ate $(-\lambda + c)$ where $c$ is some numerical correction for electronegativity and/or $\sigma$-bond energy. The resulting polynomial is a function of $\lambda$, the usual energy eigenvalue, and when equated to zero is identical with the traditionally obtained secular equation in appropriate energy units.

## SUMMARY

The physico–geometric representation of the ACMCP is described by the expression

$$ACMCP = \sum_{k=0}^{N} {}_{*} \left\{ (-1)^{\left(\sum\limits_{i=2}^{N-k} n_i(i-1)\right)} \cdot (2)^{\left(\sum\limits_{i=5}^{N-k} n_i\right)} \cdot \left[ \prod_{j=2}^{N-k} \prod_{m=0}^{n_i} \mathcal{Q}_m \right] S_k \right\}$$

where

$\sum_{*}$ = summation over all different combinations of non-adjacent rings where

$$\sum_{i=2}^{N-k} i n_i = N - k$$

"combination" = collection of nonadjacent rings whose connectivities enter into evaluation of a particular coefficient in ACMCP

$i$ = integer, $2 \leq i \leq N - k$

$j$ = integer, $2 \leq j \leq N - k$

$k$ = integer, $0 \leq k \leq N$; identifies degree of ACMCP term

$m$ = integer, $0 \leq m \leq n_j$

$n_i$ = integer, $0 \leq n_i \leq [(N - k)/i]$, where upper limit is the largest integer value of $(N - k)/i$; represents number of nonadjacent (nontouching) $i$-membered rings in a combination

$\mathcal{Q}_m$ = product of bond connectivities around (defining) the $m$th nonadjacent $j$-membered ring

$S_k$ = product of atomic symbols, $k$ in number, which are not included in the rings making up a given combination

$S_0 \equiv$ 1, for any $j$; a definition

$\mathcal{Q}_u \equiv$ 1; a definition

# Chemical Structure Fragmentation for Use in a Coordinate Index Retrieval System*

By RAY W. IHNDRIS

Cancer Chemotherapy National Service Center, National Cancer Institute,
National Institutes of Health, Bethesda, Maryland

Received January 28, 1964

This paper describes the changes that have been made in the Cancer Chemotherapy National Service Center (CCNSC) Chemical Information Retrieval System[1] established for analog searching of chemical structures.

The National Bureau of Standards Peek-a-Boo System[2] in use at the CCNSC consists of a precision-built machine which uses a microswitch to actuate the solenoid-operated punch when the plastic card is placed into punching position. The cards used for permanently recording "data" are 5 × 8 in., 8–10 mil vinyl plastic sheets. Each card is punched separately.

The Reader consists of a back-lighted opaque frame onto which the plastic cards are placed and over-layed with a grid so that the coincident holes allow light to shine through, and file numbers can be "read" from the coordinate position on the grid. The thousands and hundreds are read along the horizontal axis, and the tens and units are read on the vertical axis. The grid provides for 18,000 recording positions.
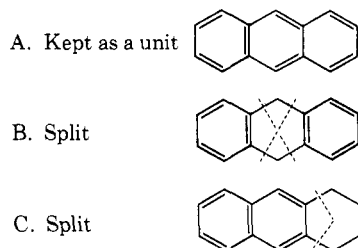
A VISIrecord file cabinet houses the plastic cards in a manner permitting easy access. The cards are shingled so that three-fourths of an inch along the top edge exposes the chemical fragment or other descriptor. The cards are notched on the right-hand edge and filed vertically so that they fit over rods in the bottom of the file cabinet. They are printed to define the space for the descriptors and also have blocks along the left-hand edge that signal the position when a card is removed. Thus it is easy to find a card for entering new data or retrieving information.

One or several descriptors are used for selection by overlaying the plastic cards, thus eliminating unwanted information. We can vary the search at any time if the desired information is not found, or conversely, if too much chaff appears.

The improvements in our system of coding chemical structures affect the ring systems, the fragment chains, and the supplier code.

**Ring Systems.**—The ring systems found in the chemicals of the CCNSC program are many and varied. Several hundred systems occur in each set of 18,000 compounds. Exactly how many appeared in the earlier sets is unknown because of the way the ring systems were recorded. Only aromatic rings (A) were kept as a unit (plus a few special heterocyclic rings). Partially hydrogenated aromatic rings (B and C) were split into fragments.

A. Kept as a unit

B. Split

C. Split

In our present system all fused and spiro rings are kept intact as indexing units.

Fragmenting a structure containing a ring system requires looking up the Revised Ring Index Number (RRI #) in "The Ring Index,"[3] which may be accomplished in several ways.

First, if one knows the name of the particular system, it can be found quickly from the alphabetic Index of Names in the back of the book. Second, if the system is not identified by name, one can search by (1) the number of rings composing the system, (2) the ring size, and (3) the ring formula. Acridine, for example, is found in the "three-ring system" in a series 6,6,6 under the ring formula $C_5N$-$C_6$-$C_6$. Each system has a serial number, which we use for coding into the Peek-a-Boo cards. In this case it is RRI-3523.

Our method for coding all possible ring systems uses a special group of 44 plastic cards which are numbered as follows.

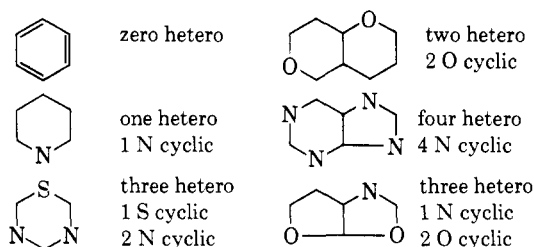| | | | | |
|---|---|---|---|---|
| 0000 | -000 | --00 | ---0 | .1 |
| 1000 | -100 | --10 | ---1 | .2 |
| 2000 | -200 | --20 | ---2 | .3 |
| 3000 | -300 | --30 | ---3 | .4 |
| 4000 | -400 | --40 | ---4 | |
| 5000 | -500 | --50 | ---5 | |
| 6000 | -600 | --60 | ---6 | |
| 7000 | -700 | --70 | ---7 | |
| 8000 | -800 | --80 | ---8 | |
| 9000 | -900 | --90 | ---9 | |

More than 10,000 ring systems can be recorded and retrieved with these cards. The four cards circled are pulled from the VISIrecord file and punched with the appropriate accession number (NSC-) in its proper coordinate position. It is necessary to punch at least four cards for each ring number.

If the particular ring system is not listed in "The Ring Index," we use the RRI # closest to where it would be placed according to *Chemical Abstracts'* rules of arranging by ring size and ring formula and add .1 to its serial number.

All single rings, fused rings, bridged rings, spiro rings, and steroid rings are included in this method of coding without regard to their degree of saturation. Previously, saturated rings were fragmented into chains or groups and were difficult to retrieve.

Occasionally, two or more ring systems may occur in the same structure. About 22 of the most common and frequently occurring rings have been entered on separate cards and not by their RRI #. This was done partly for easy retrieval and partly to reduce the need for punching two RRI #'s for the same compound which would result in false readings of numbers during retrieval. Substitution positions are also punched for benzene and cyclohexane rings, for example, mono, *ortho, para,* symmetrical, asymmetrical, and four or more substitutions.

To retrieve ring systems in a more general way, the following descriptors may be used. Hetero refers to substitutions in rings of elements other than carbon. Punch "zero hetero" if there is no hetero atom as in benzene or one or more "hetero" depending on number of noncarbon atoms in each ring system.

| | | | |
|---|---|---|---|
| | zero hetero | | two hetero<br>2 O cyclic |
| | one hetero<br>1 N cyclic | | four hetero<br>4 N cyclic |
| | three hetero<br>1 S cyclic<br>2 N cyclic | | three hetero<br>1 N cyclic<br>2 O cyclic |

Several hetero number cards may be punched depending on how many ring systems there are in one compound. The N, O, P, or S cyclic cards are also punched according to the total number of each element in each ring system.

When a double bond occurs in a nonaromatic ring, punch $>C=C<$ cyclic, or when a ketone group is part of a ring, punch $>C=O$ cyclic, or other elements such as $>C=N$- cyclic, the noncarbon element may be either in or out of the ring. Ring systems that contain fragment groups are coded by both the RRI # and the linear fragment system.

We have found this new method for coding ring systems to be very specific when using the RRI # and quite general when using the allied descriptors (*i.e.,* the hetero cards plus the cyclic designations and/or the fragment chains). Searching time is lessened and fewer false drops occur.

In the new system, two punches are required for each benzene ring, *i.e.,* the ring card plus the substitution position card, but this change solves one of our big retrieval problems. To make a survey of, say, all nitro benzenes, the $-NO_2$ card would have to be compared with each of the eight different benzene cards. With the new system we need only make one comparison, using the master benzene card and the nitro card. A more specific search can easily be made with the addition of one of the "position" cards.

The pyridine, piperidine, and partially unsaturated piperidine rings are combined on one card. This is true also for the pyrimidine ring. The N⁻ cyclic can distinguish the "onium" rings from the others. This will assure completeness in the general analog type of search. The old method would split the piperidine ring into the following fragments.

1 N cyclic
1 ring, 6 members     for
(—CH₂—)₅

or

1 N cyclic
1 ring, 6 members          for
1 —C=C—
1 —CH₂—
(—CH₂—)₂

$$>N—N—\overset{\overset{\text{O3}}{\|}}{C}—N<\quad\text{fragmented}$$
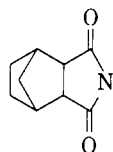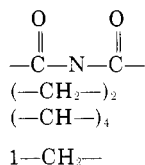$$\underset{1\quad2\quad3\quad4}{}$$

Many unwanted compounds would be retrieved when trying to reconstruct this ring. Each substitution on the ring requires a change in the search procedure without creating more specificity because these carbon fragments are used for chains as well as rings and may occur in other segments of the structure. The new system retains all rings as a unit without regard to degree of saturation.

For comparison, norbornane-2,3-dicarboximide would be entered into the system as follows.

Old

| | |
|---|---|
| Fused ring 6,5 | —C—N—C— (with =O on each C) |
| 1 N cyclic | (—CH₂—)₂ |
| 1 ring, 5 members | (—CH—)₄ |
| 1 ring, 6 members | |
| Bridge | 1—CH₂— |

New
RRI # 2242
>C=O cyclic

Nine plastic cards are required to reconstruct this formula, and this doesn't specify the exact structure or allow for substitutions on the ring. Only four RRI# cards plus a >C=O cyclic for the carbonyl group are required for this specific structure. For the norborn*ene* structure, add the >C=C< cyclic card. A ring system can be obscured by extraneous "drop-outs" if it is fragmented into its elements. We have eliminated that possibility by retaining it as a whole ring. A general search for similar ring systems can be made using combinations of any of the following headings as they fit the requirements of the search.

1 hetero atom
1 nitrogen cyclic
Bridge
>C=O cyclic

$$—\overset{\overset{\text{O1}}{\|}}{C}—N—\overset{\overset{\text{O3}}{\|}}{C}—$$
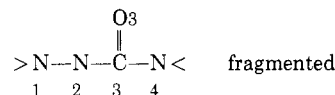$$\underset{1\quad2\quad3}{}$$

**Chain Fragmentation.**—The second part of our fragmentation system to be revised is the fragment chain of atoms, which is defined as: a series of atoms linked one to the other by single or multiple bonds in which there are no singly bonded carbon-to-carbon atoms.
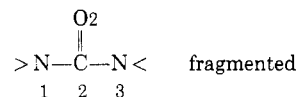
Following established rules, most organic radicals and functional groups attached to carbon chains or aromatic rings are indexed as units if there are less than four atoms in the radical or group. As in the past, separate cards are made for the following groups.[4]

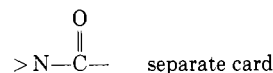| | | | |
|---|---|---|---|
| —NO | >C=O | —N=C< | —O—C—O— |
| —NO₂ | —NCO | —N=N— | C=N—N< |
| —CN | —NCS | >N—N< | —COOH |
| —CNS | —SCN | —N—C—N— | etc. |

The revision of chain fragments involves groups having four atoms in a chain or three atoms in a chain plus one one or more branches on the chain. For example

$$>N—\overset{\overset{\text{O2}}{\|}}{C}—N<\quad\text{fragmented}$$
$$\underset{1\quad2\quad3}{}$$

$$>N—\overset{\overset{\text{O}}{\|}}{C}—\quad\text{separate card}$$

A special group of sixty plastic cards is used to code the very large number of combinations possible with nitrogen, oxygen, phosphorus, sulfur, and carbon in any one of five positions in a chain. The cards are designated as follows.

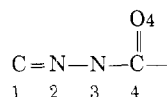| Position 1 | Position 2 | Position 3 | Position 4 | Position 5 |
|---|---|---|---|---|
| 1 N | 2 N | 3 N | 4 N | 5 N |
| 1 N= | 2 N= | 3 N= | 4 N= | 5 N= |
| 1 O | 2 O | 3 O | 4 O | 5 O |
| 1 O= | 2 O= | 3 O= | 4 O= | 5 O= |
| 1 P | 2 P | 3 P | 4 P | 5 P |
| 1 P= | 2 P= | 3 P= | 4 P= | 5 P= |
| 1 S | 2 S | 3 S | 4 S | 5 S |
| 1 S= | 2 S= | 3 S= | 4 S= | 5 S= |
| 1 C | 2 C | 3 C | 4 C | 5 C |
| 1 C= | 2 C= | 3 C= | 4 C= | 5 C= |
| 1 X | 2 X | 3 X | 4 X | 5 X |
| | | 3 END | 4 END | 5 END |
| | | | | NO END |

The "X" card can represent any other elements such as As (arsenic), B (boron), Cr (chromium), etc., when a part of the fragment chain. In case a single element to be designated by X is present, the card for this element is also punched. These cards will accommodate all combinations of elements up to five atoms in a chain plus their branchings. An "END" card is punched for the last atom of the chain. Should more than five atoms occur in a chain, the "NO END" card is punched, indicating a longer chain.

An order of preference has been established for coding these chains. If nitrogen is present, the preferred order is

Example

1. —N=N—          —N=N—C=N—
                   1   2 3   4

2. —N—N—          $$—N—N—\overset{\overset{\text{N3}}{\|}}{C}—O—$$
                   1   2   3   4

3. —N=C—          —N=C—O—C—N—
                   1   2  3  4  5

4. —N—C—          $$\overset{\cdot}{—N}—\overset{\overset{\text{O2}}{\|}}{C}—O—$$
                   1   2   3

5. —C=N—          $$—C=N—N—\overset{\overset{\text{O4}}{\|}}{C}—O—$$
                   1   2   3   4   5

After establishing ranking on the basis of nitrogen in the first two positions, further precedence is determined by oxygen, phosphorus, sulfur, and X elements. Ketone $>C=O$ and thione $>C=S$ are listed last in order of preference. Nitrogen is placed as close to the 1 position as possible. Therefore $—C=N—$ (#5) is listed with the nitrogen atoms. Over 300 different fragment chains have occurred in the first 5,000 organic compounds received this year.

The following example illustrates the procedure for coding fragment chains. First each atom is numbered according to the established order.

$$\begin{array}{ccccc} & & & & O4 \\ & & & & \| \\ C & = & N—N—C— \\ 1 & & 2 & 3 & 4 \end{array}$$

For position 1-punch # 1 C =
For position 2-punch # 2 N =
For position 3-punch # 3 N
For position 4-punch # 4 C =
For position 4-punch # 4 O =
For position 4-punch # 4 END

When double bond is indicated, the element cards for both positions concerned are punched indicating a double bond attachment.

A fragment chain may be wholly or partially included in a hetero cyclic nonaromatic ring, in which case it is coded both as a fragment chain and by the RRI #.

RRI # 0869

$$\begin{array}{ccc} & O2 & \\ & \| & \\ —O—P—O— \\ 1 & | 2 & 3 \\ & O & \\ & 2 & \text{fragment} \end{array}$$

The ring is also coded with the following descriptors.

4 hetero
3 O cyclic
1 P cyclic

A Rolodex file card is maintained on each ring system and fragment chain for quick reference by the chemist making the structural breakdown. As new arrangements of the elements in the fragment chain are received, a Rolodex card is made and filed. This file can often supply answers to requests because only ring systems and fragment chains that have been entered into the system are recorded. If the card is not there, it is not necessary to search the coordinate index file for analogs.
the coordinate index file for analogs.

**Supplier Codes.**—Suppliers of the compounds screened by the CCNSC are assigned a numerical designation; for example, number 1 is the now defunct Chemical Biological Coordination Center. Number 2 is Sterling-Winthrop Research Institute, number 3 is Johns Hopkins School of Medicine, etc. As the personnel at the various companies and universities change, each new person receives a new letter designation which is placed after the number previously assigned.

It had been the practice to punch a plastic card for the individual as well as for the company or university so that a person could be distinguished from others at the same location. About 1,000 cards were required to record

the suppliers and their representatives concerned with the first 54,000 compounds.

The method used for recording the RRI # has been adapted for use with the Supplier Code. At present, there are 650 numbers assigned to the suppliers plus the alphabetic designations for the individuals.

With the following plastic cards, all the Supplier Codes can be entered. The letter "P" is reserved for purchased compounds. The letters "I" and "O" are not used because they may be mistaken for numbers.

| | | | |
|---|---|---|---|
| 000 | -00 | --0 | A through (Z) |
| (100) | (-10) | --1 | |
| 200 | -20 | --2 | |
| 300 | -30 | (--3) | |
| 400 | -40 | --4 | |
| 500 | -50 | --5 | |
| 600 | -60 | --6 | |
| 700 | -70 | --7 | |
| 800 | -80 | --8 | |
| 900 | -90 | --9 | |

By punching the three numerical cards of the Supplier Code, the company or university submitting the compounds will be recorded. The fourth card to be punched is the alphabetical designation for the individual chemist. Cards are marked in red to avoid confusion with Ring Index numbers.

To retrieve the NSC numbers of compounds from a supplier, one has only to reassemble the code number, for example

Supplier Code Number 113-Z
Ohio State University
Dr. Robert S. Brown

Number cards 100, -10, and --3 will show all compounds submitted through the university. With the addition of the alphabetic "Z" card, only those submitted by Dr. Brown will be listed.

With this system, compounds from any one of more than a thousand suppliers can be retrieved with the use of the appropriate four of the fifty-four plastic Supplier Code Cards.

## SUMMARY

Two main divisions of organic structures—ring systems and the fragment chains—have been coded in a novel way that permits the recall of these structures with exactness without having a separate plastic card for each ring or fragment chain. Additional new ring systems and fragment chains can be recorded on the basic set of plastic cards so that the number of cards in the file will remain approximately the same. These improvements not only have reduced the file size from over 1,500 to less than 500, but also have resulted in an increased depth of entry into the structure. These improvements may be integrated with existing electronic data processing equipment should anyone wish to adapt them to his particular use.

REFERENCES

(1)  D. F. Gamble, "A Coordinate Index of Organic Compounds," Presented before the Division of Chemical Literature, 127th National Meeting of the American Chemical Society, Cincinnati, Ohio, March 21, 1955.

(2)  W. A. Wildhack, and J. Stern, "Peek-A-Boo System," in "Punched Cards—Their Applications to Science and In-dustry," R. S. Casey, J. W. Perry, A. Kent, and M. M. Berry, Ed., Reinhold Publishing Corp., New York, N. Y., 1958, Chapter 6.

(3)  A. M. Patterson, L. T. Capell, and D. F. Walker, "The Ring Index," 2d Ed., American Chemical Society, Washington, D. C., 1960.

(4)  "The Naming and Indexing of Chemical Compounds by Chemical Abstracts," Introduction to the Subject Index of Vol. 56, (Jan.–June, 1962) of Chemical Abstracts.

## AUTHOR INDEX
### VOLUME 4

## SUBJECT INDEX
### VOLUME 4