# Graphical Representation of Chemical Structures in Chemical Abstracts Service Publications

ALAN L. GOODSON

Chemical Abstracts Service, Columbus, Ohio 43210

The main features of the comprehensive formatting guidelines used for creating chemical structure diagrams for Chemical Abstracts Service (CAS) publications are described. The guidelines were designed to standardize format as much as possible and to conserve column space in publications. Methods used for reducing crowding in diagrams are described.

The chemical structure diagram (hereafter referred to as "diagram" or "graphical representation") is a two-dimensional graphical representation of the spatial arrangement of the atoms and associated bonding of a substance. Its origin and development as a universal means of conveying information between chemists have paralleled the development of atomic theory.[1] Today, the diagram is central to the application of the computer to such fields as chemical education[2] and synthesis planning.[3] For this reason, considerable effort has been expended on developing efficient methods for computer manipulation of diagrams. However, Rush[4] recently drew attention to the lack of "general rules for generating graphical representations" of chemical structures. He stated that "the only generally available rules . . . are those contained in the 'Ring Index'"[5a] and that they "are not readily amenable to computerization". The "Ring Index" rules "were established to aid in the naming (uniquely) of rings and systems of rings" so they do not cover acyclic systems, substituents, stereochemistry, etc., and are therefore not comprehensive.

To try to develop a set of hard-and-fast formatting rules for generating graphical representations of chemical structures would be impractical because there are always exceptions. This is why Rush finds "the graphic dialect . . . is at present essentially an informal, intuitive language". It is more practical to develop a set of preferred guidelines, permitting use of less-preferred guidelines when required by the nature of the chemical structures. To be of value, such guidelines should not be confined to one class of chemical structures (e.g., ring systems), but should be comprehensive, including acyclic systems, substituents, stereochemistry, etc.

CAS has been using such comprehensive formatting guidelines for nearly 10 years. Because they have been used for diagrams published in *Chemical Abstracts* (*CA*) issues (Figure 1a) and indexes (Figure 1b) and in the "Parent Compound Handbook"[6] (Figure 1c), the guidelines have also included the following basic principles: (1) diagram format should be standardized as much as possible; (2) diagrams published in *CA* issues should resemble as closely as possible the corresponding diagrams published in the original documents; (3) diagrams published in *CA* volume indexes and the "Parent Compound Handbook" should reflect their associated index names; (4) to reduce crowding, "nonessential" locants should be omitted from the diagrams appearing in *CA* volume indexes and the "Parent Compound Handbook"; (5) column space should be conserved in all publications and services.

The guidelines have been published in a very abbreviated form in the "Index Guide" since 1972.[7] They were designed for manual drawing of diagrams to be published in *CA* issues and indexes and proved to be of equal value after computer processing of graphical data was introduced at CAS in 1974.[8,9]

Conversion from manual processing to computer processing of graphical data began in 1969 as follows: The first step was the transfer of responsibility for production of publication-quality diagrams from the printer to CAS. A limited set of internal formatting guidelines was developed for use by the new CAS graphic arts staff to ensure that diagrams appearing in *CA* volume and collective indexes would conform to the existing formatting standards.

We soon recognized, however, that for the conversion to continue successfully, a more comprehensive set of formatting guidelines was needed to ensure consistency of format among diagrams appearing not only in *CA* volume and collective indexes but also in *CA* issues and future CAS services.

Diagrams drawn in accordance with the new guidelines first appeared in Volume 76 of *CA* issues and volume indexes, which were published in the first year of the 9th collective index period (1972–1976). Conversion of the index diagrams to the new format required manually redrawing the Structure Cut File, i.e., the 40 000-diagram master file from which the *CA* volume and collective index diagrams were selected. The effect of the new guidelines is exemplified in Figure 2.
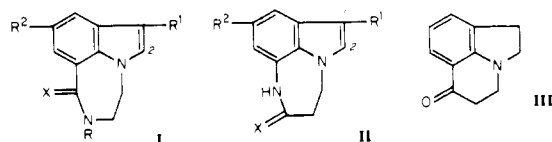
We similarly redrew the diagrams appearing in the Introduction to the "Index Guide"[7a] and a file of over 600 diagrams that occur frequently in *CA* issues (i.e., the Repetitive Structure File).

During this period, textual material for *CA* issues and indexes was being transferred initially from the CAS data base to film via an IBM 2280 photocomposition device. The IBM 2280 was later replaced by an Autologic APS-4 photocomposition device which produced full-page positives. Diagrams, however, were still hand drawn, photoreduced, and manually added to "windows" in the film or page positives already containing the text. To bring the procedure for handling diagrams into line with that for handling text, we built computer-readable files by copying the redrawn manual files with the aid of a Digital Equipment Corporation PDP-15 graphics system and storing these files on disk or tape for an IBM 370/168 host computer. We took the opportunity at this time to make minor modifications to the guidelines as dictated by experience. Diagrams to be published in *CA* issues were created similarly. When required for publication, the graphic data were selected from the data base and photocomposed, together with associated textual material, on an Autologic APS-4, producing full-page, publication-quality output. This procedure has been described more fully elsewhere.[8]

In this manner, we built computer-readable versions of (1) the Structure Cut File, which now contains over 52 000 diagrams, and (2) the Repetitive Structure File, an entry of which is illustrated in Figure 3. We also built two new reference files. One file, the Fragment Dictionary, contains over 400 graphs (in which atoms and bonds are not specified) of commonly occurring structural fragments, such as 1- through 4-atom bridges in rings of 3–8 atoms (e.g., Figure 4a), and of ring systems difficult or time consuming to construct, such as boron cages (e.g., Figure 4b). Each fragment has an associated intelligence-containing recall code as an aide-mémoire. Use of the Fragment Dictionary improves consistency of format and production rates for construction of diagrams (including the Repetitive Structure File) to appear in CAS publications and services. For example, Figures 4c and 4d,
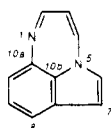
GRAPHICAL REPRESENTATION OF CHEMICAL STRUCTURES

*J. Chem. Inf. Comput. Sci., Vol. 20, No. 4, 1980* **213**

### (a) *CA* abstract

42573a  4,5 - Dihydropyrrolo[3,2,1 - jk][1,4]benzodiazepin - 7 - (6H)ones and 4,5-dihydropyrrolo[1,2,3-ef][1,5]benzodiazepin-6-(7H)ones. Hester, Jackson B., Jr. (Upjohn Co.)  U.S. 3,734,-919  (Cl. 260-239.3T;  C 07d), 22  May 1973,  Appl. 03  Nov 1969; 8 pp.  Division of U.S. 3,642,821 (C.I 76: 153805k).  Pyr-

rolo[3,2,1-jk][1,4]benzodiazepines  and  pyrrolo[1,2,3-ef][1,5]-benzodiazepines [I, II;  R  =  H,  Et₂N(CH₂)₂, Me₂N(CH₂)₃, Me; R¹ = H, Me₂NCH₂, R² = H, Cl; X = H₂, O] with tranquilizing and anticonvulsant activity were prepd.  Thus, the pyrroloquinolinone III, prepd. by known methods was treated with NaN₃ in polyphosphoric acid to give a mixt. of the 1,2-dihydro analogs of I and II (R–R² = H, X = O), which are easily sepd. Dehydrogenation of these with Pd/C gave I and II (R–R² = H, X = O).  I (R = R¹ = H, R² = Cl, X = H₂) had an ED₅₀ = 10 mg/kg (i.p., mice) in the chimney test.

### (b) *CA Chemical Substance Index* Entry

Pyrrolo[1,2,3-ef]-1,5-benzodiazepine [209-56-3]

——. 9-chloro-1,2,3,4,6,7-hexahydro- [28740-83-2].
 76: P 127029r
 prepn. and central nervous system activity of, 82: P
 43482u
——. 9-chloro-1,2,3,4,6,7-hexahydro-2-phenyl-
 monohydrochloride [28740-99-0], 77: P 5541b
——. 9-chloro-1,2,3,4-tetrahydro- [28740-90-1], 76: P
 153805k; 79: P 42573a
——. 1,2,3,4,6,7-hexahydro-
 dihydrochloride [28740-82-1]
 76: P 127029r; 80: 120720h
 prepn. and central nervous system activity of, 82:
 P 43482u
——. 1,2,3,4-tetrahydro- [27158-85-6], 76: P 153805k;
 79: P 42573a

### (c) *Parent Compound Handbook* Entry

GMKPH                                    209-56-3
  Pyrrolo[1,2,3-ef]-1,5-benzodiazepine
  8CI Pyrrolo[1,2,3-ef][1,5]benzodiazepine
    C₁₁H₈N₂
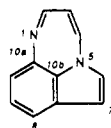    T567 1A M BN JNJ
    3-Ring System
    5,6,7
    C₄N-C₆-C₅N₂

**Figure 1.** Graphical representation of a chemical structure, drawn according to the formatting guidelines, as it appears in some CAS publications.

which appear in *CA* indexes and the "Parent Compound Handbook", were derived from Figures 4a and 4b, respectively. The second file is the Ring Image File which contains all (i.e., over 20 000) ring graphs identified by the CAS Chemical Registry System. A ring graph is illustrated in Figure 5. The Ring Image File is used to support Algorithmic Structure Display, as described elsewhere.[9]

## BASIC PRINCIPLES

We mentioned above that, in addition to being comprehensive, the guidelines include five basic principles. The first of these principles—standardization of format—was achieved by (1) portrayal of all multiple (especially double) bonds in
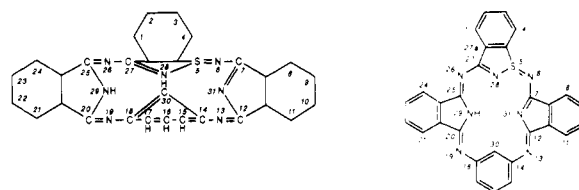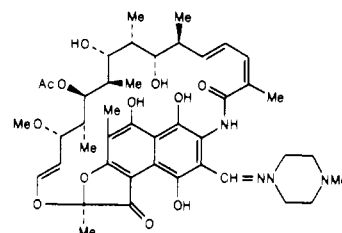
**Figure 2.** Effect of the new formatting guidelines on the appearance of chemical structure diagrams.

I. u338 Rifampicin

**Figure 3.** Repetitive Structure File entry.

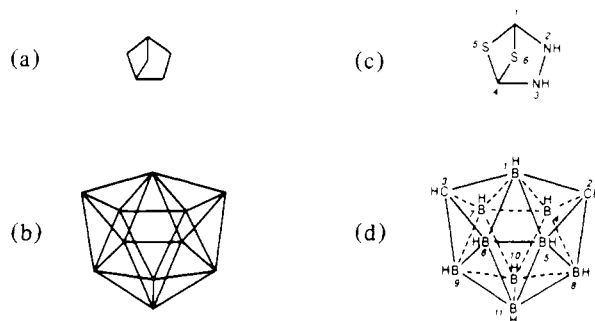Fragment Dictionary Entry          Published Diagram

(a)                    (c)

(b)                    (d)

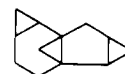**Figure 4.** Examples of Fragment Dictionary entries and their use.

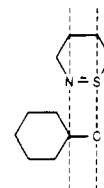**Figure 5.** Ring Image File entry.

**Figure 6.** Standard bond length.

*CA* index diagrams (see Figure 1) as well as in *CA* issue diagrams and development of the necessary rules covering placement of the multiple bonds in all diagrams; (2) use of a standard bond length for both acyclic and ring bonds, as illustrated in Figure 6; (3) standardization of ring shapes; (4) use of previously published orientation rules for fused ring systems.[5b] We standardized ring shapes by using templates for hand-drawn diagrams and, subsequently, by incorporating the standard ring shapes into the menus used for computer-assisted creation of diagrams.[8] We use combinations of the standard ring shapes, for example, in creating diagrams of fused (see Figures 3 and 5) and von Baeyer ring systems. We draw "no-atom bridges" (i.e., direct bonds) of von Baeyer ring
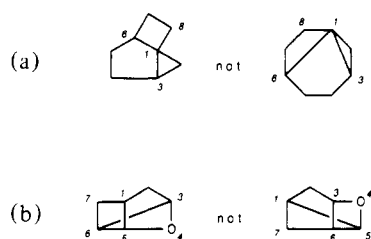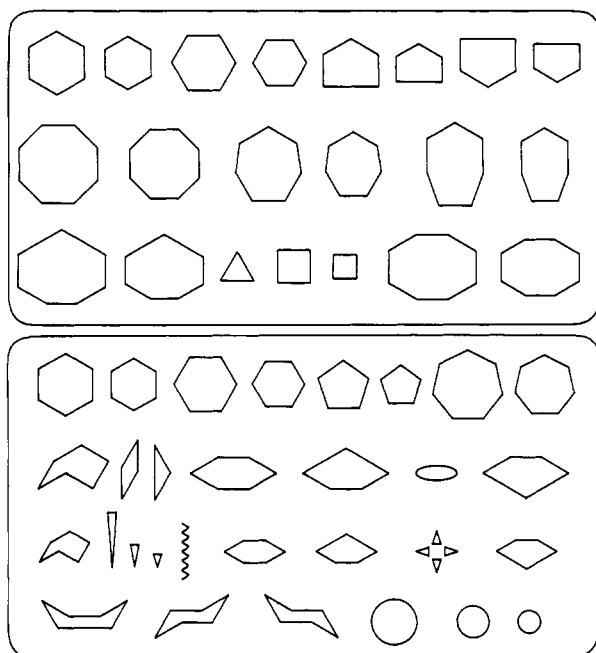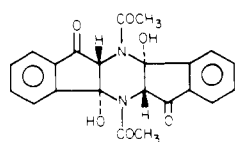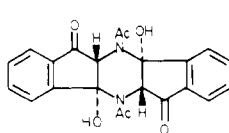
(a)  not 

(b)  not 

**Figure 7.** von Baeyer ring systems.



**Figure 8.** Templates used by graphic arts staff for drawing diagrams published in both *CA* issues and indexes.

Primary Literature Diagram          *CA* Issues Diagram



**Figure 9.** Comparison of diagrams from the primary literature and from *CA* issues.

systems as ring fusions (Figure 7a). Where there is a choice in *CA* index diagrams, we draw the bridge with the lowest numbered bridgehead as the ring fusion (Figure 7b).

We designed new templates because commercially available ones were of limited practical utility (Figure 8). We found that two sizes of template were necessary (with standard bond lengths of $^7/_{16}$ and $^5/_8$ in.), the larger being used for crowded diagrams (see "Crowded Diagrams"). Use of these templates and the new guidelines significantly increased the productivity of the graphic arts staff. The format of computer-produced diagrams was similarly standardized, the features of the templates being incorporated into the computer programming.

The second principle requires that diagrams published in *CA* issues should reflect the contents of the original documents and should not, therefore, differ radically from diagrams published in the primary literature, as illustrated in Figure 9.

Diagrams published in *CA* volume indexes and the "Parent Compound Handbook" reflect the preferred *CA* names of the substances they depict in accordance with the third principle, as illustrated in Figure 10. Here, the ring system is named as containing two bridges (etheno and methano) which are
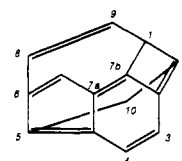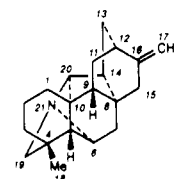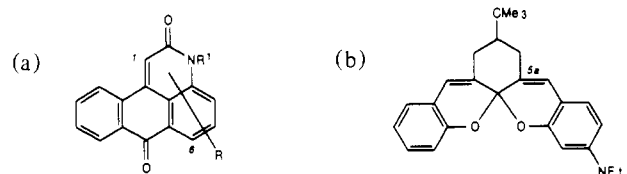
1,6-Etheno-2,5-methano-1H-cyclobut[e]indene

[18429-39-5]



**Figure 10.** Illustration of how *CA* volume index diagrams reflect names of substances and broken bond lines that do not imply stereochemistry.



**Figure 11.** Illustration of broken bond lines that do imply stereochemistry.

(a)  (b) 

**Figure 12.** *CA* issue diagrams with locants.

given a lower priority in the drawing of the diagram: where the methano bridge crosses the parent ring system (i.e., 1*H*-cyclobut[e]indene), it is the methano bridge bond lines that are broken where they cross the ring parent bond lines to avoid any implication of the presence of an atom at each crossover. Such broken bond lines convey no implication of stereochemistry; i.e., the broken bonds in Figure 10 do not imply that the methano bridge is "behind" the parent ring system. This constrasts with stereochemical diagrams (e.g., Figure 11), where broken bond lines do imply stereochemistry.

Locants are shown in *CA* issues diagrams when two or more substances differ by the position of one or more substituents, the positions of the substituents being specified in the text of the abstract. For example, in Figure 12a the variable group R was attached to position 1 or position 6, both of which were highlighted in the (general) diagram, and the positions were referred to in the accompanying text of the abstract as follows: "I (R = 1-NH₂, 6-NH₂, 1-OH; R¹ = H, Me". Such general—or Markush— structures (see also Figure 1a) are published only in *CA* issues, never in *CA* indexes.

Locants are also shown in *CA* issue diagrams when there is need to refer to a specific atom or bond whose position is not immediately apparent. For example, locant 5a of Figure 12b was referred to in the text of the abstract as follows: "which directs protonation into position 5a".

In contrast to *CA* issue diagrams, locants are assigned, with few exceptions (e.g., the two oxygen atoms of spirostan), to every nonhydrogen atom in *CA* index diagrams. Because these diagrams are two-dimensional representations of three-dimensional chemical structures, there is sometimes so much information to be displayed in a confined area that crowding occurs. A number of methods used for reducing crowding are discussed under "Crowded Diagrams". One method differs sufficiently from previous practice that it became one of the principles (the fourth) incorporated into the guidelines.

Previous practice required that the locant number of every nonhydrogen atom that was not a ring junction be shown. We felt this practice sometimes provided too much data, thus
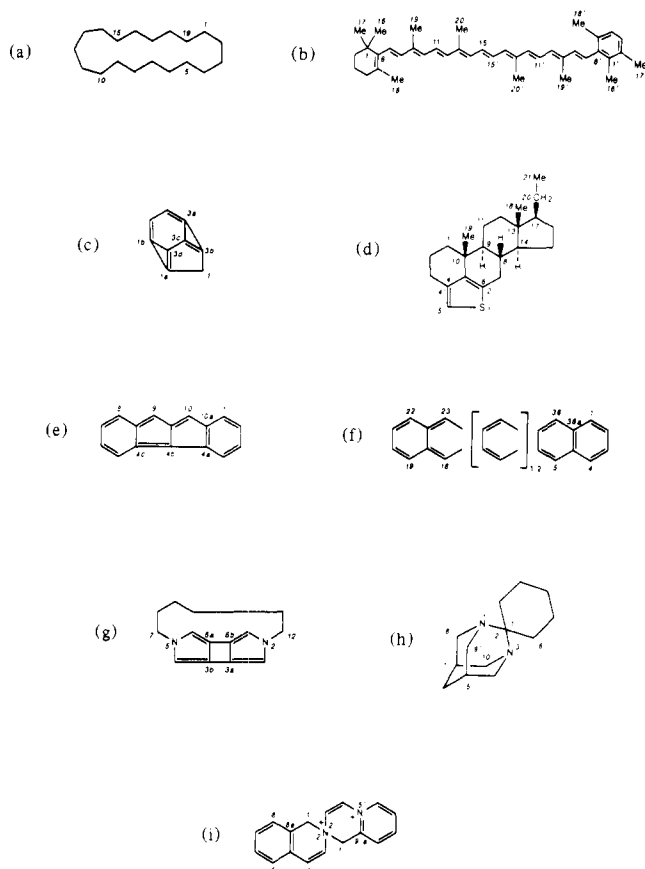
**Figure 13.** Citation of locants.



(a) $Co_5(NO_3)_2(SO_4)_4$

(b) $BOC-Gly-Asn-Leu-Ser(Bzl)-Thr(Bzl)-N(CO_2Me)(CH_2)_5CO_2H$

**Figure 14.** Examples of linear representations of chemical structures incorporated into the texts of *CA* abstracts.



**Figure 15.** Horizontal orientation of diagrams.

contributing to overcrowding, and at other times provided too little, e.g., when the numbering of ring fusions was not immediately obvious. We therefore developed the following rules that reduce crowding of diagrams but, at the same time, provide more information than previous practice.

Rule 1. Locant numbers are assigned by starting at the top or top right of a ring system and proceeding in a clockwise manner[10] (Figure 13, all diagrams except those of natural products[11]).

Rule 2. Locant numbers are shown for the first and last atoms in the numbering sequence and for all hetero atoms (Figure 13, all diagrams).

Rule 3. Where locant numbers are not shown for other reasons, for both cyclic and acyclic sequences of atoms greater than eight, every locant number that is divisible by five is shown (Figure 13a,b).

Rule 4. The locant number of every interior fusion is shown, together with the exterior ring fusion beginning the sequence (Figure 13c), except in steroid fusions (Figure 13d).

Rule 5. When more than one ring fusion occurs together in a ring periphery, each of the ring fusion locant numbers is shown (Figure 13c,e).

Rule 6. Where single ring fusions have not been shown because of the preceding rules, the locant numbers before and after the ring fusions are shown (Figure 13e,f).

Rule 7. At ring fusions with two numbering systems, the locant numbers of the ring fusions are shown (Figure 13d).

Rule 8. The locant numbers of bridgeheads and the first and last locant numbers of each bridge are shown (Figure 13g,h).

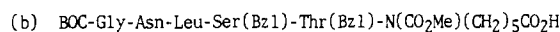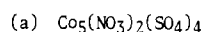Rule 9. Locant numbers of all spiro atoms are shown (Figure 13h).
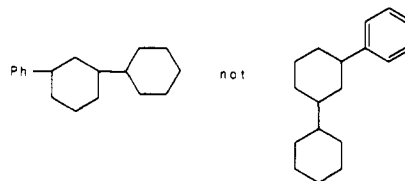
Rule 10. Where there are breaks in the sequence of locant numbers, e.g., where there are repeating units, the locant numbers immediately preceding and following the breaks are shown (Figure 13d,f).

Rule 11. When a charge sign and a locant number are located at the same position, the locant number is placed in the least congested position, usually outside the ring, and the charge sign in the next least congested position (Figure 13i).

In accordance with the fifth principle, column space is conserved in *CA* issues by incorporating as many linear representations of chemical structures as possible into the text. Such representations range widely from simple acyclic structures (e.g., Figure 14a) to quite complex structures such as polypeptide intermediates (e.g., Figure 14b, in which BOC is defined as $-CO_2CMe_3$ and Bzl as $-CH_2Ph$). Further space is conserved and uniformity improved by use of standard symbols. Some are of general use, such as Me, Et, Pr, Bu, and Ph for methyl, ethyl, *n*-propyl, *n*-butyl, and phenyl groups, respectively, while others are well established in specialized subject areas, such as the amino acid symbols in Figure 14b. Some symbols are restricted in use. Ac and Bz represent acetyl and benzoyl groups, respectively, only when they are attached to atoms other than carbon and are not numbered in *CA* index diagrams; at other times the groups are represented by MeCO and PhCO instead. Nonstandard symbols (e.g., BOC and Bzl in Figure 14b) are always defined.

Another method of conserving column space is by horizontal orientation of diagrams, as illustrated by a *CA* issue diagram in Figure 15. However, column width is limited in *CA* issues to the equivalent of about 13 fused six-membered rings of standard size, while in *CA* indexes it is equivalent only to about 8.5. Large diagrams therefore require some form of compression, which can be achieved in a variety of ways. Some structures can be reduced in size, but this is ultimately limited by loss of detail. Where an atom or group of atoms is repeated a number of times, the atoms can be represented by repeating units, as illustrated by Figure 16a, a *CA* index diagram. Where a larger group of atoms is common to more than one diagram in an abstract, the common atoms are implied by use of a wavy line, as in Figure 16b. In both *CA* issues and indexes, some diagrams are continued below when they exceed the column width (Figures 16c and 16d, respectively). Finally, large diagrams that cannot be compressed by any of the methods just described are rotated counterclockwise through 90°, as illustrated in Figure 16e, overriding the standard orientation rules.

## CROWDED DIAGRAMS

In addition to not citing "nonessential" locants, which was discussed above as the fourth principle, crowding of diagrams is also reduced as follows:

While *CA* issue diagrams resemble diagrams published in the original documents, and the need and opportunity to reduce
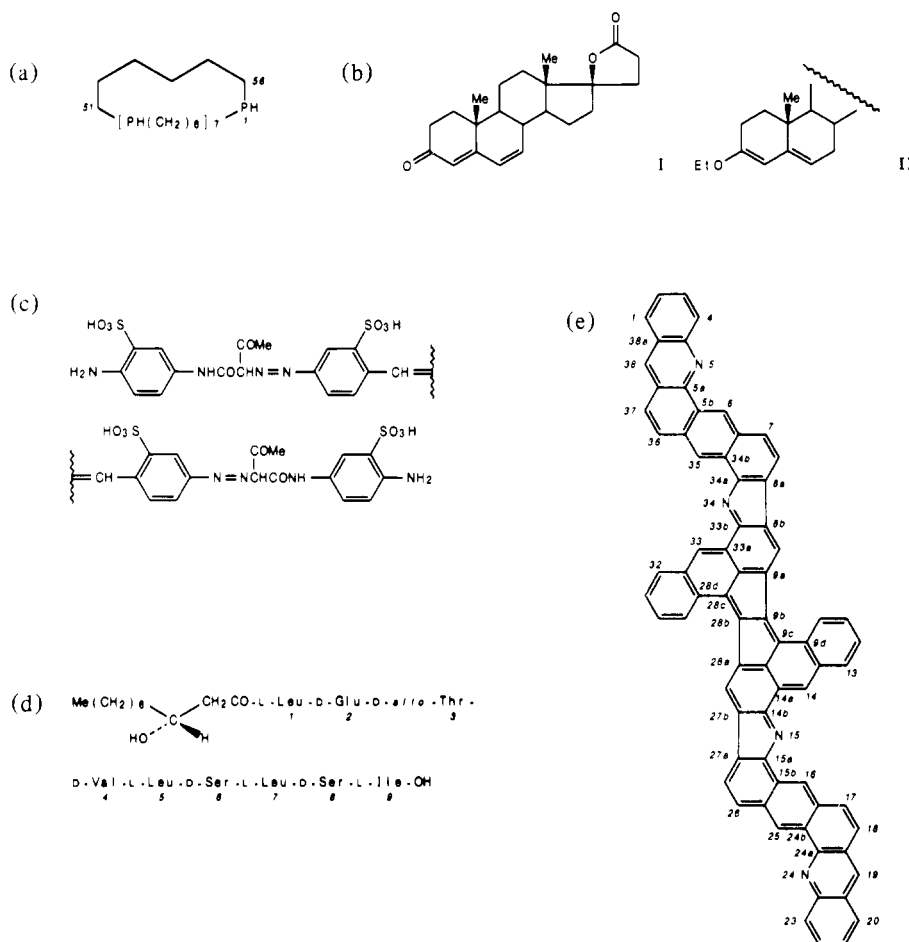
**Figure 16.** Methods employed with large diagrams to conserve column space.
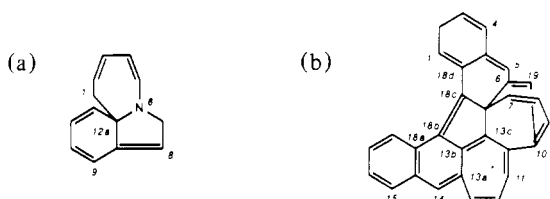


**Figure 17.** Distortion of ring systems.

crowding are therefore limited, crowding in *CA* index diagrams is reduced in a number of ways.

When, in *CA* index diagrams, an atom is common to three rings and two of the rings are not ortho-fused, as in 1*H*,7*H*-azepino[2,1-*i*]indole (Figure 17a), the component[12] (azepino) is distorted, while the base component[12] (indole) is not.
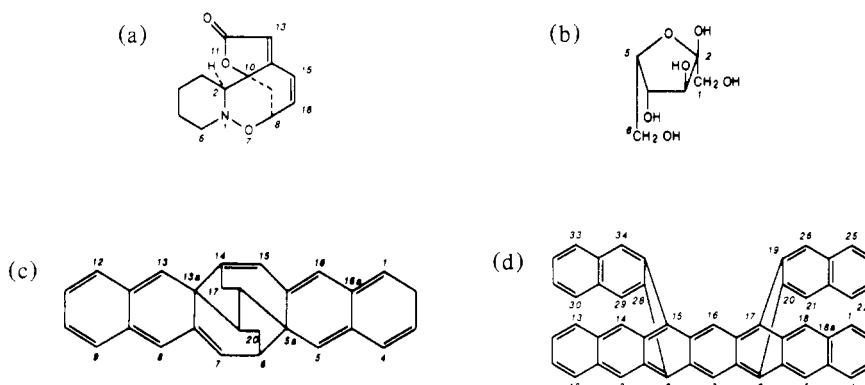
However, if distorting the component causes crowding or other problems, the base component is distorted, as in 2*H*-10,6-methenodibenzo[*c,g*]heptaleno[2,1,10-*jkl*]fluorene (Figure 17b).

When atoms or groups attached to ring systems are crowded, the atom or group is moved to a less crowded position, as in Figures 3 and 18a, or a bond is lengthened, as in Figure 18b. If the inside of a ring is crowded with, say, a complex bridge, the ring (and other associated rings, if any) is enlarged, as in Figure 18c. If, however, ring enlargment is inadequate, the bridge is drawn outside the ring, as in Figure 18d.

**Note Added in Proof.** The algorithm developed for generating chemical structure diagrams,[9] based on the guidelines described, is now being used in graphic display of structures
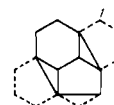


**Figure 18.** Crowded diagrams.

retrieved by CAS ONLINE, a new search service.

## SUMMARY

Because development of rigid formatting rules for generating graphical representations of chemical structures would be impractical, a set of comprehensive guidelines has been developed at CAS. The guidelines provide preferred formats for acyclic structures, rings and ring systems, and representation of stereochemistry. They also describe a number of methods for reducing crowding in diagrams. The guidelines, which have been used for CAS publications for about 10 years, were originally developed for manually drawn diagrams, but have proved to be of equal value for subsequent computer-assisted generation of diagrams. More general use of these guidelines (such as in chemical education, research and development, and information storage and retrieval) would facilitate communication of chemical structure information among scientists.

## REFERENCES AND NOTES

(1) Crosland, M. P. "Historical Studies in The Language of Chemistry"; Dover: New York, 1978.
(2) For example: Soltzberg, L. J. "Computer Graphics for Chemical Education", *J. Chem. Educ.* **1979**, *56*, 644–9.
(3) For example: Corey, E. J.; Wipke, W. T. "Computer-Assisted Design of Complex Organic Syntheses", *Science* **1969**, *166*, 178–92; "Computer-Assisted Drug Design", *ACS Symp. Ser.* **1979**, *No.* 112.
(4) Rush, J. E. "Handling Chemical Structure Information"; *Annu. Rev. Inf. Sci. Technol.* **1978**, *13*, 209–62.
(5) (a) Patterson, A. M.; Capell, L. T.; Walker, D. F. "The Ring Index", 2nd ed.; American Chemical Society: Washington, DC, 1960. (b) See pp xii–xiii. See also ref. 10.
(6) Blake, J. E.; Brown, S. M.; Ebe, T.; Goodson, A. L.; Skevington, J. H.; Watson, C. E. "Parent Compound Handbook—Successor to the Ring Index", *J. Chem. Inf. Comput. Sci.* **1980**, *20*, 162–7.
(7) "Chemical Abstracts Index Guide": (a) Vol. 76, 1972, Introduction, Paragraphs 8, 15, 140–63, and 202–12. (b) Vol. 76–85, Cumulative, 1972–6, Appendixes, Paragraph 15.
(8) Blake, J. E.; Farmer, N. A.; Haines, R. C. "An Interactive Computer Graphics System for Processing Chemical Structure Diagrams"; *J. Chem. Inf. Comput. Sci.* **1977**, *17*, 223–8.
(9) Dittmar, P. G.; Mockus, J.; Couvreur, K. M. "An Algorithmic Computer Graphics Program for Generating Chemical Structure Diagrams"; *J. Chem. Inf. Comput. Sci.* **1977**, *17*, 186–92.
(10) The wording in this rule is a brief generalization of orientation and numbering rules A-22, A-31.2, A-32.23, A-34.2, A-41.2, B-1.5, B-3.4, B-10, C-12.7, C-12.8, and C-15 in: IUPAC. "Nomenclature of Organic Chemistry: Sections A, B, C, D, E, F and H"; Pergamon: Oxford, 1979. The rules in this reference supersede those in ref 5.
(11) Figure 13c appears to violate the orientation and numbering rules of ref 10. However, this ring system is drawn for orientation purposes as though every component ring were six membered, as shown, i.e., as if the ring system were pyrene. While pyrene numbering begins at position 1 in the accompanying diagram, locant 1 of Figure 13c is assigned



to the first available position according to rule A-22 of ref 10.
(12) These terms are used as in rule A-21.5 of ref 10.

# Method for Generating a Chemical Reaction Index for Storage and Retrieval of Information

MARGARET A. MOSBY*

Tompkins-McCaw Library, Medical College of Virginia, Virginia Commonwealth University, Richmond, Virginia 23298

LEMONT B. KIER

Department of Pharmaceutical Chemistry, Medical College of Virginia, Virginia Commonwealth University, Richmond, Virginia 23298

A new method for indexing chemical reactions is described. The calculation of a reaction connectivity index results in a unique number. This number does not provide hierarchical or relational information. It encodes the concept of the reaction process in a unique identifier which is suggested for use, much as the CAS Registry number is used, to optimize ease of storage, manipulation, and retrieval from large computer files.

## INTRODUCTION

The retrieval of information in chemistry has become increasingly complex. The total amount of literature to be covered in an exhaustive search is now so vast as to preclude a systematic manual examination of the available sources. It is, however, possible to use a computer for this task. The Chemical Abstracts Service files, Derwent's files, and the American Petroleum Institute's file, as examples, are all available for computer access in some form through one or more data-base vendors. The major benefit, of course, is the speed with which a body of relevant material can be extracted from a much larger file of bibliographic citations or chemical data. The primary drawbacks with this approach are associated with the forms in which data is put into the computer files, with the processing capabilities of the computer, and with the software commands available to search the files.

The main objective of the information scientist performing a literature search is to maximize relevant retrieval without sacrificing completeness. This objective is most easily achieved when searching files that are structured by the use of a controlled vocabulary and/or other formal systems. Hence, the first recognized need in chemical information retrieval was a unique way to identify chemical compounds that would eliminate the confusion of multiple names used for the same structure. The Chemical Abstracts Service Chemical Registry System has done this since 1965, based on a system first developed by Gluck.[1] The registry number is based on the representation of the chemical structures in the form of topological tables and is unique and unambiguous. Now it is at least as important to the chemical community to have