Chemical Registries—in the Fourth Decade of Service[†]

Robert E. Buntrock‡
Buntrock Associates, Inc., 11335 300th Avenue NW, Princeton, Minnesota 55371
Received August 4, 2000

Methods for precise yet usable descriptions of virtually any topic are essential, especially for identification of chemical compounds and materials. Systematic nomenclature and precise chemical structures, if known, are the ultimate in description of chemical compounds. However, there has always been a pervasive need for brief, yet precise methods to register chemical compounds for use in both indexing and retrieval of additional information. The most predominant chemical registration system is the Chemical Abstracts Service (CAS) Registry System, begun in 1965. The history and use of this system will be described, and its importance to a number of disciplines, not just chemistry, will be described.

Definition: precise definition is at the heart of every discussion or dialogue, "learned" or otherwise (Figure 1). After all, that is the function of language—to allow interpersonal communication on an organized basis.

As observed earlier, 1 chemical information has two unique features compared to all other species of information: (1) chemical structure and (2) chemical reactions, which include chemical structures, physicochemical entities, and vector functions.

However, chemical structure is not always well defined, and even if it is, precise nomenclature is required to define chemical entities.^{2,3} It soon becomes obvious that neither structure nor precise chemical nomenclature is that utilitarian, especially for use by the nontechnical (or nonchemical) public but even for effective, precise communication between chemists.⁴⁻⁶

Examples of nomenclature difficulties in the public sphere abound, including a very prevalent confusion of "silicon" and "silicone", a problem confounded by different spellings for either entity in some languages other than English. Word processor spell checkers do not always solve the problem of poor editing, a phenomenon on the increase. For years, the author has advised educators, including high school teachers, to teach at least some chemistry from the newspapers and mass media. Hardly a day goes by without the possibility of finding some errors in reporting about science, technology, and especially chemistry. Some are quite humorous, but some are ignorant or even dangerous.

A recent news article described problems with controlling abuse of GHB, a prominent "date rape" or "party drug". GHB, or gamma hydroxybutyrate, was banned by the FDA in 1991, but related chemical ingredients—with similar hazardous effects—of party drugs are proliferating, including GBL, gamma butyrolactone, and BD, 1,4-butanediol (of course the product street names are considerably more exotic). Even if better nomenclature or CAS Registry Numbers were used for identification of ingredients in products, would it help? That is debatable.

definition 1: an act of determining

registry 3a: a place of registration

4a: an official record book

Figure 1. Webster's New Collegiate Dictionary.

24,055-9 1,4-Butanediol, 99+% [110-63-4] HO(CH2)4OH

FW 90.12 mp 16 deg. ... 2g ... ; 100 g ...

B8480-7 1,4-Butanediol, 99% [110-63-4] HO(CH2)4OH

... 1kg ... ; 3kg ... ; 18 kg ...

Figure 2. Listing for 1,4-butanediol, Aldrich Catalog, 1998–1999.

For decades, chemists—always pragmatic—have responded by developing lists of chemical compounds, with associated attributes. Such lists are chemical registries (Figure 1), and the individual compounds can be tersely and effectively described by an ID number. Although commercial chemical catalogs still list their "registries" by nomenclature—with all of the associated problems—compounds are then ordered by the catalog number, making it a registry number (Figure 2). Note that in the catalog excerpt, there are two listings for 1,4-butanediol, one for 99+% purity grade, and one for 99% purity, each with a separate catalog number. Both, of course, have the same CAS Registry Number (CASRN), shown in brackets. The Aldrich catalog may be either just a serial number or it might convey more information than simple identity.

Organizations in the business of making new chemical compounds, especially in the pharmaceutical and agricultural chemical industries, developed internal registries very early in their history. One result of these registries was that new compounds with potentially marketable activity were always referred to by the registry number—typically a two letter company code followed by 4 to 5 numerals—at least until they acquired a trade name. One somewhat amusing side effect that began appearing in the 1960s was the appearance of two sets of numbers: one, a serial number, used only internally, and the second, a randomly assigned number for use outside the organization. The reason: to not allow the competition to infer the level of your research activity.

Abstracting and indexing organizations began to see the advantages of registry numbers not only for codifying

[†] Paper CINF 40, 218th National Meeting, American Chemical Society, Aug. 24, 1999, New Orleans, LA, presented at the Skolnik Award Symposium honoring Stuart M. Kaback.

[‡] Corresponding author phone: (763)389-8370; e-mail: buntrock2@earthlink net

information in their files but also for allowing another means of searching the files. Prominent examples are the Chemical Abstracts Service (CAS) Registry System, the Derwent List of Registry Compounds, and the Beilstein file. Every compound in the Beilstein file has a Beilstein Registry Number (BRN) as well as Beilstein System Numbers and Lawson Numbers. The BRN is strictly a serial number with no information conveyed other than identity. However, the latter two numbers implicitly contain information on chemical composition and structure. The Beilstein File also contains CASRN for about 2/3 of the compounds in the file. Divergence in usage of CASRN between the Beilstein database and the CAS Registry file does occur however. Some Beilstein records have more than one CASRN per BRN, largely due to differences in policy in registering stereoisomers.

Derwent Registry Numbers (DRN) were originally assigned only to about 2000 compounds commonly encountered in the patent literature and are strictly serial numbers. For the remainder of this article, the discussion will be centered on the CAS Registry System, reflecting the widespread use and importance of CASRN.

As defined by CAS, ^{8,9} a registry system is an inventory of chemical substances with means of inputting, processing, searching, retrieving, and outputting information about the substances. Although the fundamental representation of substances is by structure, structure representation can be accomplished by the following: (1) nomenclature, (2) linear notation (e.g., Wiswesser Line Notation or WLN, SMILES), (3) fragmentation codes (e.g., GREMAS, Derwent), and (4) connection tables (e.g., Morgan, DARC). CAS chose the Morgan connection table method for the CAS Registry File.

The CAS Registry System was originally designed as a labor saving device in support of the indexing effort used to prepare the Chemical Abstracts database. Prior to the advent of CAS Registry in 1965, with the exception of a list of ca. 2500 common chemicals, there was no good way to determine if a compound had already appeared in the database. Each potentially new compound had to be drawn and named and the name compared to the list of index names. Since the advent of Registry II in 1968, compounds to be indexed undergo name and structure matching procedures against the Registry File. Those found are henceforth referred to by CAS Registry Number (CASRN). Those not found by either method are added to the file as new compounds and a Registry Number is generated.

The design, implementation, and performance of the CAS Registry System have been well documented over the years. 9,10 However, the remainder of this presentation will focus on the use of CAS Registry Numbers for searching of a number of files and use in other inventories.

The format of CAS Registry Numbers is intentionally unique. As mentioned previously, the number is a serial number; the next available number is assigned to the next new compound to be entered into the file. (As an aside, some Registry Numbers have disappeared from the file. For example, formaldehyde—CASRN 50-00-0 or compound number "5000"—is serially the second compound in the file.) The last number is an algorithmically assigned check character. The format is from two to six numerals followed by a hyphen, followed by two numerals, followed by a hyphen, and concluding with one numeral. As indicated, 50-

```
999-99-9999 = Social Security Number (of one of the patriarchs?)
```

```
1-763-389-8370 = telephone number
60563-3024 = US Postal Zip +4 Code
```

58-08-2 = CAS Registry Number (of caffeine)

Figure 3. Examples of formatted numbers.

```
L1 ANSWER 1 OF 1 REGISTRY COPYRIGHT 1999 ACS
RN 106-98-9 REGISTRY
CN 1-Butene (8CI, 9CI) (CA INDEX NAME)
OTHER NAMES:
CN
    .alpha.-Butene
    .alpha.-Butvlene
CN 1-Butylene
CN Butene-1
CN Ethylethylene
DR 1735-75-7, 54366-07-3, 33004-02-3
MF C4 H8
LC STN Files: AGRICOLA, ANABSTR, APILIT, APILIT2, APIPAT, APIPAT2,
      BEILSTEIN*, BIOBUSINESS, BIOSIS, CA, CAOLD, CAPLUS, CASREACT, CEN,
      CHEMCATS, CHEMINFORMRX, CHEMLIST, CBNB, CHEMSAFE, CIN, CSCHEM,
       CSNB, DETHERM*, DIPPR*, EMBASE, GMELIN*, HODOC*, HSDB*, IFICDB,
       IFIPAT, IFIUDB, MEDLINE, MRCK*, MSDS-OHS, NIOSHTIC, PIRA, PROMT,
      SPECINFO, TOXLINE, TOXLIT, TRCTHERMO*, TULSA, USPATFULL, VTB
        (*File contains numerically searchable property data)
     Other Sources: DSL**, EINECS**, TSCA**
         (**Enter CHEMLIST File for up-to-date regulatory information)
 H3C-CH2-CH ---- CH2
            7746 REFERENCES IN FILE CA (1967 TO DATE)
            114 REFERENCES TO NON-SPECIFIC DERIVATIVES IN FILE CA
            7758 REFERENCES IN FILE CAPLUS (1967 TO DATE)
             11 REFERENCES IN FILE CAOLD (PRIOR TO 1967)
```

Figure 4. STN REG file record for 1-butene.

00-0, formaldehyde, is compound 5000, and 1746-01-6, dioxin, is compound 174,601. Compare the format to other number formats (Figure 3).

Note that the format makes the numbers more readily identifiable even at a glance.

Note the example provided of a CAS Registry File record on STN, e.g., for 1-butene (Figure 4).

The Registry Number is shown, followed by nomenclature—first index name, then synonyms, deleted Registry Numbers (DR, if any), molecular formula, source of registration, files containing this Registry Number (locator field, LC), structure, the existence of abstract references in the CAOLD file (if any), and the current number of references in the CA and CAplus files.

There are a number of reasons for deleted Registry Numbers. One is a different source of registration. Another was an early attempt to provide additional "faceted" Registry Numbers, which would link compounds similar in stereochemistry, etc. If these were encountered earlier, especially in CAS files, they are preserved in the database, and a search for them will retrieve the current compound record in the Registry File.

| CAS Produced | Produced by others |
|--------------|--------------------|
| CA | Beilstein |
| | |
| CAplus | BIOSIS |
| Стриз | B 10313 |
| | |
| CAOLD | HSDB |
| | |
| CASREACT | MEDLINE |
| | |
| CIVI | DDC) IT |
| CIN | PROMT |

Figure 5. Examples of STN files containing CASRN.

Early in the process, a list of chemical names and Registry Numbers appeared, the TSCA Candidate List which was commonly known as the "Pre-TSCA" list. Many of the Registry Numbers therein were deleted. A successor list is the Registry Handbook—Common Names, which is published on microform. Of course, the complete list of Registry Numbers also appears in print as the Registry Handbook. Even in the current Registry System, a number of CASRN exist that have no references in bibliographic or data files. They merely register substances (usually less defined compositions or mixtures) that appear on one or more of the national regulatory lists, like the TSCA Inventory or EINECS.

Not only has the CAS Registry System become the cornerstone of the CAS indexing process, but also CAS Registry Numbers are being used for a number of other purposes. They appear in almost 60 files on the STN system, several of which are produced by CAS. A few examples are shown (Figure 5).

CASRN are used in sales (chemical catalogs, including Aldrich), transportation (including export/import), regulatory agency reporting, and disposal. In fact, they are required in many of these cases, especially the latter.

An example of use of CASRN in other files is shown for the Chemical Abstracts bibliographic files (CA files) from CAS is shown (Figure 6).

At the dawn of the Registry Systems era, CAS Registry numbers appeared in the CA file abstract as well as in the index phrases. However, that practice has been discontinued. Since 1987, author inspired names have been provided for Registry Numbers in CA File index term phrases.

Until recently, role qualifiers (preparation, uses, etc.) were only provided for ca. 1000 commonly cited chemical compounds and some controlled index terms (CT). The preparation role was indicated, as suffix "P", for all Registry Numbers, where appropriate. However, beginning in 1995, an expanded list of roles is now algorithmically applied to all Registry Numbers and CTs in the CA Files. The algorithms are good, but not perfect, as indicated in the example of sulfolane (Figure 6). Note that sulfolane was used in the preparation of filters—it was not prepared itself. Also note the other roles that were applied for sulfolane. The "use" roles are obviously accurate, but I think more information is needed to determine if sulfolane is indeed a reactant in this process rather than just a solvent.

Assignment of CASRN by the CAS Registry System is quite accurate. Less accurate assignments are usually the

```
ANSWER 17 OF 191 HCA COPYRIGHT 1999 ACS
    125:257263 HCA
TI Filters for selective separation of cells or substances from blood
IN Onodera, Hirokazu; Suemitsu, Junsuke
PA Asahi Medical Co, Japan
SO Jpn. Kokai Tokkyo Koho, 6 pp.
PI JP08196627
    Filters for selective sepn. of cells or other substances from blood
     contain arom. ring and/or olefin chain linked to OCH2NR1COCR2(R2)(R2)
     or OCH2NR1COCH(R2)(R3) [ R1 = H or alkyl; R2 = halo; R3 = halogenated
     hydrocarbons]. As an example, polystyrene nonwoven fabrics were
     treated with a mixt, contg. N-hydroxymethyltribromoaceytamide.
     sulfolan and trifluoromethanesulfonic acid and the resultant nonwoven
     fabric was then treated with anti-human CD4 for immobilization. \dots
IT 126-33-0P, Sulfolan 1493-13-6P, Trifluoromethanesulfonic acid
     17354-02-8P 91298-06-5P
     RL: NUU (Nonbiological use, unclassified); PNU (Preparation,
     unclassified); RCT (Reactant); PREP (Preparation); USES (Uses)
        (in prepn. of filters for selective sepn. of cells or other
        substances from blood)
```

Figure 6. STN HCA (Chemical Abstracts) file record.

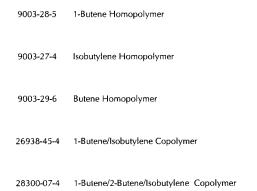


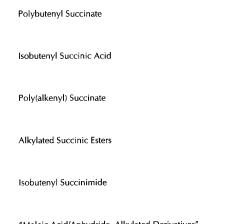
Figure 7. Examples of CASRN of "polybutenes".

result of author supplied information that is insufficiently accurate. This is often the case with patents, which are often written by nonchemists. The author previously described the difficulties that occur with C4 (or higher) compounds.¹¹ Although some misassignments were found to be made by the CA indexers, the infamous "polybutenes" reside in an information swamp (definitely not on the Information Highway). To comprehensively search these oligomers of mixed butenes, at least 40 CASRN must be used, in addition to compound names, the nomenclature for which is even less precise. Several of the CASRN are only partially correct, and some are not accurate at all. Polybutenes are oligomers of mixed butenes, primarily of isobutylene with much smaller amounts of 1-butene and 2-butene. CASRN for some of the "polybutenes" are shown are in increasing order of accuracy (Figure 7).

The first, 1-butene homopolymer, although often used for polybutenes, is not correct and should be reserved for crystalline poly 1-butene.

The situation is even more murky for derivatives of polybutenes, including the even more legendary infamous "PIBSAs" (polyisobutenyl succinates), the majority of which do not have CAS Registry Numbers. Unfortunately, the terminology is even more chaotic—some examples are shown (Figure 8).

Some PIBSAs are indexed like the last example, as derivatives of maleic or succinic acids or anhydrides, both



"Maleic Acid/Anhydride, Alkylated Derivatives"

Figure 8. Examples of PIBSAs.

by text names or by derivative CAS Registry Numbers. A good share of the blame for this indexing nightmare can be directed to sloppy nomenclature and description by the original authors. However, the polybutenes and PIBSAs represent several cases where new, more definitive CASRN should be assigned, rather than attempting to get by with established but inadequate CASRN.

Especially among those dealing with regulatory agency information, a number of myths about Registry Numbers and associated information have sprung up over the years.⁶

Myth 1: CAS Registry Numbers contain information on chemical composition.

Wrong. As stated before, Registry Numbers are serial numbers. The associated chemical information is contained in the file record, including references. Some numbers in other files do contain structural information, e.g., Lawson Numbers in the Beilstein file.

Myth 2: Toxic substances have low-numbered CAS Registry Numbers.

Not necessarily. Compounds that were described often in the literature tended to acquire Registry Numbers very early in the history of the system, including many of those that are toxic. The reverse is also not true: not all "old" compounds are toxic.

Myth 3: Per its title, the TSCA Inventory only lists toxic substances.

By any realistic definition, "toxicity" is both a relative concept, yet highly dependent on data (i.e., "how much", "toxic compared to what?"). Depending on the definition, all compounds are toxic, at least in some concentrations and environments. However, those compounds more commonly perceived as "toxic" are also commonly encountered in commerce and hence are on TSCA or other regulatory lists. Their hazardous potential must be established, if not already done.

Myth 4: All CASRN in the various sources are assigned and used accurately and precisely.

Not completely true. The primary method of assigning CAS Registry Numbers with the CAS Registry System and the CA files has already been described and is accurate to a very high degree. However, other methods are used to assign Registry Numbers to compounds appearing in other files. As previously described, 4-6 even if CAS assigns CASRN for compounds in databases other than its own, these CASRN may not have the same precision as for compound informa-

```
ANSWER 1 OF 1 REGISTRY COPYRIGHT 1999 ACS
    75-05-8 REGISTRY
CN Acetonitrile (8CI, 9CI) (CA INDEX NAME)
    Acetonitrile cluster
    Cvanomethane
    Ethanenitrile
    Ethvl nitrile
CN Methane, cyano-
CN Methanecarbonitrile
    Methyl cyanide
    Methyl cyanide (MeCN)
    3D CONCORD
    54841-72-4
MF
    C2 H3 N
CI
     STN Files: AGRICOLA, AIDSLINE, ANABSTR, APILIT, APILIT2, APIPAT,
       APIPAT2, BEILSTEIN*, BIOBUSINESS, BIOSIS, CA, CANCERLIT, CAOLD,
       CAPLUS, CASREACT, CEN, CHEMCATS, CHEMINFORMRX, CHEMLIST, CBNB,
       CHEMSAFE, CIN. CSCHEM. CSNB. DETHERM*, DDFU, DIPPR*, DRUGU, EMBASE,
       GMELIN*, HODOC*, HSDB*, IFICDB, IFIPAT, IFIUDB, IPA, MEDLINE,
       MRCK*, MSDS-OHS, NAPRALERT, NIOSHTIC, PDLCOM*, PIRA, PROMT, RTECS*,
       SPECINFO, TOXLINE, TOXLIT, TRCTHERMO*, TULSA, ULIDAT, USPATFULL, VTB
         (*File contains numerically searchable property data)
     Other Sources: DSL**, EINECS**, TSCA**
         (**Enter CHEMLIST File for up-to-date regulatory information)
 H 3 C -- C --- N
           21089 REFERENCES IN FILE CA (1967 TO DATE)
            289 REFERENCES TO NON-SPECIFIC DERIVATIVES IN FILE CA
           21141 REFERENCES IN FILE CAPLUS (1967 TO DATE)
```

Figure 9. STN file record for acetonitrile.

tion totally within the control of CAS. Due to chemical "puns" like "ether", the precision may decrease even further for CASRN assigned by algorithm. Information users should always observe *caveat emptor*, but should they expect a "CAS-housekeeping seal of approval"? That is not a bad idea.

10 REFERENCES IN FILE CAOLD (PRIOR TO 1967)

The preceding discussion hopefully made it obvious that CASRN are extremely valuable in searching for chemical substances. Should CASRN be used universally for precise identification of chemical substances, including in all published articles and patents? It is tempting to say yes. A recent submission to the Chemical Information discussion list (CHMINF-L) inquired about the identity of "MeCN". Several members responded with acetonitrile (Figure 9). One responder was reminded of an "urban myth" (probably true) that some shipper refused to accept "methyl cyanide" for shipment but would accept "acetonitrile". Would use of the CASRN have helped in this situation? Possibly, but if CASRN became a utility, who should be expected to administer the process?

It is important to note that there are costs associated with the use of CASRN outside of CAS, especially use in databases produced by other organizations. The details of such agreements are understandably proprietary. In some cases, failure to continue the agreement has led to cessation of inclusion of new CASRN in databases. The use of CASRN is becoming more and more ubiquitous, but it is only fair for CAS to recoup the distinct costs of building and maintaining such a quality system. As stated above, if

CASRN became a public utility, who would be expected to both administer and pay for the maintenance of the system?

One final admonition: when describing chemical substances, use Chemical Abstracts Registry Numbers for definition and accuracy but use them wisely.

REFERENCES AND NOTES

- (1) E.g., Buntrock, R. E.; Wolff, T. E. Information to Order. *Chemtech* **1994**, *4* (April), 8–12.
- (2) Buntrock, R. E. Jargon. DATABASE 1986, (Dec.), 110-112.
- (3) Buntrock, R. E. The Language of Chemistry. DATABASE 1989, (Aug.), 126–127.
- (4) Buntrock, R. E. Chemical Compound Registration Algorithmic or Otherwise. DATABASE 1994, (Feb.), 108–110.
- (5) Buntrock, R. E. Gold (CASRN = 7440-57-5) is Where You Find It, or Caveats on Finding Chemical Substances using CASRN. DATA-BASE 1995, (June/July), 50-55.

- (6) Buntrock, R. E. Playing the Chemical Numbers Game. CAS Registry Numbers (CASRN) Revealed. DATABASE 1996, (Aug./Sept), 72– 77
- (7) E.g., Buntrock, R. E. A Little Learning DATABASE 1990, (Oct.), 123–124.
- (8) Swartzentruber, P. E. Chemical Registry Systems, Abstracts of the ACS National Meeting, CINF 38, Aug. 30, 1984.
- (9) Weisgerber, D. W. Chemical Abstracts Service Chemical Registry System: History, Scope, and Impacts. J. Am. Soc. Inf. Sci. 1997, 48, 349–360.
- (10) Blackwood, J. E.; Blower, P. E.; Layten, P. E.; Lillie, D. H.; Lipkus, A. H.; Peer, J. P.; Qian, C.; Staggenborg, L. M.; Watson, C. E. J. Chem. Inf. Comput. Sci. 1991, 31, 204–212.
- (11) Buntrock, R. E. Documentation and Indexing of C₄ Compounds: Pathways and Pitfalls. *J. Chem. Inf. Comput. Sci.* **1989**, 29, 72–78.
- (12) Albert, A. Xenobiosis; Chapman and Hall: London, 1987; pp 1-3.

CI000109Q