routine.

### *Data Organization*

The method for organizing data records in WORDS is dictated by the analytic methods involved in the technique

which requires that the unit of information be the word itself. Thus, each word must constitute an independent and separable machine record.

The standard record format is a collection of eleven fixed length fields comprising thirty characters of information. Each of these eleven fields has a mnemonic associated with it and is referred to by that mnemonic. WORD, a field of fifteen characters, holds the actual english word comprising that record. SPKR, a one digit field, designates the speaker of that word. SPEECH, a one digit code for the part of speech, is inserted by the PARSE series. TIME, a one digit code, is not currently in use. INTV, a three digit number, allows designation of the interview in which the word was found and SEGM, a two digit number, allows referencing the particular observation within that interview. SEQ, a three digit number inserted by the SPLIT program during data input, locates the specific sequential position of the word within the interview/segment combination and SUBSEQ, a one digit number, allows insertion of up to nine words between any two originally input words. SEGTOT, a five digit number indicates total words emitted by the particular speaker in that segment. FREQ, a five digit field, is initially set to one by SPLIT and is then free for whatever use is required by the run. Finally, SPARE, a four digit field, is open to various internal uses by other WORDS programs.

All WORDS programs, other than the initial program SPLIT and the mathematical programs, are written to process records of this length and this format.

### *Program Organization*

WORDS programs belong to one of six functional types: systems control, sorting, editing, record keeping, printing and statistical. Each of these blocks will be described along with a list of contained programs.

#### *Systems Control*

These six programs all have the common function of maintaining data flow within the system and between the system and the resident monitor. The block subsumes: CHECK, COPY, FILER, MNTRA, MNTRB, and PRINT.

#### *Sorting*

All sort programs are really a third level monitor making use of the same basic applied program, IBM SM148.

Calling any sort program actually calls this monitor which brings in the main sort program segments as needed, makes modifications in the segments as required for the particular sort version called, maintains linkages between sort segments and retrieves necessary statistics before returning control to MNTRB. The block subsumes: OMITS, SORTS, and SUMMS.

#### *Editing and Format Manipulation*

The nine programs have either the function of changing the data in a file or changing the fields within records in that file or both. Such changes may be accomplished by changing fields, removing or replacing complete records. Programs included are: EDIT, FIXST, PARS1, PARS2, IDIOM, SPLIT, TEXT1, TEXT2, and STRIP.

#### *Record Keeping*

The design of WORDS makes it important to maintain a record of changes made to the data. A series of five programs, HSTR1-HSTR4 and HSTRY, process all editing changes before turning control to the actual editing programs. Since I/O scheduling for the series is complex with internally called sort modules, MNTRA accepts a call for HSTRY PKG upon receipt of which it manufactures the appropriate sequence.

#### *Printing*

The two modular printing programs, RERYT and PRISM, are used solely to produce output files on the resident monitor's print tape for later listing on an offline 1401.

#### *Statistical*

The eight programs in this functional block are designed to carry the reduced WORDS data through an intercorrelation matrix, listing of that matrix, factor analysis, varimax rotation and listing, and finally a scoring procedure utilizing the factor loadings from the varimax data. The programs are: LISTR, CORR1, DUBLR, DECOD, FACTR, VARMX, VDCOD, and SCORE. Since the first four of these programs are constant in any intercorrelational procedure, MNTRA accepts a call for COREL PKG to produce all required calls and I/O.

*CHECK.* This is the first program called in any WORDS run. Its purpose is to make an extensive series of validity checks on each of the control cards input to the run in order to catch any errors at the beginning rather than the later part of the run.

*CORRI.* The intercorrelation matrix program of the series. It will handle a matrix of up to 999 variables.

*DECOD, DUBLR.* DUBLR and DECOD function as paired programs which will almost invariably follow CORRI. DECOD is designed to produce an easily legible output from the CORRI matrix output by replacing the variable identification numbers by their English words and by allowing an ordering of obtained correlations or by screening them against a pre-set criterion level. DUBLR precedes DECOD and is used to expand the upper symmetric matrix produced by CORRI into a complete matrix (less the main diagonal).

*EDIT.* Used to make substantive changes in the data file. It is a very flexible program and allows, among other things, the change of any given word to another, the deletion of any specific word or of all occurrences of that word, the deletion of sets of interviews, speakers, segments, etc. The goal of EDIT is to reduce the total number of different words in the system; in that sense, it is the epitome of the entire system since all reduction changes idiosyncratic to the set of interviews under analysis are accomplished by EDIT.

*FACTR.* The factor-analytic program of the system. A principal-components algorithm is used to extract up to ninety-nine factors from an intercorrelation matrix of maximum order  $215 \times 215$ .

*FILER.* Allows the production of tape files, all on one tape, designed to serve as future input to the system. Allows the re-input of any of these files on future runs. The program has several safeguards in that files input to the system must be labelled, the input tape is automatically removed after use, the output tape is also removed.

*FIXST.* Designed for the merging and updating of the striplists used by the STRIP program.

*HSTR Series.* A series of five programs scheduled and called by use of the HSTRY PKG call. The programs HSTR-1, -2, -3, -4, and HSTRY maintain an accurate record of all data changes made via EDIT. The cumulated history is saved by the FILER program and allows a re-start procedure from an earlier point as well as a history of all changes made to the data.

*IDIOM.* IDIOM and its second part, PRIDE, will locate idiomatic usages which must be treated differently than regular words since the separate words within an idiom cannot be worked separately. Idiom-constructions, furnished on cards, are located and both listed and punched

in a format that is appropriate input to EDIT for any necessary changes.

*LISTR.* The physical format of records as they are kept by the WORDS system is not appropriate for that set of programs which are designed to mathematically analyze the data. It is the job of the LISTR program to re-format the data when it is finally ready for analysis.

*MNTRA, MNTRB.* These are, respectively, the administrative and executive monitors of the WORDS system. In brief, MNTRA will accept and process the run control cards which instruct the system as to what programs are needed, when, where, etc., and to construct from this list of cards a calling sequence for the program run configuration; it also makes extensive validity checks on the cards. When MNTRA has constructed this calling sequence it turns control to MNTRB which then takes over the actual calling of each program in the sequence and the task of maintaining adequate communications between programs.

*OMITS.* A sorting program with the facility for deletion of records which are equal to each other, according to furnished parameters upon which equality is to be assessed, leaving only the first of such records intact. Thus, were the user to want a list of every different word in the data, use of OMITS on the sorting parameter of the word itself will cause deletion of all words which are alike save for the first in the string; the remainder then becomes one record for every different word in the system.

*PARSI, PARS2.* These programs will insert a part of speech code into each record (word) in the data. PARS1 operates mainly on a dictionary lookup basis, although some logical manipulation is done, in order to make assignments where the grammar code is unambiguous. After resorting the output from PARS1, PARS2 takes over in order to assign the remaining codes according to fairly extensive analytic rules. The assignment of parts of speech is important in reduction since it allows reduction by rule (e.g., "delete all articles and conjunctions") via the EDIT program and in that it allows combinations of words with other words because the specific meaning of the word is defined with its part of speech, e.g., *like* = *enjoy* vs. *like* = *similar*.

*PRINT.* A dual purpose program. PRINT is called in order to terminate the WORDS system control over the computer before returning the machine to the resident monitor. Before doing so, however, PRINT will compile a record of statistics and messages produced during the actual computer run itself.

*PRISM.* This is the main printing program of WORDS. PRISM is designed to allow the printing of those files which have been selected for output in order to provide

either a record of certain of the results of that run or to provide an indication of information for planning future runs.

*RERYT*. Where PRISM is primarily designed to produce lists of any data file contents, RERYT is specifically intended to produce a printed copy of the interviews under analysis in a format similar to that of the original typescript. Thus, each speaker is separated from every other, periods are restored to the end of sentences, spacing separates segments, etc. RERYT is used when it is desirable for the user to inspect the "state" of the data after various transformations have been made. After an extensive EDIT run, for example, it is useful to be able to read the interviews in their present form to determine "clinically" just how much meaning is still being retained in the data.

*SCORE*. Using the correlation matrix input data prepared by LISTR this program accepts a deck of cards — one for each word in the matrix — punched with the varimax loadings for each word on each of the factors on which it loads highly. By using the frequencies of occurrence and the loadings as multipliers, it computes a factor score for each factor in each observation. It then produces a printing file with the factor score means and standard deviations across the entire data set and then lists, for each observation, the raw score and standard score of each of the factors.

*SORTS*. A sorting program which, unlike SUMMS or OMITS, makes no physical change to the total file in terms of deletion or summarization. SORTS will simply re-order the records in the file into whatever order is specified by the sorting parameters; other than re-ordering, no changes are made.

*SPLIT*. This program serves as the entry point for raw data into the system. When an interview is originally punched, as many words are placed on each card as is feasible. These cards are then put on tape either via an offline 1401 operation or online by COPY. SPLIT takes this card image tape as input and produces a separate WORDS record for every word on every card of the input data. SPLIT will also insert within each record all necessary data for determining the origin of the word, i.e., segment number, interview number, speaker, etc. In addition, it will also assign a sequencing number to each word as a function of its position in the segment in which it was found. Use of these identifying origin data allows SORTS to restore the data to its original order no matter how it has been re-ordered by any other program.

*STRIP*. Like EDIT, STRIP is designed to make substantive changes to the interviews under analysis. Where EDIT makes such changes on the basis of the specific interviews being handled, STRIP is designed to be applied to every set of interviews that comes along. STRIP may be considered as a de-inflection program whose task it is to place the words in the data into their root form. STRIP, unlike EDIT, has no flexibility in terms of options. It can only replace a given word with another as this replacement is specified by a deck of pre-punched cards. No deletions or other types of changes than replacement, are permitted.

*SUMMS*. Like OMITS, SUMMS deletes records equal to each other so that only one record of each type (cf. OMITS) remains. Unlike OMITS, however, SUMMS first adds the frequency data of each record deleted to the frequency data of the first record in the string. Thus applying both OMITS and SUMMS to the same data with the same sorting parameters would yield the same set of records on output but SUMMS would have accumulated within each record the summed frequency of all the records deleted.

*TEXT1, TEXT2*. A pair of programs designed to provide a KWIC type of listing with the programs producing a record for each input word which reports the two words preceding and following that specific word and reports also the interview, segment, and sequence numbers as well as the part of speech for the key-word. Unlike IDIOM which searches for specific idiomatic constructions and then reports both the idiom and the sentence which it contains, the TEXT-programs are designed more as a dictionary producer which allow the user a "random-access" approach so that *any word* can be located in context at any time.

*VARMX*. The varimax rotation program to be applied to output from FACTR. The program will rotate any set of up to thirty-three factors from the FACTR output tape to a criterion of simple structure. The set of factors to be rotated are selected from the input set by control card punching.

*VDCOD*. This program serves VARMX as does the DECOD program CORRI. It allows production of an easily legible listing from VARMX with each variable number being replaced by the corresponding English word under analysis, with all factor-loadings ordered by absolute value and with a listing of communalities and variance proportions attributable to each factor.

## APPENDIX 2

## Factors Extracted from Wizard of Oz Using Frequency Selection Choices

Factor 1		Factor 2		Factor 3		Factor 4		Factor 5	
<i>ax</i>	91	<i>Oz (lady)</i>	94	<i>farmer</i>	86	<i>Munchkins</i>	82	<i>Uncle Henry</i>	92
<i>oil</i>	91	<i>Oz (head)</i>	90	<i>brick</i>	84	<i>Witch</i>	82	<i>house</i>	90
<i>Tin Woodman</i>	86	<i>kill</i>	85	<i>Scarecrow</i>	82	<i>East</i>	80	<i>Aunt Em</i>	89
<i>tin</i>	86	<i>send</i>	83	<i>road</i>	66	<i>woman</i>	75	<i>bed</i>	74
<i>leg</i>	85	<i>lovely</i>	78	<i>stuff</i>	66	<i>old</i>	73	<i>small</i>	74
<i>right</i>	79	<i>eye</i>	67	<i>feel</i>	62	<i>little</i>	71	<i>sun</i>	72
<i>arm</i>	77	<i>do</i>	63	<i>eat</i>	60	<i>wear</i>	61	<i>door</i>	71
<i>body</i>	72	<i>help</i>	62	<i>yellow</i>	59	<i>live</i>	51	<i>laugh</i>	67
<i>soon</i>	70	<i>answer</i>	58	<i>walk</i>	58	<i>set</i>	51	<i>look</i>	52
<i>girl</i>	65	<i>throne room</i>	53	<i>hurt</i>	57	<i>face</i>	49	<i>reach</i>	51
<i>head</i>	60	<i>surprise</i>	52	<i>straw</i>	56	<i>Dorothy</i>	48	<i>Toto</i>	41
<i>once</i>	55	<i>will</i>	45	<i>place</i>	52	<i>people</i>	48	<i>run</i>	41
<i>work</i>	54	<i>no</i>	41	<i>few</i>	50	<i>Good Witch</i>	44	<i>middle</i>	39
<i>tinsmith</i>	50	<i>tell</i>	39	<i>no</i>	50	<i>can</i>	-41	<i>sit</i>	38
<i>can</i>	42	<i>Oz</i>	38	<i>brain</i>	49	<i>grow</i>	39	<i>back</i>	-33
<i>grow</i>	37	<i>West</i>	38	<i>other</i>	47	<i>land</i>	39	<i>first</i>	33
<i>far</i>	-36	<i>many</i>	36	<i>man</i>	39	<i>Cowardly Lion</i>	-38	<i>land</i>	33
<i>help</i>	33	<i>straw</i>	34	<i>Toto</i>	38	<i>dress</i>	37	<i>grass</i>	32
<i>old</i>	33	<i>die</i>	31	<i>crow</i>	37	<i>country</i>	35	<i>eye</i>	31
<i>come</i>	32	<i>great</i>	31	<i>do</i>	37	<i>great</i>	-34	<i>hand</i>	31
<i>one</i>	32	<i>grow</i>	31	<i>leave</i>	36	<i>house</i>	32	<i>one</i>	31
<i>put</i>	32	<i>return</i>	31	<i>mind</i>	35	<i>look</i>	32	<i>ask</i>	-30
<i>return</i>	31			<i>right</i>	33	<i>the group</i>	-31	<i>fall</i>	30
<i>face</i>	-30			<i>know</i>	32			<i>hard</i>	30
<i>look</i>	-30			<i>number</i>	-32				
				<i>soon</i>	32				
				<i>keep</i>	-30				
				<i>pole</i>	30				
% Variance		4.80		5.20		4.70		4.80	
Factor 6		Factor 7		Factor 8		Factor 9		Factor 10	
<i>wolf</i>	90	<i>coward</i>	88	<i>Gaylette</i>	93	<i>Winkies</i>	67	<i>mouse</i>	91
<i>lie</i>	78	<i>near</i>	75	<i>Quelala</i>	84	<i>tinsmith</i>	65	<i>Queen Mouse</i>	90
<i>crow (noun)</i>	77	<i>Cowardly Lion</i>	67	<i>time</i>	82	<i>set</i>	54	<i>safe</i>	75
<i>die</i>	77	<i>heart</i>	62	<i>Winged Monkeys</i>	81	<i>careful</i>	52	<i>turn</i>	60
<i>lay</i>	65	<i>know</i>	54	<i>Golden Cap</i>	74	<i>night</i>	51	<i>all</i>	52
<i>number</i>	65	<i>Toto</i>	41	<i>call</i>	59	<i>day</i>	50	<i>grass</i>	51
<i>come</i>	51	<i>stuff</i>	41	<i>fly</i>	58	<i>tear</i>	50	<i>try</i>	50
<i>tear</i>	51	<i>try</i>	39	<i>next</i>	51	<i>pretty</i>	48	<i>field</i>	49
<i>Wicked Witch</i>	48	<i>great</i>	37	<i>wish</i>	42	<i>leave</i>	-43	<i>run</i>	44
<i>fly</i>	48	<i>fast</i>	36	<i>the group</i>	40	<i>like</i>	-43	<i>hurt</i>	41
<i>Winkies</i>	46	<i>big</i>	35	<i>field</i>	40	<i>forest</i>	-42	<i>fast</i>	39
<i>straw</i>	43	<i>no</i>	33	<i>lose</i>	37	<i>Good Witch</i>	-40	<i>come</i>	38
<i>Golden Cap</i>	42	<i>return</i>	32	<i>together</i>	-35	<i>basket</i>	39	<i>near</i>	37
<i>one</i>	42	<i>run</i>	32	<i>glad</i>	33	<i>rule</i>	-39	<i>work</i>	36
<i>foot</i>	39	<i>beast</i>	31	<i>land</i>	32	<i>work</i>	38	<i>bring</i>	35
<i>stand</i>	36	<i>tell</i>	-31	<i>good</i>	31	<i>keep</i>	37	<i>open</i>	33
<i>next</i>	32	<i>tin</i>	31	<i>sure</i>	31	<i>last</i>	37	<i>far</i>	32
<i>ask</i>	-30	<i>reply</i>	30	<i>once</i>	30	<i>mind</i>	37	<i>live</i>	31
<i>time</i>	30					<i>yellow</i>	36	<i>speak</i>	31
						<i>few</i>	35	<i>big</i>	30
						<i>start</i>	35	<i>yellow</i>	30
						<i>bring</i>	34		
						<i>hand</i>	34		
						<i>stand</i>	-33		
						<i>friend</i>	32		
						<i>live</i>	-32		
						<i>beast</i>	-31		
						<i>ask</i>	-30		
						<i>lay</i>	30		
% Variance		3.50		4.20		3.70		3.70	

Factors Extracted from Wizard of Oz Using Frequency Selection Choices

Factor 11		Factor 12		Factor 13		Factor 14		Factor 15	
<i>spectacles</i>	87	<i>flower</i>	90	<i>terrible</i>	63	<i>Silver Shoes</i>	83	<i>room</i>	86
<i>Guardian of the Gates</i>	86	<i>stork</i>	88	<i>man</i>	58	<i>end</i>	76	<i>green</i>	85
<i>want</i>	75	<i>sleep</i>	82	<i>home</i>	-56	<i>water</i>	70	<i>soldier</i>	70
<i>Emerald City</i>	70	<i>bright</i>	56	<i>promise</i>	54	<i>Wicked Witch</i>	65	<i>dress</i>	67
<i>bright</i>	56	<i>fast</i>	56	<i>please</i>	53	<i>foot</i>	59	<i>morning</i>	65
<i>long (adj.)</i>	-45	<i>last</i>	54	<i>think</i>	53	<i>power</i>	57	<i>wait</i>	65
<i>glad</i>	41	<i>carry</i>	52	<i>stand</i>	52	<i>use</i>	51	<i>girl</i>	59
<i>first</i>	40	<i>find</i>	51	<i>voice</i>	51	<i>take</i>	46	<i>see</i>	51
<i>wish</i>	40	<i>like</i>	50	<i>little</i>	41	<i>Dorothy</i>	44	<i>bed</i>	50
<i>sun</i>	39	<i>fall</i>	45	<i>Oz</i>	40	<i>hard</i>	41	<i>throne room</i>	48
<i>surprise</i>	39	<i>let</i>	44	<i>speak</i>	37	<i>begin</i>	37	<i>night</i>	48
<i>man</i>	34	<i>hand</i>	42	<i>one</i>	36	<i>bring</i>	-35	<i>door</i>	47
<i>speak</i>	34	<i>wait</i>	40	<i>back</i>	-35	<i>beauty</i>	-30	<i>pretty</i>	43
<i>eat</i>	33	<i>must</i>	79	<i>rule</i>	-33	<i>open</i>	-30	<i>pass</i>	39
<i>open</i>	32	<i>take</i>	-38	<i>forest</i>	-32			<i>big</i>	38
<i>may</i>	31	<i>river</i>	35	<i>must</i>	32			<i>silk</i>	38
<i>night</i>	31	<i>field</i>	34	<i>Kansas</i>	-31			<i>get</i>	-37
<i>beast</i>	-30	<i>few</i>	31	<i>beauty</i>	-31			<i>course</i>	36
<i>give</i>	30	<i>lovely</i>	31	<i>first</i>	-31			<i>middle</i>	36
<i>great</i>	-30			<i>friend</i>	31			<i>one</i>	36
<i>like</i>	30			<i>wait</i>	30			<i>can</i>	-34
								<i>speak</i>	34
								<i>begin</i>	-33
								<i>cry</i>	30
								<i>eye</i>	30
								<i>many</i>	30
								<i>wear</i>	30
<b>% Variance</b>	<b>3.70</b>		<b>4.10</b>		<b>3.10</b>		<b>3.30</b>		<b>4.40</b>
Factor 16		Factor 17		Factor 18		Factor 19		Factor 20	
<i>balloon</i>	84	<i>pretty</i>	55	<i>pole</i>	76	<i>tree</i>	82	<i>courage</i>	74
<i>air</i>	71	<i>Hammer Heads</i>	54	<i>river</i>	75	<i>side</i>	63	<i>real</i>	64
<i>silk</i>	71	<i>kind</i>	53	<i>middle</i>	60	<i>branch</i>	57	<i>brain</i>	61
<i>make</i>	64	<i>dress</i>	43	<i>let</i>	48	<i>seem</i>	53	<i>very</i>	61
<i>basket</i>	63	<i>reach</i>	41	<i>get</i>	45	<i>other</i>	49	<i>many</i>	57
<i>lose</i>	55	<i>rest</i>	41	<i>water</i>	44	<i>walk</i>	45	<i>use</i>	53
<i>go</i>	52	<i>will</i>	41	<i>land</i>	40	<i>journey</i>	44	<i>sure</i>	50
<i>day</i>	45	<i>indeed</i>	-40	<i>rest</i>	40	<i>first</i>	43	<i>give</i>	49
<i>get</i>	45	<i>back</i>	39	<i>fast</i>	39	<i>real</i>	-43	<i>reply</i>	49
<i>now</i>	45	<i>grow</i>	-39	<i>West</i>	34	<i>the group</i>	42	<i>day</i>	47
<i>should</i>	45	<i>thank</i>	39	<i>leave</i>	34	<i>must</i>	42	<i>find</i>	47
<i>together</i>	45	<i>take</i>	37	<i>animal</i>	-31	<i>long (adj.)</i>	41	<i>heart</i>	46
<i>tear</i>	38	<i>field</i>	36	<i>great</i>	-31	<i>thank</i>	-41	<i>think</i>	45
<i>people</i>	37	<i>head</i>	36	<i>begin</i>	30	<i>turn</i>	41	<i>fear</i>	44
<i>find</i>	34	<i>country</i>	35	<i>may</i>	30	<i>forest</i>	39	<i>people</i>	44
<i>will</i>	34	<i>friend</i>	35			<i>Kansas</i>	-38	<i>Oz</i>	40
<i>Kansas</i>	33	<i>pass</i>	34			<i>kind</i>	-38	<i>put</i>	37
<i>thank</i>	32	<i>course</i>	-33			<i>next</i>	37	<i>big</i>	-36
<i>Oz</i>	31	<i>sit</i>	32			<i>tell</i>	-37	<i>good</i>	34
<i>last</i>	31	<i>voice</i>	31			<i>rest</i>	36	<i>face</i>	33
		<i>must</i>	30			<i>bring</i>	-35	<i>live</i>	32
		<i>open</i>	-30			<i>course</i>	-35	<i>may</i>	30
						<i>answer</i>	-34	<i>morning</i>	30
						<i>surprise</i>	34		
						<i>Good Witch</i>	-31		
						<i>beast</i>	31		
						<i>look</i>	30		
<b>% Variance</b>	<b>3.70</b>		<b>2.90</b>		<b>2.90</b>		<b>4.10</b>		<b>3.80</b>