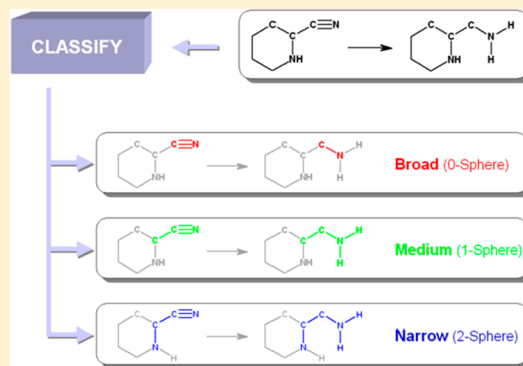


Algorithm for Reaction Classification

Hans Kraut,^{*,†} Josef Eiblmaier,[†] Guenter Grethe,[‡] Peter Löw,[†] Heinz Matuszczyk,[†] and Heinz Saller[†][†]InfoChem GmbH, Landsberger Strasse 408/V, D-81241, Munich, Bavaria, Germany[‡]352 Channing Way, Alameda, California 94502-7409, United States

S Supporting Information

ABSTRACT: Reaction classification has important applications, and many approaches to classification have been applied. Our own algorithm tests all maximum common substructures (MCS) between all reactant and product molecules in order to find an atom mapping containing the minimum chemical distance (MCD). Recent publications have concluded that new MCS algorithms need to be compared with existing methods in a reproducible environment, preferably on a generalized test set, yet the number of test sets available is small, and they are not truly representative of the range of reactions that occur in real reaction databases. We have designed a challenging test set of reactions and are making it publicly available and usable with InfoChem's software or other classification algorithms. We supply a representative set of example reactions, grouped into different levels of difficulty, from a large number of reaction databases that chemists actually encounter in practice, in order to demonstrate the basic requirements for a mapping algorithm to detect the reaction centers in a consistent way. We invite the scientific community to contribute to the future extension and improvement of this data set, to achieve the goal of a common standard.



1. INTRODUCTION

1.1. Background. Reaction classification has many different applications.^{1–3} It can be used in teaching to show the formal similarity of reactions, and it is an efficient method for systematic indexing of reactions in books and databases. Classification complements structure-based retrieval systems, and is a technique that encourages browsing approaches. In reaction retrieval systems, it is a useful tool for postsearch management of large hit lists, simplification of query generation, linking of reaction information from different sources, and access to generic types of information. It helps in deriving knowledge bases for reaction prediction and synthesis design and can be applied in prediction of new reactions. It is useful in setting up automatic procedures for analyses and correlations and in quality control and overlap studies.

Recent reviews^{2,3} divide classification techniques into two categories: “model-driven” methods and “data-driven” methods. In the former, a preconceived model is imposed, whereas in the latter, a computer automatically analyzes a set of reactions and generates a classification.

Manual reaction indexing methods have been based on the name or class of the product, or functional groups, or the reagent used, rather than on a true reaction classification. The use of such methods, or “name reactions” (e.g., Claisen condensation), or reaction mechanisms, is insufficient for a complete classification. Way back in 1938, Weygand⁴ proposed the classification of reactions on the basis of bonds broken or formed. Theilheimer⁵ developed Weygand's system into a classification based on four types of reaction change: addition, elimination, rearrangement, and exchange. Balaban⁶ devised

reaction families based on 6-electron pericycles. Hendrickson⁷ extended Balaban's analysis, and Arens^{8–10} independently derived a related but more general system.

Even before Balaban's work, Vladutz¹¹ recognized that the Theilheimer classification was too broad and lacked subclasses. (The ideal classification system should be hierarchical.) He developed the concept of the reaction center, or reaction site, and he advocated use of all the bond changes occurring during the reaction, as opposed to the single changes considered by Theilheimer.

In yet another model-driven method, Zefirov and colleagues^{12,13} identified first reaction centers and then a reduced system or reaction equation, by removing unchanged bonds. Their SYMBEQ program^{14,15} takes a so-called “formal-logical” approach based on graph theory and generates symbolic equations (SEQs) from topologies of bond changes at the reaction center. Fujita,^{16–18} like Vladutz,¹⁹ created a unitary reaction representation, which he called the imaginary transition state (ITS). Hendrickson^{20–22} built on all these earlier methods to provide a simple, general linear notation for reactions. The principle is implemented in COGNOS²³ using the large ChemReact²⁴ database from InfoChem for fast retrieval. Ugi and Dugundji's “BE-matrices”²⁵ for representing reactions have been used by their co-workers and successors in a number of computer programs,³ one of which was extended by Herges to find a new class of reactions.^{26–28}

Received: July 26, 2013

Model-driven classification methods did not consider functionality beyond the reaction center or the derivation of subclasses based on topology proximal to the center; data-driven methods have addressed this issue. Wilcox and Levinson²⁹ used a computer to derive not only a “minimum reaction concept”, similar to Zefirov’s SEQ, but also a “complete reaction concept”. Blurock^{30,31} derived an “expanded reaction center” for each product and reactant. Later he used classification in the automatic generation of detailed mechanisms in combustion of hydrocarbons.^{32,33} Gelernter and co-workers³⁴ used machine learning in a “conceptual clustering” technique. Sello^{35,36} has developed a hierarchical reaction classification using a property-based method. InfoChem’s reaction classification algorithm^{3,37} is the subject of the current paper. Christ and co-workers³⁸ have recently written what they suggest is a similar algorithm.

In the retrosynthetic analysis package Route Designer (now called ARChem) core extension to *relevant* neighboring functionality is the key feature that distinguishes the method³⁹ from the core extension found in KOSP⁴⁰ and in CLASSIFY,^{3,37} where extension is based entirely on bond distance from the reaction core. On the other hand, as the ARChem rules encode specific structural features they are less generally applicable than “shell-based” approaches such as CLASSIFY.³⁷ Ongoing improvements to ARChem include enhancing the rule extraction engine to improve treatment of regioselectivity, stereoselectivity, and interfering functional groups.

In more recent work, Varnek and his co-workers⁴¹ have addressed the problem of similarity search and classification of chemical reactions using neighborhood behavior and Condensed Graphs of Reaction (CGR) approaches. A CGR (similar to Fujita’s ITS) merges all molecules involved in a reaction into one molecular graph, allowing reactions to be considered as pseudomolecules. Descriptors can be calculated on this graph. Thus chemical processes can be treated like classical compounds and subjected to similarity-based virtual screening approaches. In a follow-up paper the authors use machine learning methods in conjunction with the CGR approach to identify errors in the atom-to-atom mapping of chemical reactions produced by an automated mapping tool by ChemAxon.⁴²

One possible drawback of these data-driven approaches is that they are based only on the topology of the extended reaction center and do not detect reactions that are mechanistically similar but different in topology. (ARChem may be an exception since it attempts to use reaction mechanisms.) Gasteiger’s HORACE system^{43,44} used inductive and resonance effects, and charge distribution, as mechanistic descriptors, and a topology based on Gelernter’s classification method³⁴ to produce a reaction hierarchy. A limitation was the fact that similarity between two reactions in different subclasses could not be measured. Chen and Gasteiger addressed that limitation in further work based on Kohonen neural networks,^{45,46} producing a two-dimensional (as opposed to one-dimensional) classification scheme. A trained Kohonen network can also be used to predict reactions. Chen and Gasteiger’s method^{45,46} has also been used for the analysis of reaction databases.³⁷ Funatsu’s team adopted the Chen-Gasteiger methodology⁴⁷ and also used Kohonen networks in a method⁴⁸ which numerically characterizes the field around molecules based on electrostatic and steric interactions with a pseudoreactant.

1.2. Reaction Centers and Atom Mapping. Molecules represented as graphs⁴⁹ can be compared using graph matching or isomorphism techniques. Graph matching can be formulated as a problem involving the maximum common subgraph (MCS) among the collection of graphs being considered. The maximum common substructure is the largest substructure common to the collection of graphs under consideration. One application of an MCS algorithm is automatic detection of a reaction center and mapping the atoms in the product with respect to those in the reactant(s): see Figure 1.

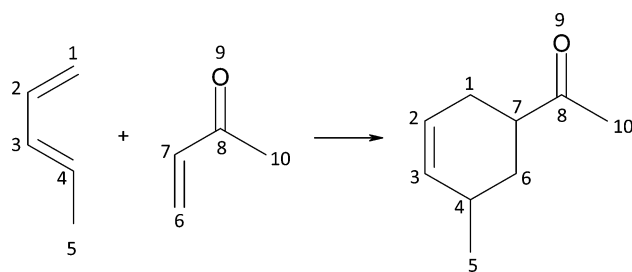


Figure 1. Atoms 1–4, 6, and 7 are in the reaction center; bonds 4–5, 7–8, 8–9, and 8–10 are unchanged.

Automatic reaction mapping and reaction center detection have been ably reviewed very recently by Chen et al.⁵⁰ Traditional reaction mapping algorithms have been based on extended-connectivity or MCS algorithms. Determining the MCS between two or more graphs is a combinatorially intractable, NP-complete problem⁵¹ but an approximate MCS method and a method using approximate reaction sites as input to an exact MCS routine were developed at Sheffield University many years ago.^{52–54} This work, and the work of Funatsu’s team,⁵⁵ was extended, in work described below, enabling the CLASSIFY algorithm to find reaction centers for more complex reactions. Other publications^{56–58} are more tangentially relevant to the current article. The MCS literature has been well reviewed.^{51,59} While an MCS algorithm is probably the most common way of mapping reactions for reaction database applications, other approaches have been used recently, especially for biochemical reactions.^{60–67}

In their review on MCS algorithms, Ehrlich and Rarey⁵⁹ concluded that new MCS algorithms need to be compared with existing methods in a reproducible environment, preferably on a generalized test set, yet the number of test sets available is small. The purpose of this publication is to supply a challenging test set of reactions, usable with InfoChem’s software or other classification algorithms, and to encourage comparisons. We believe that the reaction sets used in most previous publications are not truly representative of the range of reactions that occur in real reaction databases. We start with a historical introduction about the InfoChem CLASSIFY algorithm.

1.3. Motivation. Beginning in the late 1980s InfoChem started to develop a deep understanding of the storage and handling of chemical structure and, in particular, reaction information. The background for this activity was a major project funded by the German Bundesministerium für Forschung und Technologie (BMFT) with a view to the development and implementation of an electronic form of the printed abstracts series *ChemInform*⁶⁸ published by FIZ CHEMIE Berlin. This ChemInform project was only part of a much wider BMFT initiative which in parallel also funded the creation of both the electronic versions of *Beilsteins Handbuch*

der Organischen Chemie (organic structures)⁶⁹ and *Gmelins Handbuch der anorganischen Chemie* (inorganic structures).⁷⁰ The subject of the ChemInform project was the design and development of software modules for the capturing, handling, storage, and management of organic chemical reactions together with associated data. A data model specific to the ChemInform reactions was designed and implemented.

In 1989 InfoChem acquired an exclusive license to a reaction database which was jointly built by the All-Union Institute of Scientific and Technical Information of the Academy of Sciences of the USSR (VINITI) and the German Zentrale Informationsverarbeitung Chemie, Berlin (ZIC) since 1974. In 1989, this database, called "SPRESI", contained 1.8 million reaction records.⁷¹

InfoChem was able to read the data from the 768 magnetic data tapes provided and to convert them into an industry standard format, which theoretically would have enabled InfoChem to distribute the data to customers. Unfortunately, the reaction database management systems commercially available and used at that time (around 1990) were not able to handle more than approximately 100,000 records because of software and hardware limitations. The size of the SPRESI reaction database exceeded that of the databases commercially available at that time by several orders of magnitudes: ORAC (version 6.5) offered 25,000 reactions, SynLib (version 2.2) contained 39,000 reactions, and REACCS (version 6.1) contained 19,000 records from the Current Literature file and 47,000 from the Theilheimer database.⁷² The economic need to commercialize the SPRESI data set, however, forced InfoChem to develop ideas about how to break down the vast amount of data meaningfully and how to define and produce useful subsets. Different approaches, for example extracting only those reactions that lead to heterocyclic compounds, were discussed together with interested customers. Finally, a concept similar to the way sets of chemical reactions are published in the chemical literature was conceived.

Usually, in order to arrange reaction information clearly and to compact it in publications, only the core reaction describing the chemical transformation is graphically displayed in a general form. Variations of moieties or variable side chains are represented by a generic placeholder in the core reaction depiction and are specifically enumerated in a subsequent table along with additional information such as specific yield values (Figure 2).

This concept allows the author to represent a series of different reactions having the same transformation type using just one general representation. Using first the topological information of the reacting center and second, the structural

R1	R2	Conditions	Yield (%)
Me	Ph	aq HCl, EtOH, reflux, 30min	82
Me	t-Bu	aq HCl, EtOH, reflux, 30min	74
4-ClC ₆ H ₄	4-ClC ₆ H ₄	AcOH, EtOH, H ₂ O, 75°C, overnight	69
4-Me(CH ₂) ₇ OC ₆ H ₄	4-O ₂ NC ₆ H ₄	TsOH, EtOH, reflux, 5h	72
4-PhC ₆ H ₄	4-PhC ₆ H ₄	AcOH, 100°C, 24h	74
4-Me ₂ NC ₆ H ₄	4-Me ₂ NC ₆ H ₄	AcOH, 100°C, 24h	84

Figure 2. Generic representation for a series of different, specific reactions describing the same transformation type. Source: Thieme, Science of Synthesis,⁷³ Pyrazines (Update 2011), 16.14.5, p 261.

environment around the reaction centers, InfoChem developed an algorithm which generates a distinct numerical classification code for each different reaction type. As a result, all reactions sharing the same reaction type will have the same calculated classification code and are called "similar reactions". After the reaction class codes are calculated for all records of a reaction data set, one representative reaction can be selected from each set of similar reactions, applying certain criteria (e.g., yield, publication year, type of journal etc.), and aggregated in a newly created data collection. Although this resulting subset of reactions will be significantly smaller than the original one, the chemical essence and diversity will be retained. The application of additional criteria (e.g., number of examples per reaction type, exceeding a certain yield, etc.) finally allowed InfoChem to create a variety of even smaller, high quality reaction collection products.

2. METHODS

2.1. Atom–Atom Mappings and Reaction Centers for CLASSIFY. The initial step for the calculation of reaction class codes is the determination of the reaction centers. A bond is defined as being in a reaction center if it is made or broken. An atom is defined as being in a reaction center if it changes its number of implicit hydrogens, valency, number of π -electrons, or atomic charge, or if at least one connecting bond is in a reaction center (see Figure 3). Bond order changes are

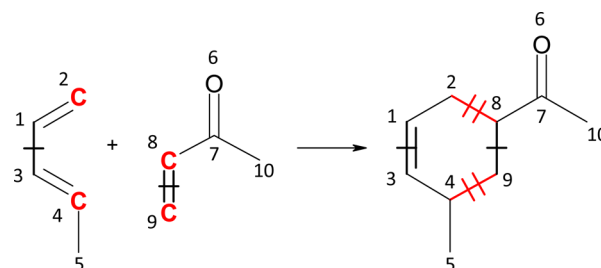


Figure 3. Reaction center definition for Diels–Alder reaction as required by CLASSIFY. Atoms 2, 4, 8, 9 and bonds 2–8 and 4–9 are the reaction center; bonds 8–9 and 1–3 have only a change of bond order and are not in the reaction site.

translated into this representation, marking the atoms of the bond where the number of hydrogens or number of π -electrons changes. This reaction center representation is specific to CLASSIFY. (In all the figures, bonds made or broken are marked with double lines in red; bonds changed are marked with a single line bold.)

The first step and the most difficult challenge for an algorithm is the detection of the atom–atom mappings between reactant and product molecules in a way that represents the chemistry of the reaction equation correctly. The reaction centers can then be derived from the atom–atom mapping information without difficulty. We have developed a program, ICMAP,⁷⁴ and optimized it to deal with the sort of "real-life" reaction data contained in reaction databases abstracted from the literature or patents (which have many unbalanced reaction equations, for example), as opposed to idealized collections of reactions that have often been used to test mapping algorithms.

2.1.1. General Approach of the Reaction Center Recognition Algorithm. The fundamental procedure of the algorithm is to test all maximum common substructures (MCS)

between all reactant and product molecules in order to find an atom mapping containing the minimum chemical distance (MCD). This means that the final solution should map the atoms of the reactants and products in such a way that a minimum number of bonds is broken in the reactants and a minimum number of bonds is made in the product. Some additional chemical rules are applied in the rating of the MCD. Thus breaking and making of bonds containing hetero atoms is preferred to breaking and making of carbon–carbon bonds. Changes in bonds connected to hydrogen atoms are rated in the same way as changes to carbon–carbon bonds are and should be avoided in the solution. Additional information is given in the program documentation.⁷⁵ In most cases, the MCD of a reaction describes the correct reaction site. Even if the reaction mechanism differs from the MCS solution the algorithm groups similar reactions very reliably in one class. This consistent detection of the reacting centers for different reactions is a basic requirement for the calculation of the reaction class codes.

2.1.2. Reaction Categorization Examples and Test Set. Recent publications^{50,59} have suggested that a generalized test set is needed for comparison of the different solutions to the atom–atom mapping problem. For this publication we created a representative set of example reactions from a large number of reaction databases that chemists actually encounter in practice, in order to demonstrate the basic requirements for a mapping algorithm to detect the reaction centers in a consistent way. The test set is grouped into different levels of difficulty which are described in the following sections. The test set (which can be extended in future) is available in Accelrys (formerly MDL) RDfile format⁷⁶ for downloading from http://www.infochem.de/content/downloads/classify_testsets.shtml, together with a description in PDF format. Additionally, InfoChem offers a service to calculate class codes online for single reactions using the current version of the InfoChem atom mapping/CLASSIFY tool; this is available at <http://www.infochem.de/products/software/icmap.shtml> under “test the software”.

Chen et al. presented a first categorization in their paper.⁵⁰ We have set up a more detailed categorization of reactions based on intensive tests and result analysis for different reaction databases in order to optimize CLASSIFY. The test set provided is a manual selection of examples from different commercial products. All these databases (except ChemInform) have been developed and produced by InfoChem, including the mapping and reaction classification as part of the delivery. The results of this reaction data processing step were subject to intensive quality assurance at InfoChem, and we were thus able to categorize the reactions as described. ChemInform, produced by FIZ CHEMIE Berlin, was also analyzed. The sources of the test set examples are listed in Table 1.

Some other examples are included from manual test input during development or from other publications.^{50,59} Since the different groups of reactions are detected automatically by ICMAP, the number of examples may in the future be extended. The different groups of examples required different solutions and extensions of the ICMAP algorithm. The numbers of examples of the groups are shown in Table 2. Note that the total number is not 104, although there are 104 reactions of the test set, because one reaction can occur in more than one group. For example the reaction in Figure 5 is an example for Group 2 “Deceptively Simple Transformations” and Group 3 “Incomplete Reaction Equations with Reduced

Table 1. Sources of Examples in Test Set

database	database vendor	no. of examples from this database in test set
SPRESI ⁷¹	InfoChem	25
Comprehensive Asymmetric Catalysis ⁷⁷	Springer	10
Glycoscience ⁷⁸	Springer	10
Science of Synthesis ⁷³	Thieme	15
Encyclopedia of Reagents for Organic Synthesis (EROS) ⁷⁹	Wiley	15
ChemInform ⁶⁸	FIZ CHEMIE Berlin	20
other examples	none	9

Table 2. Number of Examples of the Different Groups in the Test Set

group	no. of examples
Group 1	19
Group 2	16
Group 3	39
Group 4	14
Group 5	14
Group 6	19
Group 7	28
total	149

Complexity” because of the leaving group (see description below).

2.1.2.1. Group 1: Trivial Reactions. Trivial reaction equations (see, for example, Figure 4) are those where each atom on the reactant side matches to an atom in the product. For this type of reaction the best MCS is the solution. This is the only category where a simple approach for MCS searching can be applied (since every reacting atom can be expected to have a mapped partner in the reaction site on the other side of the equation), but this type of reaction is not very common in real-life reaction databases. The frequency of occurrence of trivial reactions is about 5–10% in the tested reaction databases. We included some examples in the test set that have different numberings of the atoms for the same molecule. An MCS algorithm must find the correct solution independent of the atom numbering (isomorphism) for consistency reasons.

2.1.2.2. Group 2: Deceptively Simple Transformations due to Symmetry or Similarity. Chen et al. propose this problem as a specific complexity class⁵⁰ and we agree. This group of reactions overlaps with the other six groups described in section 2. The MCS is in most cases a set of solutions. An algorithm must be able to find all solutions and test for the best according to the MCD rule. The example given and described by Chen and co-workers⁵⁰ is given in Figure 5. They explain that the correct solution is a Diels–Alder ring closure, and they also present a wrong MCS solution having a worse chemical distance (CD) because of additional hydrogen count changes for the carbon atom pairs marked 15 and 17. Hydrogen count changes are a completely broken σ bond on the reactant side and made on the product side; they cause a worse CD value.

Chen’s example is included in our test set together with an additional example of our own: see Figure 6. This example shows that a too simplistic approach to finding the MCS can miss the correct atom–atom mapping in the case of a simple symmetry. One methyl group of the *tert*-butyl MCS fragment of the reactant is assigned to the reaction bond carbon instead of

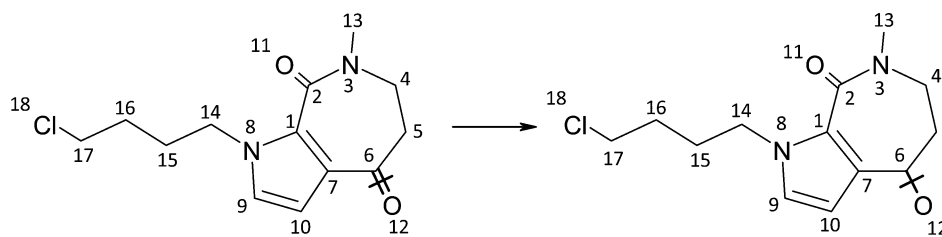
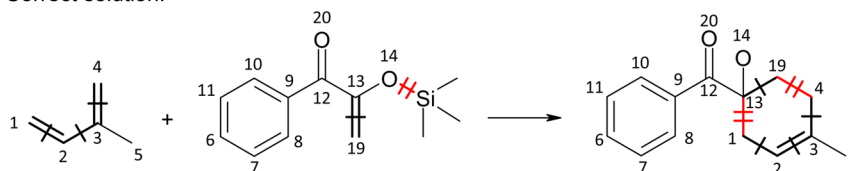


Figure 4. Complete mapping between reactant and product: only one bond order change (reduction) from ketone to alcohol group.

Correct solution:



Incorrect solution:

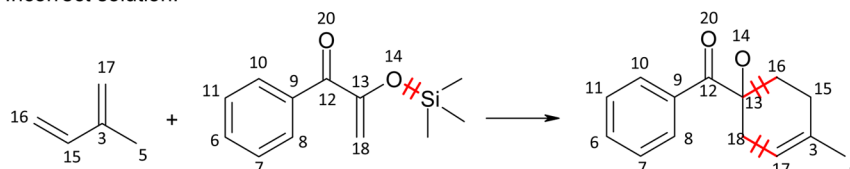
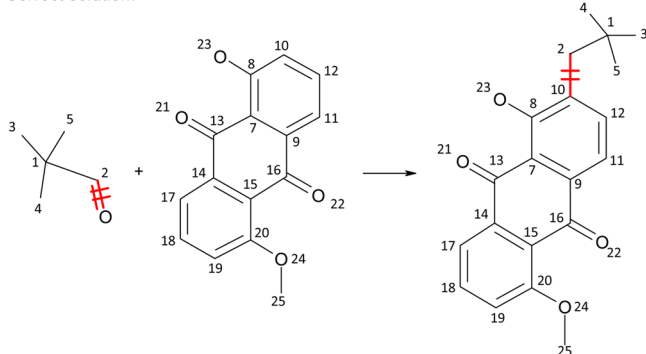


Figure 5. A “deceptively simple” transformation.

Correct solution:



Incorrect solution:

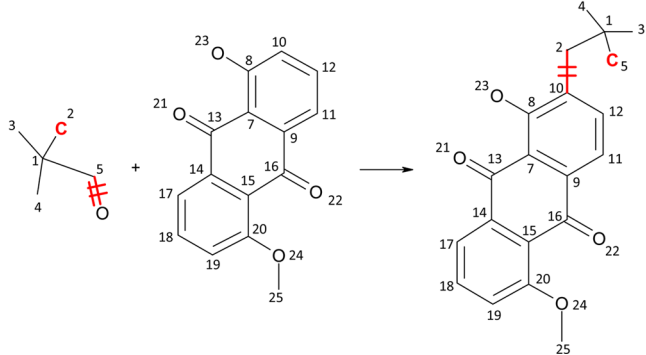


Figure 6. Atom–atom mapping in the case of a simple symmetry.

the matching methyl group on the product side. The additional hydrogen changes due to this mismatch would result in a wrong reaction site and a wrong classification code. The reason for the incorrect solution is the same as in the first example: all MCS solutions are not tested for the best solution, bearing in mind the MCD rule. If an algorithm has this type of wrong result it

tends to be dependent on the atom numbering in the molecules and to produce inconsistent results for a reaction database.

2.1.2.3. Group 3: Incomplete Reaction Equations with Reduced Complexity. Chen et al. describe this category as “Missing pieces from either side”.⁵⁰ We separated this category into two different levels of difficulty since the problem is of different complexity for a program depending on whether small or large pieces are missing, and in particular if reactant copies have to be detected as well.

Group 3 contains reactions where the items missing on either the reactant or the product side are only one or more individual single atoms or small groups of atoms. The information given is sufficient for the MCS algorithm to find the correct solution. Since the molecules’ sigma skeletons are well-defined and unambiguous, the correct reaction mapping and centers can be detected in a straightforward way. We found that about 40–50% of reactions in our databases belong to this category. An example is given in Figure 7. Here there are three different reactants and the reaction equation is incomplete, but it is nearly balanced: one oxygen atom is missing in the product. Implicit hydrogen atoms are not relevant for the algorithm to find the correct solution and thus are not counted.

A second example is given in Figure 8. The $\text{Si}(\text{CH}_3)_3$ group from the second reactant is missing completely on the product side, as is the oxygen of the first reactant.

Although reactions in Group 3 do not present a major challenge for a mapping algorithm, there are reactions in the related, more complex group where there are possible reasons for errors if a mapping algorithm is not able to analyze all possible MCS solutions. Some examples follow (Group 4).

2.1.2.4. Group 4: Complex Incomplete Reactions. In this group the number of unmapped atoms on the reactant or product side is quite high, and this increases the complexity and the requirements for the MCS algorithm.

2.1.2.4.1. Unmapped Groups in Reactants. This category is very common and represents about 30 to 40% of the reactions

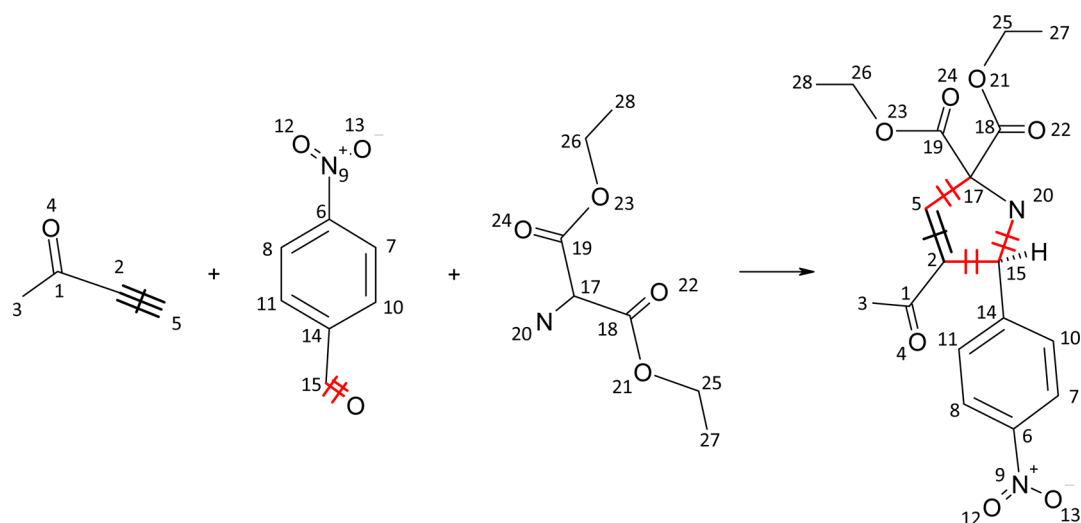


Figure 7. A noncomplex, incomplete reaction equation.

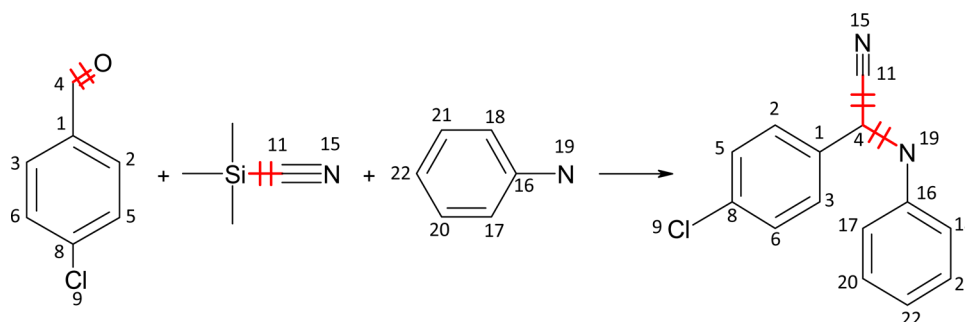


Figure 8. Another noncomplex, incomplete reaction equation.

in the databases we analyzed. Most algorithms can handle this type of reaction if the unmapped group is not too similar to the mapped part of the reactants. An example is shown in Figure 9.

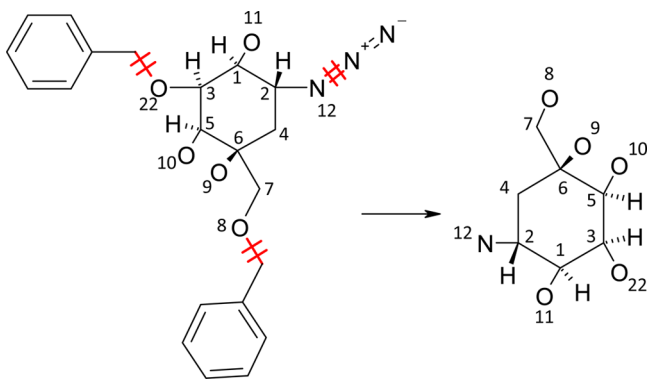


Figure 9. Two unmapped groups of seven atoms each, plus two nitrogen atoms.

2.1.2.4.2. Unmapped Groups in Reaction Products. This group of reactions is usually a big challenge for a reaction center detection algorithm since the solutions are ambiguous and not clearly defined. Reaction databases of high production quality tend to have a smaller percentage of this reaction category, but many of these reactions are correctly abstracted according to the way in which they appear in the literature. The frequency of occurrence also depends on the source documents. Since incomplete reactions frequently appear in the literature, they

have to be addressed by a mapping program. Our test databases contain about 20–30% of reactions that have an incompletely mapped product. The fraction of reactions with missing groups having more than three non-hydrogen atoms unmapped is about 2–3%. If an algorithm is not designed to test reactant copies to be mapped more than once in the product, the percentage will be even higher (*vide infra*). Figure 10 shows a highly incomplete reaction: several reactants are missing. A plausibility function can detect that less than 50% of the product atoms have been mapped to the reactant and thus report a warning for quality assurance.

2.1.2.5. Group 5: Reactions Containing More than One Product. Chen et al. describe this group as “The presence of alternative products”.⁵⁰ This group of reactions contains highly unbalanced reaction equations which require additional recognition rules. The number of reactions belonging to this group varies from database to database. In our complete tests we found 10% to 15% of multiple product reactions. This percentage is dependent on the field of chemistry covered by a database. In Figure 11 the mapping and reaction centers are identical for both products.

2.1.2.6. Group 6: Reactant Copy Symmetric. Chen et al. describe this group as “multiply used reagents”.⁵⁰ This reaction group occurs to the extent of 5 to 10% in the tested reaction databases. A reactant is abstracted only once by the database builder in the reaction equation but is mapped multiple times in the reaction product. This convention in publishing reactions causes highly unbalanced reaction equations which a mapping

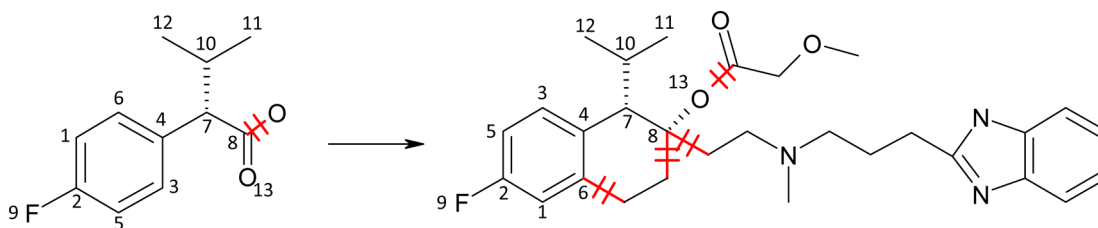


Figure 10. A highly incomplete reaction.

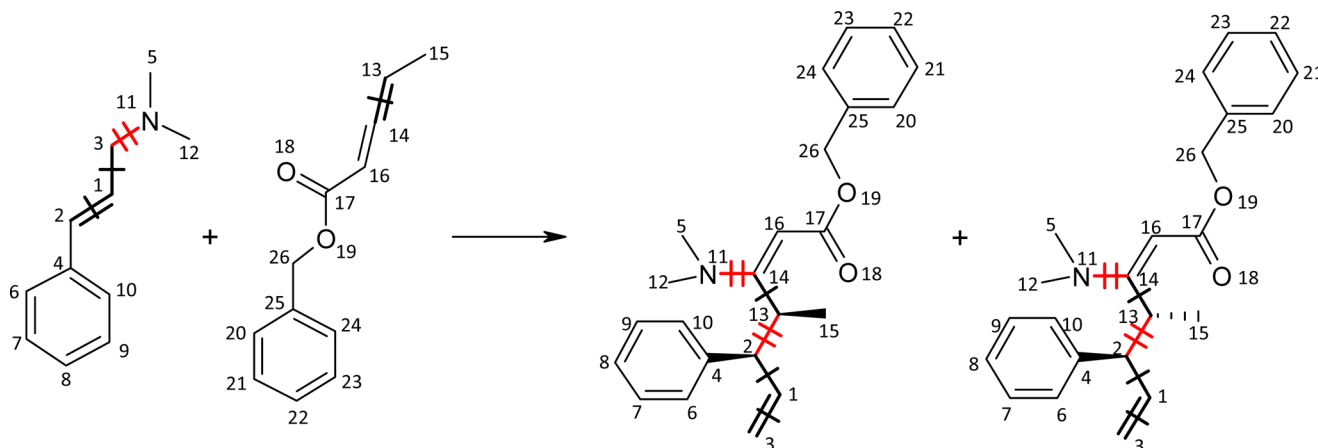


Figure 11. A highly unbalanced reaction equation.

algorithm should be able to handle. An example is given in Figure 12.

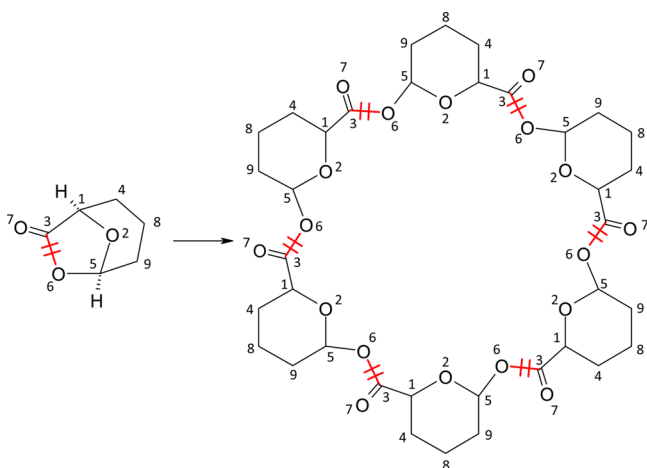


Figure 12. A highly unbalanced reaction equation where a reagent is used multiple times.

If an algorithm does not test for reactant copies, the results are still acceptable since the mapping of one copy also describes the transformation. The class codes would also be consistent, but a problem arises in chemical plausibility checks if reactant copies are not taken into consideration. If the product mapping is extremely incomplete, as it would be in the example in Figure 12 without multiple reactant copies mapped into the product, the check for the percentage of mapped atoms in the product would fail (and this is the most important check). For quality assurance purposes reactions having less than 50% of atoms mapped should be reported to the database builder for manual

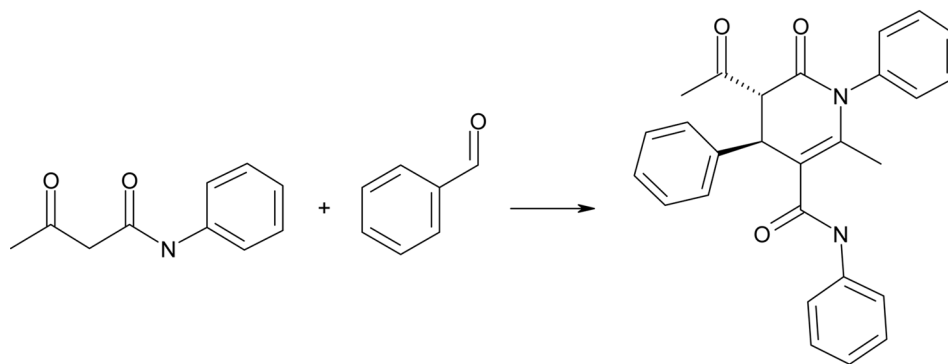
control. Usually, if so many atoms are unmapped, there is an input error, or an important reactant is missing.

2.1.2.7. Group 7: Reactant Copy with Different Reaction Centers. The most complex task for a mapping algorithm is to find copies of reactants reacting differently. The detection of these different reaction centers including the copy of the reactant is essential for the assignment of a correct class code. This type of reaction represents about 1% of the test databases. This is not a high percentage, but these reactions are chemically very interesting examples and well suited as substrates for algorithms designed for testing reactant copies. The first reaction in Figure 13 shows an example as it was stored in the database after abstraction from the literature. This reaction actually includes just two species: benzaldehyde and an anilide. ICMAP introduces a second molecule of the anilide, with the result shown in the second reaction in Figure 13. The anilide is mapped twice into the product but with different reaction centers. For the calculation of the class codes, the reaction centers of both reactant copies are required for this type of reaction. If only one copy of the reactant is included, the different reacting centers will be merged and the correct transformation is lost.

2.2. Reaction Classification. 2.2.1. Overall Class Codes.

Atom hash codes are calculated for all atoms of a reaction center cluster (where a reaction center cluster consists of reaction center atoms connected by reaction center bonds or an isolated reaction center atom) using a modified Morgan Algorithm,⁸⁰ including atom properties for the node values “atom type”, “valence state”, “total number of bonded hydrogen atoms (implicit plus explicit)”, “number of π -electrons”, “aromaticity”, and “formal charges” (see reaction center definition in section 2.1). The sum of all atomic reaction center hash codes of all reactants and one product of a reaction provides the unique overall reaction hash code, called the

Reaction stored in database:



Mapping from ICMAP, including the reactant copy required for correct reaction classification:

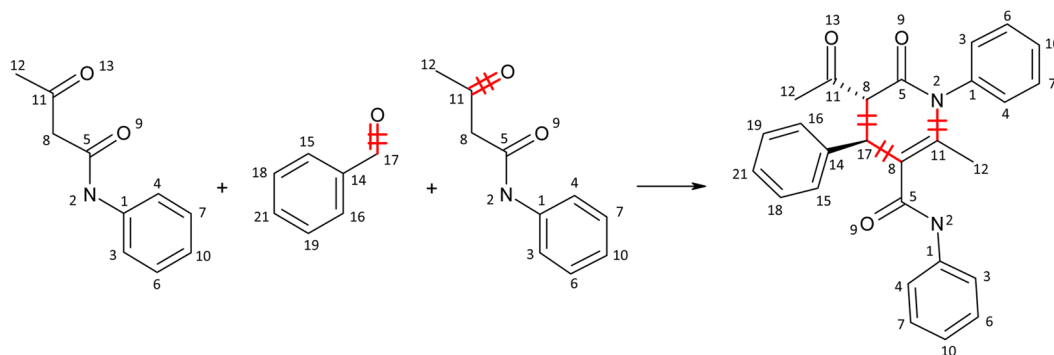


Figure 13. Two identical reactants react differently.

InfoChem Reaction Classification Code, a 15 digit numeric code.

To reflect the chemical environment of the reaction centers, the reaction center clusters may be expanded by adding bonds from reaction center atoms to atoms outside the cluster that are not in reaction centers. A broad similarity search includes reaction centers alone. Atoms in the immediate environment of the reaction center (spheres) may be included for a medium or narrow search (see Figure 14). Using reaction centers alone will

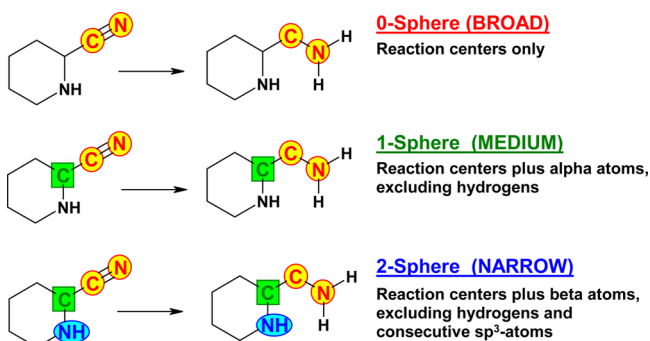


Figure 14. Inclusion of atoms in the immediate environment (spheres) of a reaction center.

give a large-sized cluster or hit list; reaction centers plus alpha atoms, excluding hydrogen atoms, will give a medium-sized cluster or hit list; and reaction centers plus beta atoms, excluding consecutive sp³ atoms, will give a small-sized cluster or hit list. Three hash coded numbers are thus assigned to each reaction (see Table 3).

Table 3. Values of Overall Class Codes for Reaction in Figure 14

type	value
BROAD	279121007655778
MEDIUM	310218305151297
NARROW	316997033077692

2.2.1.1. Generalizations. For all spheres, generalized atom types are used for the reaction classification algorithm. The following substitutions apply:

- Li, Na, K, Rb, Cs, Fr → “alkali metal atom”
- Be, Mg, Ca, Sr, Ba, Ra → “alkaline earth metal atom”
- He, Ne, Ar, Kr, Xe, Rn → “noble gas atom”
- all transition group elements → “transition metal atom”
- B to F atom types are treated as is
- all higher homologues of main groups III to VII (new IUPAC 13 to 17) are treated as “group atoms” of the respective group.

All atoms, excluding hydrogen atoms, of the first sphere around reaction center atoms are used to form the “atom clusters” of the MEDIUM level. When building the “atom clusters” for the NARROW level by adding the neighbor atoms of the second sphere to the clusters, hydrogen atoms and regular tetravalent carbon atoms (sp³-atoms) are ignored. All hetero atoms, and carbon atoms with multiple bonds or in aromatic rings, are considered to be different from regular atoms.

All the resulting atom clusters of the reactant site and one product of the product site are used to generate the *overall class codes* for the reaction type, and multiple occurrences of atom clusters are treated as only one occurrence. This means that

reactions showing several identical chemical transformations will be classified as equivalent to the reaction showing the chemical transformation only once. If a reaction contains two or more product molecules, one set of overall class codes is generated for each product molecule; the overall class codes refer to the transformation leading to this particular product.

2.2.2. Discrete Class Codes. For the calculation of *discrete class codes* the same rules apply as for calculating overall class codes, except that the resulting reaction center clusters will not be combined to form one class code for a product of a reaction, but will be treated individually, to give a set of class codes. In the example in Figure 15, a carbonyl group and a nitro group

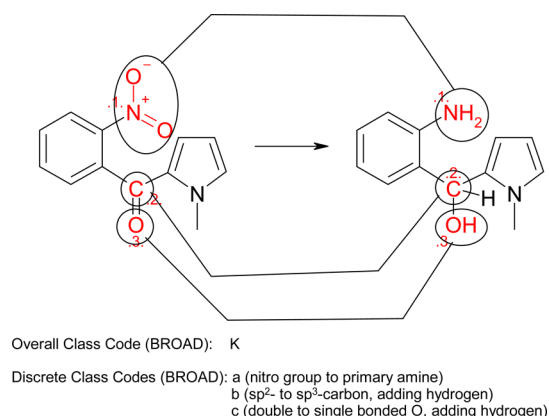


Figure 15. Reaction with two individual transformations.

are reduced simultaneously. The resulting overall class code (K) represents the combination of both transformations, by combining all reaction centers. The resulting discrete class codes (a, b, c) each represent only one reaction center cluster, consisting of one reaction center at the product site and the corresponding reaction center(s) at the reactant site (see Table 4). The correspondence is derived from the atom–atom mappings.

Table 4. Values of Overall and Discrete Class Codes for Reaction in Figure 15

type	value
overall BROAD	311052534254677 (K)
overall MEDIUM	325926056603205 (K')
overall NARROW	389072727557711 (K'')
discrete BROAD	237315105735000 (a) 242815105630000 (b) 221815105840000 (c)
discrete MEDIUM	293748155277334 (a') 316046205034775 (b') 234388286213123 (c')
discrete NARROW	325260571815478 (a'') 366741408247168 (b'') 316046181712831 (c'')

3. RESULTS AND DISCUSSION

3.1. Practical Applications. The CLASSIFY algorithm has been used in a number of applications such as searching for similar reactions (those having the same transformation as a given one), postprocessing (clustering) of reaction search

results, linking databases, database filtering, reaction type search, and analyzing the diversity of databases.

CLASSIFY is the backbone of the “Similar Reactions” feature in the Web version of SciFinder,⁸¹ and it was used until recently in reaction similarity searching in Reaxys.⁸² (Reaxys now appears to have related software that uses five rather than three spheres.) CLASSIFY is used in the functionality “transformation search”, and in postsearch hit list clustering, in Isentris 3.2.⁸³ In a project now closed, DiscoveryGate reaction databases⁸⁴ and InfoChem’s Integrated Major Reference Works⁸⁵ were linked using reaction types.

CLASSIFY has been used to produce subsets⁸⁶ of the SPRESI⁷¹ database. In 1991 a database with 400,000 reaction types (ChemReact)²⁴ was derived from the 2.3 million SPRESI reactions available at that time. InfoChem has also produced smaller subsets. ChemSynth, for example, contains over 104,000 chemical reactions, obtained as a subset of the ChemReact file. The reactions selected for ChemSynth were picked choosing only those reaction types of ChemReact (i) that have at least two example reactions and (ii) whose yield is greater than 50%, and (iii) have been mentioned in leading journals in the field of organic chemistry more than once. The 68,000 reaction types contained in ChemReact68 have been derived from ChemSynth by selecting only those with at least five examples for each reaction type.

Reaction type search (RTS) is implemented in InfoChem’s products⁸⁶ ICCARTRIDGE/ICFSE and SPRESI^{web}. (FSE stands for “fast search engine”.) ICCARTRIDGE is a software module designed to integrate chemical structure and reaction retrieval into the relational database system Oracle. It uses ICFSE to handle and search millions of structures and reactions. InfoChem is currently working on RTS, based on CLASSIFY class codes, as part of the “all in one reaction search” (RSA). RSA is intended as a Google-like search for users who do not know in advance what kind of results a certain query might yield and do not know whether to obtain those hits with an exact, a substructure, or a reaction type search. When choosing RSA the user can just draw any reaction and the system will in turn conduct an exact reaction search (XRS), an exact reaction search for reactions possibly having additional reactants or products (XRSSUB), a combined reaction type and reaction substructure search (RSS_RT MEDIUM), and a reaction type search (MEDIUM). Moreover, the results are also sorted by these criteria (which means that exact reaction search hits come first, then XRSSUB hits, etc.). The underlying idea is that even in the case where the more precise search does not find any reactions, perhaps the next, fuzzier one will.

In a scientifically and commercially vital application, CLASSIFY is used to assess the chemical diversity of a reaction database. Analyzing the number of different discrete class codes and the count of example reactions for each single class code gives the only independent, reproducible, and provable evidence of the nature and the coverage of organic chemistry included in a given reaction collection. This information can be used to judge the quality and the shortcomings of different reaction databases. As a significant example the diversity analysis of the 4.1 million reactions in SPRESI Version 2.8, using “narrow” similarity search,⁸⁷ is shown in Figure 16. The interesting regions of such a graph are the peaks indicating many reaction examples belonging to one particular type, and the singletons, which probably indicate erroneous manual data entries, inadequate or incorrect mappings (for example in a complex rearrangement), or very unusual but valid chemistry.

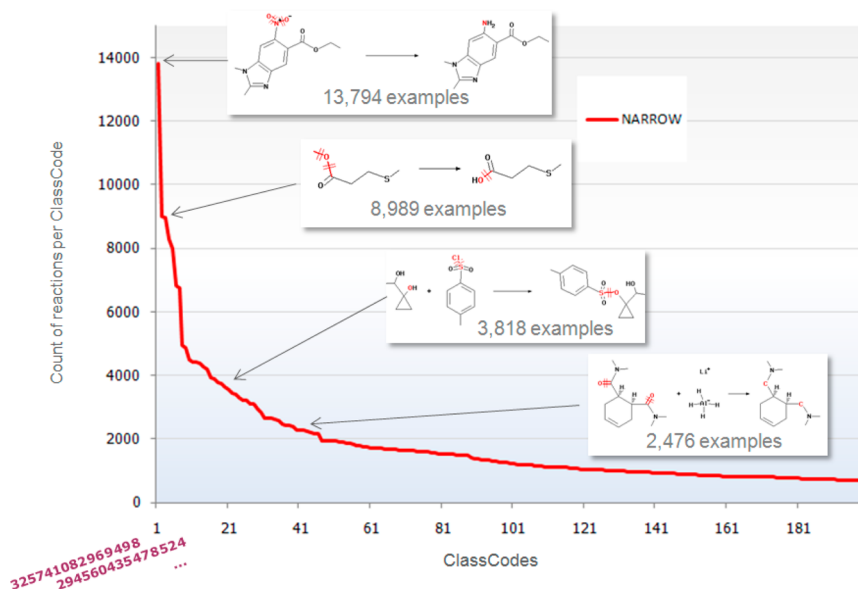


Figure 16. Diversity analysis of a reaction database (SPRESI v. 2.8).

A highly interesting conclusion of this analysis (Figure 16) is that out of the total number of 4.1 million reactions about 348,927 reactions (8.5%) are represented by just 180 different reaction types (0.015% of all different reaction types), while the full database can be reduced to 1,206,773 different reaction types. This strongly indicates that the core and the fundamentals of organic chemistry are covered by just a relatively small number of different chemical transformations.

Such an analysis is useful for quality and plausibility checking when building reaction databases. In the future, in processes which automatically extract reactions from full text or from images, this sort of analysis can be used for the automatic verification and validation of algorithmically generated reaction data.

4. CONCLUSION

Since the InfoChem CLASSIFY tool is still of high relevance, and is used in many different applications, we still maintain the software. We focus for this purpose on the mapping module as the core functionality providing the basic information for the reaction class codes. We see our test set as one of the most important steps in the future for quality assurance, and for definition of a standard for solutions dealing with reaction centers and atom mapping, but we recognize that it is an initial data set that would benefit from extension. We want to contribute to the future extension and improvement of this data set and invite suggestions and additional information from the scientific community to achieve the goal of a common standard.

■ ASSOCIATED CONTENT

Supporting Information

RDfiles for the test set and a description of the test set in PDF format. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*E-mail: hk@infochem.de.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

We are grateful to Daniel Lowe of NextMove Software for reading a draft of this paper, examining the test set, and making helpful suggestions. The authors were funded by InfoChem GmbH.

■ REFERENCES

- (1) Bawden, D. Classification of chemical reactions: potential, possibilities and continuing relevance. *J. Chem. Inf. Comput. Sci.* **1991**, *31* (2), 212–216.
- (2) Hendrickson, J. B.; Chen, L. Reaction classification. In *Encyclopedia of Computational Chemistry*; Schleyer, P. v. R., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer, H. F., III, Schreiner, P. R., Eds.; John Wiley & Sons: Chichester, UK, 1998; Vol. 4, pp 2381–2402.
- (3) Chen, L. Reaction classification and knowledge acquisition. In *Handbook of Chemoinformatics*; Gasteiger, J., Ed.; Wiley-VCH Verlag: Weinheim, Germany, 2003; pp 348–390.
- (4) Weygand, C. *Organische-chemische Experimentierkunst*; Barth: Leipzig, Germany, 1938.
- (5) *Theilheimer's Synthetic Methods of Organic Chemistry*; Tozer-Hotchkiss, G., Ed.; Karger: Basel, Switzerland.
- (6) Balaban, A. T. Chemical graphs. III. Reactions with cyclic six-membered transition states. *Rev. Roum. Chim.* **1967**, *12* (7), 875–98.
- (7) Hendrickson, J. B. Multiplicity of thermal pericyclic reactions. *Angew. Chem.* **1974**, *86* (2), 71–100.
- (8) Arens, J. F. A formalism for the classification and design of organic reactions. I. The class of (+)n reactions. *Recl. Trav. Chim. Pays-Bas* **1979**, *98* (4), 155–161.
- (9) Arens, J. F. A formalism for the classification and design of organic reactions. II. The classes of (+)n+ and (+)n- reactions. *Recl. Trav. Chim. Pays-Bas* **1979**, *98* (6), 395–399.
- (10) Arens, J. F. A formalism for the classification and design of organic reactions. III. The class of (+)nC reactions. *Recl. Trav. Chim. Pays-Bas* **1979**, *98* (9), 471–483.
- (11) Vladutz, G. Concerning one system of classification and codification of organic reactions. *Inf. Storage Retr.* **1963**, *1*, 117–146.
- (12) Zefirov, N. S.; Trach, S. S. Systematization of tautomeric processes and formal-logical approach to the search for new topological and reaction types of tautomerism. *Chem. Scr.* **1980**, *15* (1), 4–12.
- (13) Zefirov, N. S. An approach to systematization and design of organic reactions. *Acc. Chem. Res.* **1987**, *20* (7), 237–243.

- (14) Zefirov, N. S.; Trach, S. S. Symbolic equations and their applications to reaction design. *Anal. Chim. Acta* **1990**, 235 (1), 115–134.
- (15) Zefirov, N. S.; Baskin, I. I.; Palyulin, V. A. SYMBEQ program and its application in computer-assisted reaction design. *J. Chem. Inf. Comput. Sci.* **1994**, 34 (4), 994–999.
- (16) Fujita, S. Description of organic reactions based on imaginary transition structures. 1. Introduction of new concepts. *J. Chem. Inf. Comput. Sci.* **1986**, 26 (4), 205–212.
- (17) Fujita, S. Canonical numbering and coding of imaginary transition structures. A novel approach to the linear coding of individual organic reactions. *J. Chem. Inf. Comput. Sci.* **1988**, 28 (3), 128–137.
- (18) Fujita, S. Canonical numbering and coding of reaction center graphs and reduced reaction center graphs abstracted from imaginary transition structures. A novel approach to the linear coding of reaction types. *J. Chem. Inf. Comput. Sci.* **1988**, 28 (3), 137–142.
- (19) Vladutz, G. Do we still need a classification of reactions? In *Modern Approaches to Chemical Reaction Searching*; Willett, P., Ed.; Gower: Aldershot, UK, 1986; pp 202–220.
- (20) Hendrickson, J. B. Descriptions of reactions: their logic and applications. *Recl. Trav. Chim. Pays-Bas* **1992**, 111 (4), 324–335.
- (21) Hendrickson, J. B. Comprehensive system for classification and nomenclature of organic reactions. *J. Chem. Inf. Comput. Sci.* **1997**, 37 (5), 852–860.
- (22) Hendrickson, J. B. Systematic signatures for organic reactions. *J. Chem. Inf. Model.* **2010**, 50 (8), 1319–1329.
- (23) Hendrickson, J. B.; Sander, T. COGNOS: a Beilstein-type system for organizing organic reactions. *J. Chem. Inf. Comput. Sci.* **1995**, 35 (2), 251–260.
- (24) ChemReact; InfoChem: Munich, Germany, 2013. <http://www.infochem.de/products/databases/chemreact41.shtml> (accessed June 7, 2013).
- (25) Dugundji, J.; Ugi, I. Algebraic model of constitutional chemistry as a basis for chemical computer programs. *Fortschr. Chem. Forsch.* **1973**, 39, 19–64.
- (26) Herges, R.; Hoock, C. Reaction planning: computer-aided discovery of a novel elimination reaction. *Science (Washington, DC, U. S.)* **1992**, 255 (5045), 711–713.
- (27) Herges, R. Ordering principle of complex reactions and theory of contracted transition states. *Angew. Chem., Int. Ed. Engl.* **1994**, 33 (3), 255–276.
- (28) Herges, R. Coarctate transition states: the discovery of a reaction principle. *J. Chem. Inf. Comput. Sci.* **1994**, 34 (1), 91–102.
- (29) Wilcox, C. S.; Levinson, R. A. A self-organized knowledge base for recall, design, and discovery in organic chemistry. In *Artificial Intelligence Applications in Chemistry*; American Chemical Society: Washington, DC, 1986; Vol. 306, pp 209–230.
- (30) Blurock, E. S. Computer-aided synthesis design at RISC-Linz: automatic extraction and use of reaction classes. *J. Chem. Inf. Comput. Sci.* **1990**, 30 (4), 505–10.
- (31) Blurock, E. S. Reaction: system for modeling chemical reactions. *J. Chem. Inf. Comput. Sci.* **1995**, 35 (3), 607–616.
- (32) Blurock, E. S. Detailed mechanism generation. 1. Generalized reactive properties as reaction class substructures. *J. Chem. Inf. Comput. Sci.* **2004**, 44 (4), 1336–1347.
- (33) Blurock, E. S. Detailed mechanism generation. 2. Aldehydes, ketones, and olefins. *J. Chem. Inf. Comput. Sci.* **2004**, 44 (4), 1348–1357.
- (34) Gelernter, H.; Rose, J. R.; Chen, C. Building and refining a knowledge base for synthetic organic chemistry via the methodology of inductive and deductive machine learning. *J. Chem. Inf. Comput. Sci.* **1990**, 30 (4), 492–504.
- (35) Sello, G.; Termini, M. Classification of organic reactions using similarity. *Tetrahedron* **1997**, 53 (41), 14085–14106.
- (36) Sello, G. Reaction classification by similarity: the influence of steric congestion. *Tetrahedron* **1998**, 54 (21), 5731–5744.
- (37) Grethe, G. Analysis of reaction information. In *Handbook of Chemoinformatics*; Wiley-VCH Verlag: Weinheim, Germany, 2003; pp 1407–1427.
- (38) Christ, C. D.; Zentgraf, M.; Kriegl, J. M. Mining electronic laboratory notebooks: analysis, retrosynthesis, and reaction based enumeration. *J. Chem. Inf. Model.* **2012**, 52 (7), 1745–1756.
- (39) Law, J.; Zsoldos, Z.; Simon, A.; Reid, D.; Liu, Y.; Khew, S. Y.; Johnson, A. P.; Major, S.; Wade, R. A.; Ando, H. Y. Route Designer: a retrosynthetic analysis tool utilizing automated retrosynthetic rule generation. *J. Chem. Inf. Model.* **2009**, 49 (3), 593–602.
- (40) Satoh, K.; Funatsu, K. A novel approach to retrosynthetic analysis using knowledge bases derived from reaction databases. *J. Chem. Inf. Comput. Sci.* **1999**, 39 (2), 316–325.
- (41) de Luca, A.; Horvath, D.; Marcou, G.; Solov'ev, V.; Varnek, A. Mining chemical reactions using neighborhood behavior and condensed graphs of reactions approaches. *J. Chem. Inf. Model.* **2012**, 52 (9), 2325–2338.
- (42) Muller, C.; Marcou, G.; Horvath, D.; Aires-de-Sousa, J.; Varnek, A. Models for identification of erroneous atom-to-atom mapping of reactions performed by automated algorithms. *J. Chem. Inf. Model.* **2012**, 52 (12), 3116–3122.
- (43) Rose, J. R.; Gasteiger, J. HORACE: an automatic system for the hierarchical classification of chemical reactions. *J. Chem. Inf. Comput. Sci.* **1994**, 34 (1), 74–90.
- (44) Chen, L.; Gasteiger, J.; Rose, J. R. Automatic extraction of chemical knowledge from organic reaction data: addition of carbon-hydrogen bonds to carbon-carbon double bonds. *J. Org. Chem.* **1995**, 60 (24), 8002–8014.
- (45) Chen, L.; Gasteiger, J. Organic reactions classified by neural networks: Michael additions, Friedel-Crafts alkylations by alkenes, and related reactions. *Angew. Chem., Int. Ed. Engl.* **1996**, 35 (7), 763–765.
- (46) Chen, L.; Gasteiger, J. Knowledge discovery in reaction databases: landscaping organic reactions by a self-organizing neural network. *J. Am. Chem. Soc.* **1997**, 119 (17), 4033–4042.
- (47) Satoh, H.; Sacher, O.; Nakata, T.; Chen, L.; Gasteiger, J.; Funatsu, K. Classification of organic reactions: similarity of reactions based on changes in the electronic features of oxygen atoms at the reaction sites. *J. Chem. Inf. Comput. Sci.* **1998**, 38 (2), 210–219.
- (48) Satoh, H.; Itono, S.; Funatsu, K.; Takano, K.; Nakata, T. A novel method for characterization of three-dimensional reaction fields based on electrostatic and steric interactions toward the goal of quantitative analysis and understanding of organic reactions. *J. Chem. Inf. Comput. Sci.* **1999**, 39 (4), 671–678.
- (49) Warr, W. A. Representation of chemical structures. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **2011**, 1 (4), 557–579.
- (50) Chen, W. L.; Chen, D. Z.; Taylor, K. T. Automatic reaction mapping and reaction center detection. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **2013**, in press.
- (51) Raymond, J. W.; Willett, P. Maximum common subgraph isomorphism algorithms for the matching of chemical structures. *J. Comput.-Aided Mol. Des.* **2002**, 16 (7), 521–533.
- (52) Lynch, M. F.; Willett, P. The automatic detection of chemical reaction sites. *J. Chem. Inf. Comput. Sci.* **1978**, 18 (3), 154–159.
- (53) Willett, P. The evaluation of an automatically indexed, machine-readable chemical reactions file. *J. Chem. Inf. Comput. Sci.* **1980**, 20 (2), 93–96.
- (54) McGregor, J. J.; Willett, P. Use of a maximum common subgraph algorithm in the automatic identification of ostensible bond changes occurring in chemical reactions. *J. Chem. Inf. Comput. Sci.* **1981**, 21 (3), 137–140.
- (55) Funatsu, K.; Endo, T.; Kotera, N.; Sasaki, S. Automatic recognition of reaction site in organic chemical reactions. *Tetrahedron Comput. Methodol.* **1988**, 1 (1), 53–69.
- (56) Stahl, M.; Mauser, H. Database clustering with a combination of fingerprint and maximum common substructure methods. *J. Chem. Inf. Model.* **2005**, 45 (3), 542–548.
- (57) Raymond, J. W.; Kibbey, C. E. An automated method for exploring targeted substructural diversity within sets of chemical structures. *J. Chem. Inf. Model.* **2005**, 45 (5), 1195–1204.

- (58) Gardiner, E. J.; Gillet, V. J.; Willett, P.; Cosgrove, D. A. Representing clusters using a maximum common edge substructure algorithm applied to reduced graphs and molecular graphs. *J. Chem. Inf. Model.* **2007**, *47* (2), 354–366.
- (59) Ehrlich, H.-C.; Rarey, M. Maximum common subgraph isomorphism algorithms and their applications in molecular science: a review. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **2011**, *1* (1), 68–79.
- (60) Reitz, M.; Sacher, O.; Tarkhov, A.; Truembach, D.; Gasteiger, J. Enabling the exploration of biochemical pathways. *Org. Biomol. Chem.* **2004**, *2* (22), 3226–3237.
- (61) Koerner, R.; Apostolakis, J. Automatic determination of reaction mappings and reaction center information. 1. The imaginary transition state energy approach. *J. Chem. Inf. Model.* **2008**, *48* (6), 1181–1189.
- (62) Apostolakis, J.; Sacher, O.; Koerner, R.; Gasteiger, J. Automatic determination of reaction mappings and reaction center information. 2. Validation on a biochemical reaction database. *J. Chem. Inf. Model.* **2008**, *48* (6), 1190–1198.
- (63) Blum, T.; Kohlbacher, O. MetaRoute: fast search for relevant metabolic routes for interactive network navigation and visualization. *Bioinformatics* **2008**, *24* (18), 2108–2109.
- (64) Blum, T.; Kohlbacher, O. Using atom mapping rules for an improved detection of relevant routes in weighted metabolic networks. *J. Comput. Biol.* **2008**, *15* (6), 565–576.
- (65) Heinonen, M.; Lappalainen, S.; Mielikainen, T.; Rousu, J. Computing atom mappings for biochemical reactions without subgraph isomorphism. *J. Comput. Biol.* **2011**, *18* (1), 43–58.
- (66) First, E. L.; Gounaris, C. E.; Floudas, C. A. Stereochemically consistent reaction mapping and identification of multiple reaction mechanisms through integer linear optimization. *J. Chem. Inf. Model.* **2012**, *52* (1), 84–92.
- (67) Latendresse, M.; Malerich, J. P.; Travers, M.; Karp, P. D. Accurate atom-mapping computation for biochemical reactions. *J. Chem. Inf. Model.* **2012**, *52* (11), 2970–2982.
- (68) ChemInform; Wiley-VCH: Weinheim, Germany, 2013. <http://www.fiz-chemie.de/cheminform/> (accessed June 7, 2013).
- (69) Lawson, A. J. The Beilstein database. In *Handbook of Chemoinformatics*; Wiley-VCH Verlag: Weinheim, Germany, 2003; pp 608–628.
- (70) Vogt, J.; Vogt, N.; Schunk, A. Databases in inorganic chemistry. In *Handbook of Chemoinformatics*; Wiley-VCH Verlag: Weinheim, Germany, 2003; pp 629–643.
- (71) SPRESI; InfoChem: Munich, Germany. <http://spresi.de> (accessed June 7, 2013).
- (72) Borkent, J. H.; Oukes, F.; Noordik, J. H. Chemical reaction searching compared in REACCS, SYNLIB, and ORAC. *J. Chem. Inf. Comput. Sci.* **1988**, *28* (3), 148–150.
- (73) *Science of Synthesis*; Thieme: Stuttgart, Germany, 2013. <http://www.thieme-chemistry.com/en/products/reference-works/science-of-synthesis/format/electronic-edition.html> (accessed June 7, 2013).
- (74) IC_{MAP}. <http://www.infochem.de/products/software/icmap.shtml> (accessed June 7, 2013).
- (75) IC_{MAP} documentation. <http://infochem.de/downloads/documentations.shtml> (accessed June 10, 2013).
- (76) Dalby, A.; Nourse, J. G.; Hounshell, W. D.; Gushurst, A. K. I.; Grier, D. L.; Leland, B. A.; Laufer, J. Description of several chemical structure file formats used by computer programs developed at Molecular Design Limited. *J. Chem. Inf. Comput. Sci.* **1992**, *32* (3), 244–255.
- (77) *Comprehensive Asymmetric Catalysis*; Jacobsen, E. N., Pfaltz, A., Yamamoto, H., Eds.; Springer: Berlin, Germany, 1999.
- (78) *Glycoscience. Chemistry and Chemical Biology*; Fraser-Reid, B. O., Tatsuta, K., Thiem, J., Eds.; Springer: Berlin, Germany, 2008.
- (79) *e-EROS Encyclopedia of Reagents for Organic Synthesis*; John Wiley & Sons: Hoboken, NJ, 2013. <http://onlinelibrary.wiley.com/book/10.1002/047084289X> (accessed June 7, 2013).
- (80) Morgan, H. L. The generation of a unique machine description for chemical structures - a technique developed at Chemical Abstracts Service. *J. Chem. Doc.* **1965**, *5* (2), 107–113.
- (81) SciFinder; Chemical Abstracts Service: Columbus, OH. <http://www.cas.org/products/scifinder> (accessed June 7, 2013).
- (82) Reaxys; Elsevier Information Systems: Frankfurt, Germany. <https://www.reaxys.com/info> (accessed June 7, 2013).
- (83) Isentris; Accelrys: San Diego, CA. <http://accelrys.com/products/informatics/decision-support/isentris.html> (accessed June 7, 2013).
- (84) DiscoveryGate; Accelrys: San Diego, CA. <http://accelrys.com/products/databases/database-access/discovery-gate.html> (accessed June 7, 2013).
- (85) *Integrated Major Reference Works 2006*; InfoChem: Munich, Germany. <http://www.infochem.de/news/releasedisplay.shtml?release0601.shtml> (accessed June 7, 2013).
- (86) InfoChem software and databases. <http://www.infochem.de/products/index.shtml> (accessed June 7, 2013).
- (87) Eigner-Pitto, V.; Kraut, H.; Saller, H.; Matuszczyk, H.; Loew, P.; Grethe, G. Reaction classification, an enduring success story. In *Abstracts of Papers, 241st ACS National Meeting & Exposition, Anaheim, CA, United States, March 27–31, 2011*; American Chemical Society: Washington, DC; CINF-2.