

AD659057

This document has been approved  
for public release and sale; its  
distribution is unlimited.

DDC  
OCT 5 1967  
REGISTERED  
B

Approved by the  
CLEARINGHOUSE  
for the release of information  
regarding the activities of the  
Department of Defense

ANNUAL REPORT.  
AUTOMATIC INDEXING  
AND ABSTRACTING

M-21-67-1

March 1967

Annual Progress Report  
Office of Naval Research  
Contract Nonr 4440(00)

Reproduction in whole or in part is permitted for  
any purpose of the United States Government.

Electronic Sciences Laboratory  
Lockheed Palo Alto Research Laboratory  
LOCKHEED MISSILES & SPACE COMPANY  
A Group Division of Lockheed Aircraft Corporation  
Palo Alto, California

**PRÉCIS**  
**RESEARCH PROGRESS REPORT**

**Title:** "Annual Report: Automatic Indexing and Abstracting," Annual Progress Report, Office of Naval Research, Contract Nonr 4440(00).

**Background:** This investigation is concerned with the development of automatic indexing, abstracting, and extracting systems. Basic investigations in English morphology, phonetics, and syntax are pursued as necessary means to this end.

**Condensed Report Contents:** The third annual report on automatic indexing and extracting summarizes progress in three areas of investigation.

- (1) Application of English word morphology to automatic indexing and extracting
- (2) Automatic indexing and information retrieval by the thesaurus method
- (3) Studies in phonetic English

The first two areas were supported by the Office of Naval Research; the studies in phonetic English were supported in part by the ONR, and in part by grants from the Lockheed Independent Research fund.

In the first section the theory and design of the "sentence dictionary" experiment in automatic extraction is outlined. Some of the computer programs needed for this experiment have been completed; documentation for these programs is included.

The second section is a discussion of an unconventional information retrieval system which may be used to retrieve English or Russian technical literature. Specifically, there is a discussion of the query language and of the data base structure. The design of this system is complete and programming has begun.

The third section comprises a final report in an exhaustive study of the relation between the orthographic and phonetic forms of English monosyllables. This paper reports on the methods and accuracy of the algorithm to map the orthographic forms of English into the phonetic forms given by five different phonetic authorities, or to map the phonetic form given by one authority into that given by another.

**For Further Information:** The complete report is available in the major Navy technical libraries and can be obtained from the Defense Documentation Center. A few copies are available for distribution by the author.

## FOREWORD

This report marks the completion of the third year in which the Office of Naval Research has contributed support to the research in the Information Sciences at the Lockheed Palo Alto Research Laboratories of the Lockheed Missiles & Space Company. The first two sections report on work supported solely by the ONR, except for computer support contributed by LMSC; the third section reports on work supported largely by the independent research program of LMSC, with some support in the terminal phases from the ONR funds.

It is convenient to consider the work reported here as dealing with significant data bases. During the first year of the program, a major part of the effort went into establishment of a word-data base. The English Word Speculum, which has been distributed to ONR program participants, illustrates the nature of this data base. In the second yearly report, and in the present report, examples are given of exploration and application of the word data base and of a phrase data base to problems in linguistics and information analysis.

The first section of this report documents how the word data base has been utilized in the development of a computer system which creates a sentence data base which will be used for further investigation of the indexing and abstracting process. This sentence data base, scheduled for completion this year, will consist of sentences selected from chapters of various books provided with an index. It will carry structural information in the form of part-of-speech strings for each sentence, and an indication of the "significance" of each sentence as measured by the index.

The second section describes how an English/Russian phrase data base can be used in the development of a technique for obtaining English indexes from untranslated

Russian text. This section also shows how such index terms can be used to form a data base to support a retrieval system in which syntactic relationships can be utilized in the user's input descriptors.

The third section documents how a word data base giving both orthographic and phonetic forms of words has been used to develop an algorithm for converting from graphic to phonetic forms.

The group at Lockheed takes this opportunity to express its thanks for the continued support and encouragement given by the members of the Information Sciences Branch of the Office of Naval Research.

## CONTENTS

Section		Page
	FOREWORD	iii
	ILLUSTRATIONS	vi
1	AN EXPERIMENT IN THE USE OF SYNTACTIC INFORMATION IN AUTOMATIC EXTRACTING	1-1
	1.1 Introduction	1-1
	1.2 Outline of Experiment	1-3
	1.3 Experiment Documentation	1-6
2	AUTOMATIC RETRIEVAL OF ENGLISH AND RUSSIAN TECHNICAL TEXT	2-1
	2.1 Introduction	2-1
	2.2 The Query Language	2-3
	2.3 Automatic Indexing	2-9
	2.4 File Structure	2-14
	2.5 IBM 360 Configuration for Retrieval System	2-16
	2.6 Data Base	2-23
	2.7 References	2-24
3	COMPUTABLE RELATIONS BETWEEN THE ORTHOGRAPHIC AND PHONETIC FORMS OF ENGLISH MONOSYLLABLES	3-1
	3.1 Introduction	3-1
	3.2 Mechanized Converter of Orthographic to Phonetic Forms of Words	3-4
	3.3 Comparison of the Phonetic Representations of Elementary Words in Five Dictionaries	3-14
	3.4 Selection of Exception Words	3-34
	3.5 Extension of the Program To Include Monosyllables	3-36
	3.6 Summary and Conclusions	3-41
	3.7 References	3-44
Appendix		
A	RUSSIAN TEXT	A-1
B	ABSTRACTS USED AS DATA BASE	B-1

## ILLUSTRATIONS

Figure		Page
1-1	DICT Program Flow Diagram	1-9
1-2	DICTF Program Flow Diagram	1-12
1-3	RSEINTR Program Flow Diagram	1-15
1-4	Part-of-Speech Program Flow Diagram	1-19
1-5	SENDICT Program Flow Diagram	1-23
1-6	SENINT Program Flow Diagram	1-25
2-1	Simple Index	2-11
2-2	Complex Index	2-12
2-3	File Arrangement	2-17
2-4	Field and Record Codes	2-19
2-5	Basic Data File Structure for Random Access Files	2-20
2-6	Sample Data File for Item <u>energy</u>	2-21
2-7	File Characteristics	2-22
3-1	Functions Performed for Purposes of Pulse-Input, Voice-Output Communication With Computers	3-5
3-2	Outline of Algorithm for Extension to Polysyllabic Words	3-7
3-3	Outline of Algorithm for Obtaining Phonetic Form From the Orthographic Representation of Elementary Words	3-9
3-4	Operation of the Program for the Word SPICE	3-10
3-5	Operation of the Program for the Word GRUDGE	3-12
3-6	Comparison of Phonetic Data by Using the Phonetics in the Shorter Oxford Dictionary as the Reference Set	3-15
3-7	Data Organized According to the Marker-Orthographic Vowel-Marker Set and According to Phonetic Representations in Specified Dictionaries (I)	3-17
3-8	Data Organized According to the Marker-Orthographic Vowel-Marker Set and According to Phonetic Representations in Specified Dictionaries (II)	3-18

<b>Figure</b>		<b>Page</b>
3-9	Statistics on Disagreement Among Dictionaries for Phonetics of Elementary Words (I)	3-21
3-10	Statistics on Disagreement Among Dictionaries for Phonetics of Elementary Words (II)	3-22
3-11	Data Organized According to Terminal Rhyme Words According to Jones (I)	3-23
3-12	Data Organized According to Terminal Rhyme Words According to Jones (II)	3-23
3-13	Alphabetically Arranged List of Homonyms	3-26
3-14	Number of Words Forming Homonyms in Each of the Five Dictionaries	3-27
3-15	Graphic Representation of the Number of Homonym Sets Among the Elementary Words in Five Dictionaries	3-28
3-16	Graphic Representation of the Number of Homonym Sets Among the Double Standard Elementary Words in Five Dictionaries	3-30
3-17	Graphic Representation of the Number of Homonym Sets Among the Algorithmic and Elementary Words in Five Dictionaries	3-33
3-18	Modification of the Program To Include the English Monosyllables	3-35
3-19	Statistics on Disagreement Among Dictionaries for Phonetics of Additional Monosyllables (I)	3-39
3-20	Statistics on Disagreement Among Dictionaries for Phonetics of Additional Monosyllables (II)	3-39
3-21	Statistics on Polysyllabic Pronunciation of Words Ending in ES According to Five Dictionaries	3-40
3-22	Statistics on Polysyllabic Pronunciation of Words Ending in ED According to Five Dictionaries	3-40
3-23	Graphic Representation of the Number of Homonym Sets Among the English Monosyllables in Five Dictionaries	3-42



## Section 1

# AN EXPERIMENT IN THE USE OF SYNTACTIC INFORMATION IN AUTOMATIC EXTRACTING\*

### 1.1 INTRODUCTION

There is often a relationship between the orthographic form of a word and its syntactic function or its meaning. By analogy, it is possible that there is a relationship between the syntactic structure of a sentence and the significance of the sentence. It is important to establish whether or not such a relationship does exist, and if so, to discover how it can most easily be exploited.

Now that an algorithm exists for assigning parts of speech to words, it is practicable to design an experiment for correlating automatically the degree to which this structure-significance relationship exists. Part of speech can be regarded as definitive of structure; a human-compiled index can be regarded as indicative of significance. Given text for which human-compiled indexes exist, it should be possible to identify those sentences which contain the index reference phrases and those which do not. Let us designate as "indexible" and "more significant" those sentences in a text which contain a phrase listed in the index, and as "nonindexible" those not containing a phrase listed in the index. By assigning a part-of-speech category to each word of each sentence, sentences can be syntactically characterized or typed by the part-of-speech strings of their words. Then a simple sorting and comparing experiment will determine if the part-of-speech strings can be used in differentiating between "indexible" and "nonindexible" sentences.

In carrying out this experiment, a sentence file will be produced and then ordered according to (1) the part-of-speech strings of the sentences and (2) the categorization of sentences into indexible and nonindexible types. This ordered file can then be examined to see if (1) there is a manageable number of different syntactic sentence

---

\*Work on this task was conducted by L. L. Earl.

types, i. e., different part-of-speech strings; and (2) if there is a partition between sentence types categorized as indexible and those categorized as nonindexible.

If there is a manageable number of sentence types and a complete partition, by type, of indexible and nonindexible sentences, we will have achieved a sentence file which is also a sentence dictionary of indexible sentence types. This means that if the chosen corpus was sufficiently representative, a type of syntactic extracting will have been achieved, since any text can be extracted by comparing its sentence types with those in the sentence dictionary of indexible sentence types. If, on the other hand, there is clearly an unmanageable number of sentence types or clearly no relationship between sentence types and index phrases, then we will have proved that extracting by syntax alone is not practical. The actual result is likely to lie somewhere in between and is likely to suggest avenues of research by which the partition between indexible and nonindexible sentences could be improved. Some possible avenues are as follows:

- Use a system of syntactic analysis to find equivalence classes among sentence types such that the number of sentence types will be reduced without destroying the partition between indexible and nonindexible sentences.
- Find a finer (as opposed to a grosser) classification of part of speech than that used to produce more distinct sentence types with a better chance of partition between indexible and nonindexible sentences.
- Find a way in which other criteria (such as frequency counts) can be combined or superimposed on syntactical criteria to achieve a better partition between indexible and nonindexible sentences.

Once a satisfactory partition between indexible and nonindexible sentences has been achieved, the next research target will be the identification of the actual index phrase or phrases within an indexible sentence. If this identification can be achieved, not only automatic extracting but also automatic indexing will be possible. Exploration of this possibility should wait until completion of the extracting experiment.

## 1.2 OUTLINE OF EXPERIMENT

The necessary tasks in the extracting experiment include preparation of text, computer programming and checkout of text-handling routines, processing of text, and analysis of results. Without regard for how each step will be carried out, we can outline the steps of the experiment as follows.

First, a corpus of indexed text must be chosen which has sufficient variety and volume to be a statistically meaningful sample. Then each sentence must be categorized as indexible or nonindexible. The text must then be keypunched and read into the computer with this categorization. The part of speech of each word in the sentence can then be determined automatically, and a record can be produced in a storage device, giving the sentence itself, its length, its position in the input text, its index categorization, and its part-of-speech string. When a record has been formed for each sentence, we will have a sentence file whose records can be ordered by part-of-speech string and by categorization. This indexed file must then be examined and analyzed as described in the introduction to ascertain the degree to which syntax can be used in extracting and indexing text. Progress in these tasks can now be outlined in more detail.

The computer programs for the IBM 360/30 necessary for completion of the experiment can be summarized as follows:

- (1) A routine to assign parts of speech to words, by dictionary look-up for special-function words, by graphemic analysis for the bulk of the language
- (2) A routine to form the part-of-speech dictionary on disk, of all words for which the part of speech cannot be determined algorithmically
- (3) A routine to read text from cards, isolating and organizing individual sentences for processing purposes
- (4) A routine to create a sentence file on disk, with each record containing a sentence, the part-of-speech codes for the words in the sentence, and identification, such as the page and sentence number and indexible or nonindexible categorization

- (5) A routine to sort a sentence file according to the part-of-speech codes, so that sentences of like structure can be compared
- (6) A routine to print a sentence file, or a portion thereof

Of the computer programs listed, Numbers 1 through 4 have been programmed and checked out, and documentation is included in subsection 1.3.1. Program Number 5 is being worked out. A simple form of Number 6 has been programmed, checked out, and documented, but a more sophisticated program with more options will probably become necessary.

The dictionary necessary for the part-of-speech routine has been prepared and exists on cards and on tape. The random access dictionary on disk which is used by the part-of-speech routine is prepared each time, just before its use, from the tape version (approximately a 4-minute task). The compilation of the words for the dictionary is described in Section I.4 of the March 1966 annual report. There are at present 916 entries in the dictionary.

The selection and preparation of the input text is just as large a task as the compilation of the part-of-speech dictionary and programming of the text-handling routines. The index phrases which refer to the chosen text must be identified and the indexable sentences marked. The text must then be keypunched in a format which can be read by the computer program. (The keypunching rules which have been formulated are given in subsection 1.3.2.) It has been decided that 10,000 sentences, or about 500 pages of average text, should constitute a large enough sample for the indexing experiment. A lesser sample will be used initially to see if it appears that a relationship between syntax and meaning does indeed exist and if it is worthwhile pursuing the exact nature and limitations of the relationship.

Nine texts have been chosen so far for the experiment, and the indexable sentences in each have been marked. Key punching has been completed for four of the texts. Each text comprises one chapter from one of the following books.

- (1) A. Doak Barnett, Communist China & Asia, Challenge to American Policy, New York, published for the Council on Foreign Relations, Harper, 1960.
- (2) Ntel W. Chamberlain, The Firm: Micro-Economic Planning and Action, New York, McGraw-Hill, 1962.
- (3) Charles Coulston Gillispie, The Edge of Objectivity: an Essay in the History of Scientific Ideas, Princeton, N.J., Princeton University Press, 1960.
- (4) Walter James Greenleaf, Occupations; A Basic Course for Counselors, Washington, D. C., Government Printing Office, 1954.
- (5) John G. Gurley and Edward S. Shaw, Money In a Theory of Finance, Washington, Brookings Institution, 1960.
- (6) Robert L. Heilbroner, The Making of Economic Society, Englewood Cliffs, N. J., Prentice-Hall, 1962.
- (7) Mark S. Massel, Competition and Monopoly, Legal and Economic Issues, Washington, Brookings Institution, 1962.
- (8) Robert M. Palter, Whitehead's Philosophy of Science, Chicago, University of Chicago Press, 1960.
- (9) S.M. Siegel, The Plant Cell Wall - A Topical Study of Architecture, Dynamics, Comparative Chemistry, and Technology in a Biological System (International Series of Monographs of Pure and Applied Biology, Plant Physiology Division, Vol. 2), N. Y., Pergamon, 1962.

All phrases in the index which referred to the chosen chapter were extracted with the page reference. These phrases were then identified on the referenced page, and any sentence containing such a phrase was marked "indexible." Identification of the phrases was done in whatever way seemed most reasonable; however, many problems came up in such identification. Some index entries in Communist China and Asia can be used to illustrate the difficulties. One entry was "Foreign policy, Chinese Communist, pp. 77-83." These pages contained specific examples of Chinese Communist foreign policy as well as general statements about it. In sentences containing both the words "Chinese Communist" (or "Chinese") and also the words "foreign policy" (or "policy"), both phrases were marked and the sentence was called indexible. However, sentences were also called indexible when the words "foreign policy" stood for "Chinese Communist foreign policy." (For example, "National security is another basic national interest which underlies Peking's foreign policy.") Another entry was "Ideology, Chinese Communist." In this

case any sentence containing either "Chinese" or "Communist" and a form of "Ideology" (i.e., ideological) was called indexible. For the entry "Ideology of Mao," "Maoism," "thought of Mao-Tse-tung" and "Maoist ideas" were all accepted as index phrases.

Thus an attempt was made to identify as indexible all sentences whose thought was concerned with an item in the index. This does, unfortunately but necessarily, involve a subjective judgement. Often, as in the case of the entry "Ideology of Mao," the exact entry does not occur in the text at all. Even when it does, it does not seem reasonable to ignore other sentences containing concepts which the human indexer regarded as significant, albeit the form of phraseology differs slightly from that in the index. However, should the correlation between syntax and indexibility prove to be poor, one possibility will be to eliminate the "indexible" categorization from all sentences in which there is not an exact match between phrase and index entry.

### 1.3 EXPERIMENT DOCUMENTATION

In this section the details of the text-handling routines are documented. Rules for keypunching the text are given, consistent with the reading routine developed. For each checked-out computer program there is a general description, input and output format description, listing of the program's control cards or calling sequence for running under the BOS monitor, and a logic flow diagram (Figs. 1-1 through 1-6).

#### 1.3.1 Program Documentation

##### ● DICT PROGRAM

Description: DICT is a routine which reads dictionary entries of a prescribed format from the card reader and stores them in a condensed, coded form on magnetic tape. Logical IOCS is used for both reading and writing. Both input and output are double buffered. The output tape can be sorted using the standard SORT-MERGE program, and this tape then is the input to DICTF, which forms a dictionary on disk, for random access by keys. Control cards for both DICT and for the SORT-MERGE program are given below.

**Input and Output Format:** Each dictionary entry is punched on a separate card. The word occupies columns 1 through 12; it is unlikely that function words will be longer but if they are they will be truncated. The part-of-speech codes are given by two-character alphabetic codes which start in column 30 and are separated by blanks. For example, the word "after" which is a preposition, conjunction, and adverb would have a PR for preposition in columns 30 and 31, a CJ for conjunction in columns 33 and 34, and an AV for adverb in columns 36 and 37. On tape, each word occupies an 80-byte record of which the first 15 bytes are used, 12 bytes for the EBCDIC codes for the word itself, and 3 bytes for a binary representation of all the part-of-speech possibilities of the word. (Eighty-byte records were used because redundancies were encountered using 15-byte records.) Each binary bit of the 3-byte part-of-speech code represents a different part-of-speech possibility. A 1 bit in a given position indicates the presence of that particular part-of-speech possibility. The two-character codes for part of speech are given below, with the corresponding bit position in the 3-byte binary code.

Part of Speech	Alphabetic Code	Binary Bit Position	
		From Left	From Right
Ends in S	HS	24	1
Noun-Adjective	NA	23	2
Adjective	AJ	22	3
Verb	VB	21	4
Adverb	AV	20	5
Preposition	PR	19	6
Conjunction	CJ	18	7
Pronoun	PN	17	8
Interjection	IJ	16	9
Past Verb	PV	15	10
Accusative	AC	14	11
Present Participle	NG	13	12
Past Participle	PP	12	13
Negative	NV	11	14
Auxiliary	AX	10	15
Future	FT	9	16
Reflexive	RF	8	17
Noun Plural or Collective	NP	7	18
Article	AR	6	19

Thus the significant part of the tape record for the word "after" will be as follows:

bytes	1 2 3 4 5 6 7 8 9 10 11 12	13 14 15
	AFTER	00 00 70
	└──────────────────┘	└──────────┘
	EBCDIC Codes	Hexadecimal as shown

Input cards are read and condensed into the 15-byte format and written on tape unit 180, unblocked, until the IOCS end-of-file card (with slash and asterisk) is encountered, at which time an EOF is written on tape. Any cards with illegal format or codes are not written on tape; the card image is printed on the printer instead.

In forming a sentence dictionary, the first step is normally the running of DICT, followed by a DSORT run to order the entries. The resultant tape is held and can be quickly converted to a random access file (by key) on disk by the DICTF program.

Control Cards:

```

0001 // JOB DICT
0002 // DATE 66133 5/13/66
0003 // ASSGN SYS002,X'180',T1,X'90'
0004 // EXEC LOADER,R
0005 /*

0001 // JOB DSORT
0002 // DATE 66034 2/3/66
0003 // ASSGN SYS000,X'180',T1,X'90'
0004 // VOL SYS001,SORTW
0005 // DLAB 'BK BDS SORT WORK AREA' 1094737' C
0006 // XTENT 1,000,000,0199009,'094737',SYS001
0007 // ASSGN SYS002,X'180',T1,X'90'
0008 // EXEC
0009 // OPTION MESSAGES=SYSLST,PRINT,STORAGE=16000,LABEL=(U,U)
0010 // INPFI INPUT=T,VOLUME=1,BLKSIZE=(80,X)
0011 // OUTFI OUTPUT=T,BLKSIZE=80
0012 // RECORD TYPE=F,LENGTH=(80)
0013 // SORT FIELDS=(1,12,4),FORMAT=CH,SIZE=1000
0014 // END
0015 // END

```





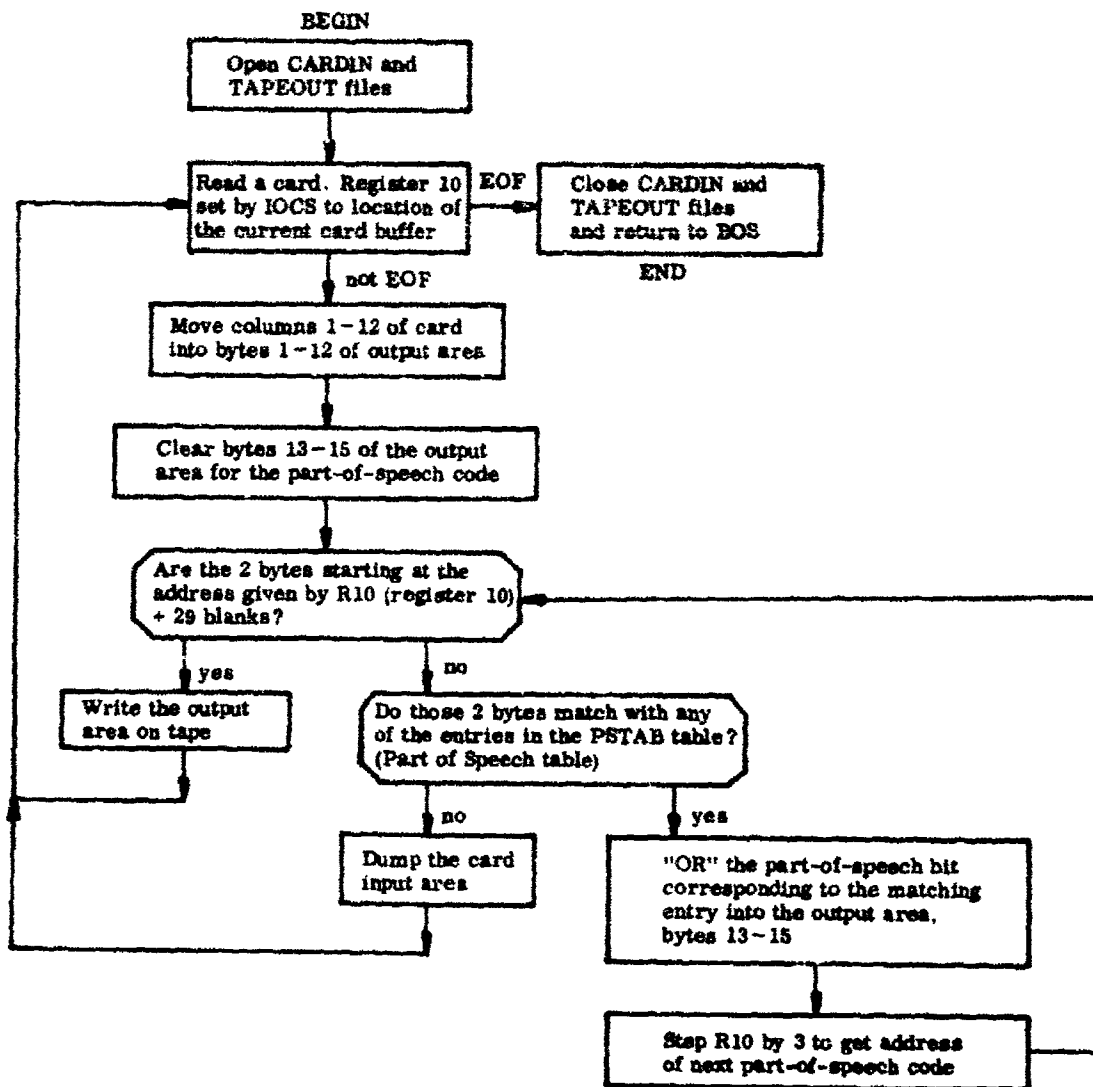


Fig. 1-1 DICT Program Flow Diagram

## • DICTF PROGRAM

Description: DICTF takes the sorted output of DICT from tape 180 and forms a random access file on disk, using logical IOCS. Both input and output are double buffered. Eighty-byte records are input from the tape; the first 15 bytes (the only significant bytes) are used unchanged to form 15-byte disk records. The first 12 bytes (the EBCDIC word) become the key by which the record can be accessed and the last 3 bytes (the part-of-speech codes) constitute the argument supplied when the record is accessed.

Error Messages: A dump is supplied if the disk file is successfully built and closed. If errors are encountered, the 15-byte error records are dumped, headed by a code to indicate the nature of the error. These codes are as follows:

INDEXFULL	cylinder index area is full
DISKERR	uncorrectable disk error
DISKFUL	prime data area on cylinder is full
DUPREC	duplicate record in input data
ORDER	sequence error in input data
LENGTH	length error in input data
TAPEERR	error in reading input tape

The 15-byte error dump is followed by a full dump and control is returned to BOS.

Input and Output Format: The input format is given under the DICT documentation, where it is the output format.

Output records are unblocked. They consist of a key area of 15 bytes, the EBCDIC coded word, and a data area of 3 bytes, comprising the part-of-speech code as described in the DICT documentation. Records are stored and indexed by logical IOCS for random access, in the disk areas specified by the XTENT cards. (See control cards below.) A track index is formed in the area given by the XTENT 4 card;

no master index is used. The XTENT 1 card gives the primary data area and the XTENT 2 card gives the overflow data area.

Control Cards: Note that for BOS version 6.2, the name in the DLAB card must be changed with each run, due to a bug in logical IOCS.

---

```
0001 // JOB DICTF
0002 // DATE 66206 7/25/66 MONDAY
0003 // VOL SY5001,SDICT
0004 // DLAB ' SDICT11 1094737' C
0005 // 0001,66100,66100,'000000000000'
0006 // XTENT 4,001,0101000,0103009,'094737',SYS001
0007 // XTENT 1,002,0104000,0115009,'094737',SYS001
0008 // XTENT 2,003,0116000,0125009,'094737',SYS001
0009 // EXEC LOADER,R
0010 /*
```

---

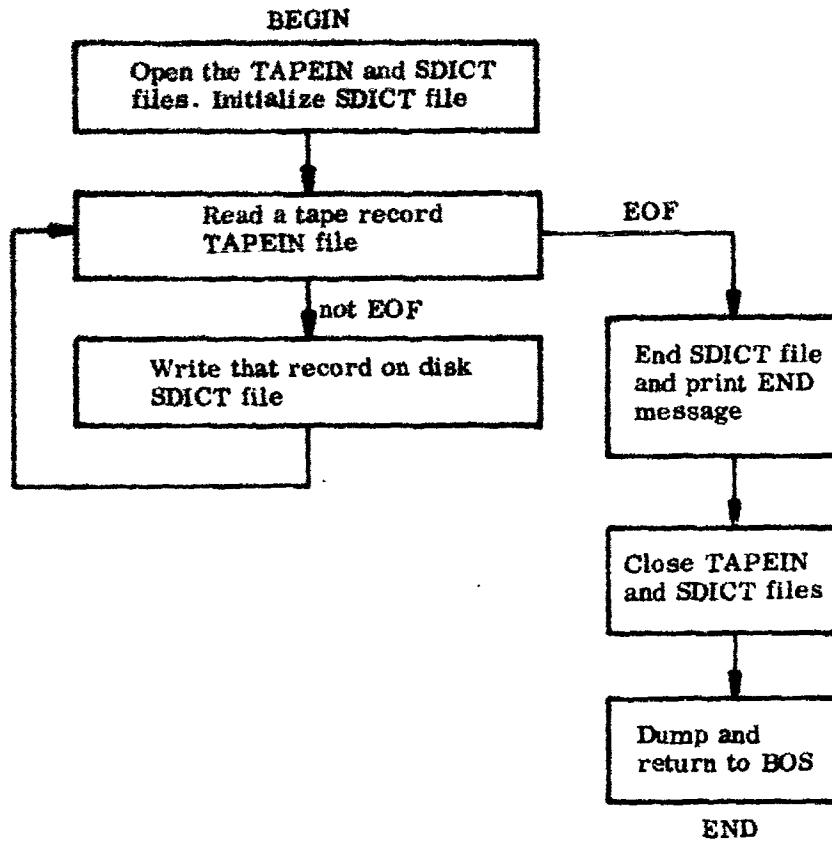


Fig. 1-2 DICTF Program Flow Diagram

● RSENTR (Read Sentence Routine) PROGRAM

Program Entry: RSENT

Description: RSENTR is a subroutine which reads text from cards, providing one sentence to the using routine each time it is called. Double buffering is used in reading text under the logical IOCS system; both buffers are set aside within the RSENTR program and are of no concern to the using routine. The sentence provided to the user is stored as designated in the calling sequence, and a "map" is provided with the location and length of each word. RSENTR also provides the using routine with the page and sentence number of each current sentence, the number of words in the sentence, and an indicator which shows whether or not the sentence was marked indexible.

Calling Sequence:

L	15 = V(RSENT)
BALR	14, 15
DC	Y(PAGE)
DC	Y(SENTBUFF)
DC	Y(SENTMAP)
DC	Y(INDIC)
DC	Y(ERROR)
DC	Y(EOF)

Normal return

**PAGE** name of the location where page number, sentence number, and number of words will be stored, 2 bytes for the page number followed by 1 byte each for sentence number and number of words.

**SENTBUFF** name of the location where the card image (EBCDIC) will be stored. A 1000-byte buffer is recommended.

- SENTMAP** location of a map giving the location and length of each word in the sentence. There is one full computer word in the map for each English word or mark of punctuation in the sentence; the first half-word contains the location of the English word or mark of punctuation and the second half-word contains the number of bytes occupied by the word. The first full word corresponds to the first English word, the second to the second English word, etc.
- INDIC** location of a byte in storage which will be set to I if the sentence is indexible, or to N if it is not.
- ERROR** an error return, which is not currently being used. The computer will stop or give a dump on error.
- EOF** the EOF return; all sentences have already been processed.

Input Format: If it is desired to set the initial page number at a value other than zero, the first card should have a \$ in column 1 and the page number in columns 2 through 6, all of which must be punched. The page number can be made to identify the text by adding a large constant, e. g., if text 1 starts at page 256 of book A and text 2 starts at page 1001 of book B, the initial card for text 1 can read \$10256 and for text 2, \$21001. The page number will be stepped by 1 and the sentence count reset to zero every time a card with a \$ in column 1 (and the rest blank) is encountered. Text is punched free form in columns 1 through 72, except that each sentence must be followed by either two blanks, or one blank plus one dollar sign. If a dollar sign is present, the sentence will be identified as indexible. (See Calling Sequence.) Detailed rules for keypunching are given in subsection 1.3.2. The EOF card has a slash in column 1 and an asterisk in column 2.

Output Format: See Calling Sequence.

Control Cards: RSENTR is a subroutine stored in the relocatable library. An INCLUDE RSENTR must follow the PHASE card of the using program.

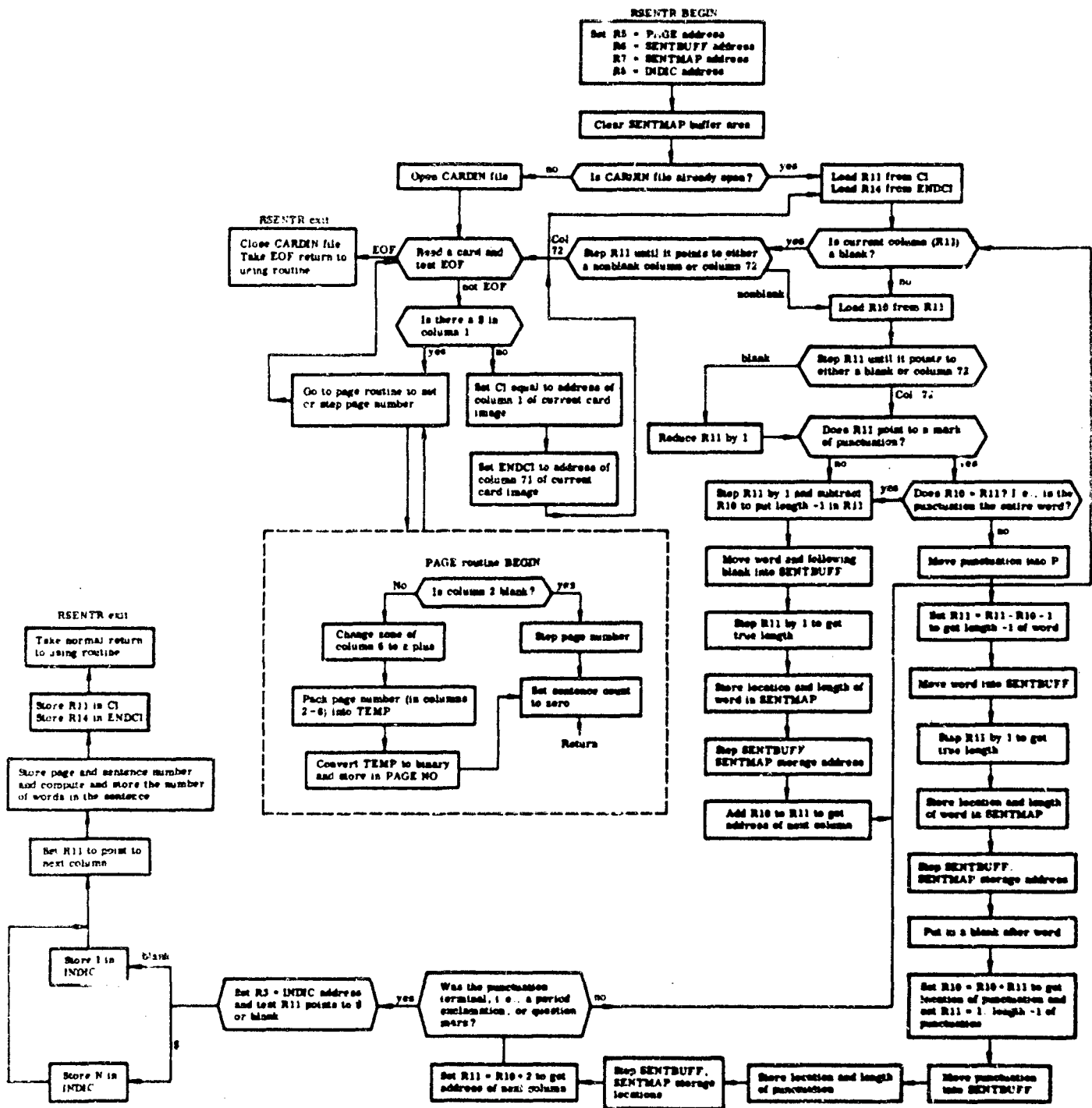


Fig. 1-3 RSENTR Program Flow Diagram

● POS (Part of Speech) PROGRAM

Description: POS is a subroutine which determines the parts of speech of the English word whose location and length are given in register 10, storing the part-of-speech code in the location specified by register 11. POS uses logical IOCS to look up each word in the indexed sequential dictionary file (on disk) produced by DICTF. For all words not in the dictionary, the parts of speech are determined by analyzing the graphemic structure of the word as described in Section I.4 of Part I of the 1966 Annual Report.

Input and Output Format: POS assumes that the address and length of the word whose parts of speech are to be determined are in a location specified by R10 (register 10). The first 2 bytes of the full word specified by R10 give the address at which the English word is stored and the second 2 bytes give the number of letters in the word.

The part of speech codes assigned by POS are the same as those described in the DICT documentation. This 3-byte code is stored into the location specified by R11.

Calling Sequence: POS is entered at POSBEGIN with a calling sequence such as:

```
L          15, = V (POSBEGIN)
BALR      14, 15
RETURN
```

Any program using POS must have the VOL, DLAB, and XTENT cards which define the dictionary file. (See DICTF and SENDICT documentation.)

Programming Notes: POS makes use of several specially defined macros, as described below.

1.  $\mathcal{L}$ LOC  $\mathcal{S}$ AFIX  $\mathcal{L}$ AFFIX,  $\mathcal{L}$ LENGTH,  $\mathcal{L}$ POS



This macro puts the letter string  $\mathcal{E}AFFIX$ , which has a length  $\mathcal{E}LENGTH$  at the location  $\mathcal{E}LOC$ . Immediately after the string it stores the 3-byte code for the parts of speech indicated by  $\mathcal{E}POS$ .  $\mathcal{E}POS$  may have one of the following values:

A	adjective
N	noun
V	verb
NV	noun and verb
PV	past verb
NG	present participle
PP	past participle
NP	noun plural
MNV	noun and verb, becoming noun only for words of 4 or more syllables
AJV	adjective and adverb

If  $\mathcal{E}POS$  is blank, a zero is stored instead of a part-of-speech code.

## 2. $\mathcal{E}LOC$ \$ENDT

This macro stores a hexadecimal FF at  $\mathcal{E}LOC$ . This is the code POS uses for ending a table.

## 3. $\mathcal{E}LOC$ \$LETR $\mathcal{E}R$ , $\mathcal{E}YES$ , $\mathcal{E}NO$ , $\mathcal{E}LEN_1$ , $\mathcal{E}STRING_1$ , $\mathcal{E}LEN_2$ , $\mathcal{E}STRING_2$ ...

This macro generates a code for checking the nth letter in the KERNEL buffer (where n is in register  $\mathcal{E}R$ ) against the letter string  $\mathcal{E}STRING$  of length  $\mathcal{E}LEN$ , with a transfer to  $\mathcal{E}YES$  if a match occurs or to  $\mathcal{E}NO$  if it does not.  $\mathcal{E}LEN$  and  $\mathcal{E}STRING$  can be repeated in pairs to check for more than one string (the OR condition). Thus if R8 = 3, \$LETR 8, YES, NO, 1, A, 2, ES, 1, B will result in a transfer to YES if the third byte in the KERNEL buffer starts the sequence A, or ES, or B. Otherwise a transfer to NO will result. If either  $\mathcal{E}YES$  or  $\mathcal{E}NO$  is left out the transfer will be to the next instruction in sequence.

4.  $\mathcal{E}LOC$   $\$LETN$   $\mathcal{E}N$ ,  $\mathcal{E}YES$ ,  $\mathcal{E}NO$ ,  $\mathcal{E}LEN_1$ ,  $\mathcal{E}STRING_1$ ,  $\mathcal{E}LEN_2$ ,  $\mathcal{E}STRING_2 \dots$

This macro is identical to  $\$LETR$  except in the designation of the letter to be checked. In  $\$LETN$ ,  $n$  is given by  $\mathcal{E}N$ ; in  $\$LETR$   $n$  is given in register  $\mathcal{E}R$ .



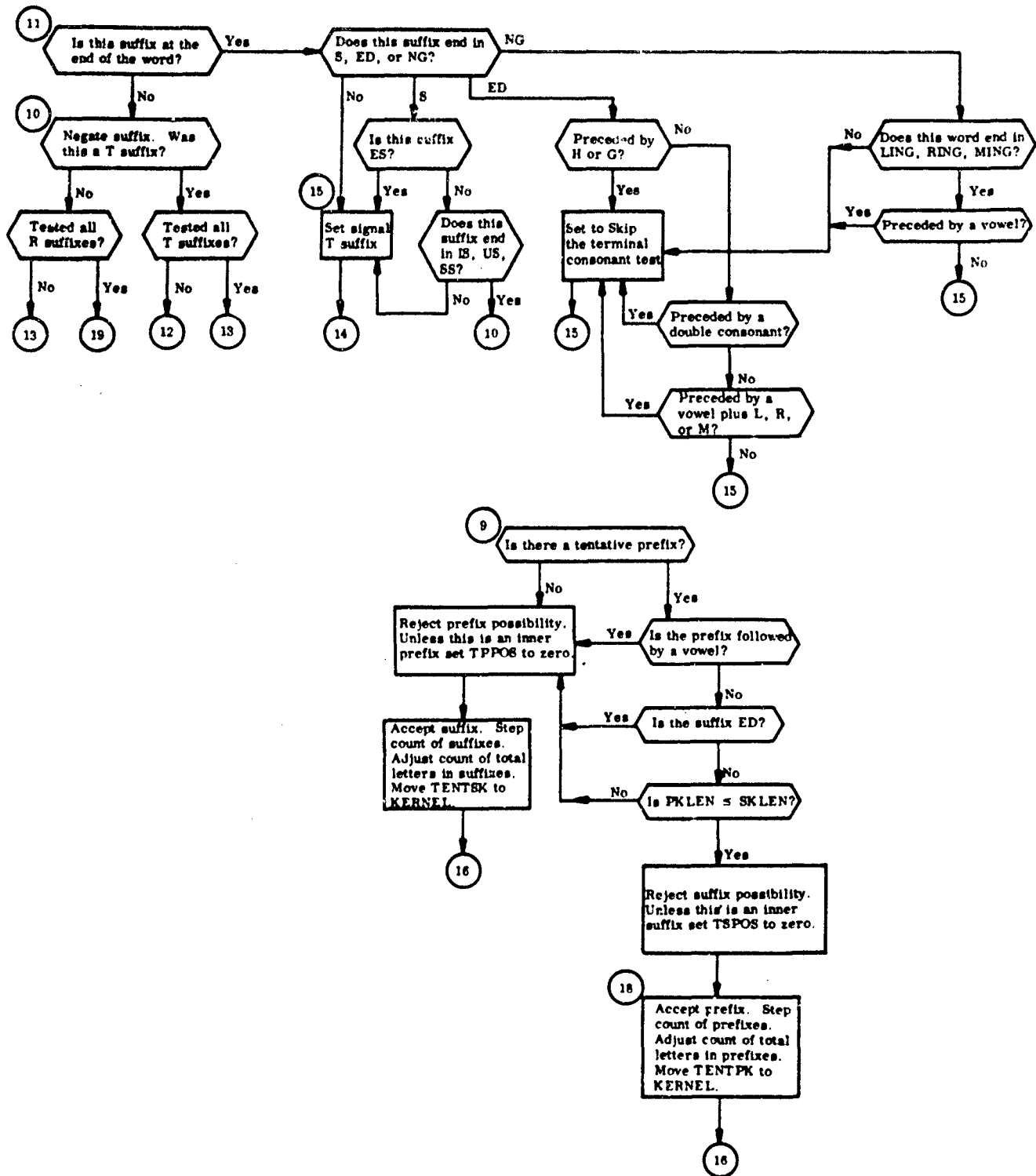


Fig. 1-4 Part-of-Speech Program Flow Diagram (Cont.)

● SENDICT (Sentence Dictionary) PROGRAM

Description: SENDIC creates a sentence file on disk, with each record of the file representing one sentence. SENDIC assembles each sentence, its part-of-speech strings, and its identification into a variable length record, and uses logical IOCS to form a blocked sequential file of these records. RSENTR is used to read the sentences from cards and POS is used to assign part-of-speech codes to the words.

Input Format: Data input is as for RSENTR.

Output Format: Each variable-length record corresponds to one sentence. The blocksize is currently set at 2728 bytes, set by consideration of the requirements for using the IBM disk-sort routine for SENSORT (not yet completed). The size of each record is calculated and stored by SENDIC; logical IOCS accumulates as many records as possible into a block. The record format is as follows:

Field Number	Number of bytes	Description
0	4	First 16 bits give the total record length in bytes, right justified; logical IOCS uses the second 16 bits
1	1	The character I or N for indexible or not indexible
2	2	The page number in binary, right justified
3	1	The sentence number in binary, right justified
4	1	The number of English words in the sentence, in binary, right justified
5	3M where M = number of English words	The part of speech codes for the words, 3 bytes each
6	Variable	The sentence in EBCDIC

Control Cards: Cards 19 through 22 are representative data cards; see the data input section of the RSENTR documentation.

```

0001 // JOB SENDIC
0002 // DATE 66172 6/21/66
0003 // VOL SYS001,SDICT
0004 // DLAB ' SDICT11                                1094737',      C
0005 //                                0001,66100,66100,'0000000000000000'
0006 // XTENT 4,001,0101000,0103009,'094737',SYS001
0007 // XTENT 1,002,0104000,0115009,'094737',SYS001
0008 // XTENT 2,003,0116000,0123009,'094737',SYS001
0009 // VOL SYS001,SENTFL
0010 // DLAB ' SENTFILE1                                1094737',      C
0011 //                                0001,66100,66100,'0000000000000000'
0012 // XTENT 1,000,0001000,0100009,'094737',SYS001
0013 // EXEC LOADER
0014 // PHASE SENDIC,S
0015 // INCLUDE RSENTR
0016 // INCLUDE SENDIC
0017 // INCLUDE POS
0018 // ENTRY SOSTART
0019 $00200
0020 CHAPTER II. ART. LIFE, AND EXPERIMENT.
0021 PHYSICS HAS BEEN THE CUTTING EDGE OF SCIENCE SINCE GALILEO, AND ITS
0022 /*

```

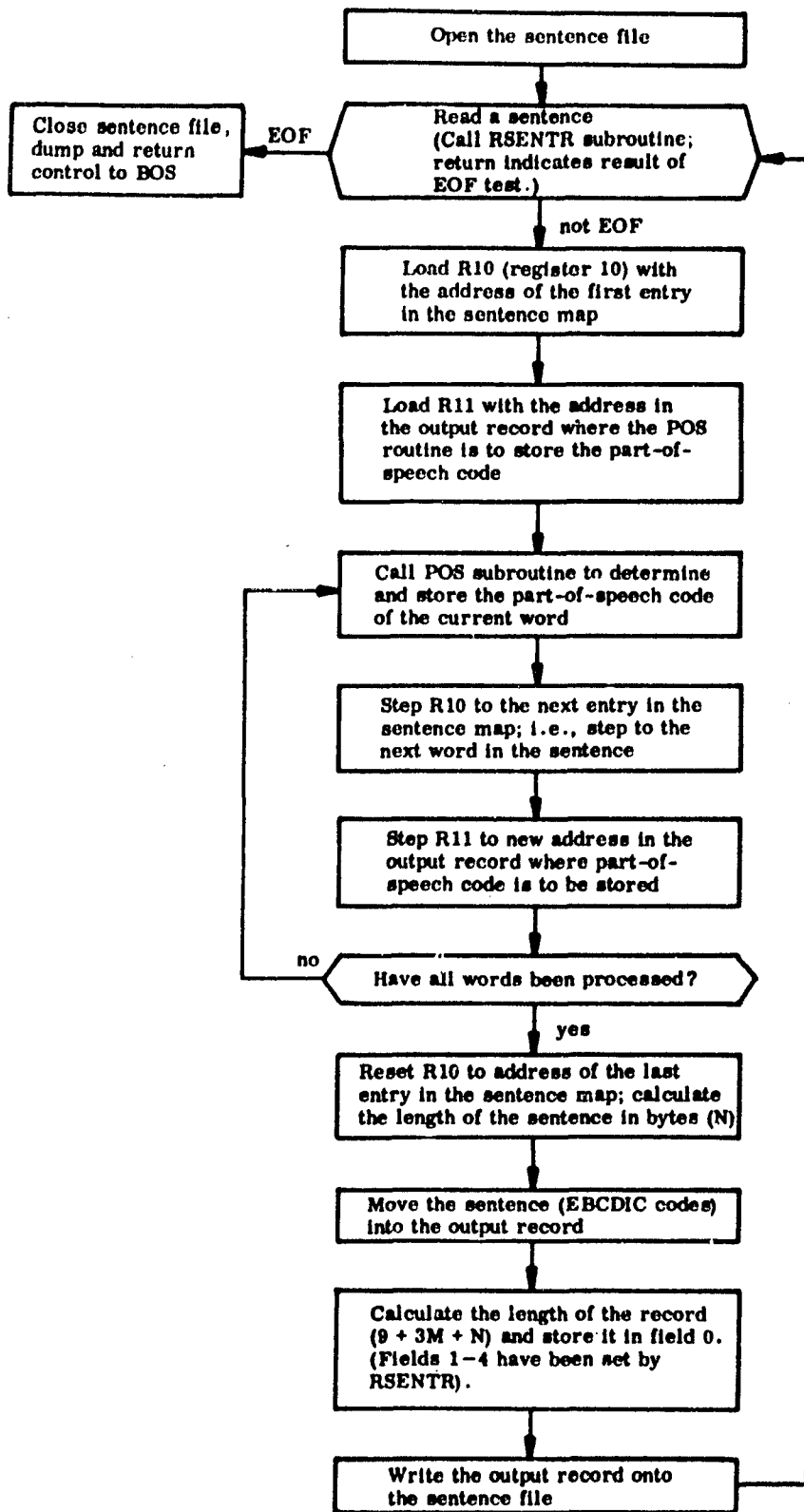


Fig. 1-5 SENDICT Program Flow Diagram

● SENINT (Sentence Interpretation) PROGRAM

Description: SENINT is a program which reads sentence records from the disk file produced by SENDIC or SENSORT, interprets them, and prints the sentences and sentence information on-line. SENDIC (see write-up) is the program which produces the original sentence file. SENSORT (not yet completed) processes the original sentence file to produce an ordered sentence file. SENINT is currently being revised to take its input from tape as well as from disk.

Input and Output Format: The input format, i.e., the format of the sentence file on disk, is identical to the output format of SENDIC. The XTENT card of the SENINT job must be the same as the XTENT card used when the file to be interpreted was produced. Normally, this XTENT card will be that used for the output file of the SENDIC or SENSORT run.

The output is on-line, on the printer. Two sentences are output on each page, one at the top and one at midpage. The first line gives the sentence category in column 2 (I for indexible, N for nonindexible), the page number in columns 9 through 16, and the sentence number in columns 25 through 32. The second line is blank. Starting with line three, the part-of-speech codes for the words appear; a slash separates the parts of speech of one word with those of the next. Following the parts of speech there are three blank lines. The next and as many following lines as necessary contain the sentence.

Control Cards:

```
0001 // JOB SENINT
0002 // DATE 66284 10/11/66
0003 // VOL SYS001,SENTFI
0004 // DLAB ' SENTFILE1 1094737', C
0005 // XTENT 1,000,0001000,0100009,'094737',SYS001
0006 // EXEC LOADER
0007 PHASE SENINT.S
0008 INCLUDE SENINT
0009 ENTRY SISTART
0010
```



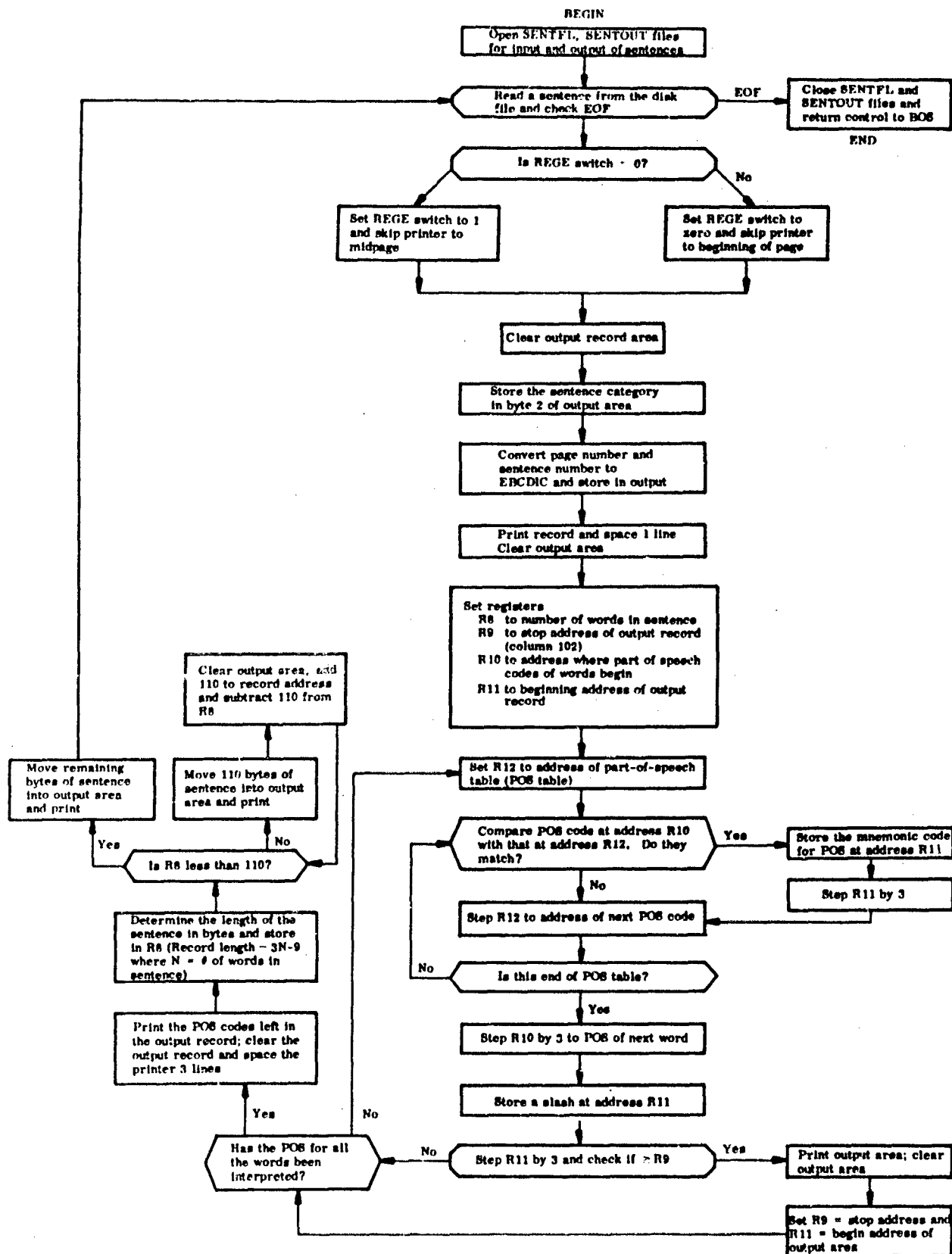


Fig. 1-6 SENINT Program Flow Diagram

### 1.3.2 Key punching Rules

1. At the end of each page, following the last sentence which begins on that page, punch a card with a \$ in column 1 and the rest of the card blank.
2. The text should be punched more or less naturally in columns 1 through 72, as in typing, except that the following rules must be observed.
  - (a) Each sentence must end with a period, exclamation, or question mark followed by either two blanks, or if there is a red \$ following the sentence, by one blank and then one \$. There may be any number of blanks before the beginning of the next sentence.
  - (b) Do not split a word between cards. Make sure that the last word terminates at or before column 72 and leave the rest of the card blank. Words must be separated by one blank. More blanks are O.K.
  - (c) Punctuation will usually be handled as in typing, punched immediately following the word preceding, but inserting a blank in front of the punctuation will be O.K.
  - (d) Any mark of punctuation other than a hyphen should be followed by a blank. Thus should there be contiguous marks of punctuation, separate them by blanks.
  - (e) Do not punch elipsis marks or use periods or question marks followed by a blank in any other capacity than to terminate a sentence.
  - (f) Quotation marks follow the rules given in (c) and (d), except read Rule 4 below.

For example:

col	1	2	3	4	5	6	7	8	9	10	11	12					18	19		
	t	h	e		w	o	r	d		"			b	i	o	l	o	g	y	"

or

col	1	2	3	4	5	6	7	8	9	10	11	12					18	19	20	
	t	h	e		w	o	r	d		"			b	i	o	l	o	g	y	"

If a comma had been included

col 1 2 3 4 5 6 7 8 9 10 11 12 13 19 20  
t h e w o r d , " b i o l o g y "

3. The codes used for special characters are as follows:

<u>Name</u>	<u>Symbol</u>	<u>Punch</u>
period	.	12-8-3
less than	<	12-8-4
left paren	(	12-8-5
plus	+	12-8-6
ampersand	&	12
exclamation	!	11-8-2
dollar sign	\$	11-8-3
asterisk	*	11-8-4
right paren	)	11-8-5
semicolon	;	11-8-6
hyphen	-	11
comma	,	0-8-3
percent	%	0-8-4
dash	-	0-8-5
greater than	>	0-8-6
question mark	?	0-8-7
colon	:	8-2
apostrophe	'	8-5
equals	=	8-6
quotation mark	"	8-7

These are standard codes for the IBM 360 computer.

4. Unfortunately, the punctuation rules (c) and (d) do not apply if a quotation mark ends the sentence. The question mark and the terminal punctuation must be punched adjacent, with the quotation first. For example: end." must be punched

col 65 66 67 68 69 70 71 72

E N D " .

5. Punch titles like sentences, supplying a period.

Section 2  
AUTOMATIC RETRIEVAL OF ENGLISH AND  
RUSSIAN TECHNICAL TEXT\*

2.1 INTRODUCTION

The following pages describe the design of a real-time information retrieval system for the IBM 360/30. The system will be capable of retrieving Russian or English technical literature in response to English queries. The system will run under the LACONIQ time-sharing monitor (Ref. 2-1). The system will utilize an IBM 2260 cathode ray tube and typewriter console for input/output. The system will be checked out with a small data base representing first 50, then 100, physics articles.

Indexing and information retrieval are intimately linked and it will be shown how the structure of the data files follows directly in conception from the Russian to English Indexing System (Ref. 2-2) described in the annual report for last year.

The Russian-to-English Indexer provides for complex cross indexing. This deep level indexing allows the establishment of an extensive network of syntactic relationships between many of the index items. This syntactic information is retained in the structure of the retrieval system's data file. This, in turn, allows retrieval of a document according to syntactic relatedness of terms within it, as opposed to mere existence of the terms within the document.

Although the retrieval system has been designed to operate upon output of the Russian-to-English Indexer, it will be seen that since this output is in fact an English index, the retrieval system could be used to retrieve English material as well as Russian material. The query language and the file structure will be discussed separately.

---

\*Work on this task was conducted by H. R. Robison.

### 2.1.1 Search Strategies

The search strategy generally used in information retrieval systems is that of coordinate or inverted file searching. In such a search, the user normally formulates his request as a Boolean string of descriptors where the descriptors are words or phrases used to index the document collection. Retrieval then consists of finding those documents, or names and sources of those documents, which satisfy the conditions of the Boolean expression.

The drawback to a Boolean search where word/phrase descriptors are used to formulate the Boolean expression is well known. If a user is interested in retrieving articles having to do, say, with the direction of proton polarization he uses as his descriptors

direction  
proton  
polarization

If the Boolean search informs him that there are, say, ten such documents, all that the user knows is that ten documents contain the words direction, proton, and polarization. Whether the words are related to one another and how remains unknown.

Attempts to solve the problem of relatedness have led to the use of various formulas which attempt to quantify the relatedness of words, generally by some type of word-frequency count of index items in the document collection.

### 2.1.2 Syntactic Relatedness

The approach used here is to attempt to establish relatedness of index items by syntactic means. Thus the phrase direction of proton polarization, if it occurs in text, is regarded as a single entity because of its syntactic structure. Proton is an adjective modifying polarization. The noun phrase proton polarization is the object

of the preposition of and the prepositional phrase of proton polarization modifies the noun direction. Text words need not be contiguous to be syntactically related for here the phenomenon of syntactic government (Refs. 2-3, 2-4) exists to distribute syntactic relationships throughout an entire sentence. For example, the government pattern for comparison is

comparison of S/ with S

which may be used to link widely separated elements

Even more convincing is the comparison of the observed number of particles which have passed through both plates, and which have retained after the second plate a momentum greater than 19.3 MeV with the expected numbers of particles of various masses.

The retrieval system under discussion has the unique feature of allowing the user to syntactically combine his input descriptors, to insert prepositions, and to search the file with this combination. Thus the input descriptor could be

direction of proton polarization

rather than

direction, proton, polarization

If the system responds that ten documents contain the descriptor in question, the user knows that they contain the descriptor as he formulated it - direction of proton polarization. This is made possible because the entire file structure is cross indexed according to the syntactic relationships of the index items contained in it.

This file structure permits a direct search for documents and also allows the user to browse through it syntactically.

## 2.2 THE QUERY LANGUAGE

The user is provided with several commands which he may use to compile his own descriptor list, then search or browse through the file.

ENTER D1, D2, ... (Descriptor 1, 2, etc.)  
LIST  
ADD D1, D2, ...  
DELETE D1, D2, ...  
CALCULATE D1, D2, ...  
DISPLAY X (X is any number of titles containing  
the descriptor or any number of index  
items containing the descriptors)  
MORE  
WRITE

The descriptor, D, may be used as the actual descriptor word or it may be used as the number that represents the word in the user's descriptor list. The commands will be explained and illustrated in the display frames that follow. The particular sequence used in the example was arbitrarily selected for illustrative purposes.

ENTER D1, D2, ...: The user must select descriptors to describe his search topic and type these words on the keyboard. This command causes the descriptors D1, D2, ... to be entered in the user's Descriptor List.

LIST: This command causes a visual display of the user's Descriptor List to be displayed on the scope of the 2260.

For each descriptor entered, the system will compute the number of documents containing the descriptor and the number of index items containing the descriptor. For example, proton might occur in 25 documents of the collection and in 40 index items in the collection. Display 1 shows user input and system response to the ENTER, LIST sequence.

It may now occur to the user that it would be worthwhile to see a display of proton spectrum rather than proton and spectrum. This is accomplished by an ADD command (Display 2). This display illustrates how the user can join his descriptors in natural syntactic units.

ENTER PROTON, SPECTRUM, DEUTERON BOMBARDMENT,  
HYDROGEN, HELIUM

LIST

<u>DESCRIPTOR</u>	<u>NUMBER OF DOCUMENTS</u>	<u>NUMBER OF ITEMS CONTAINING D</u>
1. PROTON	25	40
2. SPECTRUM	75	100
3. DEUTERON BOMBARDMENT	12	20
4. HYDROGEN	80	100
5. HELIUM	60	125

Display 1

ADD PROTON SPECTRUM

<u>DESCRIPTOR</u>	<u>NUMBER OF DOCUMENTS</u>	<u>NUMBER OF ITEMS CONTAINING D</u>
1. PROTON	25	40
2. SPECTRUM	75	100
3. DEUTERON BOMBARDMENT	12	20
4. HYDROGEN	80	100
5. HELIUM	60	125
6. PROTON SPECTRUM	10	15

Display 2

The user, if satisfied with this result, could delete the proton and spectrum entries (Display 3).



DELETE PROTON, SPECTRUM or  
DELETE D1, D2

<u>DESCRIPTOR</u>	<u>NUMBER OF DOCUMENTS</u>	<u>NUMBER OF ITEMS CONTAINING D</u>
1. DEUTERON BOMBARDMENT	12	20
2. HYDROGEN	80	100
3. HELIUM	60	125
4. PROTON SPECTRUM	10	15

Display 3

Let us assume, however, that the user is satisfied with his initial Descriptor List (Display 1). He decides, perhaps, to browse through the 40 index items which contain proton. This can be accomplished with a DISPLAY command.

DISPLAY 40 ITEMS PROTON  
or  
DISPLAY 40 ITEMS D1

INDEX ITEMS FOR PROTON

1. PROTON DECAY
2. PROTON SIMULATION
- 
- 
- 
- 
- 
- 
- 
- 
- 
12. PROTON IN ROTATING MAGNETIC FIELD

Display 4

Any list can be continued on the scope by using the MORE command.

MORE

INDEX ITEMS FOR PROTON

- 13. PROTON SPECTRUM FROM DEUTERON BOMBARDMENT OF HYDROGEN AND HELIUM
- 14. DIRECTION OF PROTON POLARIZATION

:  
:  
:

Display 5

While perusing the list of index items (Displays 4 and 5), the user may wish to see a document and index item count of an item or a portion of an item. This can be achieved by adding to the Descriptor List. Thus, the ADD command is used and the number 13 and the words proton polarization typed out.

ADD (13, PROTON POLARIZATION)

<u>DESCRIPTOR</u>	<u>NUMBER OF DOCUMENTS</u>	<u>NUMBER OF ITEMS CONTAINING D</u>
1. PROTON	25	40
2. SPECTRUM	75	100
3. DEUTERON BOMBARDMENT	12	20
4. HYDROGEN	80	100
5. HELIUM	60	125
6. PROTON SPECTRUM	10	15
7. PROTON SPECTRUM FROM DEUTERON BOMBARDMENT OF HYDROGEN AND HELIUM	4	6
8. PROTON POLARIZATION	2	4

Display 6

Now the user decides to examine the titles of the articles containing item 8, proton polarization, which was just added to the Descriptor List.

DISPLAY TITLES (D8)

LARGE-ANGLE NEUTRON-PROTON CORRELATION FUNCTION  
AT 23 MEV. MALANIFY, BENDT, ROBERTS, AND SIMMONS.

NEUTRON-PROTON ELASTIC DIFFERENTIAL CROSS SECTIONS  
FROM 1 TO 6 GEV. KREISLER, MARTIN, PERL, LONGO, POWELL.

Display 7

The DISPLAY TITLES command has an option which allows the user to examine, if he chooses, the indexes to the titles.

DISPLAY TITLES, INDEXES (D8)

LARGE-ANGLE NEUTRON PROTON CORRELATION FUNCTION  
AT 23 MEV. MALANIFY, BENDT, ROBERTS, AND SIMMONS.

BENDING MAGNET, ELASTIC SCATTERING, FOUR-MOMENTUM  
TRANSFER, INCIDENT-ENERGY INTERVAL, LIQUID HYDROGEN  
TARGET, PROTON POLARIZATION, RECOIL-PROTON MOMENTUM  
AND ANGLE, RELATIVE CROSS SECTION, SCATTERED-NEUTRON  
ANGLE, SPARK CHAMBER, STEEL-PLATE SPARK CHAMBER, WELL-  
COLLIMATED NEUTRON BEAM.

NEUTRON-PROTON ELASTIC DIFFERENTIAL CROSS SECTIONS  
FROM 1 TO 6 GEV. KREISLER, MARTIN, PERL, LONGO, POWELL.

ASYMMETRY, COUNTER TELESCOPE, DIRECTION OF PROTON  
POLARIZATION, NEUTRON-PROTON CORRELATION FUNCTION,  
PHASE SHIFT ANALYSIS, POLARIZATION, POLARIZED NEUTRON,  
POLARIZED PROTON TARGET, PREDICTION OF HAMADA-  
JOHNSON POTENTIAL MODEL, PROTON.

Display 8

Finally, the user may formulate a Boolean search expression with the CALCULATE command. If he desires a count of documents which contain, say, proton, spectrum, deuteron bombardment, but not helium, he may write as follows:

CALCULATE (D1 + D2 + D3 - D5)  
25 DOCUMENTS SATISFY SEARCH ARGUMENT

Display 9

As before, he can ask for a display of the titles and their indexes.

These commands may change somewhat in their detailed definition, though no changes in function are anticipated.

A WRITE command causes titles to be printed out on the printer thus providing the user with a hard copy.

### 2.3 AUTOMATIC INDEXING

Before discussing the file structure of the retrieval system we will summarize briefly the Russian-to-English Indexing System. There are two inputs to the system.

- Russian text on magnetic tape
- A machine dictionary on magnetic tape that is a computer representation of a standard English-Russian technical phrase dictionary

The Indexer matches Russian text phrases against Russian dictionary phrases. Dictionary phrases are in canonical form; Russian reverse inflection algorithms incorporated into the Indexer transform inflected text phrases to their canonical form. When a match is found, the English translation of the match is extracted from the dictionary. The output of the system, a cross-indexed index, is constructed from the set of such English translations according to algorithms described in Ref. 2-2.

A detailed examination of some indexed output will be helpful in understanding the file structure of the retrieval system. In last year's final report, a few paragraphs of a geological article, "Phase Transformations in the Interior of the Earth," by S. M. Stishov, Nature, September 1962, were indexed as an example. We now show the index of the entire article (Figs. 2-1 and 2-2). (The article itself is presented in Appendix A.) As before, two indexes are shown. The simple index is a listing, alphabetically arranged and with duplicate entries eliminated, of the English translations of Russian phrases in the text. The complex index is a cross-indexed version of the simple index. It is constructed using syntactic information provided by the reverse inflection algorithms incorporated into the Indexer.

There are a few errors in the complex index. Both items alphabetized under F, fourfold and fourfold, silicon in arise because of improper sequencing in the Russian reverse inflection algorithms which causes fourfold to be labeled a noun instead of an adjective. Resequencing of two subroutine calls will eliminate this error. The phrase earth in series of zone should more properly be earth in series of zones. This arises because there is no inflection algorithm in the system to convert zone to zones to correspond to the Russian plural. Even so, this problem would not have been noticeable had it not been for the fact that series is a noun whose plural and singular forms coincide. The phrase hypothesis about chemical is caused by a programming error which incorrectly assumed that a word following a preposition would be a noun or noun phrase. This error can be easily corrected. The phrase state of silica in condition is a meaningless phrase which occurs because the Russian word following condition does not happen to be a dictionary entry. It is difficult to see how to correct this problem since it is a semantic problem caused by the general meaning of condition. The mind expects a precise qualifier that does not come. One feels intuitively that this problem should not occur often and, in fact, in this article this is the only such occurrence.

One of the interesting items in the complex index is crystalline structure with silicon. It is interesting because cross indexing leads to two additional entries - structure,

FAZOVYE PREVRAWENI= V GLUBINAX ZEMLI  
 PHASE TRANSFORMATION IN GLUBINAX OF EARTH

AMORPHOUS ARRANGEMENT AVERAGE DENSITY	INCREASE INDEX OF REFRACTION INDIVIDUAL INTERMEDIATE INVESTIGATION IRON IRON METEORITE	PROPERTY PURE
CHARACTER CHEMICAL CLOSEST PACKING COATING COESITE CONDITIONS COORDINATION NUMBER CORUNDUM CRYSTALLINE CUBIC PACKING	JOINT	REGION REGISTER RESEARCHER ROCK ROENTGENOSPECTRAL RUTILE
DENSE DENSITY DIFFERENTIATED DIFFERENTIATION DISTRIBUTION	LIGHT LIMIT LOWER	SCIENTIST SERIES SIDE SILICA SILICATE SILICON SLAG SPICULAR STAGE STATE STRATUM STRUCTURAL STRUCTURE SULPHIDE ORE SULPHITE SUPPORT SURFACE SYMMETRIC SYNTHESIS
EARTH EARTH SHELL EARTHQUAKE ELASTIC ELASTICITY	MAGNESIUM MANTLE MATTER MAXIMUM MELTING METAL METALLIC METEORITE METEORITIC CRATER MIXTURE MODEL MODIFICATION MONOXIDE	
FOURFOLD	OBSERVATION	
GEOLOGIST GEOPHYSICS GRADIENT	OLIVINE ORIGIN OXYGEN	TABULAR CRYSTAL TEMPERATURE TETRAHEDRAL THEORY TRANSFORMATION TRANSITION TYPE
HARDNESS HETEROGENETIC HIGH HARDNESS HOMOGENEOUS HYPOTHESIS	PACKING PERICLASE PHASE PHASE TRANSITION PRESSURE PROCESS	WUSTITE
IMPURITY		ZONE

Fig. 2-1 Simple Index

AMORPHOUS SILICA  
ARRANGEMENT  
AVERAGE DENSITY OF EARTH

CHARACTER OF LIMIT  
CHARACTER OF METAL  
CHEMICALLY DIFFERENTIATED EARTH  
CLOSEST PACKING  
CLOSEST PACKING WITH OXYGEN  
COATING  
COATING, SILICATE  
COESITE, DENSITY OF  
CONDITIONS  
CONDITIONS, STATE OF SILICA IN  
COORDINATION NUMBER  
CORUNDUM  
CORUNDUM, HARDNESS WITH  
COSMOGONY, GEOPHYSICS AND  
CRUST  
CRYSTALLINE STRUCTURE WITH SILICON  
CUBIC PACKING

DENSE MODIFICATION  
DENSITY, GRADIENT OF  
DENSITY, INCREASE OF  
DENSITY AND ELASTICITY  
DENSITY OF COESITE  
DENSITY OF EARTH  
DIFFERENTIATION OF MATTER  
DISTRIBUTION OF DENSITY IN EARTH

EARTH  
EARTH, AVERAGE DENSITY OF  
EARTH, CHEMICALLY DIFFERENTIATED  
EARTH, DENSITY OF  
EARTH, DISTRIBUTION OF DENSITY IN  
EARTH, HYPOTHESIS ABOUT ORIGIN OF  
EARTH, IRON IN CENTER OF  
EARTH, MANTLE OF  
EARTH, MATTER OF  
EARTH, SURFACE OF  
EARTH, THEORY OF ORIGIN OF  
EARTH IN SERIES OF ZONE  
EARTH SHELL  
EARTHQUAKE  
ELASTIC PROPERTY OF LOWER MANTLE  
ELASTICITY, DENSITY AND

FOURFOLD  
FOURFOLD, SILICON IN

GEOLOGIST  
GEOPHYSICS  
GEOPHYSICS AND COSMOGONY  
GRADIENT OF DENSITY

HARDNESS WITH CORUNDUM  
HIGH HARDNESS  
HOMOGENEOUS MANTLE  
HOMOGENEOUS MATTER  
HYPOTHESIS  
HYPOTHESIS ABOUT CHEMICAL

HYPOTHESIS ABOUT ORIGIN OF EARTH

IMPURITY, IRON WITH  
INCREASE OF DENSITY  
INDEX OF REFRACTION  
INDIVIDUAL MONOXIDE, MIXTURE OF  
INTERMEDIATE REGION  
INVESTIGATION  
IRON  
IRON IN CENTER OF EARTH  
IRON METEORITE, EXISTENCE OF  
IRON WITH IMPURITY

LIGHT  
LIMIT  
LIMIT, CHARACTER OF  
LOWER MANTLE, ELASTIC PROPERTY OF  
LOWER MANTLE, STRUCTURAL MODEL OF  
MAGNESIUM  
MANTLE  
MANTLE, HOMOGENEOUS  
MANTLE, MATTER OF  
MANTLE, PROPERTY OF  
MANTLE, REGION IN  
MANTLE, STATE OF MATTER OF  
MANTLE, STRUCTURE OF  
MANTLE OF EARTH  
MATTER  
MATTER, DIFFERENTIATION  
MATTER, HOMOGENEOUS  
MATTER OF EARTH  
MATTER OF MANTLE  
MAXIMUM  
MELTING OF SULPHIDE ORE  
METAL  
METAL, CHARACTER OF  
METAL, MIXTURE OF SULPHITE AND  
METALLIC IRON  
METALLIC PROPERTY  
METEORITE  
METEORITIC CRATER  
METEORITIC CRATER, ROCK OF  
MIXTURE  
MIXTURE OF INDIVIDUAL MONOXIDE  
MIXTURE OF SULPHITE AND METAL  
MODEL  
MODEL, STRUCTURAL, OF LOWER  
MANTLE  
MODIFICATION  
MODIFICATION, DENSE  
MODIFICATION OF SILICA  
MONOXIDE  
MONOXIDE OF TYPE OF PERICLASE

OBSERVATION  
OLIVINE, PROPERTY OF ZONE WITH  
TRANSITION OF  
OLIVINE, STRUCTURE OF  
OXYGEN  
OXYGEN, CLOSEST PACKING WITH

PACKING, SYMMETRIC  
PERICLASE, MONOXIDE OF TYPE OF  
PHASE

PHASE TRANSITION

PRESSURE  
PRESSURE AND TEMPERATURE  
PROCESS  
PROPERTY  
PROPERTY, ELASTIC, OF LOWER MANTLE  
PROPERTY, METALLIC  
PROPERTY OF MANTLE  
PROPERTY OF ZONE WITH TRANSITION  
OF OLIVINE

REGION, INTERMEDIATE  
REGION IN MANTLE  
REGISTER  
RESEARCHER  
ROCK  
ROCK OF METEORITIC CRATER  
RUTILE

SIDE, SUPPORT WITH  
SILICA, AMORPHOUS  
SILICA, MODIFICATION OF  
SILICATE COATING  
SILICATE STRATUM  
SILICON  
SILICON, CRYSTALLINE STRUCTURE WITH  
SILICON, TRANSITION OF  
SILICON IN FOURFOLD  
SLAG  
SPINEL, STRUCTURE OF TYPE OF  
STAGE  
STATE, STRUCTURAL  
STATE OF MATTER OF MANTLE  
STATE OF SILICA IN CONDITIONS  
STRATUM  
STRATUM, SILICATE  
STRUCTURAL STATE  
STRUCTURAL MODEL, OF LOWER MANTLE  
STRUCTURAL TYPE  
STRUCTURE  
STRUCTURE, CRYSTALLINE, WITH SILICON  
STRUCTURE OF MANTLE  
STRUCTURE OF OLIVINE  
STRUCTURE OF TYPE OF SPINEL  
SULPHIDE ORE, MELTING OF  
SULPHITE  
SUPPORT WITH SIDE  
SURFACE OF EARTH  
SYNTHESIS

TABULAR CRYSTAL  
TEMPERATURE  
TEMPERATURE, PRESSURE AND  
THEORY  
THEORY OF ORIGIN OF EARTH  
TRANSFORMATION  
TRANSITION  
TRANSITION OF SILICON  
TYPE, STRUCTURAL

WUSTITE

ZONE  
ZONE, EARTH IN SERIES OF

Fig. 2-2 Complex Index

crystalline, with silicon and silicon, crystalline structure with. For retrieval purposes, this leads to four pointers - crystalline, crystalline structure, structure, and silicon - to the phrase crystalline structure with silicon.

As mentioned in Ref. 2-2, articles do not appear in the index. If they occur in a user's input query they will be edited out.

The simple index consists of 104 distinct entries representing the 1550-word text. Most, though not all, of the entries are single words. It should be remembered that each word/phrase entry in the simple index represents a separate distinct entry in the computer dictionary. Entries in the computer dictionary will be referred to as prime items.

The complex index (Fig. 2-2) is more interesting and more meaningful. It consists of prime items from the simple index arranged in various syntactic combinations. If a prime item is a noun and cannot be syntactically combined with another prime item, then it is listed by itself in the complex index (e. g., arrangement, corundum).

A discussion of the syntactic structures of the entries in the complex index can be found in Section 6 of Ref. 2-2. For the purposes of this discussion, however, an item in the complex index can be viewed as belonging to one of three categories. The item may be a prime, secondary, or tertiary item. Prime items have been defined already as discrete entries in the dictionary. These entries are usually nouns, adjectives, noun phrases, and noun phrases followed by prepositional phrases. A secondary is a noun or noun phrase preceded by at least one adjective. The noun (or the noun phrase) and the preceding adjectives are prime items. A tertiary is a string of prime nouns or noun phrases, and secondaries separated by prepositions and/or conjunctions. For example:

- Prime items - electron  
emission  
magnetic field (noun phrase)  
rotating  
distribution  
density  
earth



- Secondary items - electron emission  
rotating magnetic field
- Tertiary items - electron emission in rotating magnetic field  
distribution of density in earth

Secondary and tertiary items may easily be cross indexed and a glance at the complex index will show that this has been done.

In the introduction it was stated that retrieval of a document could be accomplished according to syntactic relatedness of terms within the document as opposed to the mere existence of the terms themselves within the document. We can see now that it is the secondary and tertiary items which contain the information regarding the syntactic relatedness of various prime terms.

## 2.4 FILE STRUCTURE

### 2.4.1 Inverted File

One of the file structures which can accommodate a Boolean search is the inverted file. In an inverted file system each document is identified by a number of descriptors, which could be a sequence number, a title, authors, a date, source, index terms, or some other category of descriptor. A separate ordered file exists for each descriptor category to facilitate search for such a descriptor. Each of these files is called an inverted file. For example, one file has the information ordered according to authors, another according to date, another according to index term etc. The entire file need not be examined, just those portions associated with the input descriptors.

### 2.4.2 Threaded Lists

A threaded list is a list whose elements may be located discontinuously at a series of locations throughout a file. Each element of the list is associated with a pointer which

points to the location in the file of the succeeding element. The location pointer of the final element of the list is filled, perhaps, with zeros or some other marker to indicate the end of a list. Thus the phrase electron emission in rotating magnetic field is formed by the threaded list as shown.

```
100 ELECTRON
    110
    .   .
    .   .
    .   .
110 EMISSION
    347
    .   .
    .   .
    .   .
205 FIELD
    00000
    .   .
    .   .
    .   .
347 IN
    743
    .   .
    .   .
    .   .
455 MAGNETIC
    205
    .   .
    .   .
    .   .
743 ROTATING
    455
```

Lists may be altered not by physically rearranging elements of the file, but by changing the pointers as desired.

### 2.4.3 File Structure of the Retrieval System

The file structure of this retrieval system can be described as a combination of an inverted file and a threaded list. In a previous section it was stated that an index could be viewed as being composed of prime, secondary, and tertiary elements and that the secondary and tertiary elements express the syntactic relatedness of prime elements.

The retrieval system operates upon four files -- a prime/secondary file (PSF), a tertiary file (TF), a document file (DF), and an index file (IF).

The PSF has an inverted file/threaded list structure. The tertiary, document, and index files are arranged randomly. Elements of the tertiary, document, and index files and of the secondary portion of the PSF are elements of at least one threaded list. The origin of each threaded list is one of the primary elements of the PSF. The end of each threaded list is an element of the index file. A simplified example will clarify this structure. Suppose the article "Scattering of Slow Electrons by Enhanced Ion Waves Near the Geomagnetic-Field Boundary," by A. Eviator, has an index consisting only of the following two items:

distribution of density in earth  
electron emission in rotating magnetic field

these are tertiary items which have been constructed from the following prime and secondary items:

- Prime items      - density  
                          distribution  
                          earth  
                          electron  
                          emission  
                          magnetic field  
                          rotating
- Secondary items - electron emission  
                          rotating magnetic field

The arrangement of the four files for this simple example is shown in Fig. 2-3.

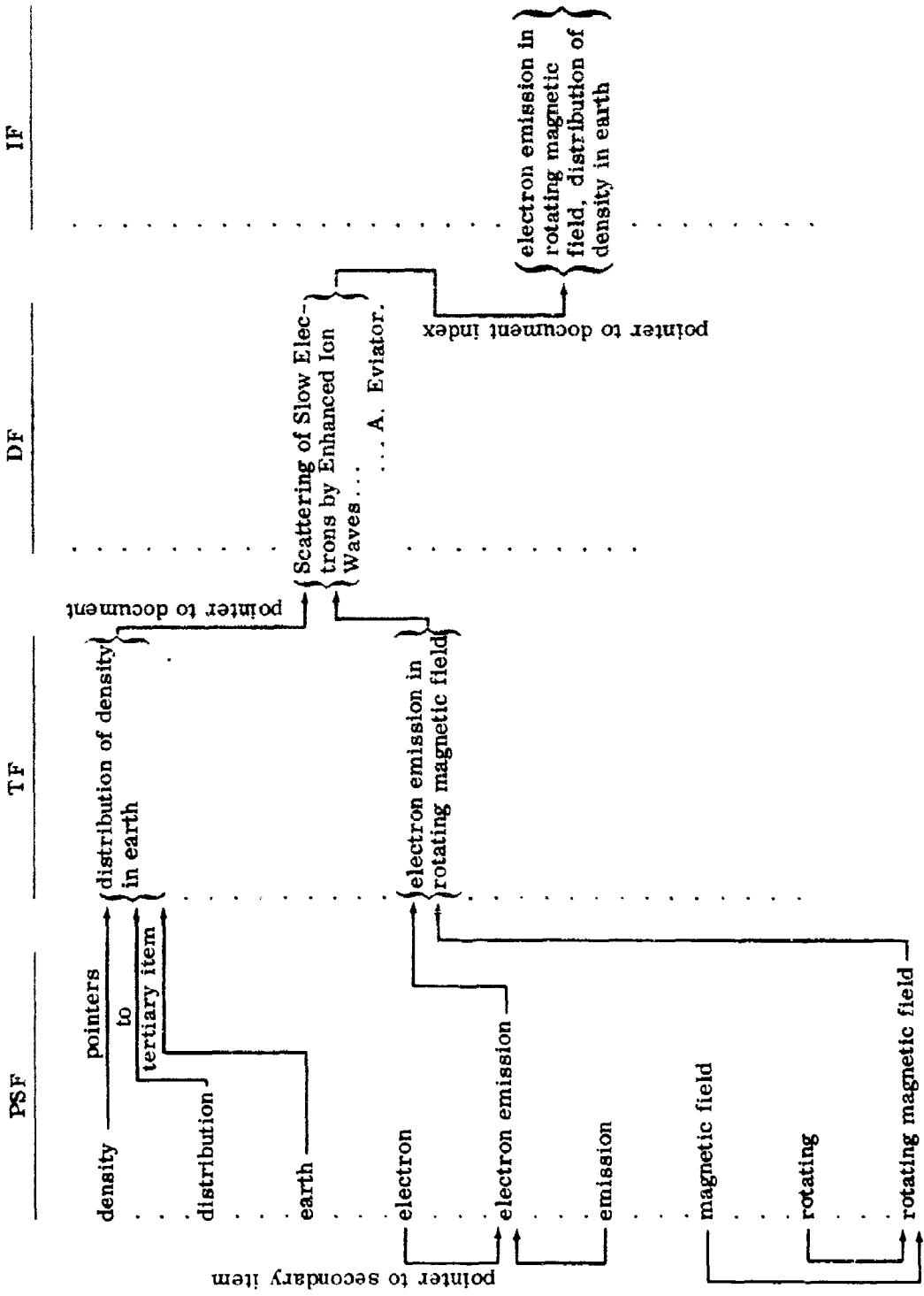


Fig. 2-3 File Arrangement

## 2.5 IBM 360 CONFIGURATION FOR RETRIEVAL SYSTEM

The detailed configuration of the data files and the details necessary for actual programming are discussed in this subsection. Though these details have not been worked out completely, it is expected that they will not differ greatly from those presented here.

This retrieval system is being designed to run under the LACONIQ monitor. LACONIQ is a time-sharing system designed to implement real-time applications. LACONIQ handles queuing, polling, chaining, and other real-time functions allowing the applications designer to concentrate on the logic of his own particular application.

Figure 2-4 is a list of various field codes used to denote the files, their possible overflows, and other information which may be required during programming.

Figure 2-5 shows a schematic diagram of the general file structure. The lower three blocks are detailed diagrams of specific fields in the top diagram. Figure 2-6 shows an example of how the data file for the prime item energy might look.

As was mentioned in the previous subsection, the data records will be kept on four different files. The files are called the primary/secondary file (PSF), the tertiary file (TF), the document file (DF), and the index file (IF). Figure 2-7 shows some of the important characteristics of these files. The overflow files are required in case data exceed the allowable record size.

The files are used as follows: the prime/secondary file (PSF) is the file against which users' descriptors are processed initially. Through Boolean operations, an attempt is made to get as close as possible to the user's request, progressing to the tertiary file (TF) if possible. From the information in the prime/secondary file (PSF) and tertiary file (TF) the program is able to tell the user how many documents in the document file (DF) contain his descriptors, and also tell him the number of "higher level" descriptors, in the PSF and TF, in which each of his descriptors are contained.

FIELD CODES

HEXADECIMAL

F2	REFERENCES TO SECONDARY INDEX RECORDS
F3	REFERENCES TO TERTIARY INDEX RECORDS
C4	REFERENCES TO DOCUMENT RECORDS
C9	ITEM
FE	PADDING
FD	CHAIN ADDRESSES

RECORD CODES

F1	PRIMARY FILE
01	PRIMARY FILE OVERFLOW
F2	SECONDARY FILE
02	SECONDARY FILE OVERFLOW
F3	TERTIARY FILE
03	TERTIARY FILE OVERFLOW
C4	DOCUMENT FILE
04	DOCUMENT FILE OVERFLOW

Fig. 2-4 Field and Record Codes

BASIC LOGICAL RECORD FORMAT

NO. OF BYTES	2	1	3 BYTES EACH	VARIABLE		
C/F	T Y P E  C O D E		ARBITRARY NUMBER OF FIELD INDICES	FIELD 1	FIELD 2	...

C/F DETAIL

NO. OF BITS	4	12
	00X0 IF X = 1, RECORD IS CONTINUED	COUNT (NO. OF BYTES IN THIS LOGICAL RECORD)

FIELD INDEX DETAIL

	LAST FIELD INDEX				
NO. OF BYTES	1	2		1	2
	FIELD IDENTIFI- CATION CODE	DISTANCE (IN BYTES) OF FIELD FROM START OF LOGICAL RECORD	...	X'FF'	NO. OF BYTES IN THIS LOGI- CAL RECORD

FIELD DETAIL

	SUBFIELD 1		LAST SUBFIELD	
	1	VARIABLE		VARIABLE
	NO. OF BYTES IN SUBFIELD	DATA	... X'00'	DATA

Fig. 2-5 Basic Data File Structure for Random Access Files

FIRST SECONDARY REFERENCE  
IS ON CYLINDER 16, TRACK 3,  
RECORD 4

5 SECONDARY (F2)  
REFERENCES  
FOLLOW

FIELD NO. 1

FIELD NO. 2

41	F1	C9	12	F2	19	FF	44	00	ENERGY	00	5	1634	2087	2091	3005	3091
----	----	----	----	----	----	----	----	----	--------	----	---	------	------	------	------	------

C/F

THIS RECORD  
IS 41 (DECIMAL)  
BYTES LONG

TYPE CODE

F1 = PRIMARY  
FILE

FIELD INDEX NO. 1

C9 = ITEM. 12 BYTES  
DISTANCE FROM  
START OF LOGICAL  
RECORD

FIELD INDEX NO. 2

F2 = SECONDARY FILE  
19 BYTES DISTANT  
FROM START OF  
LOGICAL RECORD

LAST FIELD INDEX

FF = LAST FIELD INDEX CODE  
44 BYTES IN THIS LOGICAL  
RECORD

2-21

LOCKHEED PALO ALTO RESEARCH LABORATORY  
LOCKHEED MISSILES & SPACE COMPANY  
A GROUP DIVISION OF LOCKHEED AIRCRAFT CORPORATION

Fig. 2-6 Sample Data File for Item energy



	CODE (HEX)	LENGTH (BYTES)	REFER-ENCES TO SECOND FILE	REFER-ENCES TERTI-ARY FILE	REFER-ENCES TO DOCU-MENTS	BCD DATA	PADDING	CHAIN ADDRESS
PRIMARY FILE	F1	223	0	0	0	00	0	0
PF OVERFLOW	01	223	0	0	0	-	0	0
SECONDARY FILE	F2	223	-	0	00	00	0	0
SF OVERFLOW	02	223	-	0	0	-	0	0
TERTIARY FILE	F3	223	-	-	00	00	0	0
TF OVERFLOW	03	223	-	-	00	-	0	0
DOCUMENT FILE	C4	512 MAX.	-	-	-	00	-	0
DF OVERFLOW	04	512 MAX.	-	-	-	00	-	0
INDEX FILE			-	-	-	00	0	0

NOTE: 0 = MAY HAVE  
 00 = MUST HAVE  
 - = WILL NOT HAVE

FILE	PHYSICAL RECORD SIZE (BYTES)	MAX./EXACT	LOGICAL RECORD SIZE (BYTES)	MAX. NO. OF LOG. REC./PHYS.	NO. OF PHYS. REC./TRACK	SUBJECT TO UPDATE	IN SEQUENCE
PSF	476	E	119 OR 238 OR 476	4	6	YES	YES
TF	446	E	223 OR 446	2	7	YES	NO
DF	446	M		1	VARIABLE	NO	NO

Fig. 2-7 File Characteristics

The prime/secondary file (PSF), as its name implies, contains both prime and secondary index items. Since the file is arranged alphabetically, the prime and secondary items are mixed together. Also, since secondary index items are "pointed to" by prime index items, it is possible for the PSF to point to itself. (See Fig. 2-3.)

All four files will be located on the IBM 2311 Disk Storage Drive. Such a disk has a maximum capacity of 7.25 million data bytes. Since LACONIQ is also located on disc, the retrieval system will be completely disc-based.

## 2.6 DATA BASE

A sample data base has been prepared which will be used to check out the retrieval system. The data base consists of 100 physics abstracts selected from the Bulletin of the American Physical Society, January 1966. (See Appendix B.) Of the 100 abstracts, the complex indexes of 50 have been keypunched. The rest, while indexed, have not yet been keypunched.

The abstracts, it should be pointed out, are in English. The complex indexes have been human-prepared. The procedure used to produce the indexes was to go through each abstract and to indicate primary, secondary, and tertiary items. While the primary items in this case do not correspond to any particular physics dictionary, as long as consistency in indexing is observed by the indexer, the sample data base will be an accurate simulation of a data base produced by a "real" set of primaries.

It is interesting to note that in the computer-indexed article, 104 primaries gave rise to 42 secondary and 43 tertiary items, while in the set of abstracts 613 primaries produced 168 secondary and 132 tertiary items. If the proportions of the article held, we could have expected the set of abstracts to produce around 250 secondary and tertiary items. This discrepancy is perhaps due to the terse language one would expect an abstract to be written in (as opposed to the original article). There are approximately 4000 words in the set of Physics abstracts.

In addition, since it is desirable that the retrieval system be capable of using word government information (Ref. 2-3), such patterns have been noted by the indexer and are incorporated into the data base.

## 2.7 REFERENCES

1. D.L. Drew, "LACONIQ Standards Note 4, The LACONIQ Monitor," Technical Note, Information Sciences group, Lockheed Palo Alto Research Laboratory, 22 Aug 1966
2. Annual Report: Automatic Indexing and Abstracting, Part II, English Indexing of Russian Technical Text, M-21-66-2, Lockheed Missiles & Space Company, Palo Alto, Calif., Mar 1966
3. Research in Automatic Information Abstracting and Extracting, Renewal Proposal, Office of Naval Research Contract NONr 4440(00), LMSC-894736, Lockheed Palo Alto Research Laboratory, Oct 1966
4. Ye. V. Krotevich, "Syntactic Relations Between Members of a Word Combination and Members of the Sentence," Voprosy Russkogo Yazykoznaniiya (Problems of Russian Linguistics), Vol. 2, Moscow, 1956
5. Annual Report: Automatic Indexing and Abstracting, Part I, M-21-66-1, Lockheed Missiles & Space Company, Palo Alto, Calif., Mar 1966

### Section 3

## COMPUTABLE RELATIONS BETWEEN THE ORTHOGRAPHIC AND PHONETIC FORMS OF ENGLISH MONOSYLLABLES\*

### 3.1 INTRODUCTION

Carefully checked and highly accurate sets of computable relations between the orthographic forms and the corresponding phonetic forms of English words are the basis for design of equipment such as a mechanized reader for the blind or a talking computer (Ref. 1). The former would transform into speech-like sounds, mechanically and electronically, the information on punched paper tapes (widely used for driving linotypes). Such relations are also a subject of intrinsic scientific and linguistic interest.

Although there is substantial agreement on the orthographic forms of English words, the phonetic representations corresponding to a large number of English monosyllables differ from one transcription to another (Ref. 3-1). This situation, partly as a result of different dialect patterns, is inherent in the subjectivity of phonetic transcriptions (Ref. 3-2). We do not question the phonetic representations prepared by authoritative sources, which are indeed the best efforts of skilled individuals, nor do we claim that our computer programs supplant their decision processes in any sense. Our data indicate, rather, that the seemingly subjective process of preparing phonetic transcriptions can be precisely defined by use of the phonetic information that is encoded in the orthographic forms of words. Although phonetic information is not represented explicitly in English orthography, as it is in phonetic scripts such as Sanskrit, sufficient phonetic information is available in the orthographic forms, barring the exception words.

Since the phonetic transcription of the same corpus of words differs from one reputable authority to the next, it is advisable to write separately the computable relations

---

\*Work on this task was conducted by B. V. Bhimani.

between the orthographic forms and the corresponding phonetic representations prepared by each authority, and then to arrive at an optimum set of rules that displays maximum agreement in the phonetic representations of corresponding orthographic forms. To define such an optimum set of computable relations also requires comparison of the phonetic representations of a specified and large set of words as inscribed by different authorities: such a comparison provides information on the consistent patterns among the phonetic representations studied, as well as on differences in the transcriptions. At first, such a comparison of phonetic data may seem to be needed primarily for the development of a speech recognizer, or for certain linguistic studies such as those relating the dialect patterns of English; but it is also required for ascertaining that the phonetic representations obtained from orthographic-form words, by the use of computable relations, do indeed provide data which have a high degree of agreement with authoritative transcriptions, and that the computed phonetics represent speech patterns with a low probability of confusion by the recipient. Such an evaluation of phonetic data is also an aid in the selection of exception words for increasing the accuracy of the program that maps the orthographic form of words into corresponding phonetic representations.

The procedures which will be described here do take into account the differences between the dialects studied, but not the idiosyncrasies of individual speakers; it is expected that extension of these techniques may permit treatment of the latter. The computable relations between the phonetic data representing the various dialects provide a machine with an ordered set of rules which, when properly used, can reduce the burden of storage and search for a general-purpose speech recognizer or synthesizer, as contrasted with the needs of an unordered machine which must treat each case separately.

Such a systematic model of the relation between orthographic and phonetic forms is a prerequisite for both mechanical recognition and generation of speech. For the former application, the task is the identification of acoustic characteristics with an a priori symbolic representation; for the latter, the symbols representing the phonetic form

of words must be transformed into corresponding and prespecified acoustic characteristics. This subject is discussed in recent literature (Refs. 3-3 through 3-8).

On the basis of the data discussed in subsections 3-2 and 3-3, we contend that computable relations between the orthographic forms and the corresponding phonetic representations, as well as between the various phonetic representations that correspond to any orthographic form, can be established for a group of symbols called marker-vowel-marker. This group of symbols is similar to groups described previously on the basis of logical and acoustic studies of speech (Ref. 3-9). (This paper does not present an exhaustive list of marker-vowel-marker segments or their acoustic characteristics, but a method for computing such a list is presented in subsection 2.) It is also shown in this study that the phonetic data transcribed by one authority cannot be accurately transformed into that by another authority by writing equations that map each vowel symbol of one transcription into a corresponding set of vowel symbols of another transcription or by equations that map each set of consonant symbols of one transcription into corresponding consonant symbol(s) of another transcription.

Whether the marker-vowel-marker group or symbols are in fact the units of speech as produced and/or perceived by humans is a subject for other papers. This paper discusses the evaluation of the computable relations that map into the various phonetic representations from corresponding orthographic-form English monosyllables [defined here as those words which have a single uninterrupted vowel string in their phonetic representation according to the Shorter Oxford Dictionary (SOX)(Ref. 3-10)]. The phonetic representations studied are those in five recent dictionaries (Refs. 3-10 through 3-14). In subsection 3-2 an approach is proposed for extending the work to polysyllable words of English, the validity of such an extension is considered briefly in subsection 3-5. This treatment of the phonetics of English, and of the dialect patterns of the language differs from that of the linguistic atlas (Ref. 3-15) in that the phonetic data in dictionaries are used as the basic source and computable relations are the principal objective of this study. The computer programs discussed herein are sufficiently general to assure conversion of new orthographic-form monosyllables (those not

included in the present set) with a high degree of accuracy. These computable relations require storage and operation times which will permit them to be run on a medium-size, medium-speed computer with less than 10-percent error, when computed phonetics are compared with those of five dictionaries (Refs. 3-10 through 3-14). (Error rates can be reduced to the desired level by storing a set of exception words.)

The computer programs which generate phonetic representations of orthographic-form words can readily be adapted to provide phonetic forms as transcribed by any qualified linguist. This is because the entire algorithm is related to the orthographic form of English words and the computable relations use the marker-vowel-marker symbols, which also establish computable relations between the various phonetic forms (Refs. 3-10 through 3-14). When the computation from the orthographic form into the phonetic form of words is completed, the resulting data are compared with those in other sources (Refs. 3-10 through 3-14). The computable relations are final only after a detailed manual and computer-oriented check is made of the accuracy of the resulting data.

### 3.2 MECHANIZED CONVERTER OF ORTHOGRAPHIC TO PHONETIC FORMS OF WORDS

Illustrated in Fig. 3-1 are the several functions a computer must perform to receive questions in the form of coded pulses and transmit its answers in the form of synthesized speech. The reader is referred to recent advances in computer software that promise to make it practical to write programs in a freely disposable subset of a user's natural language (Ref. 3-16), and he is assumed to be familiar with the computer functions of scanning the input statement or question, of performing the necessary searches and look-ups, and of formulating the answer. Since present-day computer output is almost exclusively orthographic, it must be transformed into appropriate phonetic form for synthesizing speech (Ref. 3-1). Similarly, input tapes which could be used in mechanized readers for the blind are also available in a form suitable only for printing. A phonetic transformation is therefore necessary for this application as well. The phonetic forms generated by such a transformation can serve as inputs to speech synthesizers. Several such synthesizers are discussed in the literature (Refs. 3-3, 3-4, 3-17, 3-18, and 3-19).

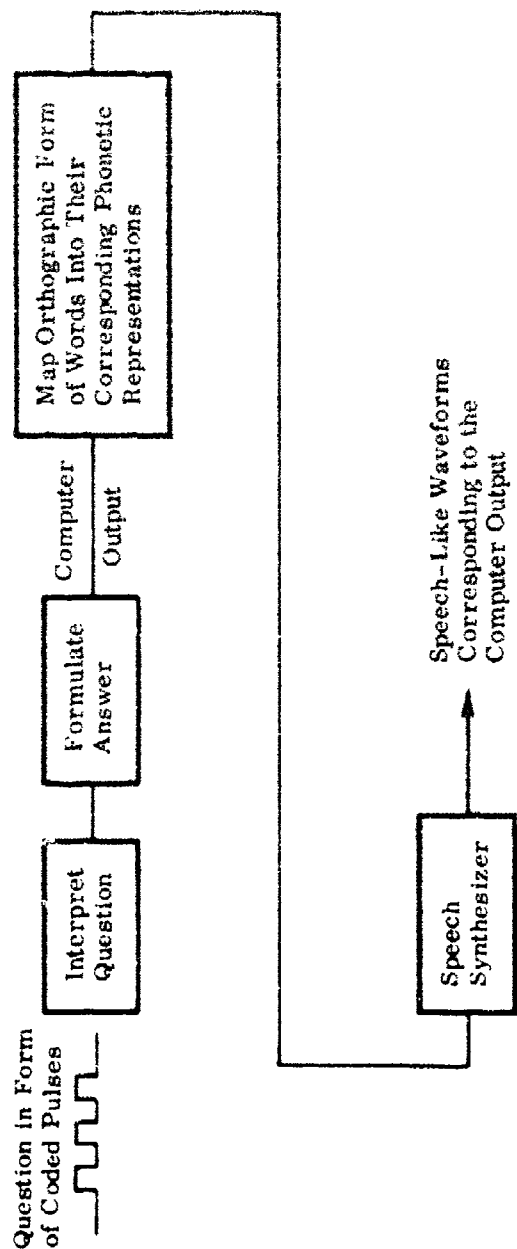


Fig. 3-1 Functions Performed for Purposes of Pulse-Input, Voice-Output Communication With Computers



The following discussion of the computable relations between the orthographic and phonetic forms of English monosyllables forms the basis of a program for more general application. Specifically, extension of the techniques for dealing with polysyllabic words requires the development of a processor for segmenting the words into their constituent syllables (Fig. 3-2). The program intended in Fig. 3-2 is beyond the capability of the hyphenation algorithms used in the printing industry, since they are required to find only some of the syllabic breaks in any polysyllabic word and at times may indicate no such breaks. For the mechanized reader, such a syllabic break must be indicated with high accuracy, and information must be provided on the stressed syllables in each word. A detailed discussion of this subject is beyond the scope of this paper. Moreover, the reading and spelling rules available in dictionaries (Ref. 3-20) or in recent papers (Ref. 3-21) have not been checked by their proponents by use of computer programs operating on exhaustive word lists and checking the results with phonetic transcriptions by more than one authority; hence, these rules may be inaccurate or incomplete and unsuitable for programming computers.

An orderly development of programs for mechanized readers requires starting with elementary words (Ref. 3-22) (essentially the one-syllable words of the language and described in detail in the references), evaluating the results obtained, and then extending the rules to include English monosyllables. (The distinctions between the elementary words and the monosyllables are considered in subsection 3.5.) Most of the elementary words and some of the additional monosyllables are kernels (for example, AC-TU-AL are three kernels in the word ACTUAL), since they contain sequences of consonant strings and vowel strings that are similar to ones observed in individual syllables of polysyllabic words. (Hence, elementary words are mostly kernel words of the language.) Since polysyllabic words are simply compoundings of the consonant strings and vowel strings that are similar to those in kernels, the accuracy with which we treat the kernel words will determine the accuracy of our work on a complete lexicon. It is also shown in subsection 3.5 that the work with elementary words can be readily applied for computing the phonetic representations of English monosyllables. It is the elementary words, therefore, for which establishing of a maximum accuracy for defining a computable set of rules is required, and the resulting phonetic data also merit the most careful check.

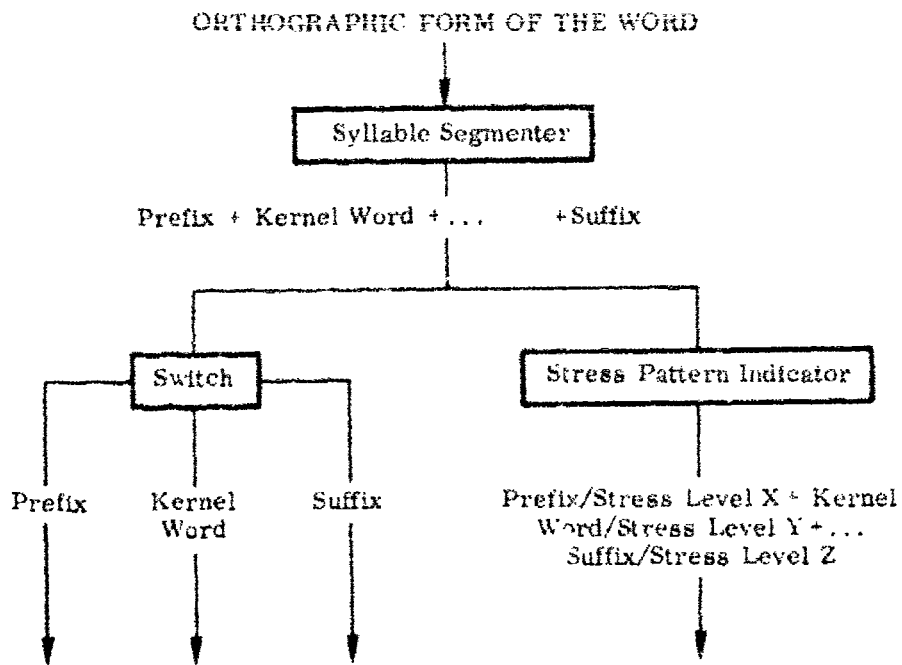


Fig. 3-2 Outline of Algorithm for Extension to Polysyllabic Words

Figure 3-3 illustrates the algorithm for the elementary words of English. The first step is the breakup of the word into vowel strings and consonant strings. For this program, the final  $\text{E}$  is classified as a marker; all other occurrences of E, as well as the letters A, I, O, U, and Y, are classified as orthographic vowels. All the remaining letters of the alphabet are classified as consonants. Words starting with vowels are classified as starting with a blank consonant  $\emptyset$ , and those ending in vowels are considered as ending in  $\emptyset$ .

The second step resolves the consonant string ambiguities which occur with orthographic letters C, G, S, TH, CH, etc. The algorithm for such a resolution of C and of G is tabulated below.

	<u>Orthographic</u>	<u>SOX Phonetic</u>
If C is followed by A, O, U:	C	k
Otherwise:	C	s
If G is followed by E or $\text{E}$ :	G	d <sub>3</sub>
Otherwise:	G	g

After these consonant strings and vowel strings have been formed, one of the special and novel features of our program is brought into play: these symbols are not transcribed directly but are first processed according to the rules of euphonic combination. [These rules, described in detail elsewhere (Ref. 3-9), define the effect of each English speech sound on each neighboring one; the referenced set is compiled specifically for English.] Such a process provides the desired phonetic transformation without regard to any particular system of phonetic symbols, and it creates an intermediate representation that is independent of any particular orthographic to phonetic transformation, representing the kernel words as consonant-marker-vowel-marker consonant strings (Fig. 3-3). Interpretation of the precise phonetic representation that corresponds to each marker-vowel-marker set provides the pronunciation for each of these words as given in the dictionary (Refs. 3-10 through 3-14). Needless to say, each dictionary requires its own set of rules for interpretation of the marking system. Figure 3-4 illustrates the operation of the program for the word SPICE, and the word GRUDGE

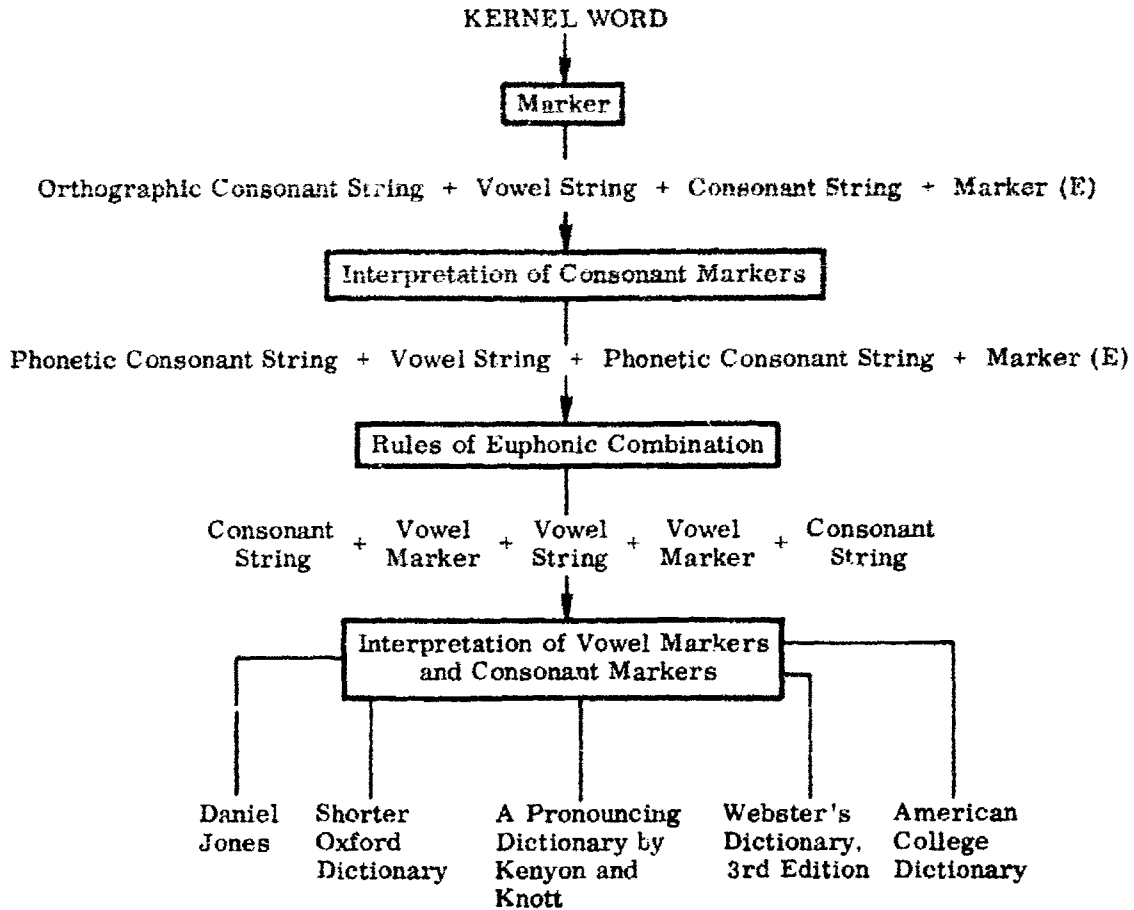


Fig. 3-3 Outline of Algorithm for Obtaining Phonetic Form From the Orthographic Representation of Elementary Words

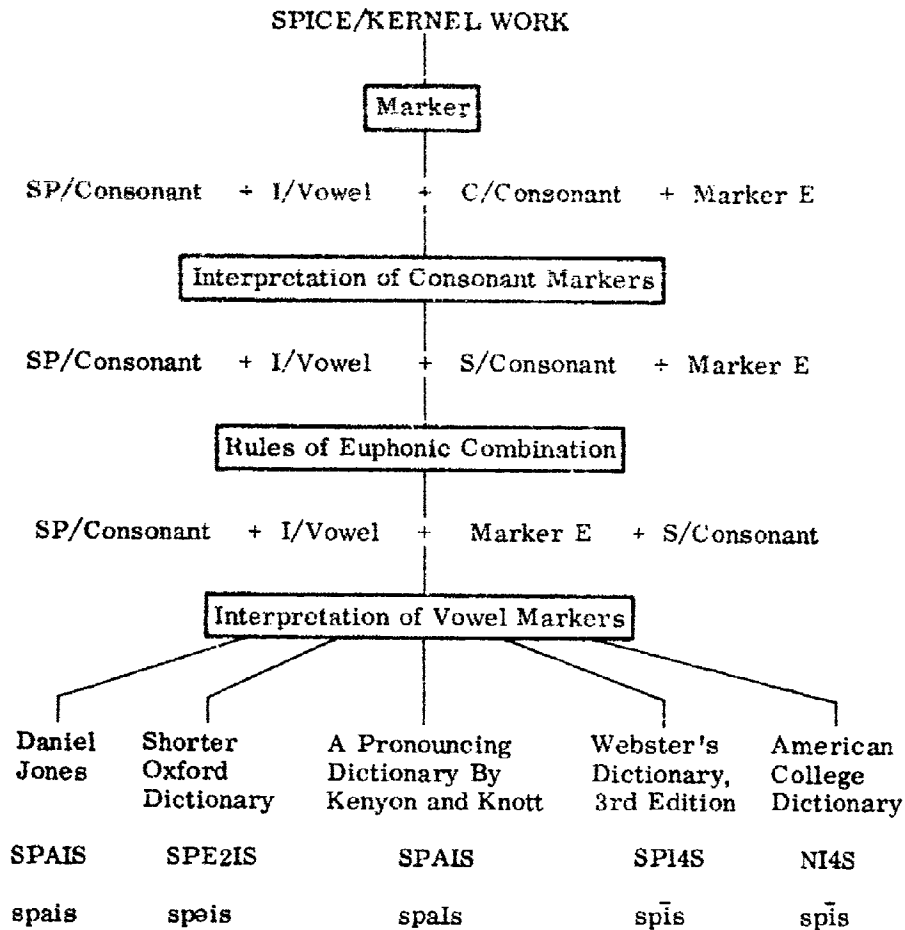


Fig. 3-4 Operation of the Program for the Word SPICE

is the subject of Fig. 3-5. Notice that the phonetic consonant string  $dd_3$  (D DZ1) is not pronounceable according to the rules of euphonic combination, and hence D is tagged as a marker. The D/marker blocks the operation of the  $\text{E}$  on the vowel U, and the pronunciation of the vowel is the same as that for the words where D follows the vowel string U, such as the word MUD.

A detailed listing of all the mapping rules needed to resolve each of the ambiguities would be of limited interest here, partly because of the unreasonably lengthy listing but primarily because of the differences in the interpretation of markers for different lexicographers. An illustration of the computable rules for mapping U into corresponding phonetic forms, as in transcriptions by Kenyon and Knott (Ref. 3-13), is presented in Table 3-1. Notice the functions of  $\text{E}$ , R, L as markers (Ref. 3-9) and the need for logical ordering of these rules. In general, the markers must be used in conjunction with one another by means of a set of precedence relations. (This may be partly responsible for the general feeling that English orthography does not present a neat phonetic representation.)

The computable rules for mapping the orthographic forms of English words into corresponding phonetic representations, as discussed herein, have been checked on an exhaustive word list (the monosyllables in the SOX). The computer programs provide phonetic representations for each of these words as recorded by five lexicographers and for each dialect studied by them (Refs. 3-10 through 3-14). Such an examination of rules and of phonetic data goes beyond work where a select set of words is compared with the phonetic transcriptions by a single authority (Ref. 3-11): the merits of our approach are obvious from the data of subsection 3.3. The rules were not checked for accuracy with published information, some of which may warrant evaluation, but they have been precisely stated. For example, the word SIGH can be analyzed as in the following tabulation.

	<u>Orthographic</u>	<u>SOX Phonetic</u>
Initial consonant string	S	s
Vowel string	I +	ei
Terminal consonant string	(GH/marker)	

GRUDGE/KERNEL WORD

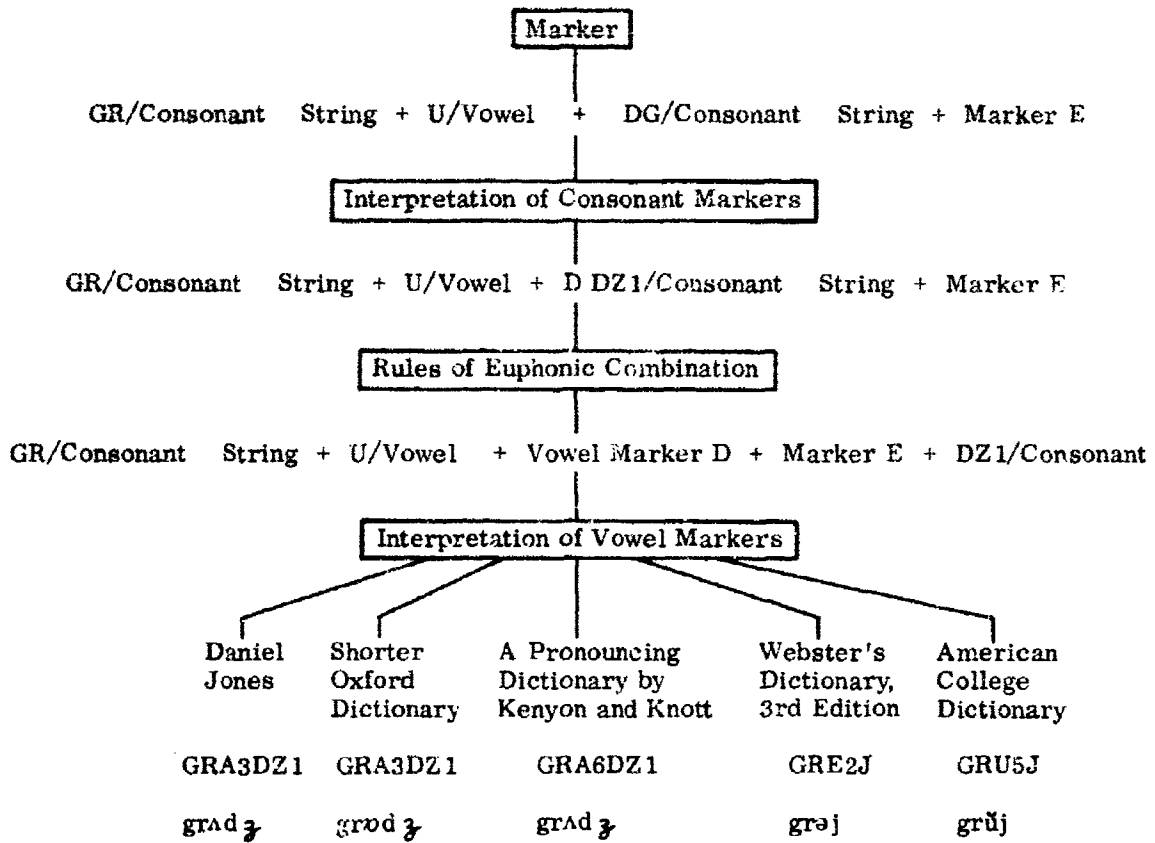


Fig. 3-5 Operation of the Program for the Word GRUDGE

Table 3-1  
 PHONETIC MAPPING OF ORTHOGRAPHIC VOWEL STRING U ACCORDING  
 TO TRANSCRIPTIONS BY KENYON AND KNOTT (Ref. 3-9)

	<u>Orthographic</u>	<u>Phonetic</u>	<u>Dialect Code</u>
If U is precluded by J or R or L, and followed by R $\bar{E}$ :	UR $\bar{E}$	UR	101
		IUr	102
		Uə(r)	103
		IUə(r)	104
Otherwise, U followed by R $\bar{E}$ :	UR $\bar{E}$	JUr	101
		IUr	102
		jUə(r)	103
		IUə(r)	104
If U is preceded by J or R or L and followed by a single consonant and $\bar{E}$ :	U Cons $\bar{E}$	{ uCons }	101
		{ IuCons }	102
Otherwise, U followed by a single consonant and $\bar{E}$ :	U Cons	{ juCons }	101
		{ IuCons }	102
Otherwise, U followed by RCons:	U RCons	3 <sup>^</sup> Cons	101
		3Cons	103
Otherwise, U followed by R:	UR	3 <sup>^</sup>	101
		3(R)	103
Otherwise, U:	U	Λ	101

The consonant string GH immediately following a vowel may seem like a marker of "long vowel sounds," but such a statement is less than 63 percent accurate; a much higher accuracy is required of the computable relations. The accuracy of our rules as well as that of the phonetic data studied is evaluated in the next section.



### 3.3 COMPARISON OF THE PHONETIC REPRESENTATIONS OF ELEMENTARY WORDS IN FIVE DICTIONARIES

Comparisons have been made of the phonetic forms of elementary words in the five dictionaries by using a set of four computer programs to test the data and the algorithm.

These programs are discussed under the following headings:

- A. Comparison of the phonetic data by using the SOX as the reference set
- B. Comparison of the phonetic data by organizing it into marker-orthographic vowel-marker sets
- C. Comparison of phonetic data organized in the terminal rhyme order of words
- D. Compilations of homonyms according to the transcriptions in each of the dictionaries studied

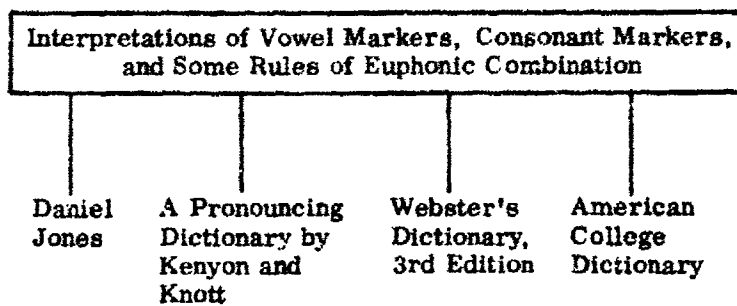
The principal objectives of such a comparison are:

- (1) To evaluate the accuracy of symbol-for-symbol comparison of phonetic data
- (2) To establish a computable set of relations between the various phonetic representations of each word
- (3) To define the smallest set of symbols over which the relations in (2) can be established
- (4) To evaluate the phonetic representations in the five dictionaries
- (5) To assign a measurement criterion to the confusion that results from the differences among the phonetic representations
- (6) To assure high accuracy for the algorithm and for the phonetic data used

#### 3.3.1 Comparison of Phonetic Data by Using the SOX as the Reference Set

The phonetic representations of elementary words, as in the SOX, were mapped into corresponding ones in the other four dictionaries (Refs. 3-11 through 3-14) for this program (Fig. 3-6). Such a mapping, while accurate [about 95 percent for data in three of the dictionaries (Refs. 3-11 through 3-13)] and indicative of the computable

**PRONUNCIATION OF KERNEL WORDS AS IN  
THE SHORTER OXFORD DICTIONARY**



**Fig. 3-6 Comparison of Phonetic Data by Using the  
Phonetics in the Shorter Oxford Dictionary  
as the Reference Set**

relations between the phonetics in the SOX and those in the other dictionaries, does not establish a complete set of simultaneously computable relations between the phonetic representations of each word in all five dictionaries. Problems also arise when a word happens to have a homonym according to the SOX but not according to one or more of the other dictionaries (Ref. 3-23). Similar limitations persist in mapping the phonetic forms in one dictionary (other than the SOX) into corresponding ones in the other four.

### 3.3.2 Comparison of the Phonetic Data by Organizing It Into Marker-Orthographic Vowel-Marker Sets

To compensate for the shortcomings of the test just described, the data were organized according to the marker-orthographic vowel-marker sets as illustrated in Fig. 3-7. (These were the groups of symbols that provided mapping from the orthographic into the phonetic forms depicted in Fig. 3-3.) Column 8 in Fig. 3-7 contains the orthographic forms of words, and the corresponding phonetic forms are presented in the same row of the matrix but arranged so that the data for each dictionary occur in a specified column, as indicated at the top of the figure. The columns numbered 2, 4, 10, 13, and 15 contain either a blank code or an asterisk; the former indicates that the phonetic representation was manually checked with that in the specified dictionary, and the latter indicates that the phonetics were calculated by the computable relations but could not be checked because the dictionary did not provide a pronunciation for the word. The orthographic form of any word is repeated in column 8 as often as distinct phonetic entries are given for it in any dictionary; rows starting with a specific word (as numbered next to the orthographic form in column 8) are arranged according to the numeric order assigned to the dialect pattern as specified in columns 12 and 15. (Notice that the row number in column 8 may differ from the dialect code, as seen for row number five of LARE in Fig. 3-8 which has a dialect code number 106.) Some dictionaries (Refs. 3-13 and 3-14) present data on a larger number of dialect patterns than do the others studied (Refs. 3-10 through 3-12). In such cases, the entries for specific dictionary columns are left blank when there is no phonetic entry corresponding to the specific row identified with a particular

Column No	British Dialects				American Dialects										
	1	2	3	4	5 6 7		8	9	10	11	12	13	14	15	16
	SOX Phonetics	Code	Jones Phonetics	Code	Status Code	Graphic Word									
	rɛʃ		rɛʃ			RARE 01	rɛʃ			rɛʃ	101		rɛʃ	101	
						RARE 02				rɛʃ	102		rɛʃ	102	
						RARE 03				rɛʃ	103		rɛʃ	103	
						RARE 04				rɛʃ	104		rɛʃ	104	
	skɛr		skɛr			SCARE 01	skɛr			skɛr	101		skɛr	101	
						SCARE 02				skɛr	102		skɛr	102	
						SCARE 03				skɛr	103		skɛr	103	
						SCARE 04				skɛr	104		skɛr	104	

Explanation of Dialect Codes

- 101 - The prevalent form as identified most frequently by the first phonetic entry in the dictionary (except R-droppers and L-droppers).
- 102 - The first variant as identified most frequently by the second phonetic entry in the dictionary (e.g., R-droppers and L-droppers).
- 103 - Eastern dialects in Keynon and Knott; R-droppers in Webster's 3rd.
- 104 - East-south entries in Keynon and Knott; dialects where V is substituted for R (in words such as BIRD) in Webster's 3rd.
- 105 - L-droppers in Webster's 3rd.

- 106 - Second variant pronunciation.
- 107 - Third variant pronunciation.

- 208 - Vowels consonants replaces the vowels in the dialect pattern, otherwise specified by the number for N in the 10N series above.
- 307 - Extreme dialects such as in Webster's 3rd where ee in the pronunciation indicated for the initial orthographic component string "ER" in dialect patterns, otherwise specified by the number for N in the 10N series above.

Fig. 3-7 Data Organized According to the Marker-Orthographic Vowel-Marker Set and According to Phonetic Representations in Specified Dictionaries (I)

Column No	British Dialects				American Dialects									
	1	2	3	4	6	8	9	10	11	12	13	14	15	16
	SOX Phonetics	Code	Jones Phonetics	Code										
	gɛr		gɛr	*	GARE 01	5X	gɛr	*	gɛr	101	*	gɛr	101	
	(gɛr)		(gɛr)	*	GARE 02	5X	lɛr	*	gɛr	103	*	lɛr	101	
					LARE 01	1W+	lɛr		lɛr	101	*	lɛr	101	
					LARE 02	1W+	lɛr		lɛr	102	*	lɛr	102	
					LARE 03	1W+	lɛr		lɛr	103	*	lɛr	103	
					LARE 04	1W+	lɛr		lɛr	104	*	lɛr	104	
					LARE 05	1W+	lɛr		lɛr	106	*	lɛr	106	

For explanation of dialect codes, see Figure 7.

Fig. 3-8 Data Organized According to the Marker-Orthographic Vowel-Marker Set and According to Phonetic Representations in Specified Dictionaries (II)

dialect pattern number. {Dialect code 101 identified the prevalent pronunciation in the United States for the American dictionaries (Refs. 3-12 through 3-14), while it represents British pronunciation patterns in the SOX and in Jones (Ref. 3-11)}.

Such an arrangement of phonetic data in rows and columns creates a pattern of phonetic vowels and consonants, orthographic vowels and consonants, and blank spaces for each of the words (Fig. 3-7). Consistent portions of such patterns for a word can be compared with those of another word in the same marker-orthographic vowel-marker group (Figs. 3-7 and 3-8). When such patterns agree with a majority of the word entries in the class of words being compared, these words are considered algorithmic and the above-mentioned pattern for such words is thus defined. Whenever the pattern of entries for a word differs from the algorithmic one, each of the graphic entries for this word is tagged with a code in columns numbered 5, 6, and 7 as illustrated in Fig. 3-8; algorithmic words do not receive such codes. The code 1W+ adjoining the word LARE, for example, indicates that all the pronunciations computed by the algorithm agree with those in the source dictionaries, but Webster's (Ref. 3-14) provides pronunciation(s) for this word in addition to the computed set corresponding to it. The code 5X adjoining the word GARE indicates that all five of these dictionaries contain phonetic forms that do not agree with those provided by the algorithm. By use of such a data comparison program, the algorithm is made to represent the phonetic data as accurately as possible. The principal changes in the algorithm and in the phonetic data resulted from the modification in the interpretation of the marking system such that the pronunciation of words with \* codes (indicating that the computed pronunciation could not be checked in the dictionary) in any one of the dictionary entries could be made to correspond to those in one or more of the other four which carries a blank code (indicating that the phonetics were actually checked) according to computable rules developed for mapping phonetics in one dictionary into that for another, thus minimizing the dependence of the algorithm on any one dictionary datum without excluding or passing judgment on actual data provided by any of the authorities.

Statistics on the codes introduced (columns 5, 6, and 7 in Fig. 3-8) by comparison of the phonetic forms, after making the above-mentioned corrections, are presented in Figs. 3-9 and 3-10. Note that the disagreement among the phonetic representations is not limited to the vowel sounds but that several consonants are also represented differently by different dictionaries. Some of these differences are indeed minor ones; e.g., Kenyon and Knott (Ref. 3-13) map the orthographic O in words such as GNOFF and KOFF as o , ɔ , ɒ , respectively, whereas the same orthographic letter in DOFF is represented as ɔ , ɒ , ɔ without a specific indication of ɒ representing the "eastern pronunciation." Another such situation is the omission of dialectal representations marked with - codes (e.g., 1W-) in column 7 of data organized according to the illustration in Fig. 3-8 and grouped accordingly in Fig. 3-9. For example, Webster's provides only "R Dropper" pronunciation for the word KERB, a word commonly used in Great Britain, but no entries are provided for the "General American Dialect(s)." (This word is among those counted as 1W- in Fig. 3-9.) Moreover, dictionaries provide occasional additional pronunciations (see discussion of LARE in Fig. 3-8) and these are grouped 1+, 2+, etc., in Figs. 3-9 and 3-10. Some words are in error because the algorithm for separating elementary words passes occasional polysyllabic ones, such as CAFE; these are identified by a P code in column 7 of the organized data (as in Fig. 3-8). Such errors should not be strictly counted against a general-purpose algorithm, and it is debatable whether the algorithm should match singular cases where only one dictionary disagrees with the computed phonetic representations. Such situations account for over 65 percent of the errors, as summarized in Figs. 3-9 and 3-10. (Excluding such singularities, the algorithm provides the phonetic transcriptions of all the words according to the corresponding representations in each of the five dictionaries for every dialect recorded therein with over 90-percent accuracy. Mapping the orthographic forms into any single specified transcription shows higher accuracy.)

### 3.3.3 Comparison of Phonetic Data Organized in the Terminal Rhyme Order of Words

The phonetic data were organized in the terminal rhyme order in a selected dictionary, as illustrated in Fig. 3-11 for the Jones dictionary (Ref. 3-11). Another portion of these

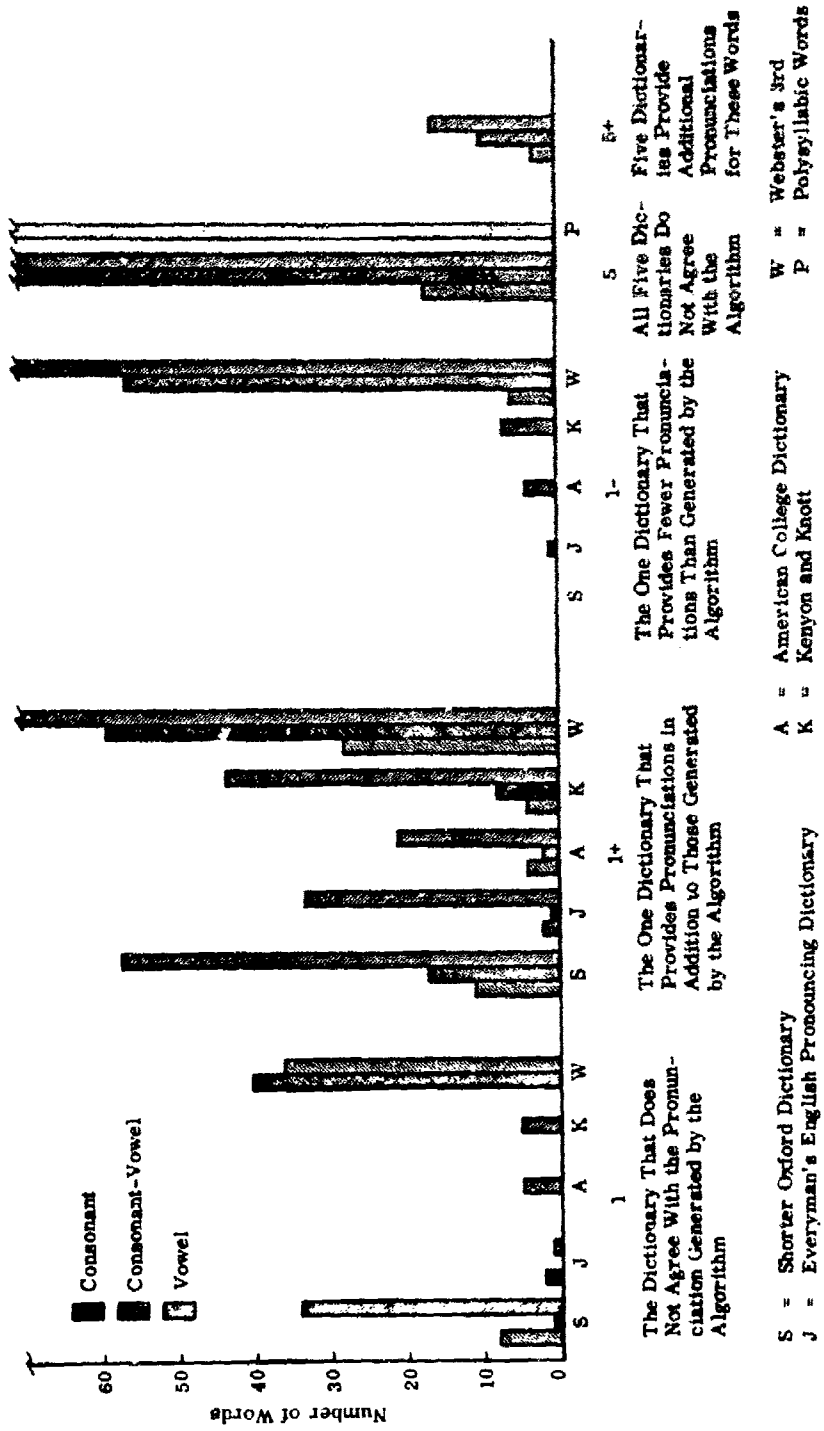
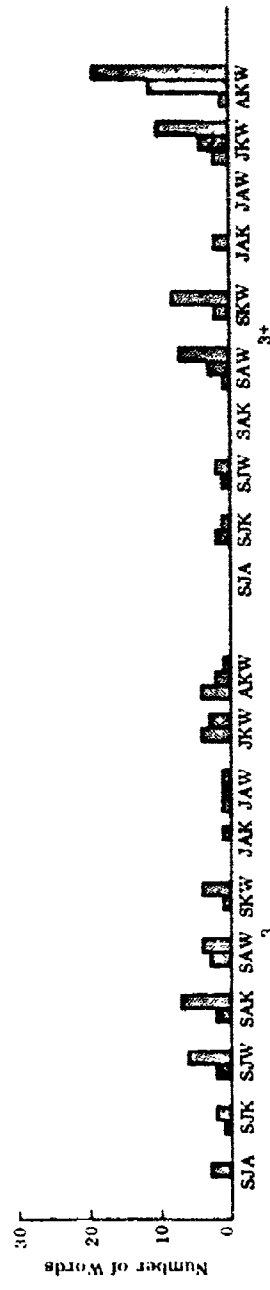


Fig. 3-9 Statistics on Disagreement Among Dictionaries for Phonetics of Elementary Words (I)

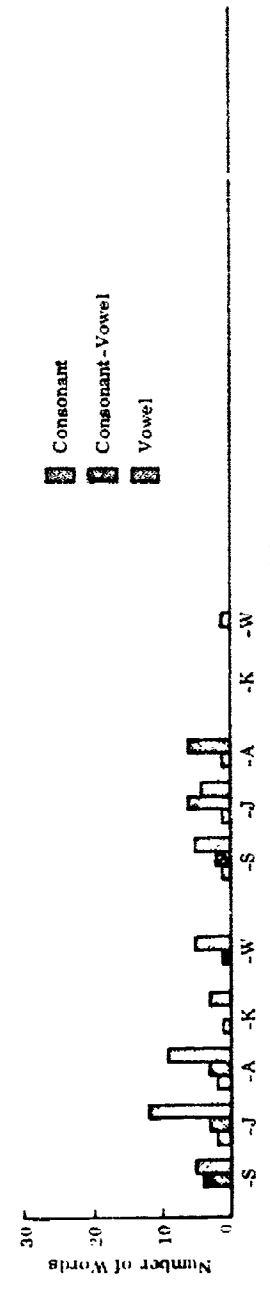




The Two Dictionaries That Disagree With the Algorithm



The Three Dictionaries That Disagree With the Algorithm



The Four Dictionaries (the Set of Five Minus the Identified Dictionary) That Disagree With the Algorithm

Consonant  
Consonant - Vowel  
Vowel

Fig. 3-10 Statistics on Disagreement Among Dictionaries for Phonetics of Elementary Words (II)

Column No.	British Dialects				5 6 7 Status Code	8 Graphic Word	American Dialects							
	1	2	3	4			9	10	11	12	13	14	15	16
	SOX Phonetics	Code	Jones Phonetics	Code			ACD Phonetics	Code	Kenyon and Knott Phonetics	Dialect Code	Code	Webster's 3rd Phonetics	Dialect Code	Code
	bə:d		bə:d		BARD 01	bə:d	101	bə:d	101	bə:d	101	bə:d	101	
	kə:d		kə:d		BARD 02		103		103		103		103	
					CARD 03	kə:d		kə:d	101	kə:d	101	kə:d	101	
					CARD 04	kə:d		kə:d	103	kə:d	103	kə:d	103	

Fig. 3-11 Data Organized According to Terminal Rhyme Words According to Jones (I)

mask		mo:sk		MASK 01	mæsk	101	mæsk	101	mæsk	101	mæsk	101
				MASK 02	mask	102	mask	102	ma(ə)sk	102	ma(ə)sk	102
				MASK 03					meisk	106	meisk	106
				MASK 04					məsk	107	məsk	107
task		to:sk		TASK 01	tæsk	101	tæsk	101	tæsk	101	tæsk	101
				TASK 02	task	102	task	102	to(ə)sk	102	to(ə)sk	102
				TASK 03					toisk	106	toisk	106
				TASK 04					təsk	107	təsk	107

For explanation of dialect codes, see Figure 7.

Fig. 3-12 Data Organized According to Terminal Rhyme Words According to Jones (II)

data from the same phonetic source and for the same phonetic vowel is presented in Fig. 3-12. Notice the differences in the mapping of the vowel in the four other dictionaries in these two figures. Such a test was repeated for each phonetic vowel data entry in each of the five dictionaries. The differences in the patterns of entries point out clearly the inaccuracy inherent in establishing a computable vowel-for-vowel relation between different dictionary entries. The inaccuracy in establishing a consonant-for-consonant relation is obvious from the discussions of Figs. 3-9 and 3-10. It can be stated that establishment of a precise relation between the phonetic forms given in five dictionaries requires the adoption of a set of symbols that is larger than a consonant string alone or a vowel string by itself (and it is not the set of syllabic nuclei as defined in recent literature (Refs. 3-5 through 3-8)). Further study of the data indicates that such a relation can be established on a marker-phonetic vowel-marker string (as was the case for mapping the orthographic forms into phonetic representations). The group of phonetic symbols over which a computable relation can be established accurately is like the segments described previously on the basis of phonetic and acoustic studies of speech (Ref. 3-9). Such an agreement may indicate that the perceived, and possibly the produced, characteristics of vowel sounds are determined by the adjoining markers in a predictable - and indeed a mechanically calculable - manner. This goes beyond results such as the concept of allophones (Ref. 3-14) or of cues between consonants and vowels (Refs. 3-3 and 3-4), by defining the precise conditions for the different speech sound representations of any word.

#### 3.3.4 Compilation of Homonyms According to the Transcriptions in Each of the Dictionaries Studied

In spite of the observed differences in the phonetic representations of elementary words, there remains the crucial test of whether they cause confusion in voice communication. Evaluation of this aspect calls for compilation of homonyms in these dictionaries (Ref. 3-23). A homonym set was defined as a set of different orthographic forms having an identical phonetic transcription according to any one dictionary of the five studied (Refs. 3-10 through 3-14). Any member of a homonym set was

called a homonym. Since each dictionary uses its own phonetic symbols, homonyms were compiled separately for each dictionary and all the elementary words were considered for this purpose. Pronunciations for words not included in a specific dictionary were generated by algorithm and checked carefully for correspondence with phonetic data for that particular word in other dictionaries. (See discussion of Figs. 3-7 and 3-8.) These homonym entries were identified with respect to dictionary and dialect patterns and the words displayed in alphabetical order (Fig. 3-13). It is readily evident that certain words are homonyms in some dictionaries but not in all of them, whereas others are homonyms in all dictionaries.

A statistical summary of the homonyms, among the 5757 elementary words, is presented in Fig. 3-14. Notice the large number of homonyms and the marked disagreement among the different lexicographers. Some normalization was performed on the data by removing dialectal pronunciations (the eastern and southern dialect patterns in the U.S., codes 103, 104, 105; all extreme dialects, code 200+) from the transcriptions as in Webster's and those in Kenyon and Knott. In spite of such normalization, there remains a marked disagreement among these transcriptions. (It may be observed that variations in dialects do not add significantly to the homonyms according to Kenyon and Knott.) Information was also compiled on the number of sets of two-word homonyms, three-word homonyms, etc., presented in Fig. 3-15. (Notice the 10 to 1 scale change between sets of three and sets of four.) Table 3-2 provides additional statistical information on the extent of disagreement among dictionary transcriptions, confirming the need for more precise representation of speech, as mentioned in the introduction.

One could object at this point that the dictionaries contain a large number of obsolete words and ones that are very infrequently used in everyday conversation. However, the decision to limit the set of words requires very careful consideration since the purpose of this study is a general-purpose system for mechanized readers, etc. Moreover the marker-orthographic vowel-marker groups studied here, while perhaps infrequent in elementary word sets, are found much more often as constituent parts

1	SOX			Jones			ACD			Kenyon & Knott			Webster's 3rd												
	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	
GRICE	02																								
GRICE	02																								
GRICE	02																								
GRID	01	B																							
GRIEVE	01	B	GREAVE	01	B	GREAVE	01	B	GREAVE	01	B	GREAVE	01	B	GREAVE	01	101	B	GREAVE	01	101	B			
GRIEVE	01	B	GREEVE	01	*	GREEVE	01	*	GREEVE	01	*	GREEVE	01	*	GREEVE	01	101	*	GREEVE	01	101	*			
GRIFF	01	X	GRIFFE	01	*W	GRIFFE	01	W	GRIFFE	01	W	GRIFFE	01	W	GRIFFE	01	101	*W	GRIFFE	01	101	*W			
GRIFF	01	X	GRIPH	01	*	GRIPH	01	*	GRIPH	01	*	GRIPH	01	*	GRIPH	01	101	*	GRIPH	01	101	*			
GRIFF	01	X	GRYPH	01	*	GRYPH	01	*	GRYPH	01	*	GRYPH	01	*	GRYPH	01	101	*	GRYPH	01	101	*			

This section of computer printout is approximately one-third of a page of a total of 386 pages.

Explanation of column numbers:

- 1 Orthographic form of a homonym, listed alphabetically (presented as in column 8 of Figures 7 and 8). This entry is repeated as often as there are homonyms in any dictionary corresponding to the phonetic representation set identified by the number in row 2.
- 2 The row number of the orthographic entry (presented as in column 8 of Figures 7 and 8) that identifies the corresponding phonetic representations of the homonym in each of the five dictionaries.
- 3 Status of the word in column 1 with respect to standard meaning as in Refs. 26 and 27.
- 4, 8, 12, 16, 21 Orthographic forms of additional homonyms in the specified dictionary and in the set containing the entry in column 1 as specified in columns 2 and 3.
- 5, 9, 13, 17, 22 The row number of the orthographic entry (presented as in Figures 7 and 8) identifying the phonetic representation in the specified dictionary.
- 6, 10, 14, 18, 24 Indicates whether the phonetic representation could actually be checked with the data in the specified dictionary (as explained in Figure 7).
- 7, 11, 15, 20, 25 The status of each entry coded for standard meaning as in Refs. 26 and 27.
- 18, 23 The dialect code associated with the phonetic representation (presented as in Figures 7 and 8).

Fig. 3-13 Alphabetically Arranged List of Homonyms

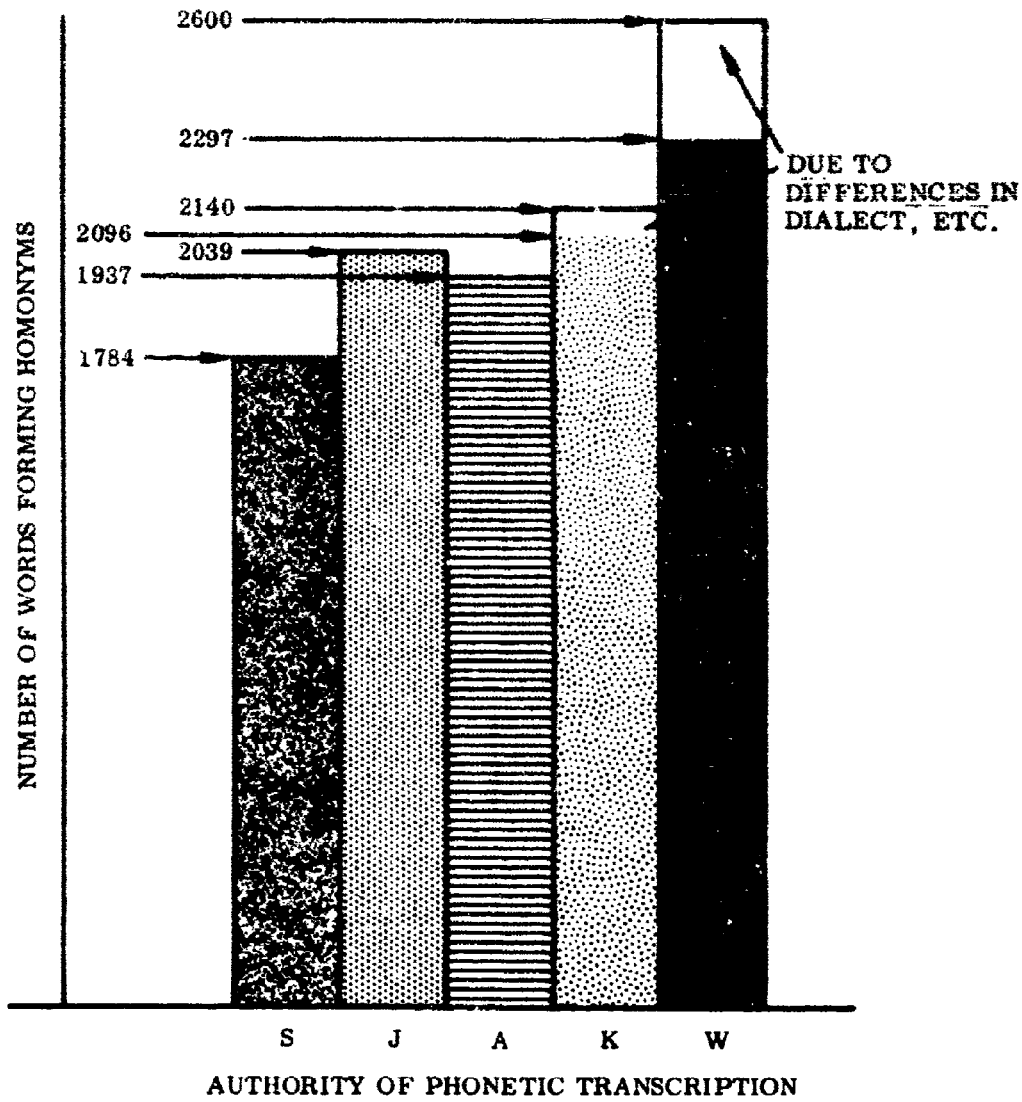


Fig. 3-14 Number of Words Forming Homonyms in Each of the Five Dictionaries

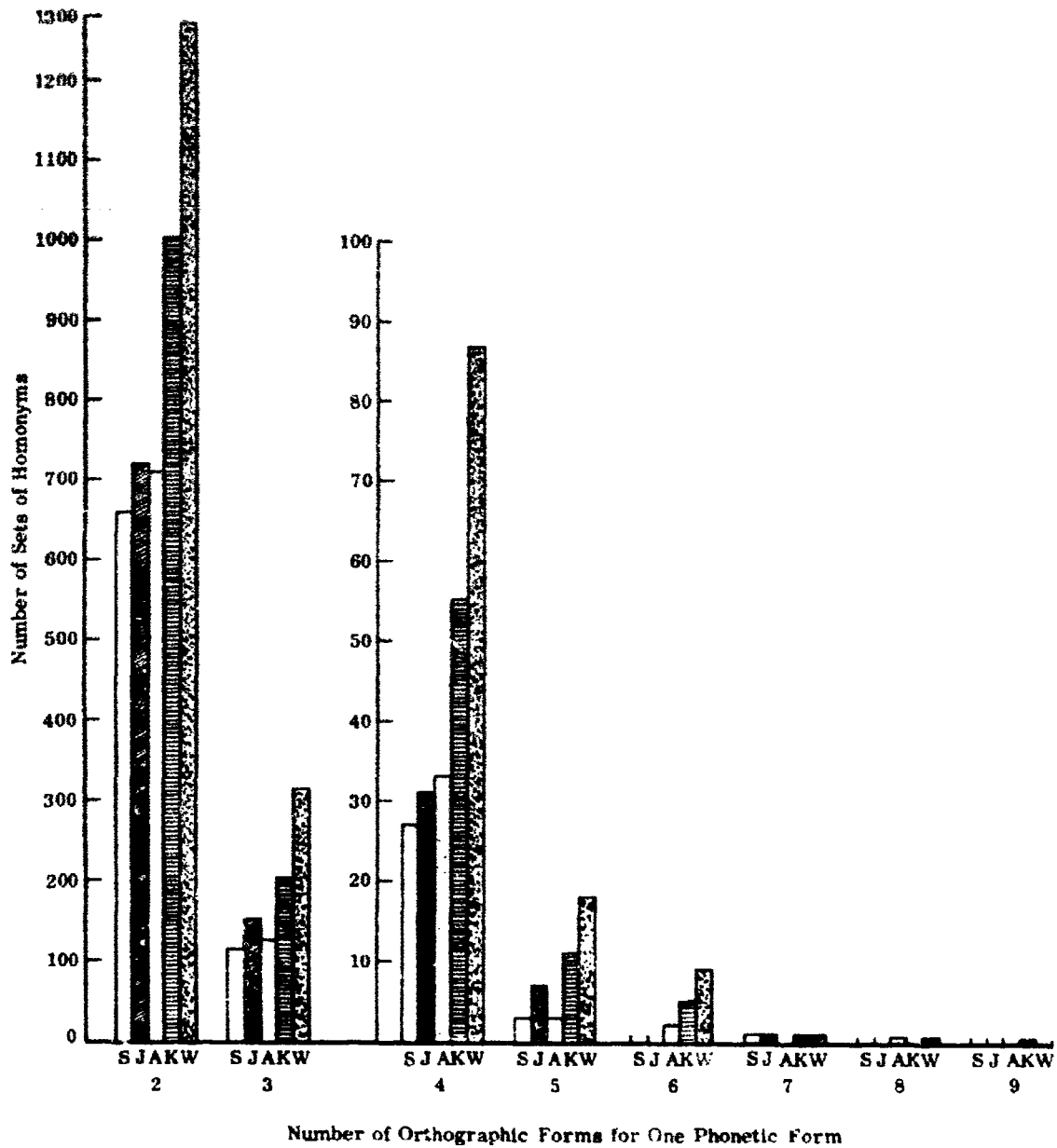


Fig. 3-15 Graphic Representation of the Number of Homonym Sets Among the Elementary Words in Five Dictionaries

Table 3-2

STATISTICAL SUMMARY OF WORDS INVOLVED IN HOMONYM SETS,  
SHOWING THE EFFECT OF DIALECT REMOVAL

Set Description	No. of Words in the Set	
	With Dialects	Without Dialects
Total Set -		
Words forming a homonym according to at least one dictionary	2966	2714
Words forming a homonym according to exactly one dictionary	746	535
Words forming a homonym according to exactly two dictionaries	236	214
Words forming a homonym according to exactly three dictionaries	189	184
Words forming a homonym according to exactly four dictionaries	290	297
Words forming a homonym according to all dictionaries	1505	1484

of polysyllabic words (Fig. 3-2) and none of these groups can be overlooked. For example, the word MOIRE is not included in the Thorndike-Lorge list (Ref. 3-25). However, the rules for mapping it are identical to those required for the second syllable of MEMOIR and the latter words appears in the referenced list (Ref. 3-25).

To evaluate the effect of the rarity of use of certain words, we incorporated into our phonetic dictionaries information about the status of each word and its parts of speech as in the SOX and in Webster's. These entries are explained in detail elsewhere (Refs. 3-26 and 3-27). With the inclusion of such information, it becomes possible to limit the word list to those words that have a standard meaning according to both SOX and Webster's, thus defining words in current usage for this discussion. The statistics on sets of homonyms compiled from the word set so defined (Fig. 3-16)



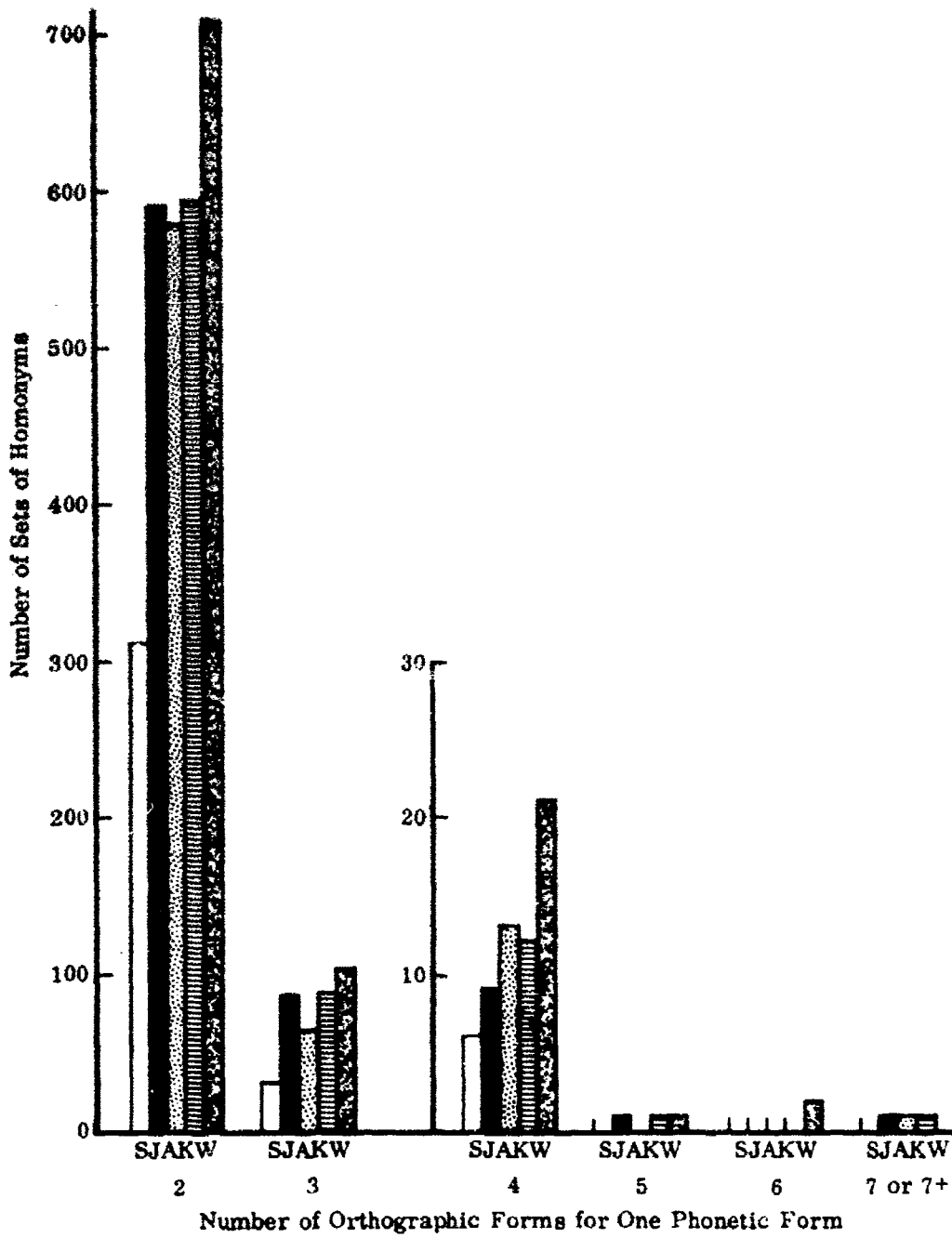


Fig. 3-16 Graphic Representation of the Number of Homonym Sets Among the Double Standard Elementary Words in Five Dictionaries

were found to some extent to resemble the distribution presented for the entire set (Fig. 3-15), but their total number was reduced approximately in proportion to the reduction in the number of words.

### 3.3.5 Comments

The results of homonym studies along with those of the second (subsection 3.3.2) and third (subsection 3.3.3) tests indicate the limitations of the symbolic representations of speech in the five dictionaries (Refs. 3-10 through 3-14). These tests also indicate the extent to which confusion exists in associating the orthographic and the phonetic forms of words. Tests one (subsection 3.3.1) and two (subsection 3.3.2) do establish the high accuracy of the computable relations for mapping the phonetics in one dictionary into that in another (Fig. 3-6), as well as of those that map the orthographic forms of words into corresponding phonetic representations (Fig. 3-3). Both of these sets of relations use the groups of symbols called marker-vowel-marker. The first test demonstrates the inadequacy of the symbol-for-symbol comparison of phonetic data, and the situation is confirmed by the results of the third test (subsection 3.3.3).

The phonetic data studied for this paper may seem to indicate that dictionary makers have not provided as precise a set of phonetic representations as desirable. However, such a conclusion is not warranted in light of five aspects discussed previously and recalled next along with additional test data.

- (1) An algorithm capable of computing phonetic forms of words from orthographic input has been demonstrated (Fig. 3-3). This is true for each of the five dictionaries studied.
- (2) An algorithm has also demonstrated the computability of the phonetic representations in four dictionaries (Refs. 3-11 through 3-14) given the phonetics in the SOX as input (Fig. 3-6).
- (3) A pattern of pronunciations can be established for each of the marker-vowel-marker groups of symbols (Figs. 3-7 through 3-10).

- (4) Although it is not possible to establish a symbol-for-symbol transformation between the phonetic representations in one dictionary and those in the remaining four (Figs. 3-11 and 3-12), such a symbolic representation would require a phonetician's decision at the time of preparation of phonetic data (Refs. 3-2 and 3-13). Although a more uniform representation may be desired, it may be the mechanism of speech perception, as well as that of speech production, that leads to an order modification of speech sounds (Ref. 3-9). (The interpretation of the linguistic aspects of this study, such as the "grapheme to phoneme" relations, the definitions of "phonemes," or the "sound changes" in the language, are beyond the scope of this paper.) It is doubtful, on the basis of data presented here, that a uniform symbolic notation can be provided for even the most well-defined dialects of English. The transcriptions in SOX and in the Jones dictionary result in a 25-percent difference in the number of homonyms. (About 10 percent of the words have different phonetic transcriptions.) Both of these represent speech patterns in the land of Professor Higgins (Ref. 3-28).
- (5) The group of symbols over which one can establish an accurate algorithmic relation for mapping the orthographic form of words into their corresponding phonetic representations are similar to those required for establishing computable relations between the phonetic data. These groups of symbols also agree with ones described on the basis of studies of the mechanism of speech production and the acoustic characteristics of speech (Ref. 3-9).

An additional check on the accuracy of the phonetic data and of the computable relations was to compile homonyms from only those words that were algorithmic (those with blank codes in columns 5, 6, and 7 of Fig. 3-8). These data (Fig. 3-17) show a marked agreement among the transcriptions studied (compare with Fig. 3-15).

Based on the preceding discussion, one can aim to provide a set of computable relations that can map the orthographic form of words into corresponding phonetic representations as specified by one or more of the authorities, and the resulting phonetic

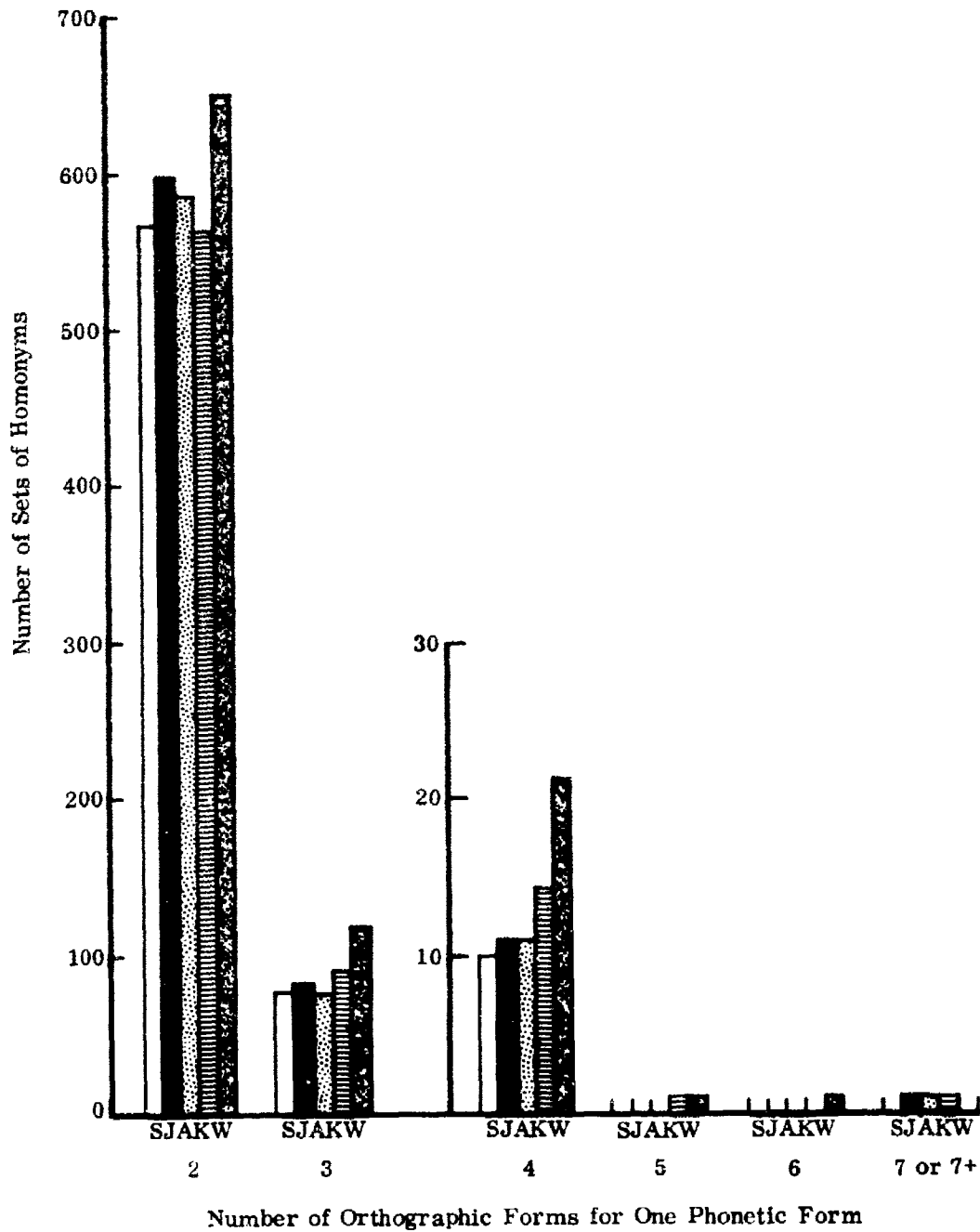


Fig. 3-17 Graphic Representation of the Number of Homonym Sets Among the Algorithmic and Elementary Words in Five Dictionaries

data can be expected to be highly accurate and useful for the operation of a mechanized reader. The high accuracy of the computations can be enhanced by incorporating certain exception words. This subject is discussed next.

### 3.4 SELECTION OF EXCEPTION WORDS

An exception word is defined as one for which the phonetic representations as computed by the rules do not agree precisely with those in each of the dictionaries (Refs. 3-10 through 3-14) for phonetic data on all the dialect and speech patterns of the language. When the algorithm is limited to providing representations in a specific phonetic transcription, an exception word is one for which the computed representations differ from the corresponding phonetic data in the authoritative source.

A set of exception words must be added to the program such that their phonetic representations can be obtained by look-up rather than by computation. Such a requirement increases the storage required and decreases the speed of operation of the program; hence, words to be included in the set of exception words must be selected with care.

It is indeed mandatory to include words such as ARE for which the algorithm is in error according to each of the five dictionaries and which is indeed an important structure word of the language. The decision may be different in the case of words that are obsolete or of little interest (but care must be taken in this decision as discussed for the word MOIRE, which will be considered again later). Moreover, it is necessary to decide on the pronunciation of words for which there is disagreement among the transcriptions (Figs. 3-8 through 3-10). The coding of words in our phonetic dictionary (columns 5, 6, and 7 in Fig. 3-8) and the incorporation of data on the dictionary status of words as well as their parts of speech (Refs. 3-26 and 3-27) makes it possible to specify the types of words one wishes to include in the class of exception words and to check the words for which questions arise, thus making the program flexible for the definition of phonetics of elementary words. Six generalized classifications of exception words are described next.

The first group contains words such as ARE for which the transcriptions in all five dictionaries do not agree with the computed values, and the method for making decisions about these is mentioned in the above paragraph.

The second group of words is that for which only one or two dictionaries disagree with the computed values, as discussed in subsection 3.3 for the words GNOFF, KOFF, DOFF, and KERB as well as those classified as 1, 1+, or 1- in Fig. 3-9. The inclusion of some such words seems justified only in special cases, and most of these represent minor problems as mentioned previously.

The third group of exception words represents inaccuracies in the algorithm that separates the elementary words from the total dictionary, resulting in polysyllabic forms. Some commonly used words (e.g., CAFE) merit inclusion in this category, but with the stipulation that they will not be effective when the algorithm operates on individual syllables of polysyllabic words. The obsolete words in this group may be ignored for most applications.

The fourth group contains words that require more than one rule for their mapping. These are not really exception words, but they may improve the efficiency of the algorithm. An example of such words as MOIRE, which calls for inclusion of a w sound between the first consonant and the vowel string. Such occurrences tend to repeat in polysyllables, as evident in the pronunciation of MEMOIR.

The fifth group contains words that have multiple phonetic representations for the same orthographic form, such as HOUSE and BOW. Such cases may be resolved only in context or by use of parts of speech information, a subject which is beyond the scope of this paper.

The sixth group contains singular words that do not fit a general pronunciation pattern, such as the differences in the pronunciation of LASS and BRASS in the SOX. Some of these words are singular in one or more dictionaries (as presented in Figs. 3-9 and 3-10), and some such as CHEF and CHERE retain their foreign

pronunciation. The decision on inclusion of these is a matter of preference by the computer programmer or the specifying linguist or by both of these parties.

With the use of codes (columns 5, 6, and 7 of Fig. 3-8) and the information on the status of words (Refs. 3-26 and 3-27) the precise criteria for the program can be specified. The data in Figs. 3-9 and 3-10 also aid in evaluation of the accuracy of a resulting program when the criteria have been decided upon. Should one wish to work with phonetic data that are different from the set studied, it may be possible to focus attention on words that are known to have disagreements with the present algorithm, hence making it possible for the programmer to ask the consulting linguist about phonetic representations of a specific set of words. With such information and the coded phonetic dictionary as aids, programs can be written to meet the exacting requirements of many groups without sacrificing unduly the operating efficiency of a program by inclusion of an unnecessarily large number of exception words, or the accuracy of the program by not incorporating the correct and complete set of rules and the requisite set of exception words.

### 3.5 EXTENSION OF THE PROGRAM TO INCLUDE MONOSYLLABLES

The English monosyllables contain all the elementary words of the language; most of their plural forms (obtained by the addition of ES) and past tense or adjectival forms (as obtained by the addition of ED); words ending in CUE, GUE, QUE, GNE, DME, etc., which have characteristics that are similar to those of elementary words; common words such as HE, SHE, ME, etc., grouped into a zero vowel category because the final E is tagged as a marker (Fig. 3-3); the letters of the alphabet by themselves; and words such as TABLE, BOTTLE, and their plural and adjectival forms, etc. The last of these sets of words is polysyllabic according to most dictionaries, but these words contain a single interrupted phonetic vowel string according to the SOX.

For extension of the algorithm to include the above-mentioned words, it was necessary to implement a portion of the program illustrated in Fig. 3-2 to separate the appropriate ES and ED endings, etc. (Fig. 3-18). The E in these endings was treated as a marker ( $\bar{E}$ ), and the algorithm for the elementary words was used for computing the phonetic forms of words. The words ending in CUE, GUE, etc., required a minor modification of the algorithm; namely, the U had to be treated as a marker that blocked the operation of  $\bar{E}$ . Words such as HE, ME, and SHE were mapped by first adding  $\bar{E}$  to their orthographic form and then computing the phonetic representations. The letters of the alphabet required a special set of rules. Most of the vowels (A, E, I, O, U) could be mapped by adding an  $\bar{E}$  to their orthographic form; consonants such as Z, B, C, and D required the addition of  $E\bar{E}$  to their orthographic form; some consonants such as F and H required (E consonant) representation in their orthographic form; and so forth. The pronunciation of words such as TABLE and BOTTLE was computed by permitting the  $\bar{E}$  to operate as a vowel marker across a (single consonant + L) group only.

By incorporating modifications such as those mentioned above, the revised algorithm was used for computing the pronunciations of English monosyllables. The second test, described in subsection 3.3.2, was applied to the completed data for just those monosyllables that were not in the set of elementary words. The results of these tests are summarized in Figs. 3-19 through 3-22. Notice that the test results are grouped for monosyllables (Figs. 3-19 and 3-20) which can be compared with data in Figs. 3-9 and 3-10; for occurrences of polysyllabic pronunciations for words ending in ES (Fig. 3-21); and for words ending in ED (Fig. 3-22). Notice also the lack of agreement among the lexicographers about the polysyllabic pronunciation of words ending in ES and ED, which makes it difficult to specify a simple algorithm. However, this subject concerns polysyllabic words and hence is beyond the scope of this paper.

Comparing the data in Figs. 3-19 and 3-20 with the data in Figs. 3-9 and 3-10, and considering that the number of words in the former case is smaller than that in the latter, one can conclude that the algorithm is just as accurate on the added set of words as it was for the elementary words (for which computable relations were defined).



INPUT NOT AN EXCEPTION -  
 SPLIT GRAPHIC  
 $C_1 V_1 C_2 V_2 C_3$

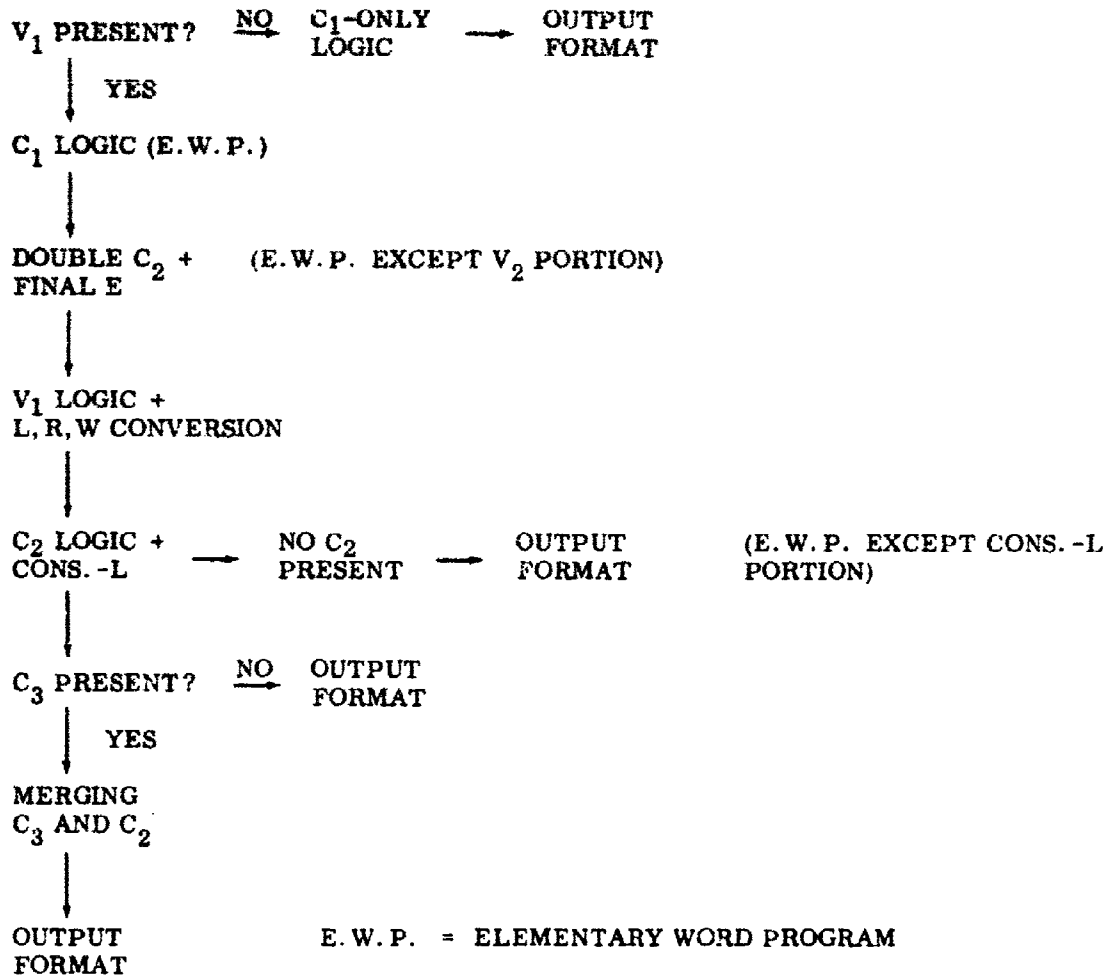


Fig. 3-18 Modification of the Program To Include the English Monosyllables

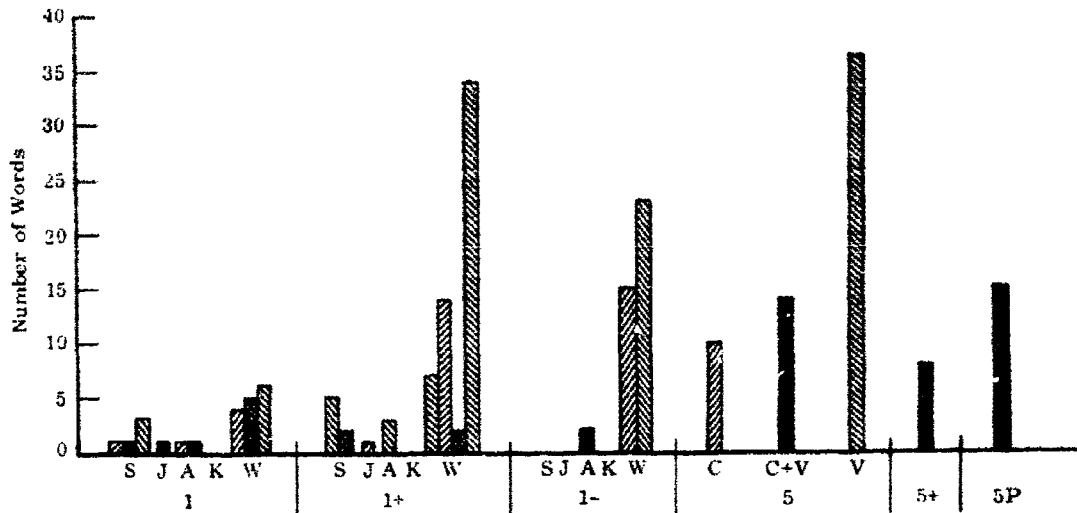


Fig. 3-19 Statistics on Disagreement Among Dictionaries for Phonetics of Additional Monosyllables (Except Those Ending in ES or ED Which Have Polysyllabic Pronunciations). Total of 1147 words represented (I)

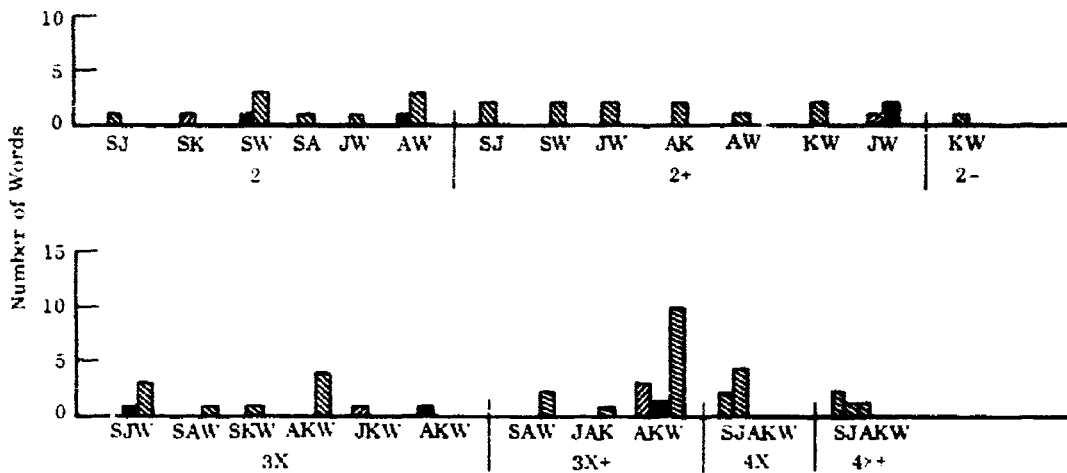


Fig. 3-20 Statistics on Disagreement Among Dictionaries for Phonetics of Additional Monosyllables (Except Those Ending in ES or ED Which Have Polysyllabic Pronunciations). Total of 1147 words represented (II)

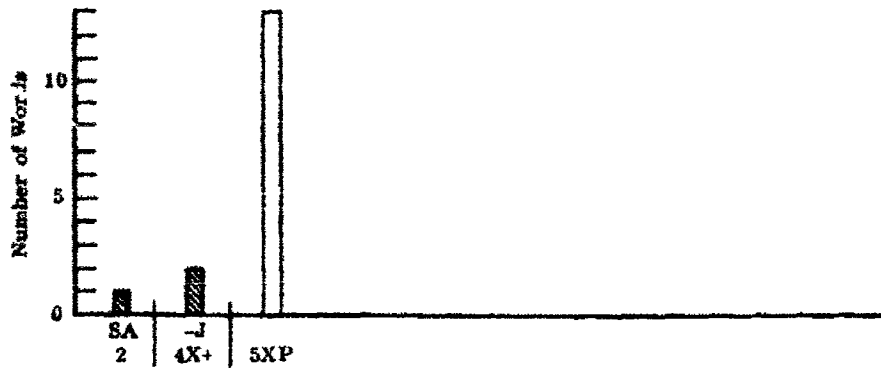


Fig. 3-21 Statistics on Polysyllabic Pronunciation of Words Ending in ES According to Five Dictionaries

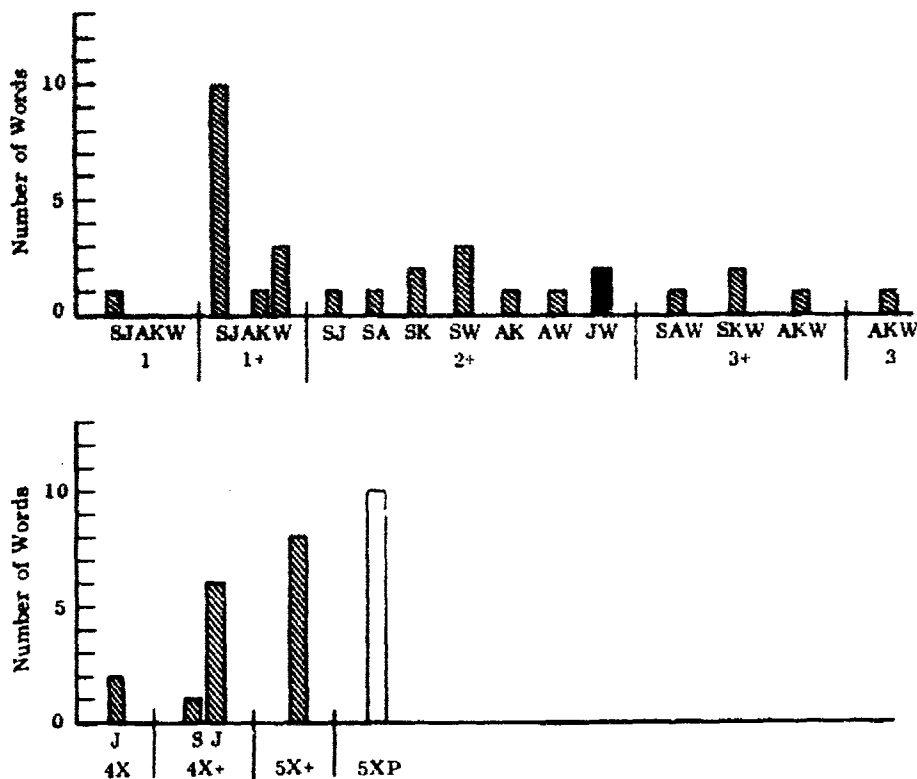


Fig. 3-22 Statistics on Polysyllabic Pronunciation of Words Ending in ED According to Five Dictionaries

The properly corrected (by procedures described in subsection 3.3.2) phonetic data on the English monosyllables were used for compiling new and complete sets of homonyms (as described in subsection 3.3.4). The data are presented in Fig. 3-23. Notice the limited effect of the addition of words to the 5757 elementary words. (See Fig. 3-15.)

The results of computation of phonetic forms of additional words indicate that the computable relations described for elementary words can indeed be extended, and the accuracy obtained with the new set of words is at least as high as that for the elementary words. Moreover, the number of and distribution of sets of homonyms are not altered significantly by the addition of words to the elementary word set. The phonetic systems and the computable relations, as discussed in this paper, remain highly accurate and suitable for working with polysyllabic words.

### 3.6 SUMMARY AND CONCLUSIONS

We have discussed the existence and accuracy of relations between the orthographic forms of English monosyllables and their corresponding phonetic representations for the various dialects and transcriptions (Refs. 3-10 through 3-14) as well as those between the various phonetic representations that correspond to any of these orthographic forms and as transcribed by five lexicographers (Refs. 3-10 through 3-14). The description starts with the computable relation for the set of elementary words and extends their applicability to English monosyllables.

Some of the difficulties in determining the phonetic representations of corresponding orthographic forms and the disagreements among lexicographers on such representations are discussed. Limitations to the specification (in a technical paper) of a set of rules that map the orthographic form of English words into corresponding phonetic representation are indicated; a more general method for computing such relations is described and evaluated; different types of errors are studied; and an approach is presented for the selection of exception words by use of codes incorporated in the

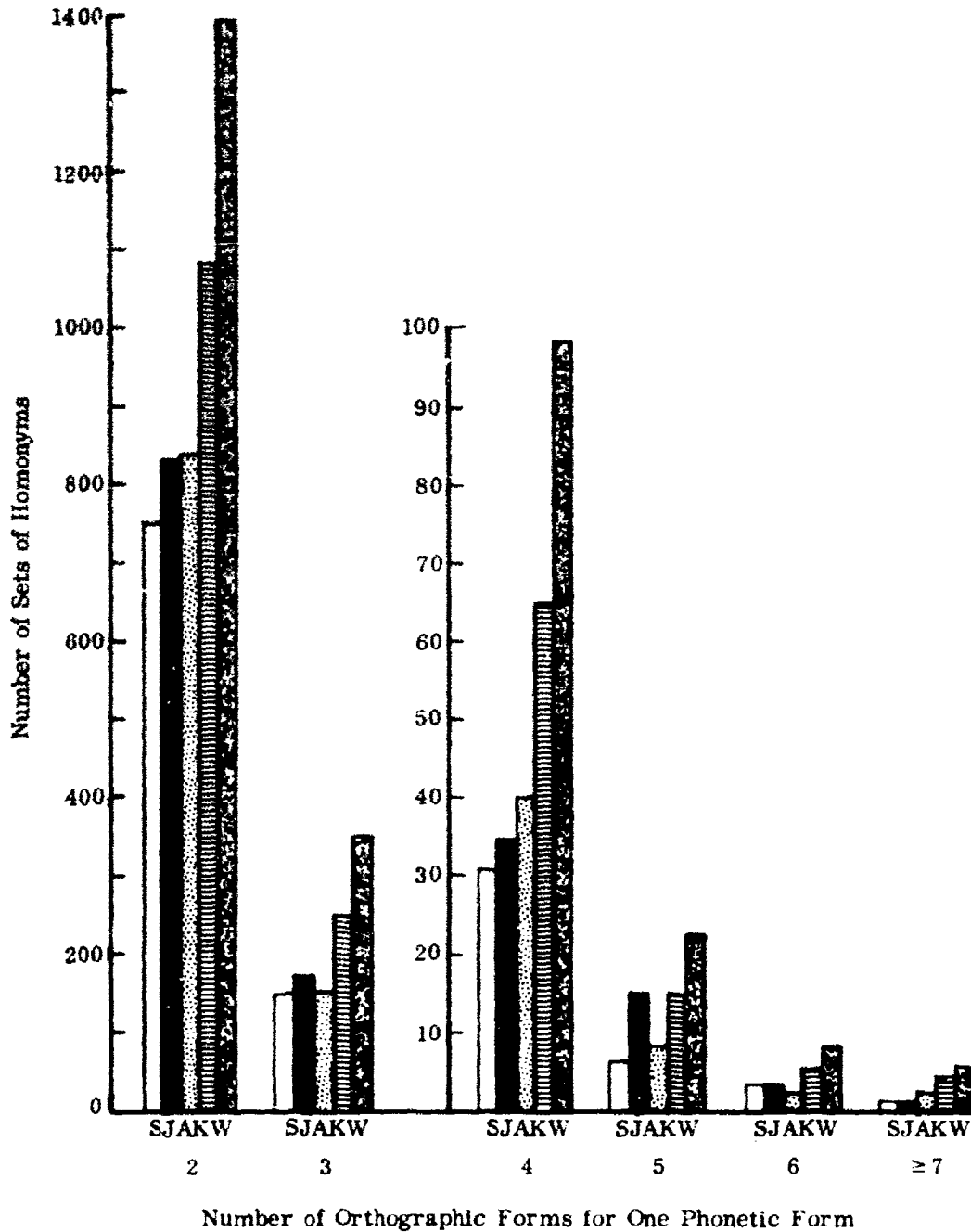


Fig. 3-23 Graphic Representation of the Number of Homonym Sets Among the English Monosyllables in Five Dictionaries. Total of 6904 words represented

unique phonetic dictionary (Fig. 3-8, and data on status of words, Refs. 3-26 and 3-27). The exception words are used for increasing the accuracy of computations, and the coded information is used for increasing the efficiency of the algorithm by proper selection of exception words.

The smallest group of symbols (called marker-vowel-marker) over which the computable relations between the orthographic and the corresponding phonetic forms of English monosyllables can be established is mentioned and evaluated. Comparison of data for such a group of symbols is shown to provide an accurate set of computable relations for mapping the orthographic forms of words into corresponding phonetic representations, as well as to provide an accurate set of computable relations between the various phonetic representations studied (Refs. 3-10 through 3-14).

Evaluation of the phonetic data indicates a marked disagreement in the symbol-for-symbol relation between the various phonetic transcriptions. The confusions that result from such discrepancies are evaluated by a compilation of homonyms. An evaluation is made of the extent to which such confusion can be attributed to differences in regional pronunciation patterns and also the extent to which the confusion can be reduced by restricting the number of words as either double standard in their syntactic status (Refs. 3-26 and 3-27) or algorithmic in their phonetic status.

This entire treatment is aimed at computer programmers and engineers; hence, an operational point of view is taken throughout. This may differ from some linguistic approaches and concepts, but consideration of these is beyond the scope of this paper. This work describes the computable relations between the orthographic and the phonetic forms of English; it provides a computable definition of the dialect patterns studied; and it defines the group of symbols, called marker-vowel-marker, for which orthographic and phonetic forms can be related as well as ones for which phonetic data can be accurately compared and related. Thus, it describes the computer programs and definitions of phonetics of English words for use with speech recognition equipment and for machines that can read English.

### 3.7 REFERENCES

- 3-1 B. V. Bhimani, R. D. Merrill, R. P. Mitchell, and M. R. Stark, "An Approach to Speech Synthesis and Recognition on a Digital Computer," Proc. of 21st National Conference of Association for Computing Machinery, A. C. M. Publication P-66, Thompson Book Co., Washington, D. C., 1966, pp. 275 - 296
- 3-2 A. H. Roberts, A Statistical Linguistic Analysis of American English, Mouton & Co., Hague, 1965
- 3-3 F. S. Cooper, P. C. Delatere, A. M. Liberman, J. M. Borst, and L. J. Gerstman, "Some Experiments on the Perception of Synthetic Speech," J. Acoust. Soc. Am., Vol. 24, 1952, pp. 597 - 606
- 3-4 F. S. Cooper, A. M. Liberman, and J. M. Borst, "The Interconversion of Audible and Visible Patterns as a Basis for Research in the Perception of Speech," Proc. Nat. Acad. Sci. (U. S.), Vol. 37, 1951, pp. 318 - 325
- 3-5 Ilse Lehiste and G. E. Peterson, "Duration of Syllable Nuclei in English," Studies in Syllable Nuclei 2, Speech Research Laboratory, University of Michigan, Ann Arbor, Mich., 1960
- 3-6 Ilse Lehiste and G. E. Peterson, "Transitions, Glides and Diphthongs," Studies in Syllable Nuclei 2, Speech Research Laboratory, University of Michigan, Ann Arbor, Mich., 1960
- 3-7 B. Lindblöm, "Spectrographic Study of Vowel Reduction," J. Acoust. Soc. Am., Vol. 35, 1963, pp. 1773 - 1781
- 3-8 K. N. Stevens, A. S. House, and A. P. Paul, "Acoustical Description of Syllable Nuclei: An Interpretation in Terms of Dynamic Model of Articulation," J. Acoust. Soc. Am., Vol. 40, 1966, pp. 123 - 132
- 3-9 B. V. Bhimani, A Multidimensional Model for Automatic Speech Recognition, Final Report, Contract AF 19(628)-2766, DDC Doc. No. AD-437-324, 1964
- 3-10 The Shorter Oxford Dictionary on Historical Principles, Oxford at the Clarendon Press, 1959, 3rd ed., revised with addenda

- 3-11 J. D. Jones, Everyman's English Pronouncing Dictionary, Dutton, New York, 12th ed., 1963
- 3-12 The American College Dictionary, Random House, 1962, New York
- 3-13 J. S. Kenyon and T. A. Knott, A Pronouncing Dictionary of American English, G. C. Merriam Co., Springfield, Mass., 1958
- 3-14 Webster's New International Dictionary of the English Language, G. C. Merriam Co., Springfield, Mass., 1961
- 3-15 Hans Kurath et al., Linguistic Atlas of the United States and Canada, Brown University, Providence, R. I., 1939-1941
- 3-16 M. I. Halpern, "XPOP: A Meta-Language Without Metaphysics," Proc. of Fall Joint Computer Conference, 1964, pp. 57-58
- 3-17 J. L. Flanagan, Speech Analysis, Synthesis, and Perception, Academic Press, 1965
- 3-18 C. G. M. Fant, Acoustic Theory of Speech Production, Mouton & Co., Gravenhage, 1960
- 3-19 K. N. Stevens, "Studies of Formant Transition Using a Vocal Tract," J. Acoust. Soc. Am., Vol. 28, 1958, p. 578
- 3-20 Funk and Wagnalls, New Practical Standard Dictionary, J. G. Ferguson Publishing Co., Chicago, Ill., 1956
- 3-21 E. M. Higginbottom, "A Study of Representation of English Vowel Phonemes in Orthography," Language and Speech, Vol. 5, 1962, pp. 67-118
- 3-22 J. L. Dolby and H. L. Resnikoff, "On the Structure of Written English Words," Language, Vol. 40, 1964, pp. 167-196
- 3-23 L. L. Earl, B. V. Bhimani, and R. P. Mitchell, "Statistics of Operationally Defined Homonyms of Elementary Words," Mechanical Translation (in press)
- 3-24 R. Bloch and G. L. Trapen, Outline of Linguistic Analysis, Linguistic Society of America, Waverly Press, Baltimore, 1942



- 3-25 E. L. Thorndike and C. Lorge, The Teacher's Word Book of 30,000 Words, Bureau of Publications, Teacher's College, Columbia University, New York, 3rd printing, 1959
- 3-26 J. L. Dolby and H. L. Resnikoff, The English Word Speculum, Vols. I through V, Lockheed Missiles & Space Company, Palo Alto, Calif., 1964
- 3-27 J. L. Dolby, H. L. Resnikoff, and MacMurray, "A Tape Dictionary for Linguistic Experiments," Proc. of Fall Joint Computer Conference, 1963, pp. 419 - 423
- 3-28 G. B. Shaw, Pygmalion, Dodd, 1939

Some of the illustrations and related commentary have already been published by the Association for Computing Machinery (item 1 in the list of references) and by the journal Machine Translation; such material is repeated here for the sake of continuity and completeness.

**Appendix A**  
**RUSSIAN TEXT**

**A-1**

**LOCKHEED PALO ALTO RESEARCH LABORATORY**  
**LOCKHEED MISSILES & SPACE COMPANY**  
**A GROUP DIVISION OF LOCKHEED AIRCRAFT CORPORATION**

# Гипотезы

## ФАЗОВЫЕ ПРЕВРАЩЕНИЯ В ГЛУБИНАХ ЗЕМЛИ

Достижения механики и физики в XVII—XVIII вв. позволили определить массу и среднюю плотность Земли. Последняя оказалась равной  $5,5 \text{ г/см}^3$ . А так как плотность наиболее тяжелых пород на поверхности Земли не превышает  $3,3 \text{ г/см}^3$ , то, естественно, возникло представление, что плотность Земли увеличивается с глубиной.

Факты существования железных метеоритов, а также в прошлом популярная теория происхождения Земли из горючего вещества Солнца привели многих ученых к мысли о концентрации железа в центре Земли. Примечательно, что уже вполне определенные высказывания французского геолога А. Дюбуа в 1866 г. о железном ядре Земли вскоре получили поддержку со стороны сейсмологов, которым в конце XIX и начале XX в. удалось установить наличие в Земле ядра.

В 20-х годах текущего столетия В. М. Гольдшмидт (Норвегия) и немецкий физико-химик Г. Тамман развили представление о том, что в первоначально расплавленной Земле происходило разделение (дифференциация) вещества по их плотности, аналогично тому, что мы имеем, например, при плавке сульфидных руд. При этом процессе появляются три слоя: шлак (силикатный слой), штейн (смесь сульфидов и металла) и собственно металл. Согласно этой гипотезе, в Земле выделялись следующие слои: силикатный и сульфидный (оболочка Земли) и металлический, состоящий из железа с примесью никеля (ядро Земли).

Американские ученые Ф. Кларк, Г. Вашингтон, Л. Адамс и др. не выделяли сульфидный слой; они полагали, что между железным ядром и силикатной оболочкой находится промежуточная область, состоящая из смеси силикатов и железа.

Теория слоистой, химически дифференцированной Земли, во многом подкреплена данными сейсмологов, которые первоначально считали, что в мантии (оболочке) Земли, т. е. в той ее части, которая расположена между земной корой и ядром, существует много границ раздела.

В дальнейшем успехи геофизики и космогонии, связанные главным образом с именами

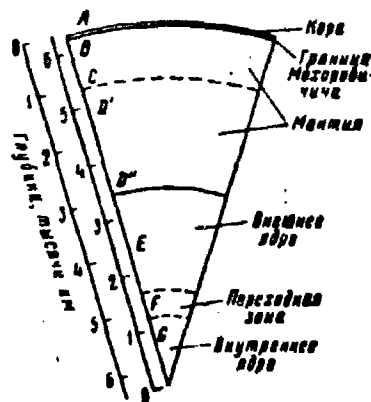


Рис. 1. Зоны в Земле по К. Буллену



Рис. 2. Пластинчатый кристалл новой фазы (микрофотография)

ми англичанина Г. Джоффриса, немца Б. Гутенберга, новозеландца К. Буллена, советского ученого О. Ю. Шмидта, американца Дж. Койшера и др., заставили внести существенные коррективы в эти представления. Были уточнены количество и характер границ раздела в Земле, изучены упругие свойства ее глубины, получены данные о распределении плотности в Земле. Одновременно в свете новых гипотез о происхождении Земли и солнечной системы аргументы в пользу огненно-жидкой стадии в развитии Земли потеряли свою силу. Геохимические наблюдения также мало согласовались с представлением о некогда расплавленной Земле.

Обобщив результаты многих исследований, К. Буллен выделил в Земле ряд зон, отличающихся друг от друга по физическим свойствам (рис. 1). Оказалось, что упругие свойства оболочки (мантии) Земли таковы, что сульфиды тяжелых металлов или металлическое железо не могут находиться в ней в заметных количествах. Однако предположительно физически и химически однородной (гомогенной) силикатной, преимущественно оливиновой<sup>1</sup>, мантии Земли вступает в противоречие с известными геофизически-

ми данными. Оказалось, что в зоне «С» резко возрастает скорость сейсмических волн (см. рис. 1) и растет электропроводность. На зону «С» приходится максимум глубоководных землетрясений. Предположение о гомогенной мантии приводит к вероятно большому моменту инерции ядра. Поэтому К. Буллен был вынужден предположить, что внутри зоны «С» также резко увеличивается градиент плотности.

Еще в 1936 г. известный английский ученый Дж. Бернал пытался объяснить аномальные свойства зоны «С» переходом оливи-

на в более плотную модификацию, имеющую структуру типа шпиннели. Структурный тип шпиннели представляет собой плотнейшую кубическую упаковку ионов кислорода с катионами, расположенными в октаэдрических и тетраэдрических положениях.

Структура оливины, хотя и обладает плотнейшей упаковкой, но сильно искажена и потому переход в более симметричную упаковку шпиннели будет сопровождаться увеличением плотности.

В настоящее время получены некоторые косвенные доказательства возможности подобного перехода. Эта гипотеза, хотя и не может полностью объяснить все особенности строения мантии, но тем не менее ее значение заключается в том, что она возбудила большой интерес к проблеме физического изменения состояния вещества мантии под действием давления.

Американский геофизик Ф. Борч, детально изучивший свойства мантии, пришел к выводу, что, по всей вероятности, единственной неоднородной (гетерогенной) областью в мантии является зона «С», в которой можно ожидать физических и химических изменений. Одновременно им было показано, что упругие свойства нижней мантии или зоны «D» слишком высоки для обычных силикатов с кремнием в четвертой координа-

<sup>1</sup> Оливины — силикат магния и железа  $(MgFe)_2 \cdot SiO_4$ .

ции<sup>1</sup>, но вполне соответствуют плотноупакованным оксидам типа периклаза MgO, рутила TiO<sub>2</sub> и корунда Al<sub>2</sub>O<sub>3</sub>. На этом основании Ф. Берч предположил, что в зоне «С» происходит перестройка ферромагнитных силикатов в плотно упакованные структуры простых и сложных оксидов.

Гипотеза Ф. Берча требует перехода кремния из четверной координации в шестерную, т. е. такое расположение, где вокруг кремния располагается шесть атомов кислорода.

Однако возможность подобного изменения координации у кремния подвергалась сомнению. В связи с вышеказанным, нами, совместно с научным сотрудником Института физики высоких давлений АН СССР С. В. Провой, были поставлены опыты по изучению состояния кремнезема в условиях сверхвысоких давлений и высоких температур. Работа проводилась при помощи установки, созданной в ИФВД АН СССР и способной генерировать давление до 200 тыс. ат и сочетания с высокими температурами.

В качестве исходных веществ употреблены кварц и аморфный кремнезем. Опыты проводились при давлениях от 36 до 145 килобар и температурах от 1200 до 2000 °С. При давлениях 115–145 килобар и температурах около 1500 °С удалось обнаружить неизвестную фазу с высоким показателем преломления, в виде игольчатых и пластинчатых кристаллов размером до 0,5 мк (рис. 2).

Новая фаза имела высокую твердость, близкую к твердости корунда, и плотность 4,36 г/см<sup>3</sup>, в то время как наиболее плотная из всех известных до сих пор модификаций кремнезема — коэсит — имеет плотность 3,91 г/см<sup>3</sup>. При помощи химических и спектральных исследований удалось установить, что новая фаза состоит из чистого кремнезема. Следовательно, нами была получена новая модификация кремнезема с очень высокой плотностью, превышающей на 64% плотность кварца и на 45% плотность коэсита.

Рентгеноструктурное исследование, проведенное совместно с акад. Н. В. Беловым,

<sup>1</sup> При этом кремне-кислородные соединения содержат кремний в окружении четырех атомов кислорода, расположенных по вершинам правильного тетраэдра (в терминологии кристаллохимиков это расположение называется тетраэдрической координацией).

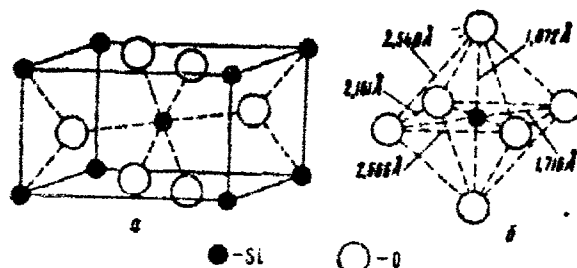


Рис. 3. Общий вид элементарной ячейки новой модификации SiO<sub>2</sub> (а) и координационный октаэдр кремния (б)

показало, что новая модификация кремнезема имеет кристаллическую структуру с кремнием в шестерной координации (рис. 3). Необходимо заметить, что условия получения новой плотной модификации SiO<sub>2</sub> совпадают с условиями, предполагаемыми для верхних частей слоя «С», где как раз и ожидалось фазовые переходы.

Приближенные подсчеты показывают, что модель нижней мантии (зона «D»), состоящая из механической смеси MgO (периклаза), FeO (вюстит) и новой модификации SiO<sub>2</sub> будет иметь плотность и упругость, вполне согласующиеся со свойствами, известными для нижней мантии из экспериментальных и теоретических данных. Но можно предполагать, что нижняя мантия не представляет собой механической смеси индивидуальных оксидов. В этом случае нижняя мантия может рассматриваться как гомогенное вещество с плотнейшей упаковкой кислорода, октаэдрические положения которой заняты магнием, железом и кремнием с неупорядоченным расположением последних. Подсчеты плотности состава, соответствующего веществу мантии, предполагая, что она находится в вышеуказанном структурном состоянии, также согласуются с известными данными.

Необходимо подчеркнуть, что предполагаемые структурные модели нижней мантии допускают возможность дальнейших трансформаций, например, в структуры с координационным числом катионов, равным восьми. Переход такого типа может объяснить резкое изменение свойств на границе ядра, что регистрируется соответствующими изменениями в скоростях сейсмических волн. Структуры с координационным числом 8 характерны для металлов и интерметаллических соединений, и весьма вероятно, что вещество Земли, имея подобную координацию, будет обладать металлическими свойствами.

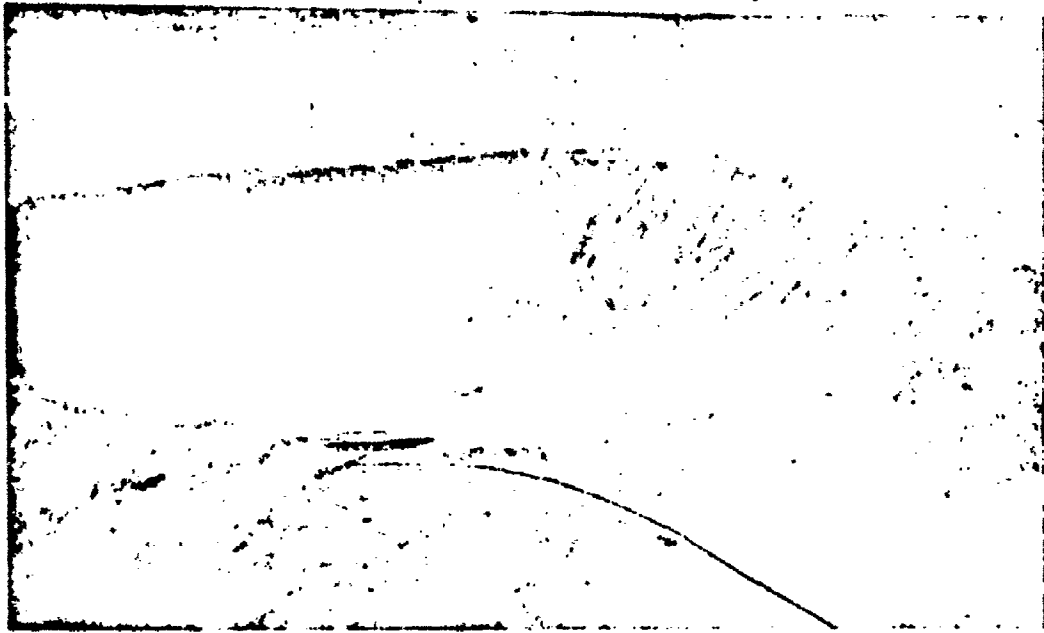


Рис. 4. Общий вид Арizonского кратера

В настоящее время нельзя решить, насколько справедливы высказанные здесь предположения. Ясно лишь одно, что синтез новой, плотной модификации  $\text{SiO}_2$  дает возможность объяснить свойства нижней мантии без привлечения гипотезы о химическом изменении вещества Земли с глубиной.

В заключение хотелось бы добавить, что группа американских исследователей (Э. Чоу, Дж. Фахи, Дж. Литтлер и Ц. Милтон) обнаружила открытую нами рутилоподобную модификацию  $\text{SiO}_2$  в породах Арizonского метеоритного кратера (рис. 4),

где она образовалась под действием высоких давлений и температур, возникших при падении метеорита. Возможно, что при падении метеоритов возникают более высокие давления, чем получаемые пока в лабораториях. Поэтому необходимо тщательно изучать породы метеоритных кратеров; может быть, будут открыты еще более интересные минералы, которые помогут нам в познании земных глубин.

С. М. С т и ш о в  
Московский государственный университет  
им. М. В. Ломоносова

Appendix B  
ABSTRACTS USED AS DATA BASE

**AA1. Neutron-Proton Elastic Differential Cross Sections from 1 to 6 GeV.\*** M. N. KRZEMIEC† (introduced by M. L. Perl), F. MARTIN, M. L. PERL, *Stanford University*, M. LONGO, AND S. T. POWELL, III, *The University of Michigan*.—Measurements of *np* elastic scattering from 1 to 6 GeV were made using spark chambers. A well-collimated neutron beam produced by the external beam of the bevatron interacted in a liquid-hydrogen target. The recoil-proton momentum and angle were measured with a bending magnet and spark chambers, and the scattered-neutron angle was measured using an array of steel-plate spark chambers. The apparatus was on rails to allow coverage of c.m. angles from about 10° to 150°. Preliminary results covering the entire angular range are presented in the form of relative cross sections vs 4-momentum transfer for each incident-energy interval. Corrections for detection efficiency and background contamination have been made.

\* Work supported in part by the U. S. Atomic Energy Commission and the U. S. Office of Naval Research.  
† National Science Foundation Predoctoral Fellow.

**AA2. Large-Angle Neutron-Proton Correlation Function at 23 MeV.\*** J. J. MALANIFY, P. J. BENDT, T. R. ROBERTS, AND J. E. SIMMONS, *Los Alamos Scientific Laboratory*.—Measurements are presented for the neutron-proton correlation function  $C_{np}$  at  $E_n = 23$  MeV and  $\theta_{c.m.} = 150^\circ$ , in addition to further data at  $\theta_{c.m.} = 180^\circ$ . Polarized neutrons from the  $T(d,n)^3\text{He}$  reaction were utilized, together with a polarized proton target (LMN) in which the polarization averaged approximately 30%. The recoiling protons were detected by a counter telescope. The value of  $C_{np}$  was derived from an asymmetry that was induced by cycling the direction of the proton polarization. The results to date lie close to the predictions of the Hamada-Johnston potential model. The effect of these data on the recent Livermore phase-shift analysis<sup>1</sup> is discussed.

\* Work performed under the auspices of the U. S. Atomic Energy Commission.  
<sup>1</sup> T. Hamada and J. D. Johnston, *Nucl. Phys.* 24, 383 (1962). These predictions were kindly calculated for us by Prof. P. Signad.  
<sup>2</sup> H. P. Noyes, D. E. Hader, R. A. Arndt, and M. H. MacGregor, *Phys. Rev.* B130, 380 (1963).

**AA3. Search for a Narrow Resonance in  $p\text{-He}^4$  Elastic Scattering.\*** D. BOYD (introduced by J. V. Kane), *Rutgers University*, P. F. DONOVAN, J. V. KANE, J. F. MOLLENAUER, *Bell Telephone Laboratories*, AND P. PARKER, *Brookhaven National Laboratory*.—Evidence for a sharp state in  $\text{Li}^6$  has been reported by Benison *et al.*<sup>1</sup> in the decay of  $\text{AHe}^4$ . Dangle *et al.*<sup>2</sup> searched for this resonance in  $p\text{-He}^4$  scattering, from a  $Q$  of 10.3–10.8 MeV in 25-keV steps on a less than 15-keV-thick target. We have repeated the experiment of Dangle *et al.*, since it has been estimated that, if the resonance is a  $T=2$  state, it may be narrower than 15 keV owing to isospin conversion and 4-particle phase-space considerations.<sup>3</sup> We also wished to extend the energy range, the statistical accuracy, and the number of angles of the Dangle group. No deviation exceeding 0.3% from a smooth cross section was found over the range  $10.35 \leq Q \leq 11.05$  using 10-keV overlapping steps, and

$9.3 \leq Q \leq 11.6$  MeV using 100-keV steps, at laboratory angles of 150° and 120°.

\* Work supported in part by the National Science Foundation.  
<sup>1</sup> M. J. Benison, B. Krishnamurthy, R. Levi-Strauss, and M. Raymond, *Phys. Letters* 13, 333 (1964).  
<sup>2</sup> R. L. Dangle, J. Jobst, T. L. Essner, *Bull. Am. Phys. Soc.* 10, 433 (1963).  
<sup>3</sup> C. Zupancic, private communication.

**AA4. Study of the Mass 4 System by the  $\text{H}^2(d,p)$ ,  $\text{H}^2(d,p)$ , and  $\text{H}^2(t,p)$  Reactions.\*** G. G. OHLSEN, R. W. NEWSOME, JR., AND R. H. STOKES, *Los Alamos Scientific Laboratory*.—Proton spectra from the deuteron bombardment of  $\text{H}^2$  and  $\text{He}^3$  have been obtained at 6 energies between 10 and 16 MeV. Reaction protons from the gas target were identified with a semiconductor  $\Delta E\text{-}E$  telescope. At a laboratory angle of 15°, the peak cross section for formation of the  $\sim 20\text{-MeV}$  state of  $\text{He}^4$  decreases smoothly from about 30 to about 6 mb/r·MeV as the energy is increased from 10 to 16 MeV. No evidence for a similar state appears in either the  $\text{H}^2(d,p)$  or the  $\text{H}^2(t,p)$  reactions. This confirms the accepted isospin of zero for this state. Spectra from all three reactions show a broad anomaly (width 2–3 MeV) at an excitation energy about 2 MeV above the  $t+\pi$  or  $t+p$  mass.

\* Work performed under the auspices of the U. S. Atomic Energy Commission.

**AA5. Study of the Mass 4 System by the  $\text{H}^2(t,d)$  and  $\text{He}^3(t,d)$  Reactions.\*** R. H. STOKES, NELSON JARMIS, R. W. NEWSOME, JR., AND G. G. OHLSEN, *Los Alamos Scientific Laboratory*.—Deuteron spectra from the bombardment of  $\text{H}^2$  and  $\text{He}^3$  with a 21.8-MeV triton beam have been obtained at laboratory scattering angles in the range 10°–20°. In contrast to the  $\text{He}^4(d,p)$  reaction,<sup>1</sup> breakup particles from the excited  $\text{He}^4$  nuclei cannot contribute to the observed spectra. The virtual state of  $\text{He}^4$  at  $\sim 20\text{-MeV}$  excitation energy is clearly observed in the  $\text{He}^3(t,d)$  reaction, but no corresponding state is observed in the  $\text{H}^2(t,d)$  reaction. In the deuteron spectra from both reactions, a broad peak occurs at about 1.7 MeV above the  $p+t$  or  $n+t$  mass. In general, the cross sections are much lower and the structure more pronounced than in the corresponding  $(d,p)$  reactions. The deuteron spectra are discussed in terms of virtual states of the  $\alpha$  particle and possible effects of reaction mechanisms.

\* Work performed under the auspices of the U. S. Atomic Energy Commission.  
<sup>1</sup> G. G. Ohlsen, R. W. Newsome, Jr., and R. H. Stokes, Paper A44, this meeting.

**AA6. Evidence for a  $2^+ T=0, S=2$  State in  $\text{He}^4$ .** Y. C. TANG, *University of Minnesota*.—Measurements from  $\text{Li}^6(\pi^+, 2p)\text{He}^4$  showed that, in the summed energy spectra of the protons, there are two strong peaks corresponding to 0 and 30 MeV excitation in  $\text{He}^4$ , and a weak peak at 20 MeV excitation.<sup>1</sup> These features can be understood by considering  $\text{Li}^6$  as  $\text{He}^4+d$ . The ground-state peak is due to absorption of the pion by the deuteron cluster. To produce the levels at about 20 MeV in  $\text{He}^4$ , the pion needs to interact with a nucleon from each of the two clusters, since these levels have cluster

structures of a nucleon plus a  $A=3$  cluster. This process has, however, only a small probability, since the clusters are well separated and pion is absorbed only when the two nucleons are in close proximity. For the peak around 30 MeV, it is produced by the pion interacting with a deuteron cluster in the  $\alpha$  cluster. The levels at this energy have, therefore,  $d+d$  cluster structure and the lowest one has  $S=2$ ,  $L=0$ , and  $T=0$ . Using similar consideration for  $\text{Li}^3(\pi^+, 2p)\text{He}^4$ , agreement with experiment can also be obtained.

<sup>1</sup>G. Charpak *et al.*, *Phys. Letters* 16, 54 (1965).

**AA7. Optical Parameters for the Elastic Scattering of  $\text{He}^3$  by  $\text{He}^4$ .** N. R. FLETCHER, F. DUNHILL, T. GRAY, AND H. T. FORTUNE, *Florida State University*.—Optical-model parameters have been determined for the elastic scattering of  $\text{He}^3$  by  $\text{He}^4$  for use in a DWBA direct reaction analysis of the reaction  $\text{Li}^3(\pi, \text{He}^3)\alpha$  observed by Heydenburg and Han.<sup>1</sup> The yield of elastically scattered  $\text{He}^3$  measured over a bombarding energy range of 12–19 MeV shows no prominent resonances. Angular distributions were measured from  $\theta_{\text{c.m.}} = 30^\circ$ – $150^\circ$  at  $E_{\text{He}^3} = 12.0, 13.5, 15.0, 16.5, 18.0,$  and  $19.0$  MeV. At 15 MeV and above, the optical potentials and radius parameters found vary only a few percent from the values:  $U=100$  MeV,  $W=10$  MeV,  $r_1=r_2=r_3=1.6$  F,  $a_1=a_2=0.6$  F. A spin-orbit interaction of about 20 MeV improves the fit at back angles; however, it is evident that mechanisms other than this are needed to account for much of the deviation of the data from the simple optical-model calculation.

<sup>1</sup>Work supported in part by the U. S. Air Force Office of Scientific Research.  
<sup>1</sup>N. P. Heydenburg and I. G. Han, *Bull. Am. Phys. Soc.* 7, 34 (1962).

**AA8. Reactions  $\text{He}^4(\text{He}^3, p_1)\text{Li}^{3,4}$  and  $\text{Li}^3(p, p_1)\text{Li}^{3,4}$ .** W. D. HARRISON, *California Institute of Technology*.—The total cross section for the reaction  $\text{He}^4(\text{He}^3, p_1)\text{Li}^{3,4}$  has been measured by observing the yield of deexcitation  $\gamma$ -rays from  $\text{Li}^{3,4}$ . Measurements cover the excitation range from 9.5 to 12.1 MeV in the compound nucleus  $\text{Be}^7$ . A study of the  $\text{Li}^3(p, p_1)\text{Li}^{3,4}$  reaction and its possible interpretation were reported previously.<sup>1</sup> Similar features are observed in both reactions: a broad maximum at about 10.7 MeV excitation and some kind of narrower anomaly at about 11.1 MeV. These features have been fitted with 2-level formulas. Combining the information so obtained with that from the  $\text{Li}^3(p, p_1)\text{Li}^3$  and  $\text{Li}^3(p, p_1)\text{Li}^4$  reactions, the following assignments are made:  $J^\pi = \frac{1}{2}^-$ ,  $T = \frac{1}{2}$  for the broader level and  $J^\pi = \frac{3}{2}^-$ ,  $T = \frac{1}{2}$  for the narrower. The  $T = \frac{1}{2}$  assignment is based on the fact that the narrower level is not observed in the  $\text{Li}^3(p, p_1)\text{Li}^3$  and  $\text{Li}^3(p, p_1)\text{Li}^4$  reactions, where both its formation and decay are forbidden. In the other reactions, where only its formation is forbidden, it appears through its interference with the  $T = \frac{3}{2}$  level. The  $T = \frac{1}{2}$  level is probably the same as that recently observed at Berkeley.<sup>2</sup>

<sup>1</sup>Work supported by the U. S. Office of Naval Research.  
<sup>1</sup>W. D. Harrison, *Bull. Am. Phys. Soc.* 9, 703 (1964).  
<sup>2</sup>C. D'Arraz, J. Ceray, and R. H. Pebl, *Phys. Rev. Letters* 14, 708 (1965).

**AA9. Phototriton Cross Section of  $\text{Li}^6$ .** N. K. SHERMAN, *McGill University*, JOHN E. E. BAGLIN, AND R. O. OWENS, *Yale University*.—The cross section for the reaction  $\text{Li}^6(\gamma, t)$  recently observed<sup>1</sup> has been measured for photon energies between 18 and 30 MeV. Bremsstrahlung of end-point energy 40 MeV produced by the Yale electron linac was used to irradiate  $\text{Li}^6$  evaporated onto Formvar films. The target thicknesses were about  $270 \mu\text{g}/\text{cm}^2$ . Tritons were identified by a

quadrupole triplet magnet and were stopped in a silicon-barrier detector. Energy resolution was 130 keV. Tritons from the  $(\gamma, t)$  reaction were gated out. Absolute cross section was obtained by normalizing the  $\text{Li}^6(\gamma, p)$  cross section via photoprotons from a deuterium gas cell, and also to the  $\text{Li}^6(\gamma, p)$  cross section<sup>2</sup> via photoprotons identified along with the phototritons. Two-body breakup was assumed. The  $(\gamma, t)$  cross section found in both ways is 0.7 mb at 21 MeV, almost as large as the photonucleon cross section at this energy.

<sup>1</sup>Work supported by the U. S. Atomic Energy Commission.  
<sup>1</sup>N. K. Sherman, R. C. Morrison, and J. R. Stewart, *Bull. Am. Phys. Soc.* 10, 341 (1965).  
<sup>2</sup>B. L. Berman, R. L. Bramblett, J. T. Caldwell, R. R. Harvey, *et al.*, UCRL-12376 (1965).

**AA10. Charge Symmetry in the Mirror Reactions  $\text{Li}^3(d, p)\text{Li}^3$  and  $\text{Li}^3(d, n)\text{Be}^3$ .** S. M. AUSTIN,† P. PAUL, A. CHEUNG, S. S. HANNA, AND W. E. MEYERROFF, *Stanford University*.—It has been pointed out<sup>1</sup> that a comparison of mirror reactions can provide a test of charge symmetry in nuclear reactions. For  $\text{Li}^3(d, p)\text{Li}^3$  and  $\text{Li}^3(d, n)\text{Be}^3$ , this comparison was made<sup>2</sup> by observing the isotropic  $\gamma$ -rays from  $\text{Li}^3$  and  $\text{Be}^3$ . This experiment has been repeated and extended with increased precision by use of a Ge(Li) detector of sensitive volume  $3 \times 2 \times 1 \text{ cm}^3$ , from  $E_d = 0.15$  to 7.2 MeV with the Stanford 3-MV Van de Graaff and a newly installed FN tandem. Although the energy range covers at least two known resonances in  $\text{Be}^3$ , no definite evidence of resonant structure was observed. Above  $E_d = 0.5$  MeV, the neutron-to-proton ratio lies between 1.14 and 1.20, with some evidence for mild variations. In the region of overlap, below  $E_d = 3$  MeV, there is qualitative agreement with the earlier observations.<sup>1</sup> Above  $E_d = 1$  MeV, where the cross sections show pronounced stripping patterns,<sup>3</sup> the neutron-to-proton ratio is in good agreement with stripping calculations. Below  $E_d = 0.7$  MeV, an observed decrease in the ratio is suggestive of the Oppenheimer-Phillips process.

<sup>1</sup>Work supported in part by the National Science Foundation and the U. S. Army Research Office.  
<sup>1</sup>D. H. Wilkinson, *Phil. Mag.* 2, 83 (1957).

**AA11. Comparison of the Reactions  $\text{Li}^3(p, n_1)\text{Be}^3$  and  $\text{Li}^3(p, p')\text{Li}^3$ .** S. S. HANNA, P. PAUL, S. M. AUSTIN,† AND W. E. MEYERROFF, *Stanford University*.—A comparison of total cross sections in the reactions  $\text{Li}^3(p, n_1)\text{Be}^3$  and  $\text{Li}^3(p, p')\text{Li}^3$  has been made by measuring the isotropic  $\gamma$ -rays emitted from  $\text{Be}^3$  and  $\text{Li}^3$  with a Ge(Li) detector of sensitive volume  $2 \times 3 \times 1 \text{ cm}^3$ . The range  $3 \text{ MeV} < E_p < 10 \text{ MeV}$  was covered with the newly installed Stanford FN tandem. In sharp contrast to the ratio determined in the  $\text{Li}^3+d$  reactions,<sup>1</sup> the ratio of cross sections in the  $\text{Li}^3+p$  reactions deviates markedly from unity and varies strongly with energy. The neutron-to-proton ratio rises to a value of 0.54 at  $E_p = 3.3$  MeV, falls smoothly to 0.21 at 5.5 MeV, rises again to 0.33 at 8.6 MeV, and then declines. The over-all dominance of the  $(p, p')$  reaction suggests the presence of a direct process in this reaction. It is noted that the maximum in the ratio at 3.3 MeV occurs at an energy where the neutron yield attains a maximum, while the minimum at 5.5 MeV can be attributed to a strong resonance in the proton yield not observed in the neutron yield. These observations are discussed in terms of isotopic spin and known levels of  $\text{Be}^3$ .

<sup>1</sup>Work supported in part by the National Science Foundation and the U. S. Army Research Office.  
<sup>1</sup>Alfred P. Sitar, Fellow. Present address: Michigan State Univ.  
<sup>2</sup>Provided by Lawrence Radiation Laboratory, Livermore.  
<sup>3</sup>S. M. Austin, P. Paul, A. Cheung, S. S. Hanna, *et al.* Paper AA10, this meeting.



**AA12. ( $\text{He}^4, \text{He}^3$ ) and ( $\text{He}^4, \alpha$ ) Reactions on  $\text{Li}^6$  and  $\text{Li}^7$ .** VAN BLUMMEL AND MORROW K. BRUSSEL, *University of Illinois*.—Targets of  $\text{Li}^6$  and  $\text{Li}^7$  were bombarded by  $\text{He}^4$  particles from the University of Illinois cyclotron, at incident beam energies of 21, 24, and 27 MeV. From the  $\text{Li}^6$  target, we obtained angular distributions of the scattered  $\text{He}^3$  particles corresponding to the ground and 2.18-MeV states. From the  $\text{Li}^7$  target, we obtained angular distributions of the scattered  $\text{He}^3$  particles corresponding to the ground, 0.48-, and 4.63-MeV states, and of  $\alpha$  particles leaving the residual  $\text{Li}^6$  nucleus in the ground 2.18- and 3.56-MeV states. Spectra were recorded in the angular range from  $30^\circ$  to  $140^\circ$  in the c.m. system. Detection was by a surface-barrier silicon detector. Targets were of enriched lithium hydride, evaporated onto thin nickel backings.

\* Work supported in part by the U. S. Office of Naval Research.

**AA13. Charged-Particle Reactions from  $\text{Li}^7 + p$  at 30 MeV.** H. H. FORSTER, D. W. DEVINS, AND C. C. KIM, *University of Southern California*.—Self-supporting foils of  $\text{Li}^7$  (purity 99.6%) were bombarded with 30.3-MeV protons from the USC linear accelerator. Angular distributions were obtained for the  $\text{Li}^7(p, \alpha)\text{Li}^4$ ,  $\text{Li}^7(p, \alpha)\text{Li}^3$ , and  $\text{Li}^7(p, pd)\text{He}^3$  reactions at laboratory angles from  $\sim 10^\circ$ – $100^\circ$ . Particle identification was achieved by getting a 2-dimensional analyzer with coincidences between the signals from a  $(dE/dx) - E$  counter system consisting of a 190- $\mu$  fully depleted silicon and a  $\frac{1}{2} \times 1$ -in. NaI(Tl) detector. The resolving time of the fast-coincidence circuit used for the angular-correlation experiment was  $\sim 6$  nsec. In the  $\text{Li}^7(p, \alpha)\text{Li}^3$  reaction, deuterons leaving  $\text{Li}^3$  in the ground state or first two excited states could be resolved and angular distributions were obtained for each deuteron group. In the  $(p, pd)$  angular-correlation experiment, two detectors were placed at equal angles with the incident beam; particle identification was used in one arm only. The resultant angular identification is discussed and compared with that obtained for the  $\text{Li}^7(p, pd)\text{He}^3$  reaction.

\* Work supported in part by the U. S. Atomic Energy Commission.

**AA14. Electric Quadrupole Moments of  $\text{Li}^7$  and  $\text{Cl}^{35}$  Nuclei.** PAUL E. CADE, *The University of Chicago*.—The electric field gradient  $q$  at the Li or Cl nucleus has been calculated for LiH, LiI, LiF, LiCl, and  $\text{LiCl}_2$ , and certain other LiX diatomic molecules using Hartree-Fock-Roothaan SCF wavefunctions obtained by Cade, Wahl, Huo, and Salez. These wavefunctions, except for that of LiCl, are believed to be very close approximations to the true Hartree-Fock wavefunctions. Using the experimental values of  $eqQ$  for either  $\text{Li}^7$  or  $\text{Cl}^{35}$  in these molecules, the electric-quadrupole moments  $Q$  of  $\text{Li}^7$  and  $\text{Cl}^{35}$  are calculated and the internal consistency of these results is considered. The variation of  $q$ , and hence  $eqQ$ , with vibrational state is also discussed on the basis of calculated results.

**AB1. Mean Life of the 1.9-MeV  $0^+$  Level in  $\text{Ce}^{138}$ .** ROBERT L. GRANAM AND JAMES S. GEIGER, *CRNL-Atomic Energy of Canada Limited*.—The mean life of the 1.9-MeV  $0^+$  state in  $\text{Ce}^{138}$  has been obtained by measuring the time correlation between x-rays from K-capture and K-1900 conversion electrons. The conversion electrons were selected using a 7-gap orange  $\beta$ -ray spectrometer located on a beam line of the Chalk River tandem accelerator. The 1.9-MeV  $\text{Ce}^{138}$  level was populated by the decay of 3.4-min  $\text{Pr}^{138}$ , which was produced by periodically bombarding a 5-mg/cm<sup>2</sup> metallic Ce target with 9-MeV protons for 30-sec intervals. The time-

correlation data were accumulated between irradiations. The detectors consisted of Naton plastic scintillators on XP-1020 photomultipliers. The timesorter was calibrated to  $\sim 1\%$  accuracy using an air-cored helical delay line.<sup>1</sup> The mean life obtained for the 1.9-MeV level in  $\text{Ce}^{138}$  is  $\tau = 0.57 \pm 0.04$  nsec. The E2  $\gamma$ -transition rate for the 300-keV  $0^+ \rightarrow 2^+ E2$  transition deduced from this lifetime and the E0, E2 branching ratio of Hisatake et al.<sup>2</sup> is 6 Weisskopf units.

<sup>1</sup>R. L. Graham, J. S. Geiger, R. E. Bell, and R. Barton, *Nucl. Instr. Methods* 13, 40 (1962).  
<sup>2</sup>K. Hisatake, Y. Yoshida, K. Etoh, and T. Murata, *Nucl. Phys.* 68, 520 (1963).

**AB2. Decay of  $\text{Eu}^{145}$ .** JOHN C. HILL, *Purdue University*.— $\text{Eu}^{145}$  sources were obtained from the  $\text{Sm}^{144}(d, n)\text{Eu}^{145}$  reaction and purified chemically. The half-life of  $\text{Eu}^{145}$  was determined to be  $5.93 \pm 0.1$  days. The  $\gamma$  spectrum was studied, using Ge-Li detectors as well as NaI(Tl) detectors. The conversion electron and positron spectra were investigated, using a 180° magnetic spectrometer.  $\gamma$ - $\gamma$  coincidences were observed using both NaI and Ge-Li detectors. A total of 23  $\gamma$  lines were observed; 13 of these have not been reported before. With the help of coincidence and intensity measurements, the transitions were incorporated into a consistent decay scheme. Two low-intensity positron groups, leading to the ground and first excited states, respectively, were observed. The levels of  $\text{Eu}^{145}$ , populated through the decay of  $\text{Eu}^{145}$ , are compared to levels of  $\text{Sm}^{145}$ , populated through the reaction  $\text{Sm}^{144}(d, p)\text{Sm}^{145}$ , observed by Kenefick and Sheline.<sup>1</sup>

\* Work supported by the U. S. Atomic Energy Commission.  
<sup>1</sup>R. A. Kenefick and R. K. Sheline, *Phys. Rev.* 139, 1479 (1963).

**AC1. Conduction-Electron Spin Resonance.** MARTIN LAMPE\* (introduced by P. M. Platzman) AND P. M. PLATZMAN, *Bell Telephone Laboratories*.—We have calculated (using a simple theory) the paramagnetic resonance absorption by conduction electrons in a thin metallic sample, extending the earlier work of Dyson.<sup>1</sup> Our results are valid for a metal sample of arbitrary thickness, with a static magnetic field  $H_0$  at arbitrary angle with respect to the sample surface, and under either classical or anomalous skin-effect conditions. The electromagnetic field is assumed to be incident normally, but not necessarily symmetrically, on both sides of the samples. Under anomalous skin-effect conditions, as the field  $H_0$  is rotated from parallel to normal, the shape of the spin resonance line is modified. For asymmetric boundary conditions on the EM field, the line decreases in intensity (to zero under certain conditions) and broadens drastically. For a symmetric or 1-sided EM field, the line shows a slight decrease in intensity and a slight narrowing. We see no evidence for any structure of the kind found by Walsh<sup>2</sup> in his electron-spin-resonance experiments on potassium. Numerical results are presented.

\* Permanent address: St. Mary's College, Calif.  
<sup>1</sup>P. J. Dyson, *Phys. Rev.* 96, 349 (1953).  
<sup>2</sup>W. M. Walsh, Jr., *Phys. Rev.* (to be published).

**AC2. Pressure Dependence of Knight Shift in Pt.** T. KUSUDA AND L. RIMAL, *Ford Motor Scientific Laboratory*.—The Knight shift  $K$  of  $\text{Pt}^{195}$  in platinum metal was measured under hydrostatic pressure up to about 8000 kg/cm<sup>2</sup> at three temperatures, 64.8°, 0°, and -78.0°C.  $|K|$  decreases with increasing pressure.  $\delta K/\delta P$  is markedly temperature-dependent and changes as  $T^3$  with temperature. The temperature dependence of  $K$  at constant volume, which has been derived from the previously obtained  $K(T)$  at constant pressure and the present pressure-dependence data, obeys a quadratic law,  $K = K_0 + BT^2$

$+O[(T/T_2)^4]$ , where  $T_2$  is the degeneracy temperature of the  $d$ -band holes. The observed volume dependence of  $K_0$  and  $B$  is analyzed using a standard band model. The volume dependence of the band parameters thus obtained is as follows: (i)  $d \log T_2 / \log V = -3.5 \pm 0.9 \sim d \log T_2 / d \log V$ ; (ii)  $d \log T_{ex} / d \log V = -0.9 \pm 1.9$ , where  $T_{ex}$  is an exchange parameter; and (iii)  $d \log n / d \log V = 4.3 \pm 1.3$ , where  $n$  is the number of the  $d$  holes per atom.

**AC3. Spin-Orbit Coupling, Magnetic- and Electric-Field Interactions of Iron ( $3d^6$ ) in Complexes of Tetragonal Symmetry.** GILDA M. HARRIS, *Stanford University*.—Ferric ion ( $3d^5$ ) has a  $^6S$  ground state. To 1st order, this multiplet does not interact under spin-orbit coupling nor is it split by any crystal field. Yet, there is experimental evidence from electron spin and Mössbauer spectra that there is a zero-field splitting of the ground state and a net electric-field gradient at the iron nucleus in complexes of lower-than-octahedral symmetry. In order to possibly explain these observations, a calculation of the energy eigenfunctions and eigenvalues of ferric ion was made in a strong octahedral field with tetragonal distortions and spin-orbit coupling among the 33 lowest-energy states. Two different zero-order basis sets were used. Ground-state splittings, low-spin conditions, electric-field gradients, magnetic susceptibilities, and field energies were obtained. A systematic study was made of the effect on all the above properties of which and how many excited states were included in the calculation, starting with just one of the degenerate partners of the lowest-lying quartet state ( $^4T_1$ ), adding its partners, then the lowest-lying doublet state, and, finally, the rest of the electronic states. This effect was marked and exceeded the difference between the 2nd order and exact calculation made in each of these approximations.

**AD1. Theory of Lattice Vibrations in Wurtzite and Zincblende.** MICHEL A. NUSIMOVIC (introduced by J. L. Birman) AND JOSEPH L. BIRMAN, *New York University*.—Previously,<sup>1</sup> the frequencies and normal modes for CdS in  $C_{6v}$  structure were calculated by solving the dynamical matrix for propagation along three directions of the wurtzite Brillouin zone. A more detailed study of this problem was made using five models of varying degrees of elaboration: (1) central forces between 1st neighbors; (2) central forces between 1st and 2nd neighbors; (3) central forces between 1st, 2nd, and 3rd neighbors; (4) central and noncentral forces between 1st, 2nd, and 3rd neighbors; (5) central and noncentral forces between 1st, 2nd, 3rd neighbors, and long-range Coulomb forces. Using appropriate coordinate transformations,<sup>2</sup> frequencies and normal modes of zincblende of  $T_d$  structure were calculated. Dispersion curves have been calculated for  $\epsilon$  (hypothetical) cubic CdS. Some comparison of the present model and the shell model is given.

\* Work partially supported by the U. S. Army Research Office (Durham) and the Aerospace Research Laboratories, Office of Aerospace Research, Wright-Patterson AFB.  
† Permanent address: Lab. de Physique de l'École Normale Supérieure, Paris.  
‡ M. Nusimovic and J. L. Birman, *Bull. Am. Phys. Soc.* 10, 616 (1965).  
§ J. L. Birman, *Phys. Rev.* 113, 1093 (1959).

**AD4. Electronic Transport in Graded Heterojunctions.** L. J. VAN RUYVEN AND F. E. WILLIAMS, *University of Delaware*.—We consider a single crystal of a semiconductor whose composition gradually changes from that of a small-bandgap material at one end to that of a large-bandgap material at the other. Although the quantum-mechanical problem of a charged

particle in a varying periodic potential has not been solved rigorously, this graded heterojunction has, in an approximate sense, a graded bandgap. In order to separate the effect of the electric field due to space charges from the effect of the graded bandgap on electronic charge transport, we consider a graded heterojunction that is doped inhomogeneously so that there is no space charge. The motion of electrons at the conduction-band edge and the motion of holes at the valence-band edge are analyzed, and in addition to the normal diffusion term an additional term is obtained that is linear in the gradient of the band edge for each type of carrier. Steady-state photoexcitation in one region leads under certain conditions to both types of carriers moving in the same direction, i.e., to the region of smaller bandgap. In addition, the motion of excitons in graded heterojunctions is considered.

\* Work supported by U. S. Army Engineer Research and Development Laboratories.

**AD5. Minority-Carrier Mobility in  $p$ -Type Germanium under High Uniaxial Stress at Room Temperature.** M. W. CRESSWELL AND J. P. MCKELVEY, *Pennsylvania State University*.—The minority-carrier mobility in a single crystal of germanium containing  $3 \times 10^{18}$  indium atoms per cc has been measured as a function of uniaxial compression in a lattice (111) direction at room temperature by observing the drift of an excess carrier concentration under the influence of an applied electric field. The maximum fractional change in length of the crystal, cut in the shape of a rod, was 0.99%, corresponding to a compressive pressure of approximately  $1.5 \times 10^8$  dyn/cm<sup>2</sup>, sufficiently high to transfer substantially all conduction electrons to a single valley. A comparison of the experimental data with theoretical estimates of mobility variation due to strain-induced population transfer in the conduction band yielded a value for the scattering anisotropy of  $\mu_r = 3.7 \pm 0.3$ . This value is considerably higher than that obtained from studies of magnetoresistance of  $n$ -type germanium, the difference being attributed to electron-hole recombination.

\* Work supported by the U. S. Office of Scientific Research.  
† Present address: Westinghouse Res. Labs., Pittsburgh.

**AD6. Strain-Induced Freeze-Out Effect in  $n$ -GaSb.** A. SAGAR, W. LEHMANN, AND M. POLLAK, *Westinghouse Research Laboratories*.—The effect of hydrostatic pressure on resistance and Hall coefficient of "undoped"  $n$ -ZnSe (Ref. 1) was measured at 195° and 300°K. The effect of uniaxial stress on resistance was also measured between 77° and 30°K. Temperature dependence of Hall coefficient and resistance of the samples (between 77° and 300°K) was similar to the "undoped" samples in Ref. 1. The resistance, Hall coefficient, and their derivatives w.r.t. pressure increased with increasing pressure so that  $R_H/\rho = \text{const.}$  At 195°K,  $R_H$  and  $\rho$  increased by 80% from one to 10<sup>4</sup> atm for a sample with  $R_H(77^\circ\text{K}) \approx 5 \times 10^4$  cm<sup>2</sup>/C. The piezoresistance value due to uniaxial stress was  $\frac{1}{2}$  of the piezoresistance value owing to hydrostatic pressure at low pressures. This indicates that the effect is a purely hydrostatic one and suggests that the mechanism responsible is a pressure-induced variation of the ionization energy. Similar observation has been made in  $n$ -GaAs by Sladek.<sup>2</sup> Our data give the unusual result that the value of the ionization energy ( $E_i \approx 0.008$  eV at normal pressure) is increased by a factor of two at 10<sup>4</sup> atm.

† M. Aven and B. Segal, *Phys. Rev.* 130, 81 (1963).  
‡ R. J. Sladek, *Phys. Rev.* (to be published).

**AD7. Avalanche Breakdown in *p*-Type GaAs.** KURT WEISER, ROBERT E. FERN, AND JOSEPH F. WOODS, *IBM Watson Research Center.*—We have continued the study of avalanching and light emission of thin high-resistivity layers<sup>1</sup> in *p*-type GaAs structures produced by diffusing Zn into Mn-doped material. Such layers, with a resistivity of typically 5000  $\Omega$ -cm, are formed at the boundary of the zinc-dominated surface layer. Capacitance measurements, as well as examination of transmitted light by phase-contrast microscopy,<sup>2</sup> indicate that the width of these layers is of the order of 1  $\mu$ . A drastic increase in current is observed at a field of approximately  $2 \times 10^5$  V/cm. Light emission is then proportional to  $(i - i_0)$ , where  $i$  is the experimental current and  $i_0$  is the Ohmic component as extrapolated from low field values. The electrons generated by the avalanching process are swept out of the layer and recombine in the zinc-rich region or the manganese region to either side of it, depending on the current direction. The spectral distribution of the emitted light differs accordingly for the two cases. The quantum efficiency of the light emission is comparable to that of ordinary GaAs diode though the power efficiency is much lower because of the high voltage (typically 16 V) needed to produce the light.

<sup>1</sup> K. Weiser and J. F. Woods, *Appl. Phys. Letters* 7, 225 (1965).  
<sup>2</sup> M. Drougard, private communication.

**AD8. EPR and Electrical Properties of the Dominant Defect in Electron-Irradiated *p*-Type Silicon.** NISSIM ALKELER AND BERNARD GOLDSTEIN, *RCA Laboratories.*—When *p*-type silicon is bombarded with electrons, the dominant defect formed is the *K* center.<sup>1</sup> We have studied the paramagnetic properties and growth of this center as functions of electron flux and bombardment energy under conditions of different resistivities, impurity dopants, Fermi level, and illumination. Introduction rates,  $g$  values, and symmetry properties are presented and discussed. The *K* center is independent of the *p*-type dopant. It is not a primary defect, but requires oxygen. At high integrated electron fluxes, the EPR-measured *K*-center concentration decreases; however, illumination and annealing experiments have established that the defects are still present but have a different charge state because they have trapped an electron. We have associated the *K* center with a previously reported 0.3-eV defect level<sup>2</sup> based on the facts that both require oxygen, both have about the same introduction rates and bombardment energy dependence, and that the value of the Fermi level at which the *K*-center EPR absorption decreases sharply is about 0.3 eV.

<sup>1</sup> Work sponsored by NASA-Goddard Space Flight Center.  
<sup>2</sup> V. Shahan *et al.*, *Intern. Conf. Phys. Semiconductors, Paris, 1964*.  
<sup>3</sup> *Space Technol. Lab. Rept. MR-32, contract NAS 5-1851.*

**AD9. Photoconductive Properties of High-Resistivity Gallium Phosphide.** BERNARD GOLDSTEIN AND S. S. PERLMAN, *RCA Laboratories.*—Static and dynamic photoconductive properties of single-crystal high-resistivity (compensated) GaP have been studied in the intrinsic and near-infrared spectral region at 300°, 77°, and 27°K. Room-temperature resistivities in excess of  $10^4$   $\Omega$ -cm have been produced by copper diffusion into either *n*- or *p*-type GaP. Photoconductivity of the high-resistivity material is strongly influenced by traps. In *n*-type material, at 300°K, infrared radiation stimulates the dark conductivity and quenches the intrinsic photoconductivity at the same photon energies; thermal-probe measurements indicate that the stimulation is due to increased hole current. At 27°K, only stimulation is observed regardless of the level of intrinsic photoconductivity, but decay characteristics show that this latter response is due to increased electron current.

At the intermediate temperature of 77°K, stimulation (electron and hole) and quenching are present and the over-all behavior is more complex. In *p*-type material, infrared radiation produces only stimulation. An energy-level diagram is presented, which can explain these observations: principal features include electron traps about 0.6 eV below the conduction band, recombination levels near the center of the bandgap, and "sensitizing" hole traps about 0.7-0.8 eV above the valence band.

**AD10. Photovoltaic and Photocapacitive Properties of Surface-Barrier Junctions of High-Resistivity Gallium Phosphide.** S. S. PERLMAN (introduced by Bernard Goldstein) AND BERNARD GOLDSTEIN, *RCA Laboratories.*—Measurements of photovoltaic currents and photocapacitance of surface-barrier junctions in high resistivity *n*- and *p*-type GaP are presented and shown to be of unique value in the study of photoelectronic processes. In particular, they provide means for the direct determination of the polarity of charge carriers released by an infrared transition. This follows from an extension of the accepted model of intrinsic and extrinsic photovoltaic and photocapacitive response, which considers effects such as infrared quenching and infrared stimulation of intrinsic response. In *n*-type material, measurements indicate that intrinsic photocapacitance is quenched by infrared radiation. At the same time, extrinsic photovoltaic response is stimulated by intrinsic radiation. In *p*-type material, intrinsic photovoltaic response is quenched by infrared radiation. All the spectral dependence curves of these effects are compared for optical-threshold energy with similar curves of photoconductivity quenching<sup>1</sup> and are found, as predicted by the model, to represent the same basic electronic transition. The unambiguous determination of photovoltaic current and photocapacitance supplies independent verification that this transition does, indeed, create positive charge carriers (holes) and involves trapping levels located approximately 0.8 eV above the valence band.

<sup>1</sup> B. Goldstein and S. S. Perlman, Paper AD9, this meeting.

**AD12. Dislocation Degeneracy in Heavily Doped Germanium and Silicon Single Crystals.** G. H. SCHWUTTER,\* *IBM East Fishkill*, AND R. GERSTE, *IBM Watson Research Center.*—Recent investigations of dislocation structures in heavily doped silicon single crystals report definite ranges of donor-impurity concentrations in which dislocations become degenerate.<sup>1</sup> Consequently, it should be possible to grow heavily doped *n*-type silicon crystals dislocation-free without any special precautions. This conclusion is not consistent with the model outlined by Dash for dislocation propagation in high-resistivity silicon crystals and also not consistent with the experimental findings of Patel *et al.* for heavily doped germanium crystals. In view of the important consequences of Mil'vidskii's work, the perfection of heavily melt-doped *n*- and *p*-type silicon crystals<sup>2</sup> is investigated through x-ray-diffraction microscopy. Test samples are cut parallel to the pull axis and large-area x-ray topographs are recorded. Our measurements confirm for silicon Patel's work. *p*-type silicon crystals are essentially free of dislocations, while in *n*-type crystals dislocations are present. The topographs also reveal impurity striations and strong impurity cores inside the crystals.

\* Work supported by the U. S. Air Force.  
<sup>1</sup> M. G. Mil'vidskii *et al.*, *Soviet Phys.-Solid State* 8, 2606 (1965).  
<sup>2</sup> Crystals kindly supplied by the Shockley Laboratory, Palo Alto, Calif.

**AE9. Comparison of the Scattering of 1200-MeV Electrons and Positrons from Protons.\*** R. L. ANDERSON, BRUNO BORGIA, AND J. W. DEWIRE, *Cornell University*.—An electron or positron beam with a 10% momentum spread is obtained from the bremsstrahlung beam of the Cornell synchrotron by using a system of magnets to select and focus members of electron pairs. The beam passes through a thin-wall 45-cm liquid-hydrogen target to a quantameter for monitoring. Scattered electrons and recoil protons in coincidence are observed in spark chambers placed on either side of the beam. The polar angles of pairs of coplanar tracks are fitted to the scattering kinematics to select events of elastic scattering. The system is designed to eliminate insofar as possible any differences in the detection of positrons and electrons. In particular, the problem of distinguishing between scattered positrons and positive pions, which was present in the Stanford experiments,<sup>1</sup> has been avoided. Analysis of a run yielding 5200 elastic scatterings in the interval  $10 \leq q^2 \leq 20 \text{ fm}^{-2}$  shows the two cross sections to be equal within statistics. Additional data are being obtained.

\* Work supported in part by the National Science Foundation.  
<sup>1</sup> A. Browman, F. Liu, and C. Schaefer, *Phys. Rev. B139*, 1079 (1965).

**AE10. Comparison of Wide-Angle Electron Pair Production with the Predictions of Quantum Electrodynamics.\*** E. EISENHANDLER, J. FEIGENBAUM, N. B. MISTRY, P. J. MOSTEK, D. R. RUST, A. SILVERMAN, C. K. SINCLAIR, AND R. M. TALMAN, *Cornell University*.—The photoproduction of electron pairs has been studied at several photon energies ranging from 500 to 1800 MeV, using a large uniform-field magnet and spark chambers. The  $e^+$  and  $e^-$  production angles accepted by the apparatus ranged from  $6^\circ$  to  $10^\circ$ . The momentum acceptance of the magnet and the bremsstrahlung peak energy were varied in such a way as to keep their ratio constant. A comparison is made with quantum electrodynamics. The deviation from Q.E.D. observed by Blumenthal *et al.*<sup>1</sup> would lead to an approximate discrepancy of 15% in the energy range of this experiment. Analysis of the data is in progress and results are presented.

\* Work supported in part by the National Science Foundation.  
<sup>1</sup> R. B. Blumenthal *et al.*, *Phys. Rev. Letters* **14**, 662 (1965).

**AE11. Muon-Proton Elastic Scattering.\*** H. VON BRIESEN, JR. (introduced by T. Yamanouchi), R. W. ELLSWORTH,† A. C. MELISSINOS, J. H. TINLOT,‡ T. YAMANOUCI,§ *University of Rochester*, L. M. LEDERMAN, M. J. TANNENBAUM,§ *Columbia University*, R. L. COOL, and A. W. MASCHKE, *Brookhaven National Laboratory*.—Preliminary results of the Brookhaven-Columbia-Rochester muon-proton elastic-scattering experiment<sup>1</sup> have been reported for 4-momentum transfers of  $q^2 = 12 - 31 \text{ fm}^{-2}$ . We present here the results for the low-momentum transfer region ( $q^2 = 8 - 19 \text{ fm}^{-2}$ ). The purified muon beam from the AGS having momenta  $p = 1.5 - 6.0 \text{ BeV}/c$  was incident on a 6-ft liquid-hydrogen target. Both incident and scattered muon tracks and the recoil-proton track were photographed in spark chambers. The recoil-proton energy was measured by its range in a heavy-plate aluminum spark chamber. Of about 28 000 triggering events, we have identified about 900 elastic events. The differential cross-section  $d\sigma/dq$  is presented. It is compared with the results of the high-momentum-transfer run as well as with  $e-p$  scattering data.

\* Work supported in part by the U. S. Atomic Energy Commission.  
 † Present address: Univ. Washington.  
 ‡ Deceased.  
 § Present address: CERN, Geneva.  
<sup>1</sup> R. Cool *et al.*, *Phys. Rev. Letters* **14**, 724 (1965).

**AE12. Pion Form Factor from Electroproduction.\*** W. W. ASH (introduced by K. Berkelman), C. W. AKERLOF, K. BERKELMAN, AND C. A. LICHTENSTEIN, *Cornell University*.—

We have been using the circulating electron beam of the Cornell 2-GeV synchrotron and an internal liquid-hydrogen target to study the reaction  $e + p \rightarrow e + \pi + p$ . The inelastically scattered electrons are momentum-analyzed in a quadrupole magnet and detected by a telescope of scintillators followed by a lead-glass Cerenkov counter. Electroproduced pions emitted along the electron-momentum-transfer direction are analyzed in a similar magnet and detected in coincidence by a scintillator telescope. The transverse and longitudinal virtual photon contributions to the electroproduction yield are separated by taking data for corresponding 4-momenta transfer  $k^2$  and pion-nucleon c.m. energy  $W$  at several different electron-scattering angles between  $15^\circ$  and  $55^\circ$  lab. Data have been taken at  $-k^2 = 3.0 \text{ F}^{-2}$  and  $W = 1200$  and  $1300 \text{ MeV}$ . The longitudinal contribution is interpreted so as to place experimental limits on the pion electromagnetic form-factor.

\* Work supported in part by the National Science Foundation.

**AE14. Electron-Induced Cascade Showers in Copper and Lead at 1 BeV.\*** WALTER RALPH NELSON, RICHARD C. MCCALL, JOSEPH K. COBB, AND THEODORE M. JENKINS, *Stanford Linear Accelerator Center*.—The longitudinal and radial development of electron-photon showers has been measured in copper and lead at 1 BeV. A new technique, described in an earlier paper,<sup>1</sup> using the thermoluminescent property of LiF has been employed to measure energy deposition. The resultant radial distributions and transition curves are compared with Monte Carlo calculations and other experiments. The fraction of incident energy that leaks out of a cylinder of radius  $r$  (radiation lengths) is plotted against  $r/\epsilon_0$ , where  $\epsilon_0$  is the critical energy of the absorbing material, and it is observed that most of the existing statistical and experimental data, including this experiment, fall on an empirical curve—regardless of incident energy or choice of absorber.

\* Work sponsored by the U. S. Atomic Energy Commission.  
<sup>1</sup> T. M. Jenkins *et al.*, *Nucl. Instr. Methods* (to be published).

**AE15. Search for Dirac Monopoles Produced by the Cosmic Radiation.** W. C. CARITHERS (introduced by R. K. Adair) AND R. J. STEFANSKI, *Yale University* and *Brookhaven National Laboratory*.—The simplest model of an electric charge and a magnetic monopole violates time-reversal invariance. Since the Fitch-Cronin-Turley effect suggests that time-reversal invariance is violated, it appeared desirable to conduct a further search for heavy magnetic monopoles. The negative results of such a search allow us to place an upper limit for the flux of monopoles in the atmosphere at  $R \leq 3 \times 10^{-16}$  monopoles/cm<sup>2</sup>-sec. The experimental design is similar to that of Malkus. Monopoles created by cosmic rays high in the atmosphere diffuse along the geomagnetic field lines and are accelerated and focused onto a nuclear-emulsion rack by a large solenoid magnet (peak field = 13 kG, magnetic moment =  $1.3 \times 10^8 \text{ G/cm}^3$ ). The final trajectory is defined by coincidence counters and a spark chamber. Assuming a specific model for monopole production by very high energy nucleon-nucleon collisions, we convert our rate into a cross section. For a 15-BeV/c<sup>2</sup> monopole,  $\sigma \leq 2 \cdot 10^{-26} \text{ cm}^2$ , or  $10^{-3} \text{ (A/Mcp)}^2$ .

**AH1. Spectroscopy of Interstellar Grains.** F. M. JOHNSON, *Electro-Optical Systems*.—Distinctive but previously unrecognized patterns corresponding to vibrational-energy repartitions within electronic bands have been identified among 18 diffuse interstellar lines. Such patterns will facilitate the search for the chemical origin of the lines. Since there is a known correlation between the strength of these absorption lines and stellar reddening, it is believed that they have their origin in interstellar grains. Thus, assignments of the lines may be tantamount to a partial chemical identification of these grains. Each of the 18 lines falls into one of three groups. The 1st and 2nd groups occur in the vicinities of 4600 and 6100 Å, respectively, with the former the more diffuse. Both of these groups are comprised of a set of 3 lines with wavenumber separations of 556 and 1568  $\text{cm}^{-1}$ . The remaining 12 lines fall within a 3rd group, also near 6100 Å, whose mean wavenumber separations can be sorted into 226  $\text{cm}^{-1}$  or multiples thereof. A discussion is given of various types of chemical species that appear most likely to be identified with these correlations.

**AH2. Inhomogeneous Cosmological Expansion, Quasistellar Sources, and Quasistellar Galaxies.** Y. NE'EMAN, *Tel-Aviv University*.—We discuss the hypothesis that quasistellar radio sources are fed by the decay of high-energy particles, produced in superdense conditions corresponding to the cosmological preexpansion stages. These are independent cores whose expansion has lagged behind. Some comments are made with respect to the possible rôle of quasistellar galaxies and the hypothesis of an oscillating model.

<sup>1</sup> I. Novikov, *Astronom. J. (USSR)* 41, 1075 (1964).

**AH3. Scattering of Slow Electrons by Enhanced Ion Waves near the Geomagnetic-Field Boundary.** AHARON EVIATAR, *University of Maryland*.—The scattering of low-energy electrons by ion-plasma oscillations in a stable plasma containing both suprathermal particles and a current is considered using the linearized Balescu-Lenard equation. The time required for such waves to scatter particles through 90° is estimated and compared to the Coulomb scattering time. This time will be short as compared to the Coulomb time or other wave-particle scattering times for slow particles, if conditions exist that inhibit Landau damping of ion waves. This can explain enhanced diffusion of slow electrons across surfaces of discontinuity in the Earth's magnetic field. Serbu<sup>1</sup> has observed 1- to 2-eV electrons whose density distribution shows no marked variation at either the magnetopause or the shock wave. We suggest that this is a result of scattering by ion waves excited by the observed fast particles coexisting with the plasma. The electrostatic oscillation spectrum has resonances at the electron and ion plasma frequencies. Observations of temperatures and flow velocities in the transition zone indicate that the conditions required by this theory are satisfied without attaining the extreme values required for instability of ion modes.

\* Work supported in part by the National Aeronautics and Space Administration.  
<sup>1</sup> G. P. Serbu, *Space Research* (North Holland Publ. Co., Amsterdam, 1963), Vol. 3.

**AH4. Distribution of Neutral Hydrogen above 120 km.** MORDEHAI LIWISITZ, *NASA-Goddard Space Flight Center*.—Results of a new study of the hydrogen distribution in the thermosphere are presented. These are based on a solution of the diffusion equation for a minor constituent through an ambient stationary atmosphere, taking into account the effect of evaporative loss in slowing down the effusion at the

base of the exosphere. Results of the calculation reveal that the "total" hydrogen content of the atmosphere above 120 km is increased by a factor of ~1.1 to ~2.0 in the range of temperatures considered (1000°–2500°K). The higher hydrogen abundance, as well as the greater amplitude of variation with temperature appear to come closer to an explanation of the observed Lyman  $\alpha$  radiation.

\* National Academy of Science—National Research Council Resident Research Associate.

**AH5. Hit Probability of Interplanetary Dust Particles in the Vicinity of Earth.** J. WILLIAMS RUSSE, *Ross Polytechnic Institute*.—The advent of vehicles into space has brought into sharp focus some of the hazard problems to be surmounted. One among these many problems is interplanetary dust particles and the hazard that it presents to spacecraft. In this regard, the question always arises as to the probability that the body will be hit by interplanetary dust particles. An attempt has been made here to answer some of these questions, at least in a preliminary way. The probability  $p(x)$  that a hit will occur exactly  $x$  times in time interval  $t$  is given by  $p(x) = (\bar{N})^x e^{-\bar{N}} / x!$ , where  $\bar{N}$  is the average number of hits. General expression of  $\bar{N}$  is  $\bar{N} = \sigma t I$ , where  $\sigma$  is the cross-sectional area of the vehicle,  $t$  the time interval, and  $I$  is the omnidirectional intensity of the particle. A preliminary computation shows that these hit probabilities are very small unless it is encountered by a stream of dust particles.

**EB1. Binding Energy of the Nucleus.** OLIVER K. MAMUEL (introduced by James Paul Wesley), *University of Missouri, Rolla*.—The nuclear binding energy of the nuclide ( $Z, A$ ) is presently defined as its stability relative to ( $A-Z$ ) neutrons and  $Z$  hydrogen atoms. This definition results from a model where the Yukawa exchange interactions are ignored and the nucleus is assumed to be composed of  $Z$  protons and ( $A-Z$ ) neutrons. The resulting binding energies of different nuclides are calculated relative to different standards, except in the special case where their assumed neutron-proton ratios are identical. Thus a comparison of the binding energies of the different nuclides has little or no meaning. Since the Weizsäcker equation considers only the intrinsic stability of the nucleus due to forces acting between the nucleons, this cannot be used to calculate the binding energy as presently defined. It is therefore proposed to define the binding energy of the nuclide ( $Z, A$ ) as its stability relative to  $A$  neutrons. The results of the two definitions of binding energy are discussed.

**EB2. Comparison of Theoretical Internal-Conversion Coefficients for Magnetic Multipoles ( $Z=39$ ).** C. P. BHALLA, *University of Alabama, Huntsville*.—New calculations of internal-conversion coefficients for  $k=0.15 \text{ mc}^2$  have been completed for the magnetic multipoles for  $Z=39$ . The present calculations are based on the following realistic model.<sup>1</sup> Atomic-screening effects both for continuum and the bound states are included by self-consistent Hartree-Fock treatment and the finite nuclear size effects are also included. The results of present calculations are compared with those of Rose<sup>2</sup> and of Siv and Band.<sup>3</sup>

\* Work supported in part by the National Aeronautics and Space Administration.  
<sup>1</sup> C. P. Bhalla, in *Internal Conversion Process*, J. H. Hamilton, Ed. (to be published), C. P. Bhalla, Conf. Nucl. Particle Phys., Sept. 1963, Liverpool.  
<sup>2</sup> M. E. Rose, *Internal Conversion Coefficients* (1955).  
<sup>3</sup> L. A. Siv and I. M. Band, in *Alpha, Beta, and Gamma Ray Spectroscopy* (1963).

**EB3. Theory of Low-Lying Spectra of Odd-Mass Nuclei.** A. I. SHERWOOD, *University of California, La Jolla*, AND A. GOSWAMI, *Western Reserve University*.—The low-lying states of an odd-mass nuclei are considered to be states of an odd (quasi) particle strongly coupled to an even-even core. The quasiparticle Hamiltonian is obtained by performing the Bogoliubov transformation. The coupled linearized equation for the amplitude of excitation of the states of the odd-mass nuclei are obtained by considering the equation of motion of a quasiparticle employing an expansion in terms of a complete set of states of the even-even core. This procedure yields a Hermitian matrix in contrast to the higher random-phase approximation, which gives non-Hermitian matrices. The states of the core-even core are treated by the (quasi) boson approximation. Numerical calculations are presented and discussed.

\* Work partially supported by the National Science Foundation.

**EB4. Theoretical Investigation of the Nuclear Properties of the Odd-Mass Pm Nuclei.** T. F. O'DWYER AND D. C. CHOUDHURY, *Polytechnic Institute of Brooklyn*.—The theory of the intermediate coupling approach in the unified nuclear model<sup>1</sup> is applied to analyze the low-energy nuclear properties of the odd-mass Pm nuclei. For this purpose, it is assumed that the last odd proton, having available the  $1g_{7/2}$  and  $2d_{5/2}$  states, is coupled to the collective surface vibrations of the even-even core. The resulting Hamiltonian of the coupled system is diagonalized including all states with up to 3 phonons of the quadrupole vibrations. With reasonable values for the coupling strength and for the "effective" spacing in energy between the  $g_{7/2}$  and  $d_{5/2}$  states, the calculated energy levels are in good agreement with the recent experimental data.<sup>2</sup> Other nuclear properties are also calculated and compared with the available experimental data.

<sup>1</sup> A. Bohr and B. R. Mottelson, *Mat. Fys. Medd. Dan Vid. Selskab.* 27 No. 16 (1953); D. C. Choudhury, *ibid.*, No. 4 (1954).  
<sup>2</sup> K. P. Gopinathan and M. C. Joshi, *Phys. Rev.* B134, 297 (1964); D. B. Fossan *et al.*, *ibid.* B140, 1 (1965); W. M. Currie and P. W. Dougan, *Nucl. Phys.* 61, 561 (1965); C. H. Chen and R. G. Arns, *ibid.* 63, 233 (1965).

**EB5. Zero-Range Surface Interaction for Closed-Shell Nuclei.** J. LETOURNEUX AND J. M. EISENBERG, *University of Virginia*.—Recently, it has been suggested<sup>1</sup> that a zero-range surface interaction may be appropriate for the description of low-energy spectrum of nuclei. Such a force is used here, together with the particle-hole formalism, to discuss excitations in closed-shell nuclei. This leads to considerable simplifications in comparison with conventional treatments, since the diagonalization of large secular matrices is replaced by the solution of dispersion relations. When isospin is a good quantum number, the formalism is essentially the same as that developed by Goswami and Pal.<sup>2</sup> The model is tested by applying it to the light nuclei, and, in addition, a detailed study of  $Pb^{208}$  is presented. The results compare quite favorably with those obtained<sup>3</sup> using a more realistic interaction. In particular, the positions of calculated levels below 4 MeV agree to within better than 0.1 MeV.

\* Work supported in part by the U. S. Atomic Energy Commission. Grants from the National Science Foundation and from Research Corporation for computer time are acknowledged.

<sup>1</sup> I. M. Green and S. A. Moszkowski, *Phys. Rev.* B139, 1790 (1965).

<sup>2</sup> A. Goswami and M. F. Pal, *Nucl. Phys.* 35, 544 (1962).

<sup>3</sup> Gillet, Green, and Sanderson, *Phys. Letters* 11, 44 (1964).

**EB6. Applications of the Neutron-Proton Pairing Theory.\*** H. T. CHEN (introduced by L. S. Kisslinger) AND A. GOSWAMI, *Western Reserve University*.—The Pal-Goswami method is applied for the calculation of neutron-proton ( $n-p$ ) pair correlation effects in the nuclei of Ni-Zn region. It is shown that the  $n-p$  pair correlation effects are important only for nuclei of  $A \leq 70$ . Interesting results are obtained for odd-mass nuclei, where the features of the low-lying states can be explained as an interplay of the relative Fermi energies of the neutron and proton and the  $n-p$  pair correlation effect. The most spectacular results are achieved for the odd-mass Zn and Ga isotopes where the theory predicts very-near-lying levels of same  $j$  as observed experimentally. For the even-even nuclei, the effect of quadrupole force is also included and good agreement is obtained with experimental data as regards the 1st  $2^+$  excited states.

\* Work supported in part by the National Science Foundation.

<sup>1</sup> M. K. Pal, in *Proceedings of Low Energy Conference, Bombay, 1963* (Atomic Energy Commission, Government of India 1963); A. Goswami, *Nucl. Phys.* 60, 228 (1964).

**EB7. Nucleon-Nucleon Force, Single-Particle Energies and Effective Interaction for  $O^{16}$  and  $F^{19}$ .** T. T. S. KEO AND G. E. BROWN, *Princeton University*.—Using the reaction-matrix theory, the Hamada-Johnston potential<sup>1</sup> was used in calculating the properties of  $O^{16}$  and  $F^{19}$ . The reaction matrix elements were computed using the separation method<sup>2</sup> for SE and TE potentials, and reference spectrum method<sup>3</sup> for SO and TO potentials. The  $V_i(Q/e)V_i$  contributions for TE tensor potential were computed using the closure approximation.<sup>4</sup> Dispersion and Pauli correction terms were investigated, but not included in the calculation. Using the linked cluster perturbation formalism, single-particle energies of  $1d_{5/2}$ ,  $2s_{1/2}$ , and  $1d_{3/2}$  were obtained by letting the valence neutron interact with the  $O^{16}$  core. Results were very encouraging. The doublet splitting was found to come overwhelmingly from the TO's force through the 1st order H-F process. The spectra of  $O^{16}$  and  $F^{19}$  were obtained by diagonalizing the effective interaction  $G_M (=G\Omega_M)$  in the  $n-d$  shell. Among the extra configurations taken care of by the wave operator  $\Omega_M$ , the  $3p-1h$  core-polarization processes were the most important. Resulting spectra were fairly satisfactory for  $O^{16}$  but less satisfactory for  $F^{19}$ .

<sup>1</sup> T. T. S. Kuo and G. E. Brown, *Phys. Letters* 10, 54 (1965).