

Experience with the Mechanized Chemical and Biological Information Retrieval System

D. P. JACOBUS, D. E. DAVIDSON, A. P. FELDMAN, AND J. A. SCHAFER
Division of Medical Chemistry, Walter Reed Army Institute of Research,
Walter Reed Army Medical Center, Washington, D. C. 20012

Received July 2, 1969

New computer methods have been developed in association with the drug development programs of the Walter Reed Army Institute of Research. Experiences with these systems are recounted. Special input devices and computer programming have been developed for the input and retrieval of conventional chemical structural diagrams. The costs, operation, and the advantages of this system are discussed. Associated files of biological properties and inventory control information have been created, which are searchable. The methods used in creating consolidated listings of selected chemical compounds and associated biological data are discussed.

A mechanized chemical and biological information retrieval system has been developed for monitoring drug development programs sponsored by the Walter Reed Army Institute of Research.

The system retrieves chemical structures according to either their chemical or their biological properties. These two requirements are met in separate ways. The criteria for chemical retrieval are structural identity, inclusion of a specified fragment, or the presence of given screens. Specialized programs have been developed for retrieval by these chemical criteria. Searches from the biological point of view are similarly guided by their own criteria, such as the kind of test, the type and degree of activity, the origin of compounds, etc. These biological and administrative criteria are handled by more conventional computer methods, as are the correlation of files by means of accession numbers.

Based on our own early experience, the experience of Bello,¹ and a general impression founded upon the "Survey of the Chemical Notation Systems,"² we estimated that the cost of handling chemical information would amount to 5% of the total drug development program. Our later experience has shown that this figure is accurate or reasonable only if the figure is averaged over a four-year span. The programming costs amount to about 10% of the first year's effort for the initial phase of the program, followed by 5% the second year, and then by 4% and 2% for the next two years. If this initial heavy commitment in information retrieval is not made, then test systems tend to become separated from the retrieval programs, and the retrieval programs are essentially a mere recording of summarized information. We have also found that a detailed inventory-control program must be coordinated with the chemical and biological programs to insure proper management of the entire drug development effort.

In a conglomerate system of this nature, the utility of each component system depends on its own performance, as well as on the smoothness of its integration with the other programs. In this context, the five basic problem areas, in our opinion, are the following:

- I. Inexpensive input
- II. Interpretation of input and file maintenance
- III. Synonym files
- IV. Questions
- V. Output

INEXPENSIVE INPUT

Although chemical input is discussed in greater detail below, the input method for biological and conventional data deserves brief mention. The methodology we use is somewhat different from that reported by Waldo,³ who keypunches laboratory results directly from data sheets. Our philosophy, however, is the same as his. Our original investigations and observations are recorded in a format suitable for direct input. The objective is to minimize human error, to reduce cost, and to provide the original investigator with summarized data which he can evaluate for accuracy before the information is transferred to the central files. The direct recording of all observations also makes it possible to send the submitter of the chemical a report of all the details of the data, rather than just a summarization.

Inexpensive chemical input has been achieved by the use of the chemical typewriter.⁴ This machine, by virtue of recording coordinates, provides discrete locations for the characters typed, and therefore digitizes a chemical structure. Important to its operation is the production of hard copy (sheets or 3- × 5-inch cards) for the typist to inspect and to make corrections on. This hard copy is also used by the person responsible for the input, for checking errors discovered in machine-processing. Such compounds are then retyped. Since some incorrectly drawn structures require the staff to consult with the submitter of the chemical, we cannot imagine a system operating without the production of hard copy at the time of input. The typewriters, manufactured by the Mergenthaler Co., have been relatively satisfactory. Our typists produce approximately 1000 correct structures per 40-hour week, or 50,000 per machine per year per shift, at the cost of approximately 10 cents per entry. Overhead figures

vary according to the scheme of operation and raise the cost to 15-25 cents per structure. Of the total input typed, approximately 15% of the structures are rejected because of typing errors. One and one-half per cent of the total material is rejected by the system because of incorrect chemistry. These rejection rates and the machine downtime are not considered in the over-all input capability of 1000 structures per week.

For a typist to type a structure, the chemical diagram must first be made available. If only names are available, a manual conversion to a chemical diagram is required. Considering overhead rates, the conversion of a name to a chemical diagram has cost between 50 cents and one dollar over a broad range of compounds.

We have developed rules for the typing of chemical structures. While we believe most of these rules represent elementary precepts of chemistry, some of these rules may be peculiar to our system. It is our opinion that differences in the interpretation of conventionally-used chemical diagrams will cause much more serious problems of compatibility between systems than differences due to the nature of the various systems. This is because many conventional diagrams are, in a wider context, ambiguous. We feel, therefore, that these rules represent mostly "good" chemistry, and need to be understood and agreed upon. With experience and advanced programming, the number of rules for proper machine analysis is decreasing, although some rules are maintained to keep the hard copy looking "attractive." Examples of typing rules include the prohibition of typing fractions in the molecular formula or structural formula, the requirement of one bond per substituent to show attachments to an atom, and the typing of hydrogen atoms when necessary to satisfy valency requirements for hetero atoms in a ring system.

INTERPRETATION OF INPUT AND FILE MAINTENANCE

The tape from the chemical typewriter must be interpreted in order to produce a conventional connection-table suitable for machine processing (Figure 1). This interpretation problem exists not only for typewriters, but also for the interpretation of input from optical scanners or from cathode ray consoles. We are not aware of any other production system interpreting chemical diagrams which are written in a conventional manner. We have emphasized conventional diagrams because we feel that the value of a chemist's time is such that we should put a premium upon easily recognized output.

Chemists do not generally write chemical structures in which each atom and each bond is shown explicitly. Instead, there is an abundance of shorthand conventions, such as $-\text{COOH}$, $-\text{C}_6\text{H}_5$, $-\text{C}_6\text{H}_{11}$, $-\text{C}(\text{CH}_3)_2(\text{CH}_2)_3$, which must be correctly interpreted by the computer. For the last four years we have been involved with two philosophies for interpreting such conventions. The first philosophy, described by Gould, Gasser, and Ryan,⁵ interprets the chemical diagrams by looking up each abbreviation in a table in order to find explicit bond-attachments. Programming for these chemical abbreviations (or glyphs) is time-consuming and requires a large computer core, but has the advantage that the interpretation of each glyph can be specified in advance. Since table of look-up may be easily changed and enlarged,

this was a natural beginning point for our efforts in interpreting conventional chemical diagrams. Confusion in the interpretation of glyphs may result, however, from their juxtaposition. The program utilizing these glyphs has been working since 1962, and although we are not using it now, the availability of third-generation hardware with large memories permits reconsideration of glyphs as a way of handling chemical diagrams. The second philosophy for handling abbreviations does this by means of rules. Rules have the advantage of high-speed processing and small core requirements; the disadvantage is that they are not easily changed, since a change in one rule may change the interpretation of other chemical structures. These rules, like chemical glyphs, also need to be understood and agreed upon in order to determine if they, in fact, represent "good" chemistry. For example, we allow $\text{CH}_3\text{SO}_2\text{CH}_3$ to be typed on input, but the rules interpret it only as dimethylsulfone. While it is not necessary that chemists learn these rules, it is necessary that they agree on the interpretation of the diagrams. Since chemistry is a diagrammatic language, there are a fair number of special cases usually associated with multivalent atoms, many of which we have taken care of in our existing programs.

The cost of running a compound through this phase of the program is approximately 13 cents if the operation is carried out at commercial daytime rates, or 1.3 cent if the cost is calculated using government equipment. The output of this first phase is a connection table similar to those used as input in other systems. When we used such tables as input, we found them time-consuming to prepare and to check, and consequently expensive. We feel that on grounds of cost alone, we have established the advantage of typewriter input.

The connection table generated from the structure might be left unchanged, but most other systems process it. This processing, however, produces additional levels of incompatibility between systems. To our knowledge, our method is the only one which does not attempt to produce a condensed connection table; nor do we attempt to generate a unique canonical notation for each structure. We feel that a critical problem in a canonical notation is that, in the absence of a rigorous mathematical proof, the algorithm used may establish false identities, which would be extremely difficult to discover. We felt, therefore, that in renouncing canonical notations, we would keep the system clean. More important, when searching for fragments, canonical notations are useless. For establishing the identity of actual compounds, our program contains all applicable screens and the molecular formula. Only if these match will it trace through the connection table, atom-by-atom, to determine identity.

To reduce the time-consuming atom-by-atom matching, we compute screen bits which are assigned to each molecule. A screen bit is a simple yes-no flag reflecting the presence or absence of a particular chemical fragment—e.g., a carboxyl group. Our screens, however, can be relatively exotic—e.g., NH_2 separated by any four atoms from OH without regard to whether the intervening atoms are aromatic, aliphatic, etc. Screen bits are generated automatically, become a permanent part of the record, and are used as an extension of the molecular formula on file.

CN(C)C1=CC=C(C=C1)C2=CC=C(C=C2)C3=CC=C(C=C3)C4=CC=C(C=C4)C5=CC=C(C=C5)C6=CC=C(C=C6)C7=CC=C(C=C7)C8=CC=C(C=C8)C9=CC=C(C=C9)C10=CC=C(C=C10)C11=CC=C(C=C11)C12=CC=C(C=C12)C13=CC=C(C=C13)C14=CC=C(C=C14)C15=CC=C(C=C15)C16=CC=C(C=C16)C17=CC=C(C=C17)C18=CC=C(C=C18)C19=CC=C(C=C19)C20=CC=C(C=C20)C21=CC=C(C=C21)C22=CC=C(C=C22)C23=CC=C(C=C23)C24=CC=C(C=C24)C25=CC=C(C=C25)C26=CC=C(C=C26)C27=CC=C(C=C27)C28=CC=C(C=C28)C29=CC=C(C=C29)C30=CC=C(C=C30)C31=CC=C(C=C31)C32=CC=C(C=C32)C33=CC=C(C=C33)C34=CC=C(C=C34)C35=CC=C(C=C35)C36=CC=C(C=C36)C37=CC=C(C=C37)C38=CC=C(C=C38)C39=CC=C(C=C39)C40=CC=C(C=C40)C41=CC=C(C=C41)C42=CC=C(C=C42)C43=CC=C(C=C43)C44=CC=C(C=C44)C45=CC=C(C=C45)C46=CC=C(C=C46)C47=CC=C(C=C47)C48=CC=C(C=C48)C49=CC=C(C=C49)C50=CC=C(C=C50)C51=CC=C(C=C51)C52=CC=C(C=C52)C53=CC=C(C=C53)C54=CC=C(C=C54)C55=CC=C(C=C55)C56=CC=C(C=C56)C57=CC=C(C=C57)C58=CC=C(C=C58)C59=CC=C(C=C59)C60=CC=C(C=C60)C61=CC=C(C=C61)C62=CC=C(C=C62)C63=CC=C(C=C63)C64=CC=C(C=C64)C65=CC=C(C=C65)C66=CC=C(C=C66)C67=CC=C(C=C67)C68=CC=C(C=C68)C69=CC=C(C=C69)C70=CC=C(C=C70)C71=CC=C(C=C71)C72=CC=C(C=C72)C73=CC=C(C=C73)C74=CC=C(C=C74)C75=CC=C(C=C75)C76=CC=C(C=C76)C77=CC=C(C=C77)C78=CC=C(C=C78)C79=CC=C(C=C79)C80=CC=C(C=C80)C81=CC=C(C=C81)C82=CC=C(C=C82)C83=CC=C(C=C83)C84=CC=C(C=C84)C85=CC=C(C=C85)C86=CC=C(C=C86)C87=CC=C(C=C87)C88=CC=C(C=C88)C89=CC=C(C=C89)C90=CC=C(C=C90)C91=CC=C(C=C91)C92=CC=C(C=C92)C93=CC=C(C=C93)C94=CC=C(C=C94)C95=CC=C(C=C95)C96=CC=C(C=C96)C97=CC=C(C=C97)C98=CC=C(C=C98)C99=CC=C(C=C99)C100=CC=C(C=C100)C101=CC=C(C=C101)C102=CC=C(C=C102)C103=CC=C(C=C103)C104=CC=C(C=C104)C105=CC=C(C=C105)C106=CC=C(C=C106)C107=CC=C(C=C107)C108=CC=C(C=C108)C109=CC=C(C=C109)C110=CC=C(C=C110)C111=CC=C(C=C111)C112=CC=C(C=C112)C113=CC=C(C=C113)C114=CC=C(C=C114)C115=CC=C(C=C115)C116=CC=C(C=C116)C117=CC=C(C=C117)C118=CC=C(C=C118)C119=CC=C(C=C119)C120=CC=C(C=C120)C121=CC=C(C=C121)C122=CC=C(C=C122)C123=CC=C(C=C123)C124=CC=C(C=C124)C125=CC=C(C=C125)C126=CC=C(C=C126)C127=CC=C(C=C127)C128=CC=C(C=C128)C129=CC=C(C=C129)C130=CC=C(C=C130)C131=CC=C(C=C131)C132=CC=C(C=C132)C133=CC=C(C=C133)C134=CC=C(C=C134)C135=CC=C(C=C135)C136=CC=C(C=C136)C137=CC=C(C=C137)C138=CC=C(C=C138)C139=CC=C(C=C139)C140=CC=C(C=C140)C141=CC=C(C=C141)C142=CC=C(C=C142)C143=CC=C(C=C143)C144=CC=C(C=C144)C145=CC=C(C=C145)C146=CC=C(C=C146)C147=CC=C(C=C147)C148=CC=C(C=C148)C149=CC=C(C=C149)C150=CC=C(C=C150)C151=CC=C(C=C151)C152=CC=C(C=C152)C153=CC=C(C=C153)C154=CC=C(C=C154)C155=CC=C(C=C155)C156=CC=C(C=C156)C157=CC=C(C=C157)C158=CC=C(C=C158)C159=CC=C(C=C159)C160=CC=C(C=C160)C161=CC=C(C=C161)C162=CC=C(C=C162)C163=CC=C(C=C163)C164=CC=C(C=C164)C165=CC=C(C=C165)C166=CC=C(C=C166)C167=CC=C(C=C167)C168=CC=C(C=C168)C169=CC=C(C=C169)C170=CC=C(C=C170)C171=CC=C(C=C171)C172=CC=C(C=C172)C173=CC=C(C=C173)C174=CC=C(C=C174)C175=CC=C(C=C175)C176=CC=C(C=C176)C177=CC=C(C=C177)C178=CC=C(C=C178)C179=CC=C(C=C179)C180=CC=C(C=C180)C181=CC=C(C=C181)C182=CC=C(C=C182)C183=CC=C(C=C183)C184=CC=C(C=C184)C185=CC=C(C=C185)C186=CC=C(C=C186)C187=CC=C(C=C187)C188=CC=C(C=C188)C189=CC=C(C=C189)C190=CC=C(C=C190)C191=CC=C(C=C191)C192=CC=C(C=C192)C193=CC=C(C=C193)C194=CC=C(C=C194)C195=CC=C(C=C195)C196=CC=C(C=C196)C197=CC=C(C=C197)C198=CC=C(C=C198)C199=CC=C(C=C199)C200=CC=C(C=C200)C201=CC=C(C=C201)C202=CC=C(C=C202)C203=CC=C(C=C203)C204=CC=C(C=C204)C205=CC=C(C=C205)C206=CC=C(C=C206)C207=CC=C(C=C207)C208=CC=C(C=C208)C209=CC=C(C=C209)C210=CC=C(C=C210)C211=CC=C(C=C211)C212=CC=C(C=C212)C213=CC=C(C=C213)C214=CC=C(C=C214)C215=CC=C(C=C215)C216=CC=C(C=C216)C217=CC=C(C=C217)C218=CC=C(C=C218)C219=CC=C(C=C219)C220=CC=C(C=C220)C221=CC=C(C=C221)C222=CC=C(C=C222)C223=CC=C(C=C223)C224=CC=C(C=C224)C225=CC=C(C=C225)C226=CC=C(C=C226)C227=CC=C(C=C227)C228=CC=C(C=C228)C229=CC=C(C=C229)C230=CC=C(C=C230)C231=CC=C(C=C231)C232=CC=C(C=C232)C233=CC=C(C=C233)C234=CC=C(C=C234)C235=CC=C(C=C235)C236=CC=C(C=C236)C237=CC=C(C=C237)C238=CC=C(C=C238)C239=CC=C(C=C239)C240=CC=C(C=C240)C241=CC=C(C=C241)C242=CC=C(C=C242)C243=CC=C(C=C243)C244=CC=C(C=C244)C245=CC=C(C=C245)C246=CC=C(C=C246)C247=CC=C(C=C247)C248=CC=C(C=C248)C249=CC=C(C=C249)C250=CC=C(C=C250)C251=CC=C(C=C251)C252=CC=C(C=C252)C253=CC=C(C=C253)C254=CC=C(C=C254)C255=CC=C(C=C255)C256=CC=C(C=C256)C257=CC=C(C=C257)C258=CC=C(C=C258)C259=CC=C(C=C259)C260=CC=C(C=C260)C261=CC=C(C=C261)C262=CC=C(C=C262)C263=CC=C(C=C263)C264=CC=C(C=C264)C265=CC=C(C=C265)C266=CC=C(C=C266)C267=CC=C(C=C267)C268=CC=C(C=C268)C269=CC=C(C=C269)C270=CC=C(C=C270)C271=CC=C(C=C271)C272=CC=C(C=C272)C273=CC=C(C=C273)C274=CC=C(C=C274)C275=CC=C(C=C275)C276=CC=C(C=C276)C277=CC=C(C=C277)C278=CC=C(C=C278)C279=CC=C(C=C279)C280=CC=C(C=C280)C281=CC=C(C=C281)C282=CC=C(C=C282)C283=CC=C(C=C283)C284=CC=C(C=C284)C285=CC=C(C=C285)C286=CC=C(C=C286)C287=CC=C(C=C287)C288=CC=C(C=C288)C289=CC=C(C=C289)C290=CC=C(C=C290)C291=CC=C(C=C291)C292=CC=C(C=C292)C293=CC=C(C=C293)C294=CC=C(C=C294)C295=CC=C(C=C295)C296=CC=C(C=C296)C297=CC=C(C=C297)C298=CC=C(C=C298)C299=CC=C(C=C299)C300=CC=C(C=C300)C301=CC=C(C=C301)C302=CC=C(C=C302)C303=CC=C(C=C303)C304=CC=C(C=C304)C305=CC=C(C=C305)C306=CC=C(C=C306)C307=CC=C(C=C307)C308=CC=C(C=C308)C309=CC=C(C=C309)C310=CC=C(C=C310)C311=CC=C(C=C311)C312=CC=C(C=C312)C313=CC=C(C=C313)C314=CC=C(C=C314)C315=CC=C(C=C315)C316=CC=C(C=C316)C317=CC=C(C=C317)C318=CC=C(C=C318)C319=CC=C(C=C319)C320=CC=C(C=C320)C321=CC=C(C=C321)C322=CC=C(C=C322)C323=CC=C(C=C323)C324=CC=C(C=C324)C325=CC=C(C=C325)C326=CC=C(C=C326)C327=CC=C(C=C327)C328=CC=C(C=C328)C329=CC=C(C=C329)C330=CC=C(C=C330)C331=CC=C(C=C331)C33

00000000000000000000000000000000 CND 3N4V78299
003N-4V-78299
C 4 V 2 0 S

O=C1C(=O)N(C(=O)O)C(=O)N1C(=O)O

```

*****WARNING - OXYGEN PRECEDES VARIABLE VALENCE ELEMENT      X=11, Y=03
*****WARNING - OXYGEN GROUP PROCESSED                          X=11, Y=03
*****WARNING - OXYGEN FOLLOWS VARIABLE VALENCE ELEMENT        X=17, Y=12
*****WARNING - OXYGEN GROUP PROCESSED                          X=17, Y=12
*****WARNING - PEROXIDE MAY HAVE BEEN CREATED OR LOST         X=17, Y=12
*****WARNING MESSAGES ONLY - COMPOUND ACCEPTED

```

```

STRUC
01      O      003
02      O      003
03      N      D01,D02,S04
04      C      S03,D05,S14
05      C      D04,S06,S07
06      O      SH,S05
07      C      S03,S08,D12
08      N      S07,D09
09      C      D09,SH,S10
10      C      S09,SH,D11
11      C      D10,SH,S12
12      C      S11,D07,S13
13      C      S12,D14,S15
14      C      D13,S04,SH
15      S      D15,D17,S13,S18
16      O      D15
17      O      D15
18      O      SH,S15

```

We have emphasized screens because of our interest in searches for inclusion. We felt that complex inclusion-questions otherwise would require multiple backtracking on large portions of the file. The screens also represent an attempt to reduce the recomputation associated with multiple questions. Of necessity we arrived at our pattern of screen bits in an arbitrary manner; we are now able to evaluate these from the point of view of information theory, since we have now accumulated both a large data base and a large series of questions. At present, our guess is that we are carrying more screen bits (500) than would result if they were selected in a systematic manner. However, since screen production is only a one-time computation, and since screens occupy relatively little

Updating and sorting files is part of the interpretation of input. A test for identity determines whether a compound will be added to the master file or be treated as a duplicate. To do this operation efficiently, the chemical file is kept in some kind of chemical order, and indeed, all existing chemical systems maintain their files in some order. Our file keeps the compounds in molecular formula order, followed by their screen-bit configurations. The order in which the biological files usually are maintained

is accession number order, thus necessitating sorting whenever chemical and biological files are correlated. Even in random access systems, corresponding indexes will have to be sorted. With large files, large sorts are a difficult and expensive operation.

SYNONYM FILES

Synonyms—not necessarily names, but synonymous numbers—arise because submitters of chemicals, testing laboratories, etc., insist on retaining their own identification numbers, and also because in our system a permanent accession number cannot be assigned to a compound when it is received, but only after its record has proceeded to a certain stage in our system.

The permanent accession number cannot be assigned to a compound by our system until it is determined that the new entry does not duplicate a previous entry on the file. All processing up to this final step carries the uncertainty that the incoming compound may be a duplicate. The assignment of an accession number means that all the previously available but less favored numbers must be changed, or else that a synonym file must be maintained so that the preferred accession number can be coupled with records using older identification numbers.

The main problem up to this point is keeping the system clean—i.e., eliminating bad input, being certain that compounds which are rejected are properly re-entered, correcting incorrect chemistry, and correcting typing errors which are not chemical errors and which therefore are accepted by the system. The main problem from this point on is the maintenance of the synonym file and the production of output in a format and sequence suitable for integration with the synonym, biology, and inventory files. We have just indicated how late in the system a permanent accession number is assigned. In both our older manually operated systems and in our computer-operated system, we would prefer a single master or registry number, but we have found it impossible to operate without synonyms. There are at least five major sources of synonyms: 1) old test systems have different accession numbers, such as the serial numbers of the World War II malarial program; 2) companies and private organizations sending in compounds use their own accession numbers for inventory control; 3) our biological testers assign their own serial number to find their data. While testers often keep track of our number, their own accession number improves the operation of their system; 4) blind accession numbers are assigned to check on the reproducibility of biological tests; 5) multiple samples of the same chemical from the same source requires lot control, which of necessity means at least a subdivision of the previous accession number. We feel this problem can only get worse, especially for systems dealing with multiple unpublished private files.

QUESTIONS

The format of our chemical questions, and some of the chemical answers, are shown in Figures 2 and 3. The chemical question for identity is asked by typing the actual structure. Similarly, a substructure search is specified by typing the actual fragment. The letter Z,

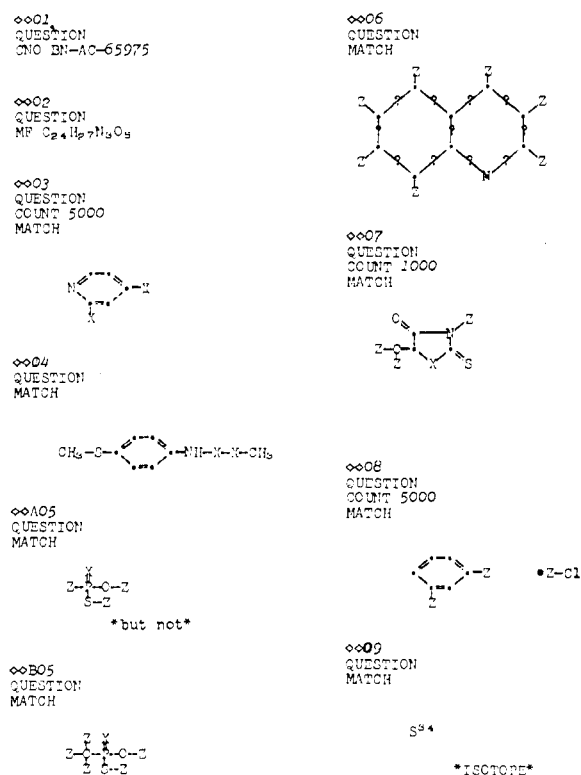


Figure 2. Examples of queries

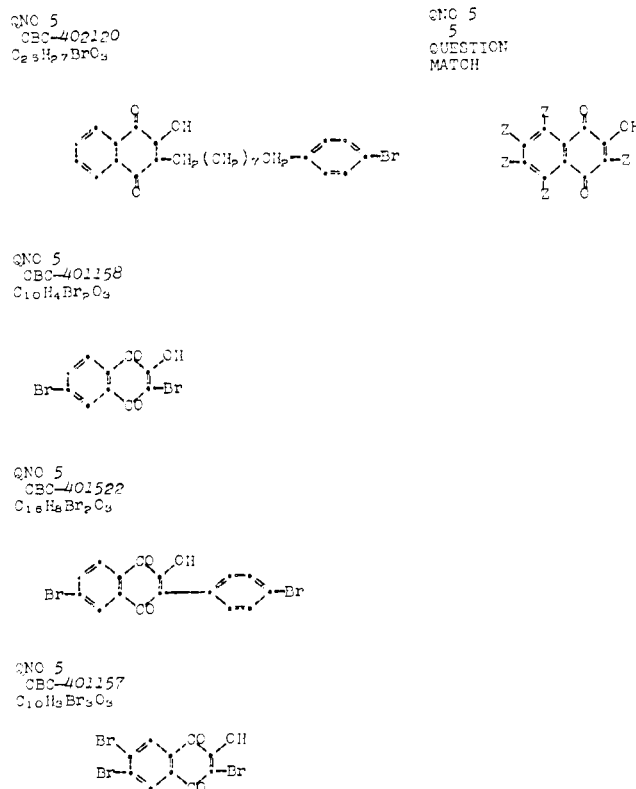


Figure 3. A chemical question and some answers

CHEMICAL AND BIOLOGICAL INFORMATION RETRIEVAL SYSTEM

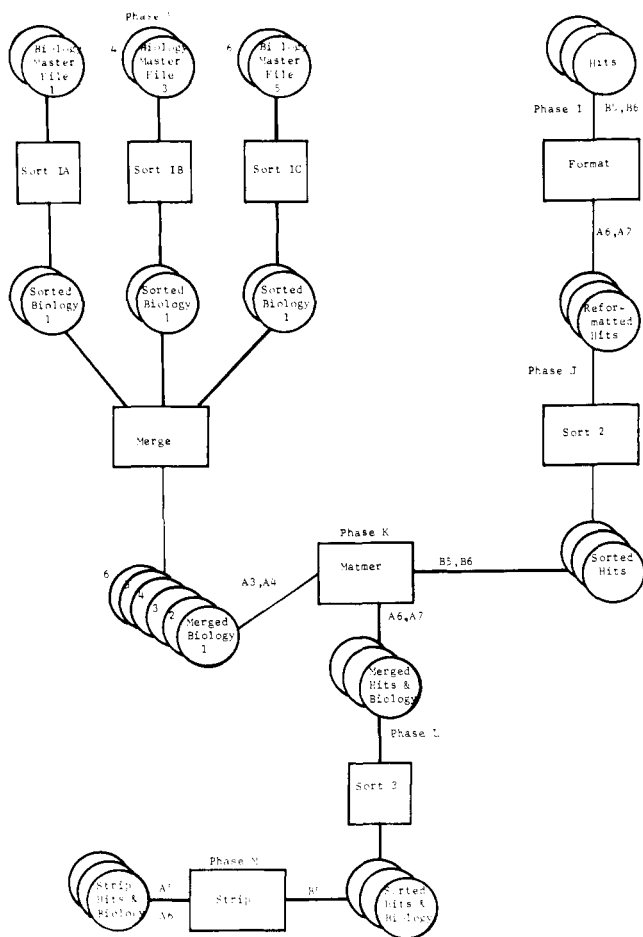


Figure 4. Merged biology and chemstructs text data

inserted at appropriate locations on the fragment, indicates locations where substituents are allowed. Following the Gould convention,⁵ the letter X, typed similarly, indicates that a nonhydrogen substituent is required at that location. Questions including "but not" conditions can be specified by typing first the fragment desired, then the fragment to be excluded.

Except as a tour-de-force, we have found that chemical output without associated data is of little interest. To obtain the chemical diagrams independently from one file and the appropriate biological information from another is impractical. With small files, this has been done for many years; with large mechanical systems, the necessary manual handling becomes an impossible task. It is the machine which must present the desired chemical, biological, and associated data integrally combined.

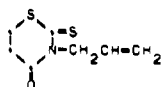
Because of the synonym problem and because the chemical order is different from the accession number order, much sorting and merging is necessary to achieve this combined printout. These data files may always outgrow on-line memory systems. As implied by Figure 4, the conventional data processing problems are large.

INEXPENSIVE OUTPUT

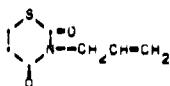
The last critical area is the development of inexpensive output. At present, the cheapest output is produced by the high-speed printers associated with digital computers. We use a Data Products printer off line. We originally chose this machine because its drum was easily engraved with the characters constituting a chemical font. This printer has since proved very reliable in production use.

Our machine prints 600 half lines per minute when printing chemical diagrams or 600 full lines per minute

4444-2583
C₇H₉N₂S₂



C₇H₉N₂S₂



RODENT RADIATION MORTALITY DATA									
SOURCE 0039					SUBMITTER N-173				
JULIAN DATE	PERCENT SOL.	PH	VEH	CHEM LD50	ROUTE ADM.	DRUG DOSE	RADIATION TYPE DOSE	NO. TIME MICE	MORTALITY BY DAYS
03/175	0.00	0.0	CMCTW	0800	IP	000.00	G 1000	30 10	000000/00010/23121
03/175	3.00	7.7	CMCTW	0800	IP	300.00	G 1000	30 15	000000/01000/14301/01001/00010/00100
03/175	3.00	7.7	CMCTW	0800	IP	600.00	G 1000	30 15	041000/00022/01002/00000/00000/10100

RODENT RADIATION MORTALITY DATA									
SOURCE 0039					SUBMITTER N-226				
JULIAN DATE	PERCENT SOL.	PH	VEH	CHEM LD50	ROUTE ADM.	DRUG DOSE	RADIATION TYPE DOSE	NO. TIME MICE	MORTALITY BY DAYS
03/176	0.00	0.0	CMCTW	0600	IP	000.00	G 1000	15 10	000000/00012/01201/3
03/176	2.00	6.2	CMCTW	0600	IP	400.00	G 1000	15 15	000000/00120/27201
03/176	2.00	6.2	CMCTW	0600	IP	200.00	G 1000	15 15	000000/00212/21101/10000/00000/00000

Figure 5. Printout of biological data and associated chemical structures

when printing biological data (Figure 5); a newer version is available which operates at 1200 lines per minute.

From the description of the above system, which handles large volumes of data, it is obvious that a number of problems had to be solved, independently of the ones presented by chemistry. Problems were encountered in the areas of systems analysis, programming, electrical and mechanical engineering, topology, information theory, and human-factor analysis. We believe that the growth and effectiveness of our system is due in large part to our attention to these areas and by our use of experts and consultants therein. The growth of the system was stimulated by the need to depend upon it for the control of a large medicinal chemistry program.

LITERATURE CITED

- (1) Bello, F., "How to Cope With Information," *Fortune*, pp. 2-10 (Sept. 1960).
- (2) "Survey of Chemical Notation Systems," National Academy of Sciences, National Research Council, Publication 1150, 1964.
- (3) Waldo, W. H., "Searching Two-Dimensional Structures By Computer," *J. CHEM. DOC.* 2, 1-2 (1962).
- (4) Feldman, A., D. B. Holland and D. P. Jacobus, "The Automatic Encoding of Chemical Structures," *J. CHEM. DOC.* 3, 187 (1963).
- (5) Gould, D., E. B. Gasser and J. F. Ryan, "Chemical Search—An Operating Computer System For Retrieving Chemicals Selected For Equal, Analogous or Related Character," *J. CHEM. DOC.* 5, 24 (1965).