

# Over 20 Years of Reaction Access Systems from MDL: A Novel Reaction Substructure Search Algorithm

Lingran Chen,\* James G. Nourse, Bradley D. Christie, Burton A. Leland, and David L. Grier

MDL Information Systems, Inc., 14600 Catalina Street, San Leandro, California 94577

Received April 24, 2002

From REACCS, to MDL ISIS/Host Reaction Gateway, and most recently to MDL Relational Chemistry Server, a new product based on Oracle data cartridge technology, MDL's reaction database management and retrieval systems have undergone great changes. The evolution of the system architecture is briefly discussed. The evolution of MDL reaction substructure search (RSS) algorithms is detailed. This article mainly describes a novel RSS algorithm. This algorithm is based on a depth-first search approach and is able to fully and prospectively use reaction specific information, such as reacting center and atom–atom mapping (AAM) information. The new algorithm has been used in the recently released MDL Relational Chemistry Server and allows the user to precisely find reaction instances in databases while minimizing unrelated hits. Finally, the existing and new RSS algorithms are compared with several examples.

## 1. INTRODUCTION

REACCS, the first REAction ACCess System from MDL, was demonstrated publicly in 1981.<sup>1</sup> In the past 20 years, MDL's reaction database management and retrieval systems have undergone three generations of development. REACCS was a mainframe application used to build, maintain, and search proprietary and commercial databases of chemical reactions and related data. It was one of the earliest and most widely used reaction database management and retrieval systems in the 1980s. MDL ISIS (Integrated Scientific Information System),<sup>2,3</sup> the second generation of MDL database management and retrieval systems, was launched in 1991. MDL ISIS is based on a client-server architecture, which allows users to run MDL ISIS/Draw (for structure input) and MDL ISIS/Base (for database searching) on a desktop computer (client) to access reaction and molecule databases managed and searched by MDL ISIS/Host running on a mainframe computer (server). With the shift of the methods from the product based combinatorial synthesis to the reaction based combinatorial synthesis over the past a few years, reaction databases have become a more valuable source for the design of new combinatorial libraries. This has sparked the innovation of the new reaction searching and management system. In 2001, MDL launched the third generation of reaction database management and retrieval systems called MDL Relational Chemistry Server.<sup>3</sup> There are two major differences between MDL ISIS and MDL Relational Chemistry Server. First, in MDL ISIS/Host reactions are stored and searched using a proprietary database engine. In contrast, MDL Relational Chemistry Server is based on a new technology provided by Oracle, the data cartridge, which allows custom search methods such as RSS, to be called directly by the Oracle search engine. In other words, the chemical representation routines have fully been integrated with Oracle, which allows end users to use the

standard RDBMS search interface (SQL) to perform reaction searching. Second, in MDL ISIS the RSS method is based on the same RSS algorithm developed for REACCS, while the RSS approach employed in MDL Relational Chemistry Server is based on a completely new RSS algorithm, which will be described in this paper.

In this article, the evolution of the RSS algorithms used in the different MDL reaction database management systems will be outlined. The main focus is a new RSS algorithm developed for the reaction cartridge. A detailed comparison of this algorithm with the existing MDL RSS methods will be discussed using several examples.

## 2. CLASSIFICATION OF RSS QUERIES

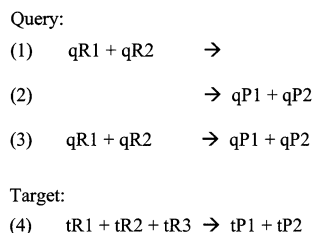
RSS is a centrally important search method for any reaction retrieval system. Similar to SSS (SubStructure Search),<sup>4</sup> RSS is designed to retrieve reaction instances from reaction databases that contain the entire query reaction. An RSS query consists of one or more substructures that represent reactants and/or products. Therefore, RSS queries can have all the features of SSS queries. RSS queries can also have reaction-specific features that an SSS query does not have. RSS queries can be classified according to different criteria.

**2.1. Classification Based on Query Components.** RSS queries can take three different forms as shown in Figure 1.

Query 1 consists of only reactants and can be used to retrieve all reactions that include both substructures qR1 and qR2 in two of their reactants from the database. Query 2 consists of only products and is designed to find all reaction instances that contain substructures qP1 and qP2 in two of their product structures. Most RSS queries take the form of query 3 that consists of both reactants and products and can be used to retrieve all reactions that contain all four substructures.

**2.2. Classification Based on Reacting Center Information.** The reacting center of a reaction is the collection of

\* Corresponding author phone: (510)895-1313; e-mail: L.Chen@mdl.com.



**Figure 1.** Three different forms of RSS queries and a multiple component target reaction: (1) An RSS query consisting of only reactants. (2) An RSS query consisting of only products. (3) An RSS query consisting of both reactants and products. (4) A target reaction with three reactants and two products.

atoms and bonds that are changed during a chemical reaction. RSS queries can be grouped into two categories based on the presence or absence of the reacting center criterion.

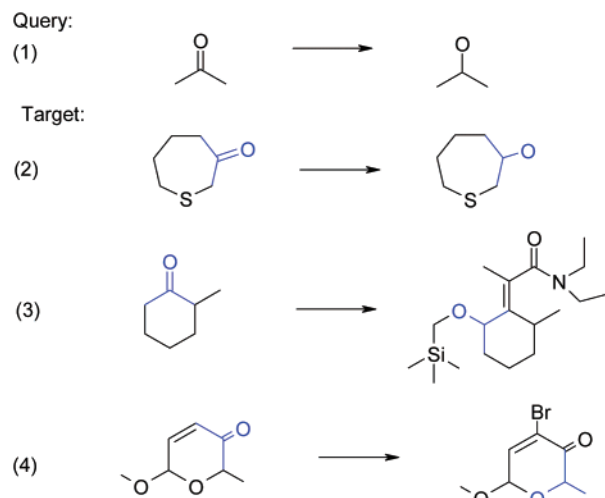
**2.3. Classification Based on Atom-Atom Mapping (AAM) Information.** AAM information defines the correspondence between reactant and product atoms of an RSS query or a target reaction. There are different types of RSS queries based on the presence or absence of the AAM relationship between reactants and products.

Typical, well-defined, modern RSS queries are those that consist of both reactants and products and contain both reacting center and AAM information. A good RSS algorithm should be able to handle such queries correctly and efficiently. It should also be able to deal with any other types of RSS queries.

### 3. THE BRUTE-FORCE RSS ALGORITHM

At the simplest level, an RSS approach does not take reaction-specific properties, such as reacting center and AAM information, into account. Thus, the simplest RSS algorithm can use an exhaustive approach. First, all combinations of query-target structure pairs on each side of the reaction are determined, and then the molecular substructure search (SSS) algorithm is used to compare each query-target structure pair. For example, to determine whether query 3 can be matched to the target reaction 4 (see Figure 1), a total of 10 query-target molecule pairs must be compared: six molecule pairs are required for comparing the reactant side (2 query reactants  $\times$  3 target reactants = 6 possibilities), and four molecule pairs are required for comparing the product side (2 query products  $\times$  2 target products = 4 possibilities). Finally, all the SSS hits are combined using a logical AND operation to determine the final RSS mapping result.

The main advantage of this algorithm is the simplicity of the procedure itself. Such an RSS algorithm is, however, inefficient because it spends a large amount of time in the time-consuming SSS analysis. Another problem is that it retrieves many inaccurate hits for a given RSS query. For example, using query 1 of Figure 2 for retrieving reaction instances in the database that convert a carbonyl to a carbon-oxygen single bond, the brute-force RSS algorithm can indeed retrieve all such transformations, such as reaction 2 of Figure 2, but it can also retrieve other reaction examples unrelated to the desired transformation, such as reactions 3 and 4 of Figure 2. This is because the information about which atoms and bonds are in the transformation is not used in the search. Any reactions whose reactant and product each contain carbonyl and hydroxyl groups respectively will be



**Figure 2.** Using query 1 of this figure, a brute-force RSS algorithm based on pure SSS retrieves a correct hit (target 2) but also false hits (targets 3 and 4). The mappings of the query onto the target reactions are highlighted in blue.

retrieved. Because of the above problems, this exhaustive approach was never used in MDL products.

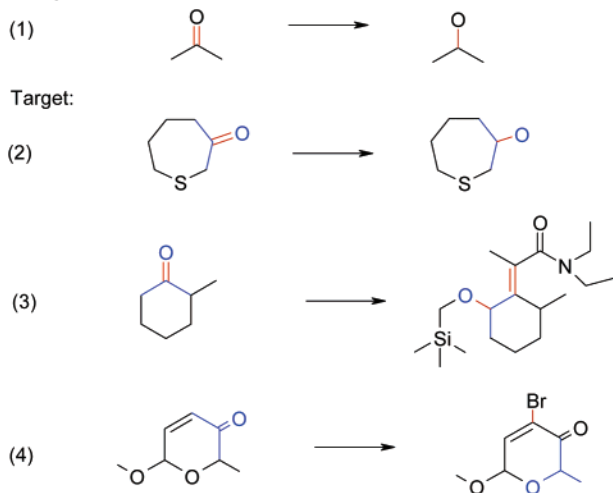
### 4. THE FIRST GENERATION OF THE RSS ALGORITHM

A more efficient RSS algorithm can be designed by taking into account the results of the SSS matches carried out earlier. This is done using a breadth-first search procedure. It also consists of sequential calls to the SSS algorithm for each substructure in an RSS query. The major difference between this algorithm and the brute-force search method is that if a query structure cannot be matched to any target structures on the same side of the reaction, the RSS process terminates immediately. Consider the same example of comparing query 3 with target 4 (see Figure 1). In the best case, only three SSS calls are needed to determine that the query 3 cannot be matched to target 4: the first query reactant, qR1, cannot be matched to any of three target reactants, tR1, tR2, and tR3, and thus the RSS process terminates immediately. Only in the worst cases does the algorithm need to carry out the comparison for all 10 query-target molecule pairs. This algorithm was implemented in the REACCS system as the first generation of the RSS method (RSS1).

It is obvious that the RSS algorithm described above is much more efficient than the brute-force search method. However, the second problem encountered by the brute-force technique still exists: no improvement of the quality of the hitlists is achieved. That is, many inaccurate hits are still retrieved from the database for a given RSS query because in this RSS algorithm, the matching between query and target structure pairs is still completely done by the pure SSS algorithm.

Because the key part of a reaction that chemists are usually most interested in is the reacting center where the changes are occurring, the reacting center bonds of most reactions in a database have been explicitly marked as shown in reactions 2, 3, and 4 of Figure 3 (in red lines). If the reacting bonds of a query are also marked, as shown in query 1 of Figure 3, the differentiating ability of the RSS algorithm can be improved by incorporating additional checking of reacting center matching into the algorithm. For example, this refined

Query:



**Figure 3.** The query and the target reactions shown here are identical with those shown in Figure 2, except that reacting bonds are marked (in red). The reacting center information helps the RSS algorithm to exclude target 4 as a hit, but it still matches this query to target 3 (false hit) as well as the correct hit (target 2).

RSS algorithm will no longer match query 1 to target reaction 4 because in reaction 4, all the C=O and C–O bonds are not reacting bonds and thus cannot be matched to the reacting C=O and C–O bonds of the query. However, the algorithm still matches query 1 to target reaction 3. This is because reaction 3 meets all the requirements of the query 1, including a reacting C=O bond in the reactant and a reacting C–O bond in the product. Therefore, we need to find some new strategy to exclude the remaining false hits.

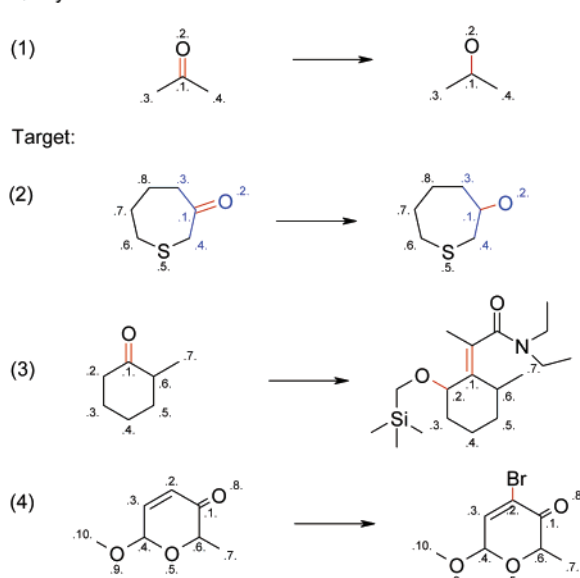
## 5. THE SECOND GENERATION OF THE RSS ALGORITHM

A simple solution to avoid retrieving unrelated hits, such as reaction 3 of Figure 3, is to establish the AAM relationships between reactant and product atoms for both query and database reactions. For RSS queries, the atom–atom correspondences between reactants and products can be assigned manually by users. To allow chemists to manually assign AAMs for RSS queries, a graphical editor was created within REACCS.

However, for target reactions the situation is more complicated. In the late 1980s, four major reaction databases were developed for REACCS with a total of about 90 000 reactions.<sup>5</sup> During that time, many companies had also built their own private reaction databases for use with REACCS. The reactions in these databases usually contained reacting center information but lacked AAM relationships between reactant and product atoms. It would have been prohibitively expensive to manually assign AAMs for all the reactions in the existing databases. Therefore, in 1988 a program called ARCP (Automatic Reacting Center Perception), later renamed as the Automapper, was developed for automatic reacting center perception and AAM assignment.<sup>5</sup>

Using the Automapper, databases can be updated by applying the Automapper to all of the reactions and then automatically re-registering the reactions with the newly perceived reacting center and AAM information. All new reactions added to existing or new databases are handled in a similar manner. The Automapper can also be used to preprocess RSS queries. For example, in the MDL ISIS

Query:



**Figure 4.** The query and the target reactions shown here are identical with those shown in Figure 3, except that atom–atom mappings between reactant and product have been assigned.

system, suppose a user draws an RSS query using MDL ISIS/Draw but does not manually specify any reacting center bonds and AAMs. This query is then transferred to MDL ISIS/Base and the Automapper can automatically perceive the reacting center and assign the AAMs for the query. The Automapper can also be explicitly invoked by users within MDL ISIS/Base.

Using the available AAM information for both queries and database reactions makes it possible for the RSS algorithm to be significantly improved.

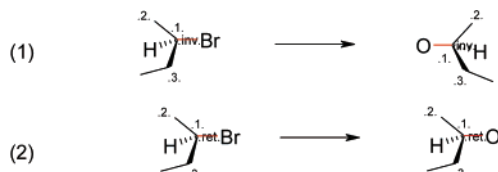
**5.1. Description of the RSS Algorithm.** The RSS algorithm that will be described in this section is similar to its predecessor (RSS1) described in section 4, with several major modifications. This algorithm is also based on a breadth-first search procedure. No query–target product pairs will be compared until all the query–target reactant pairs have been successfully matched. During the matching, all the overlapping SSS mappings between each query–target molecule pair must be recorded. After all the query–target molecule pairs have been handled, the SSS mappings are then used to perform a complete AAM checking to ensure that the atom–atom matchings between query–target reactant pairs and product pairs are consistent with AAMs within the RSS query and the target reaction. However, storing all the overlapping SSS mappings for all query–target molecule pairs of a query–target reaction pair is very memory intensive in many cases. To reduce memory usage, the RSS algorithm reduces the SSS mappings to a more compact form that can be stored in a bit set. This reduced form may lead to false equivalent mappings that result in extra hits.

With the above improvement of the RSS algorithm, most of the ambiguities involving query 1 and the target reactions in Figure 3 have been resolved by using AAM information as shown in Figure 4. Target reaction 2 can still be successfully matched by query 1, but now both target reactions 3 and 4 are no longer matched with query 1.

Furthermore, additional RSS query features have been developed based on the AAM information.

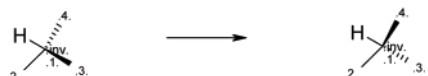


**Figure 5.** Explicit hydrogen atoms on mapped atoms between reactants and products can be used to search for reactions involving reacting hydrogens.

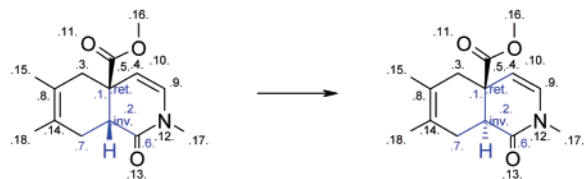


**Figure 6.** RSS queries with Inversion (query 1) and Retention (query 2) properties.

Query:



Target:



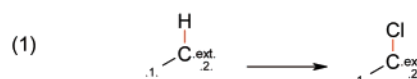
**Figure 7.** With the help of atom-atom mapping and Inversion/Retention properties, it is possible to search for stereochemically correct reactions.

**5.2. Handling of Reacting Hydrogen Atoms.** Without AAM information, searching for reactions involving reacting hydrogen atoms is difficult because hydrogen atoms are treated implicitly in MDL software. For example, searching for reduction of olefins will also lead to finding addition reactions. This problem can now be solved by comparing the hydrogens on the pair of atoms in reactants and products that have the same AAM value. In the hydrogenolysis of an allylic bromine (Figure 5), the reaction to be matched to this query not only must lose a bromine on the C.1. atom in the reactant but also must replace the bromine with a hydrogen on the C.1. atom in the product.

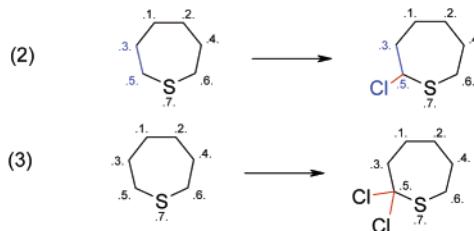
**5.3. Handling of Stereochemistry.** Transformations involving stereochemistry are important reactions in databases. Therefore, two new atom properties, inversion and retention, associated with AAM were created. Inversion and retention marks specify how the stereochemistry of an asymmetric center in a reactant or product changes (or does not change) in a reaction. The markings are as follows: (a) "Inv." which specifies that an asymmetric center inverts its stereochemical configuration during the reaction (see Figure 6(1)). (b) "Ret." which specifies that an asymmetric center retains its stereochemical configuration during the reaction (see Figure 6(2)). The basis for the determination of retention and inversion properties is the ligand atom numbering and atom mapping numbers.

The retention and inversion properties can be applied to both queries and database reactions, indicating that the transformation must proceed with the correct stereochemistry. This property is particularly useful in the search for epimerisation reactions (target reaction of Figure 7), where there are no bonds to mark as reacting centers. For example, using

Query:



Target:



**Figure 8.** The application of the ExactChange property to a pair of reactant and product atoms that have the same atom-atom mapping completely describes the changes allowed at the atom in an RSS query. In the example shown in this figure, query 1 will be matched to target reaction 2, but not reaction 3.

the query of Figure 7, which requires that the stereoconfiguration of the C.1. atom be inverted, RSS is able to find the target reaction of Figure 7 in the database.

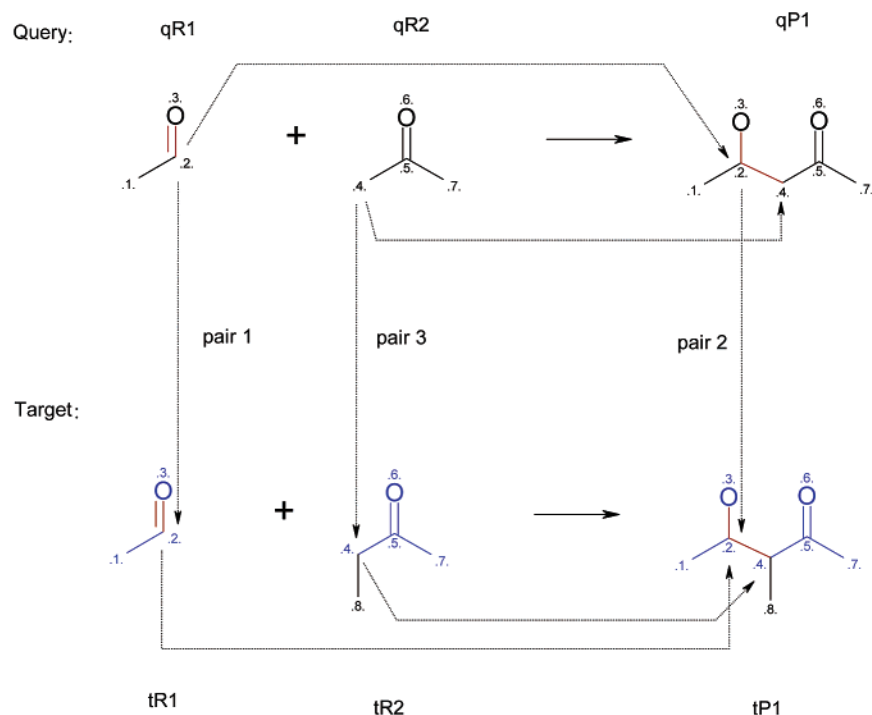
**5.4. ExactChange Match.** With the help of AAM information, it is also possible to search for very specific reactions. This is achieved by using a query property called ExactChange (its marking is .ext.). The application of the ExactChange property to a pair of reactant and product atoms that have the same AAM completely describes the bond changes allowed at the atom in an RSS query. Consider the example shown in Figure 8. The C.2. atoms of both the reactant and product of query 1 have an ExactChange flag. This specifies that one and only one hydrogen of the C.2. atom must be replaced by one chlorine during the reaction. Thus, the C.2. atom of the query reactant can only be matched to a target reactant carbon atom which loses exactly one hydrogen. Similarly, the C.2. atom of the query product can only be matched to a target product carbon atom which obtains exactly one chlorine. These two matched target atoms cannot involve any other reacting bonds, and furthermore, they must have the same AAM value. Thus, query 1 of Figure 8 will be matched to target reaction 2 but not reaction 3. This is because in the latter reaction, two hydrogens attached to C.5. atom are replaced by two chlorine atoms.

The significantly improved RSS algorithm based on the AAM information described above can be regarded as the second generation of the RSS method (RSS2) and has been used in both the REACCS and MDL ISIS systems.

## 6. THE THIRD GENERATION OF THE RSS ALGORITHM

RSS2 described in the previous section has two major problems. (1) As discussed in the previous section, RSS2 used a string with limited size to store hashed SSS mappings of all query-target molecules pairs for a given query-target reaction pair. In some cases only partial "hashed" information about all the reactant mappings could be stored. This occasionally would lead to extra hits when there was a hash collision in RSS2. This problem was not as significant in the past when reaction databases were relatively small, and returning extra incorrect hits is generally regarded as preferable to missing hits. However, in the past decade or so, many changes have taken place. Large reaction databases containing hundreds of thousands, even millions, of reactions have become widely available.<sup>6</sup> With the dramatic increase





**Figure 9.** RSS3 algorithm prospectively uses three types of information during the matching process: (1) reacting centers, (2) AAM information, and (3) the matching between query-target molecule pairs performed previously.

in the number of reactions in modern reaction databases, the major problem of reaction searching is no longer getting too few hits. Rather, today the problem has become that a single RSS search usually yields far too many hits. Because of these new challenges, elimination of unrelated hits has become one of the most important requirements in designing a new generation of the RSS algorithm.

(2) Although the incorporation of AAM information in RSS2 significantly increases its flexibility and specificity, the use of AAM information in RSS2 is not optimized. The AAM information is used in a retrospective manner. That is, it is used to ensure the correct relationships between query and target atom pairs with regard to the AAM relationships within the query and the target reactions only after all the possible, time-intensive SSS processes have been completed. In this section, we will describe a novel RSS algorithm (RSS3) based on the prospective use of both the reacting center and AAM information. First, the algorithm for matching well-defined RSS queries with typical target reactions will be described. Then, the handling of other types of RSS queries will follow. The detailed comparison of RSS3 with RSS2 will be given in the last section.

**6.1. Handling Typical RSS Queries.** As described in section 2, a well-defined RSS query consists of both reactants and products and contains both reacting center and AAM information. The same is true for a target reaction. The procedure of a typical RSS process for matching well-defined RSS queries with typical target reactions is as follows.

**(a) Creation of the Starting Atom Table for the RSS Query.** The atom-by-atom matching of two given structures begins with the selection of one atom of the query molecule and tries to match it with an atom of the target molecule. These starting point atoms are called starting atoms. During the preparation of an RSS query, the starting atoms are chosen for each query molecule. The starting atoms must be reacting atoms for any query molecule that contains

**Table 1.** Starting Atom Table Showing the Relationship between Primary Starting Atoms (PSA) and Secondary Starting Atoms (SSA) of the Molecules of the RSS Query of Figure 9

no.	qR1	qR2	qP1
1	PSA: 2		SSA: 2
2		PSA: 4	SSA: 4
3	SSA: 2	SSA: 4	PSA: 2

reacting centers. There are two types of starting atoms: primary starting atoms (PSA) and secondary starting atoms (SSA). One primary starting atom and one or more secondary starting atoms will be chosen for each query molecule. The primary starting atom must be a reacting atom that has the maximum number of neighbor atoms among all reacting atoms of a component molecule. The secondary starting atoms are chosen based on their AAM relationship with the primary starting atoms of the query molecules on the other side of the reaction. Let us use an example to explain this. The primary and secondary starting atoms for the three molecules of the RSS query of Figure 9 are summarized in Table 1. First, the reacting atom 2 of the first query reactant (qR1) is chosen as qR1's primary starting atom because this reacting atom has two neighbor atoms, while the other reacting atom 3 of qR1 has only one neighbor atom. Then, the reacting atom 2 of the query product (qP1) is chosen as qP1's secondary starting atom because this atom has the same AAM value as the corresponding primary starting atom of qR1 (see row 1 of Table 1). Similarly, the only reacting atom 4 of the second query reactant (qR2) is chosen as the primary starting atom for qR2. Then, the reacting atom 4 of the product is chosen as another secondary starting atom for qP1, because qP1's atom 4 has the same AAM value as that of the primary starting atom 4 of qR2 (see row 2 of Table 1). Next, the reacting atom 2 of the product, qP1, is chosen as qP1's primary starting atom, and then the reacting atom 2 of qR1 as its secondary starting atom because this atom has

the same AAM value as qP1's primary starting atom 2. The second query reactant, however, has no reacting atom whose AAM is equal to 2. In this case, the reacting atom 4 is chosen as qR2's secondary starting atom because this atom has the same AAM value as that of qP1's atom 4 (see row 3 of Table 1). The relationship between the primary and secondary starting atoms of the molecules of the RSS query of Figure 9 is shown in the "starting atom table" (Table 1). These starting atoms will be used in stage (b).

If the RSS algorithm starts comparison from the reactant side, then the primary starting atoms of the query reactants will be used as starting atoms for the query reactants, and the secondary starting atoms of the query products will be selected as starting atoms for the query products. On the other hand, if the RSS algorithm starts comparison from the product side, then the primary starting atoms of the query products will be used as starting atoms for the query products, and the secondary starting atoms of the query reactants will be selected as starting atoms for the query reactants.

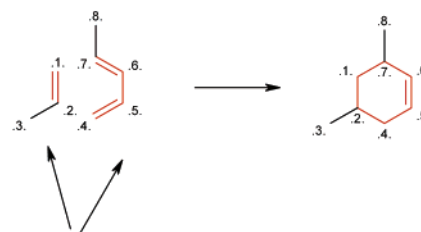
**(b) Main Stage: Matching Query-Target Molecule Pairs.** RSS3 uses a depth-first search method. The first stage involves calling the backtrack algorithm<sup>7</sup> to determine a combination for comparing query-target molecule pairs of the first reaction side, which is chosen using the method described in section (d) below. The handling of the first reaction side, however, depends directly upon the results from comparing the corresponding query-target molecule pairs of the second side. In detail, the comparison of each query-target molecule pair of the first reaction side begins by matching the query's primary starting atom to a reacting atom of the target molecule. If the matching of this molecule pair succeeds, the next molecule pair to be compared is not picked up from the same reaction side, but from the second side, as guided by the query's starting atom table. Take the query and the target reaction shown in Figure 9 as an example.

Step 1: The algorithm starts by comparing the query reactant qR1 with the target reactant tR1 using qR1's primary starting atom 2 (see Table 1) and tR1's reacting atom 2 as a starting atom pair, because tR1's reacting atom 2 is also a carbon atom and its number of neighbor atoms matches that of qR1's primary starting atom 2. The substructure matching between qR1 and tR1 is successful.

Step 2: The next step is not to compare qR2 with tR2. Rather the algorithm turns to the product side and begins to compare qP1 with tP1. This comparison starts with the query-target starting atom pair (qP1's atom 2, tP1's atom 2). Here, qP1's atom 2 is its secondary starting atom (see Table 1) corresponding to qR1's primary starting atom 2. The atom 2 of tP1 is chosen as the target starting atom because this atom has the same AAM value as atom 2 of tR1 which has been matched to qR1's primary starting atom 2.

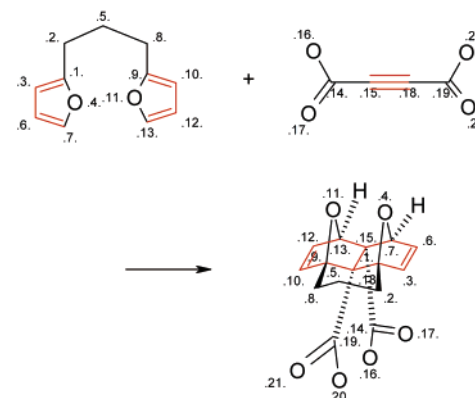
After successfully matching qP1 with tP1, the algorithm then returns to the first reaction side and starts to compare qR2 with tR2 by using the starting atom pair (qR2's atom 4, tR2's atom 4), because both qR2 and tR2 have only one reacting atom. Since the matching between qR2 and tR2 is successful, the algorithm turns to the product side and attempts to compare qP1 with tP1 again, but this time it uses another starting atom pair (qP1's atom 4, tP1's atom 4). However, a simple check of the previous matching record shows that qP1's secondary starting atom 4 has been successfully matched to tP1's atom 4. (This check can be

Query:



The query reactant consists of two disconnected fragments

Target:



**Figure 10.** A strict AAM checking is needed to avoid the query involving an intra Diels–Alder reaction being matched to a target undergoing an inter-Diels–Alder reaction.

avoided using another strategy described in section (d) below.) Therefore, no SSS call is executed. The RSS algorithm thus terminates with the conclusion that the overall matching between the query and the target reaction has succeeded.

It should be noticed that if qP1 cannot be matched with tP1 using the starting atom pair (qP1's atom 2, tP1's atom 2) in Step 2, the algorithm will not try any other starting atom pairs of qP1 and tP1, because this starting atom pair is chosen based on the matching between qR1 and tR1 and on the AAM relationships between qR1 and qP1 and between tR1 and tP1, respectively. However, in this case, the algorithm will try to match qR1's primary starting atom 2 to tR1's each remaining reacting atoms until qR1 is successfully matched to tR1 again. Then the algorithm turns to match qP1 and tP1 with a new starting atom pair. This process is repeated until either both qR1–tR1 and qP1–tP1 pairs are successfully matched, or all tR1's reacting atoms have been tried. In the later case, the algorithm will now try to match qR1 with tR2.

It should be pointed out that although the prospective use of the AAM information described above can lead to correct matching in most cases, it might result in false hits in some special cases. For example, consider the RSS query and the target reaction shown in Figure 10. The query is an intramolecular Diels–Alder reaction whose reactant consists of two disconnected substructures. The target reaction, on the other hand, is an intermolecular Diels–Alder reaction between two reactants. When comparing this query-target reaction pair, the query product is successfully matched to the target product, and the query reactant is also successfully matched to the first target reactant. This result is not correct.

To guarantee correct matching results, RSS3 uses the following procedure to solve the above problem. First, the algorithm collects all the overlapping mappings of the query-target molecule pair on the first reaction side and all the overlapping mappings of the corresponding query-target molecule pair on the second reaction side. As soon as the matching between the latter molecule pair is completed and the substructure mappings have been collected, the algorithm performs a strict check to ensure that each target reactant atom matched to the query reactant atom has the correct AAM relationship with the corresponding target product atom that matches the query product atom. This check quickly detects that the matching between the first query reactant and the first target reactant in Figure 10 is not valid because two target product atoms (C.15. and C.18.) that match query product atoms (C.1. and C.2.) do not come from the first target reactant but from the second target reactant.

As can be seen from the above description, the new algorithm is able to use the following two types of information prospectively: (1) reacting centers and (2) AAM information. This strategy allows the new algorithm not only to find a successful matching of a given query-target reaction pair faster but also to terminate sooner for a query-target reaction pair that cannot be matched, when compared with RSS2. For example, when comparing query 1 with target 3 of Figure 4 (see also section 5), RSS2 first successfully matches the query's reactant and product with target 3's reactant and product, respectively, and then it uses the collected SSS mappings to perform the AAM checking and detects that the query does not match the target. RSS3, on the other hand, terminates the searching immediately after successfully matching the query's reactant with the target's reactant using the starting atom pair (query reactant C.1., target reactant C.1.), but not matching the query's product with the target's product using that starting atom pair (query product C.1., target product C.1.), and not matching query reactant atom C.1. with the target reactant's other reacting atom: O atom.

The performance of the new algorithm can be further improved using some simple strategies described below.

**(c) Check if an RSS Query is Bigger Than a Target Reaction.** If a query has more reactants and/or products than those of a target reaction, then the query cannot be matched to the target reaction, and RSS terminates immediately. This stage is done after stage (a) but before stage (b).

**(d) Selection of the First Reaction Side.** The reaction side where the RSS algorithm starts the comparison is called the first reaction side, and the other side is called the second reaction side. The first reaction side can be either reactant or product side. It is chosen using the following procedure

```

if (qNr x tNr) <= (qNp x tNp) then
    select the reactant side as the first reaction side,
    and the product side as the second reaction side;
else
    select the product side as the first reaction side,
    and the reactant side as the second reaction side.
end if

```

where qNr, qNp, tNr, and tNp are the numbers of query

**Table 2.** There Are a Total of Six Possibilities to Compare an RSS Query with Two Reactants and a Target Reaction with Three Reactants

combination	qR1	qR2
1	tR1	tR2
2	tR1	tR3
3	tR2	tR1
4	tR2	tR3
5	tR3	tR1
6	tR3	tR2

reactants, query products, target reactants, and target products, respectively. This means that the first reaction side has a smaller number of query-target molecule pairs. The strategy to do the query-target comparisons with the smaller side first can usually lead to quicker termination of the RSS algorithm for those query-reaction pairs that cannot be matched. For example, when comparing a query-target reaction pair, in which both the query and the target reaction have two reactants and one product, if the reactant side is handled first, then in the worst case, the RSS algorithm will not terminate until all four query-target reactant pairs have been compared. However, if starting with the product side first, then the algorithm can terminate immediately after the failure of the matching between the query and target products.

Even for a query-target reaction pair that can be matched successfully, the above strategy can lead to finding a successful matching faster. For example, according to the discussion above, for the query-target reaction pair shown in Figure 9, RSS3 will choose the product side as the first reaction side because the product side has only one query-target molecule pair (qP1-tP2), while the reactant side has four possible molecule pairs (qR1-tR1, qR1-tR2, qR2-tR1, qR2-tR2). After successful matching between qP1 and tP1, the algorithm will quickly find the successful matchings of the qR1-tR1 pair using the starting atom pair (qR1's atom 2, tR1's atom 2) and of the qR2-tR2 pair using their starting atom pair (qR2's atom 4, tR2's atom 4). No extra checking is needed (cf. step 2 of stage (b) above).

The above stage (d), the selection of the first reaction side, is done after stage (c) but before stage (b).

**6.2. Handling Other Types of RSS Queries.** Other types of RSS queries (see section 2) can be treated as special cases of the well-defined RSS queries that consist of both reactants and products and that contain both the reacting center and AAM information. Therefore, the RSS process involving these queries can be handled quite easily. For example, for RSS queries consisting of only reactant(s), the algorithm is as follows.

The algorithm uses a backtrack algorithm to calculate all possible ways to compare molecules of the reactant sides of the query and target reaction. For example, when comparing query 1 of Figure 1 that consists of two reactants (qR1, qR2) and target reaction 4 of Figure 1 that contains three reactants (tR1, tR2, tR3), there are a total of six possible "combinations" as shown in Table 2.

After a choice is made, the query molecule qR1 is compared with target reactant molecule tR1. If qR1 is found to be a substructure of tR1, then qR2 is compared with tR2. If the comparison of any molecule pair fails, the current combination will be abandoned, and the next combination is tried. If the comparison of all molecule pairs returns



**Table 3.** Comparison of Three Generations of MDL RSS Algorithms

RSS version	algorithm feature	product	note
RSS1	breadth-first search based algorithm, not using AAM	REACCS (original version)	section 4
RSS2	breadth-first search based algorithm, using AAM retrospectively	REACCS (Version 7 & up), MDL ISIS/RCG	section 5
RSS3	depth-first search based algorithm, using AAM prospectively	MDL Relational Chemistry Server	section 6

positive results, that is, if the query molecules are substructures of the corresponding target reaction molecules in all molecule pairs of a combination, then the matching of the query with the target reaction is successful; otherwise, it fails. Each matching result is recorded to avoid duplicate SSS processes. For example, if qR1 fails to be matched to tR1 in combination 1, combination 2 will be skipped, and qR1 will be tried to match tR2 in combination 3.

The RSS queries consisting of only product(s) can be handled using the same algorithm described above for dealing with the queries consisting of only reactant(s).

The above approach for handling reactant-only and product-only queries is similar to that of RSS2.

In the real situation, some reactions may be missing some or all AAM values and/or reacting center information. Special methods have been designed for dealing with such special cases. For example, in some cases nonreacting atoms of a query molecule have AAM information, while all reacting atoms of the same molecule possess no AAM at all. In this case, the atom selected as a starting atom for this query molecule is an atom that has AAM and that also has the largest number of attachments among all atoms that have AAM in the query molecule. This starting atom is thus not a reacting atom.

RSS3 described above can be regarded as the third generation of the RSS approach and has been used in MDL Relational Chemistry Server,<sup>3</sup> a new MDL product released recently.

## 7. COMPARISON OF THE EXISTING AND NEW RSS ALGORITHMS

In this section we will present the comparison of different MDL RSS algorithms. In particular we will describe how RSS3 used in MDL Relational Chemistry Server differs from RSS2 in the older MDL ISIS products such as MDL ISIS/Host Reaction Gateway and MDL ISIS/Base using several concrete examples.

**7.1. Major Differences of MDL RSS Algorithms.** The major differences of the three generations of MDL RSS algorithms are summarized in Table 3. RSS1 was a breadth-first search based RSS algorithm. It was an uncorrelated RSS algorithm because it did not use AAM information, which was not available in its time. The main problem of RSS1 is that it usually retrieves too many false hits from databases. Also, the search features are quite limited. RSS1 was used in the first version of REACCS.

RSS2 is also a breadth-first search based RSS algorithm. The major difference between RSS2 and RSS1 is in that the former uses AAM information but the latter does not. Using AAM data allows RSS2 to be more flexible and specific than RSS1. RSS2 also supports more search features, such as ExactChange and Retention/Inversion features. However, the use of the AAM information in RSS2 is quite limited. This important reaction specific information is used only in the last step for checking the consistency between query-

target atom-atom matchings and reactant-product AAMs within the query and the target. To do this requires that all the overlapping SSS mappings of all the query-target molecule pairs must be recorded first. This is not always possible even using a hash technique described in section 5. It is not unusual that the hash storage overflows, and thus only part of the SSS mappings of a given query-target reaction pair can be stored, resulting in false hits. Our experience shows that RSS2 has trouble dealing with some queries, particularly those that involve explicit hydrogen atoms and the ExactChange property (see examples below.).

RSS3, the third generation of MDL RSS method, is based on a depth-first search algorithm. It is able to use AAM data fully and prospectively. In fact, the entire RSS process for handling well-defined query-target reaction pairs is guided by the AAM information. This allows RSS3 to limit the time-consuming substructure search phase for each component molecule as much as possible. Furthermore, in RSS3 only the overlapping SSS mappings of two query-target molecule pairs (one from each reaction side) need to be recorded for AAM checking. Also, for any RSS queries that do not contain two or more disconnected fragments in one component molecule, only the overlapping SSS mappings that are generated by using two starting atom pairs (one for each reaction side) need to be stored. Only for those special RSS queries that contain two or more disconnected fragments in one component molecule, is the algorithm required to collect all the overlapping SSS mappings, which are generated starting with all the possible starting atom pairs of two query-target molecule pairs (one from each reaction side). Therefore, the new algorithm is able to perform a complete AAM check to ensure the matching between the given query-target reaction pair is correct in almost all cases. The main differences between RSS2 and RSS3 are summarized in Table 4.

**7.2. ExactChange Requirements on the Query.** As described in section 5.4, an ExactChange flag allows one to mark all atoms that one wants to change exactly as one specifies them in the reactants and products of the query. However, the ExactChange convention implemented in RSS3 is not the same as that used in RSS2. In the latter case, hydrogen atoms attached to an atom with ExactChange flag can be either implicit or explicit ones. However, RSS3 requires that all reacting bonds must be explicitly drawn and marked, including bonds to hydrogen atoms.

MDL ISIS searches using RSS2 and MDL Relational Chemistry Server searches based on RSS3 might lead to quite different hits for the same query using the same reaction database. For example, MDL ISIS retrieved 47 hits from the ChemInform Reaction Library (98.1 release) using query 1 shown in Figure 11. One hit is shown in Figure 11 as Target 3.

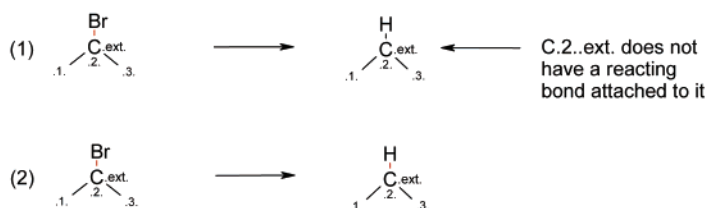
RSS3 hits nothing using the same query 1 against the above hitlist. This is due to the new ExactChange convention,



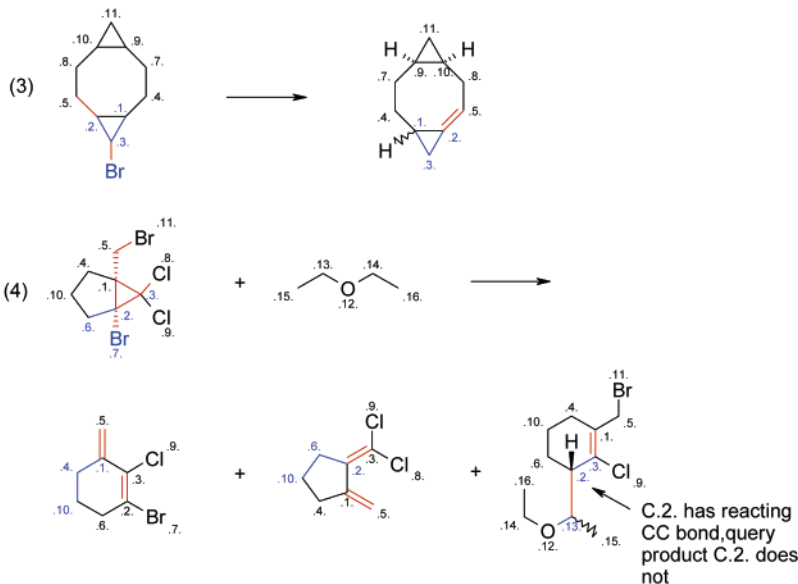
**Table 4.** Major Differences between RSS2 and RSS3

item	RSS2	RSS3	note
algorithm	see Table 3	see Table 3	sections 5, 6, & 7.1
Exact Change convention	reacting hydrogens attached to an atom with ExactChange flag can be either implicit or explicit	all reacting bonds attached to the atom with ExactChange flag must be explicitly drawn and marked, including bonds to hydrogens	section 7.2
correct AAM relationships between query and target	not always	yes	section 7.3
correct matching of queries with explicit hydrogens	not always.	yes	section 7.4
reactions related to ether and ester structures	AAM and reacting center information of O-atom in O-C may be ignored to return more hits. May return false hits	more flexible, allowing users to retrieve specific reactions	section 7.5
can detect reacting center errors?	no	yes	section 7.6
queries with conflicting reacting center settings	may return a lot of false hits	reject to return any hits for such invalid RSS queries	section 7.7
AAM conventions	see Table 5	see Table 5	section 7.8

Query:



Target:

**Figure 11.** RSS3 matches query 2 with target 3 but not query 1. The blue highlight in target 3 shows query 2. RSS3 does not match query 2 with target 4, which is a false hit returned by RSS2 using query 2.

since the C.2. atom of the query product has no reacting bonds. RSS3 requires that the C-H bond at the product side must be marked as a reacting bond because there is an ExactChange flag on the C.2. atom of C-H, such as in query 2 of Figure 11. By using this new query 2, MDL ISIS retrieved the same hitlist (47 hits) from the ChemInform Reaction Library (98.1 release) as it did using query 1. One hit is shown in Figure 11 as target 4.

RSS3 found 46 target reactions using the same query 2 (see Figure 11) against the above hitlist this time. The missing hit (target 4) turned out to be a false hit by MDL ISIS. The reason target 4 is a false hit is that the C.2. atom

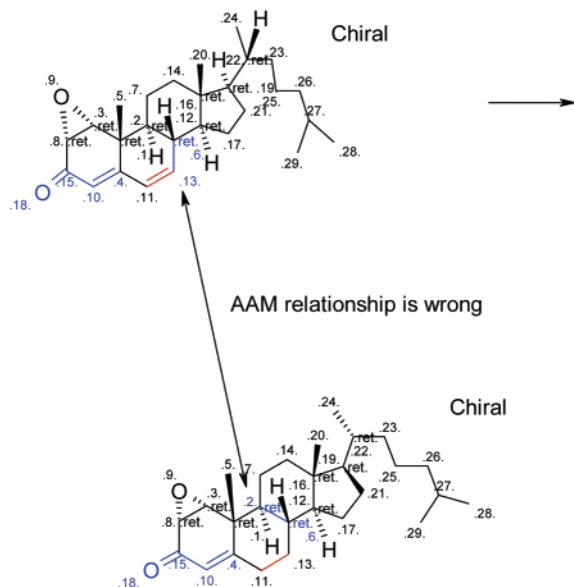
of its third product has a C-C reacting bond attaching to it while the corresponding query C.2..ext. atom does not.

**7.3. Correct Atom-Atom Mapping Relationships between Query and Target.** RSS3 strictly checks the relationships between the query/target atom-atom matchings and the atom-atom mappings within the query and within the target reaction. For this reason, RSS3 can guarantee that such relationships are correct. This was not the case in RSS2. For example, RSS2 matches the query in Figure 12 to 30 hits from the Theilheimer database. However, only 4 hits are correct among the 30 hits. In this case, RSS2 returns extra hits; one of the false hits is shown in Figure 12 as the target.

Query:



Target:



**Figure 12.** RSS2 incorrectly matches the query with the target reaction. The blue highlight shows the query in the target. RSS3 does not match this query-target pair.

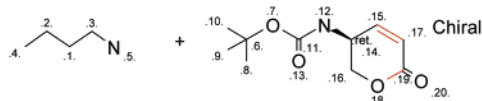
In this example RSS2 matches C.1.-C.2. in the query reactant to C.6.-C.13. in the target reactant and matches C.1.-C.2. in the query product to C.2.-C.6. in the target product. The AAM between query and target is highlighted in blue in Figure 12. This example shows the incorrect matching between the query and the target by RSS2.

As a special case, consider the following example. RSS2 matches the query with the target shown in Figure 13. However, RSS3 does not match this query-target pair, taking into account all overlapping mappings in the substructure search. This is because in the query product, two N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> groups came from the same reactant. Thus, their corresponding atoms have the same atom-atom mapping values. However, in the target product, two N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> groups came from two different reactants. Thus, the atoms of these two N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> groups have different atom-atom mapping values. For example, two N atoms have the atom-atom mapping value of 5 and 10, respectively.

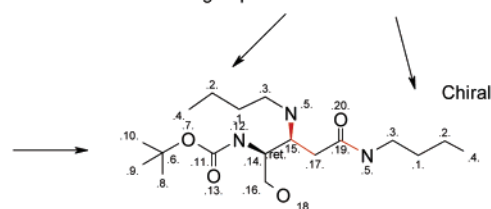
However, a simple change will allow RSS3 to match a query similar to the query in Figure 13 with target reactions similar to the target in Figure 13. One can do any of the following:

(1) Remove the atom-atom mapping values from either N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> group of the query product, as shown in queries 1 and 2 in Figure 14.

Query:



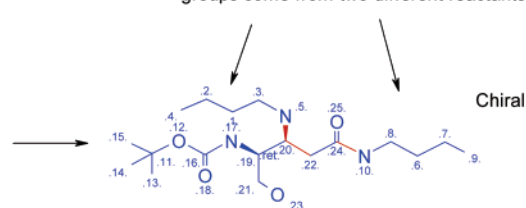
These two N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> groups come from the same reactant



Target:



These two N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> groups come from two different reactants



**Figure 13.** RSS2 matches the query with the target shown in this figure. RSS3 does not match them.

(2) Remove the atom-atom mapping values from all the N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> groups of both reactant and product of the query, as shown in the query 3 in Figure 14.

Now let us consider the second special case. RSS3 does not match the query with the target in Figure 15. This is because in the reactant and product of the query, the oxygen atoms in O-C have atom-atom mappings. However the corresponding target oxygen atoms do not have atom-atom mappings.

However, there is a simple way to allow RSS3 to match a query similar to the query in Figure 15 with target reactions similar to the target in Figure 15. One can do either of the following:

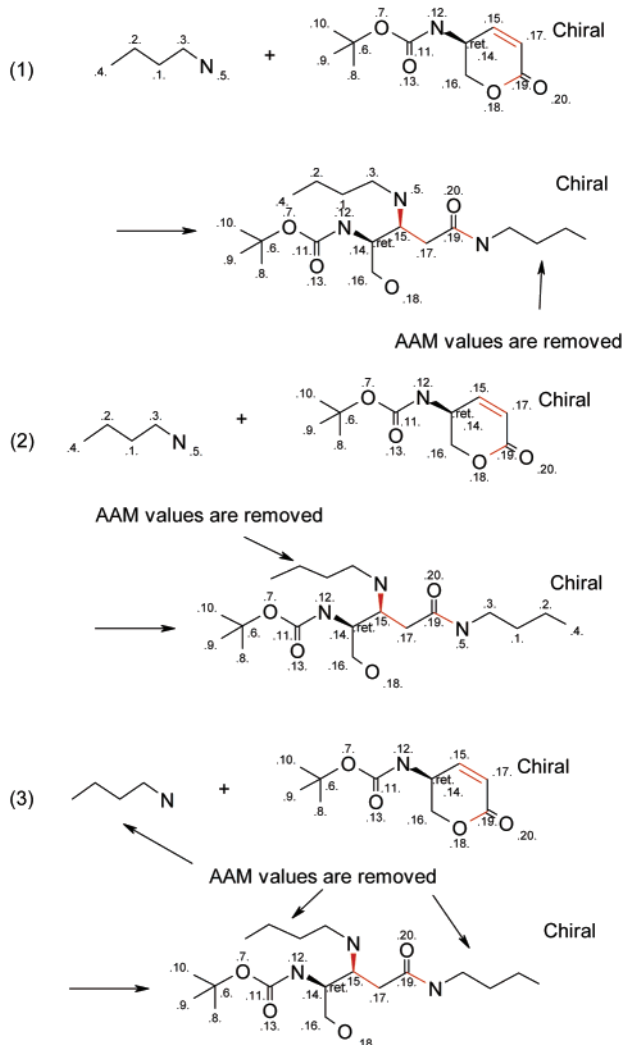
(1) Remove the AAM values from two oxygen atoms in O-C of the query, as shown in the query 1 in Figure 16.

(2) Remove the AAM values from the two oxygen atoms in O-C of the query, and also remove the reacting bond marks attached to them, as shown in the query 2 in Figure 16.

RSS3 matches both queries 1 and 2 of Figure 16 with the target in Figure 15. It should be noted that the queries such as query 2 of Figure 16 can be matched to those target reactions in which their C-O bonds in O=C-O and/or O=C-O-C=O can be either reacting bonds or unchanged bonds. (This problem is most common in ether and ester reactions and is discussed further in section 7.5.)

**7.4. Correct Matching of Queries with Explicit Hydrogens.** RSS3 is more rigorous in its handling of explicit hydrogens than was RSS2. For example, as mentioned

Query:



**Figure 14.** Three queries in this figure are identical with the query in Figure 13 except that the AAM values have been removed from the first N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> group of the product (query 1), or from the second N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> group of the product (query 2), or from all three N-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>2</sub>-CH<sub>3</sub> groups of both reactant and product (query 3). RSS3 matches any of these queries with the target in Figure 13.

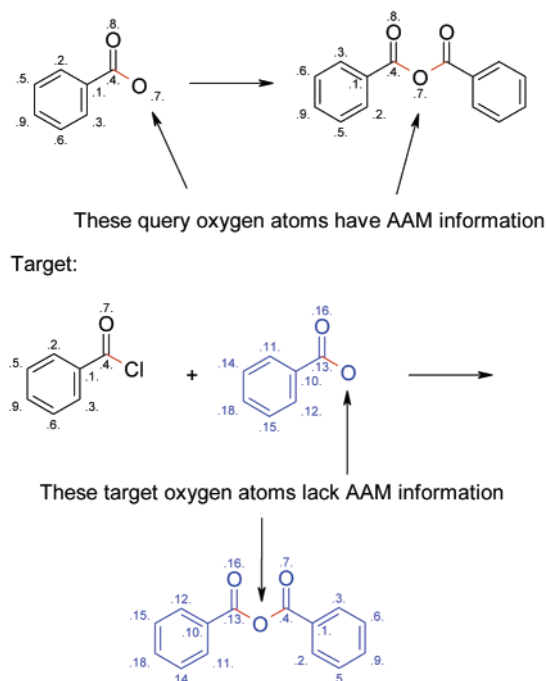
previously (section 7.2), RSS2 matches the query shown in Figure 12 to 30 hits from the Theilheimer database. However, only four hits are correct among the 30 hits. Two of them are shown in Figure 17: Target 1 is a correct hit, but target 2 is a false hit. This is because C.1.-H and C.2.-H bonds of target 2 in Figure 17 are not reacting bonds, while the query product's C.1.-H and C.2.-H bonds are reacting bonds (see Figure 12).

RSS3 matched the query in Figure 12 with only the four correct hits, such as target 1 of Figure 17, but did not match this query with any other incorrect hits returned by RSS2, such as target 2 of Figure 17.

### 7.5. Reactions Related to Ether and Ester Structures.

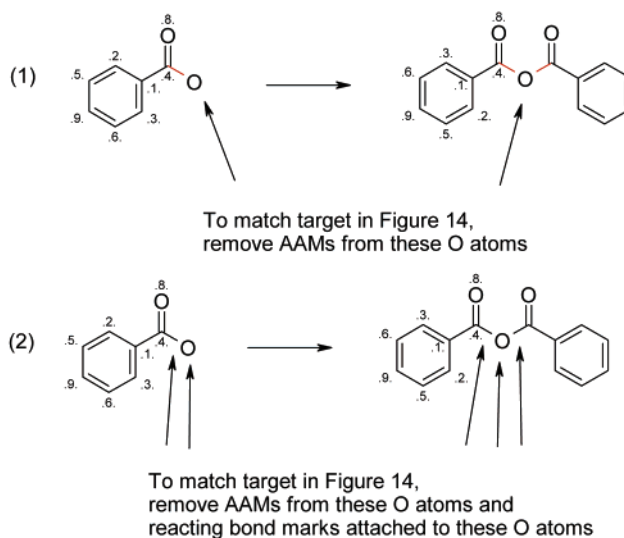
In some organic reactions, it is almost impossible to determine from which reactant an oxygen atom came. Two types of such reactions are the ether-related (C-O-C) reactions and ester-related (O=C-O) reactions. RSS3 returns correct reactions that are related to these structures. However, RSS2 did not always return the correct hits. For example,

Query:



**Figure 15.** RSS3 does not match the query with the target in this figure. This is because in the reactant and product of the query, the O atoms in O-C have atom-atom mappings. However the corresponding target O atoms do not have atom-atom mappings.

Query:



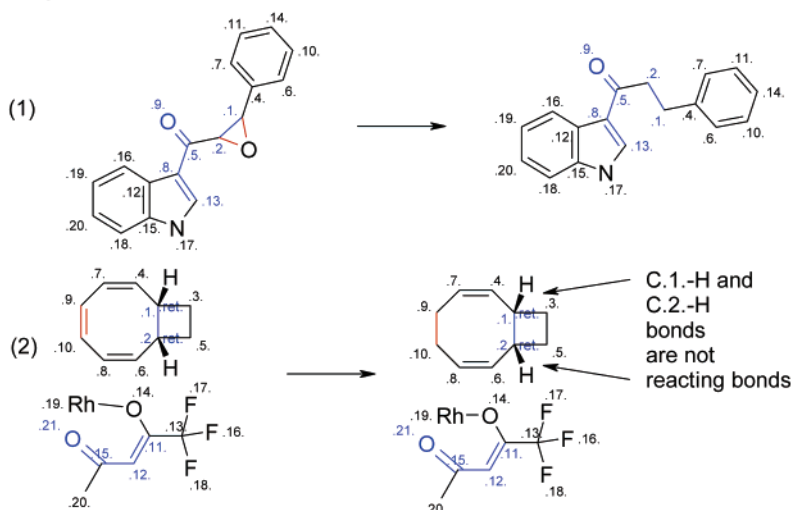
**Figure 16.** Query 1 is identical with the query in Figure 15 but without the AAM values from two oxygen atoms in O-C of the query. Query 2 is also like the query in Figure 15 but without the AAM values from two oxygen atoms in O-C of the query and without the reacting bond marks attached to them. RSS3 matches both queries 1 and 2 with the target in Figure 15.

RSS2 matches the query with the target shown in Figure 18. However, RSS3 does not match this query with the target. This is because the C.2.ext. atom of the query reactant requires that there must be no reacting bonds attached to it, but the atoms C.12., C.14., C.15., C.31., and C.32. of the target reactants all have a reacting bond attached to them.

As another example, RSS2 retrieved 81 hits from twelve MDL reaction databases<sup>8</sup> by using query 1 shown in Figure



Target:



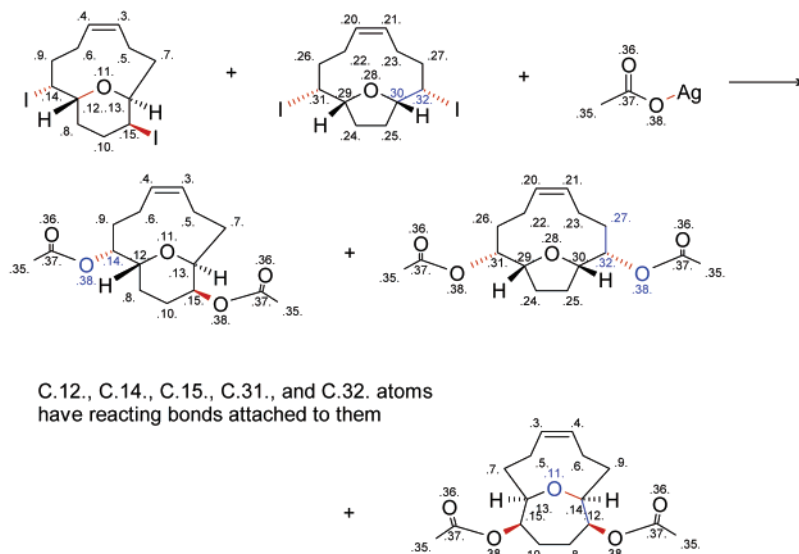
**Figure 17.** RSS3 matches the query in Figure 12 with target 1 but not target 2. RSS2 incorrectly matches the query in Figure 12 with target 2. The blue highlight shows the query in Figure 12.

Query:



C.2. doesn't have reacting bonds attached to it

Target:



C.12., C.14., C.15., C.31., and C.32. atoms have reacting bonds attached to them

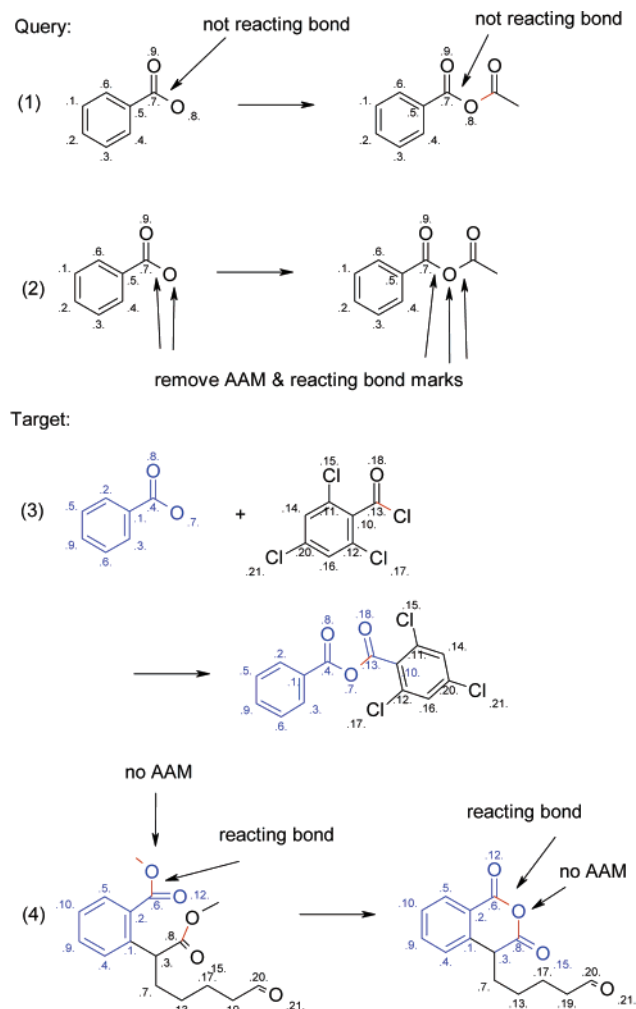
**Figure 18.** RSS2 wrongly matches the query with the target shown in this figure. RSS3 does not match this query with the target.

19. RSS3 matched this query to only 41 reactions among the 81 hits. One reaction that was hit by both RSS2 and RSS3 is shown in Figure 19 as target 3. The other reaction that was hit by RSS2 but not by RSS3 is shown in Figure 19 as target 4. This is because O—C bond in O—C=O of the query reactant is not a reacting bond, but in target 4, both O—C bonds in two O—C=O are reacting bonds.

In order for RSS3 to match query 1 with both the targets in Figure 19 and all other 81 hits returned by RSS2, query 1 in Figure 19 must be modified by changing all C—O bonds into unmarked ones, and also remove the atom—atom mapping value 8 from O—C.7., as shown in the query 2 in Figure 19.

The above example shows that RSS3 allows the user to search and retrieve more specific reactions by simply tailoring the RSS queries, if needed.

**7.6. Detection of Reacting Center Errors.** In the MDL products reacting centers of chemical reactions are represented using a bond property called “reacting center status”. This property can have the following values: 0 = unmarked, 1 = a reacting center, -1 = not a center, 2 = no change, 4 = bond made/broken, and 8 = bond order changes. Combinations of these basic property values, such as 5 = (4 + 1), 9 = (8 + 1), and 13 = (12 + 1), are also possible.<sup>3,9</sup> The latter cases occur when a reaction involves two or more alternative products. A bond of a reactant may be broken to



**Figure 19.** RSS2 matches query 1 with both targets 3 and 4. But RSS3 matches query 1 with only target 3, not target 4. After removing the AAM value 8 from oxygen atoms in O-C.7. and reacting bond marks of C.7.-O.8. of query 1, as shown in query 2, RSS3 matches query 2 with both targets 3 and 4.

form the first product, thus this bond should be marked as “bond broken” and should be assigned a value of 4. On the other hand, the same bond of the same reactant may be unchanged (but other bond(s) may be changed) to form the second product, thus this bond should be assigned a reaction center value of 2 (unchanged). The code used to represent these two different situations is 6 (= 4 + 2).

In some cases, RSS3 cannot match a reaction to itself because it contains errors in some reacting centers. Although RSS3 cannot guarantee that it will not match a reaction to itself if it has any reacting center errors, RSS3 can still detect some of the reactions that contain such errors. Therefore, if one finds that RSS3 cannot match a reaction to itself, it is

most likely because the reaction contains reacting center errors.

On the other hand, RSS2 is not able to detect such errors. RSS2 does not strictly check the relationships between the query/target atom-atom matchings and the atom-atom mappings within the query and within the target.

For example, RSS3 does not match the reaction shown in Figure 20 to itself. This is because it contains three reacting center status errors. The reacting center status values of 4 (bond made/broken) and 8 (bond order change) for the reactant bonds C.14. = O.5. and C.14.-C.12. and 2 (no change) for the product bond C.12.-C.14., respectively, are incorrect. The correct values should be 6 (=4 + 2, which means bond made/broken or unchanged), 10 (=8 + 2, which means bond order change or unchanged), and 4 (bond made/broken), respectively (see Figure 20). If all these errors are corrected, then RSS3 will match the reaction to itself.

### 7.7. Queries with Conflicting Reacting Center Settings.

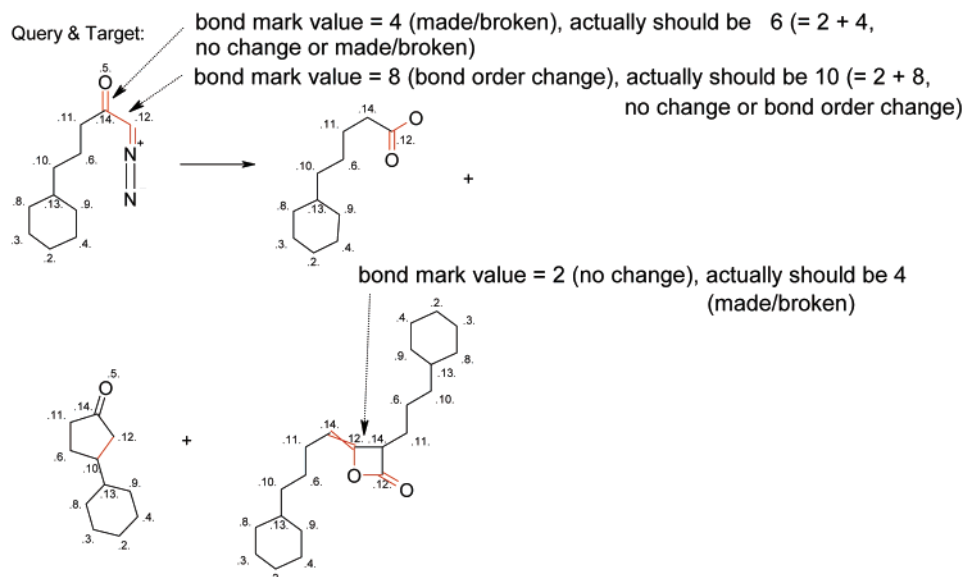
The RSS query shown in Figure 21 has a conflicting reacting center setting. This is because, on one hand, the drawing shows that a CC triple bond must be changed into a double bond or an aromatic bond in a ring (marked using the flag ‘Rn’), on the other hand, both the reactant and product bonds are marked as “unchanged”, which is indicated using the bond status value of 2 in the corresponding rxnfile<sup>3,9</sup> (which is not shown in Figure 21). With the query of Figure 21, RSS2 retrieved 936 hits from a test database containing 50 000 reactions sampled from CIRX reaction databases. Most of these hits are, however, unrelated reactions. On the other hand, RSS3 did not retrieve any hits for this invalid RSS query! If both CC bonds of the query are changed into “bond order change”, then RSS3 finds 862 hits from the same database. All of these are correct hits.

**7.8. RSS Conventions.** There are three special RSS conventions used in RSS2 (see Table 5). The first convention is that if an RSS query contains both reaction center and AAM information, while a target reaction contains reaction center data and partial AAM information, then RSS2 will take reaction center data into account but check only part of the query’s AAM information which correspond to the target reaction’s AAM and ignore query’s remaining AAM data. In contrast, RSS3 will perform a full check of all query’s reaction center and AAM data.

The second convention states that if an RSS query contains both reaction center and AAM information, while a target reaction contains only reaction center data and does not contain AAM information, then RSS will take reaction center data into account but ignore the query’s AAM information. Similarly, the third convention states that if an RSS query contains only AAM but no reaction center information, while a target reaction lacks AAM (irrespective of reaction center

**Table 5.** RSS Conventions

convention	query	target	RSS2	RSS3
1	has reacting center has AAM	has reacting center has partial AAM	checks reacting center, uses partial query AAM	checks reacting center, checks all query AAM
2	has reacting center has AAM	has reacting center no AAM	checks reacting center, ignores query AAM	checks reacting center, ignores query AAM (can turn on AAM checking)
3	no reacting center has AAM	(irrespective of reacting center) no AAM	ignores target reacting center (if exists), ignores query AAM	ignores target reacting center (if exists) ignores query AAM (can turn on AAM checking)



**Figure 20.** RSS3 does not match this reaction to itself because it contains three reacting center status errors.

Query:



**Figure 21.** This RSS query has a conflicting reacting center setting, because the drawing shows that a CC triple bond must be changed into a double bond or an aromatic bond in a ring (marked using the flag 'Rn'). However, at the same time both the reactant and product bonds are marked as "unchanged", which is indicated using the bond status value of 2 in the corresponding rxnfile but which is not shown in this drawing.

information), the RSS matching will ignore both reaction center and AAM information.

RSS2 uses the above conventions to retrieve more hits from reaction databases for a given query. For example, RSS2 matches the query with the target reaction shown in Figures 22 and 23, respectively.

For backward compatibility, the second and third conventions have also been implemented in RSS3. However, in RSS3 these conventions can easily be turned off by setting the option *checkAllQueryFeatures* to TRUE, forcing RSS to check all query features. In this case, RSS3 will not match both the query-target pairs shown in Figures 22 and 23. It should be pointed out that this new feature is not available in the current release of MDL Relational Chemistry Server (Version 2.0).

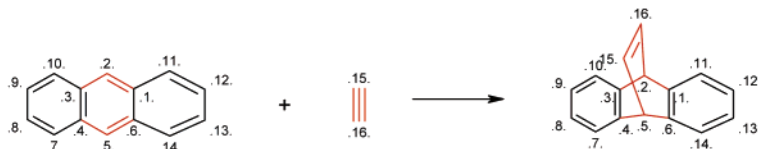
## 8. CONCLUSION

The reaction substructure searching method has become a key functionality of any reaction database management and

retrieval system. The first generation of the RSS algorithm used in REACCS mainly consisted of sequential calls to a molecule substructure search algorithm and often led to too many unrelated hits. The performance of the RSS was improved by implementing an existing molecular key screening technique. The second generation of the RSS algorithm implemented in both the REACCS and MDL ISIS systems is a significantly improved version of the previous RSS algorithm that uses AAM information. The introduction of the AAM information into the RSS algorithm not only allows it to reduce the number of unrelated hits but also makes it possible to implement several new features, such as Exact-Change searching. The performance of RSS searching was enhanced by both the molecular key and reaction key screening techniques. However, the use of AAM information was not optimized in this algorithm. Furthermore, because of inherent limitations in the implementation of this algorithm, it is not unusual that false hits are sometimes returned by this RSS method for a variety of RSS queries, especially those involving explicit hydrogen atoms and ExactChange flag.

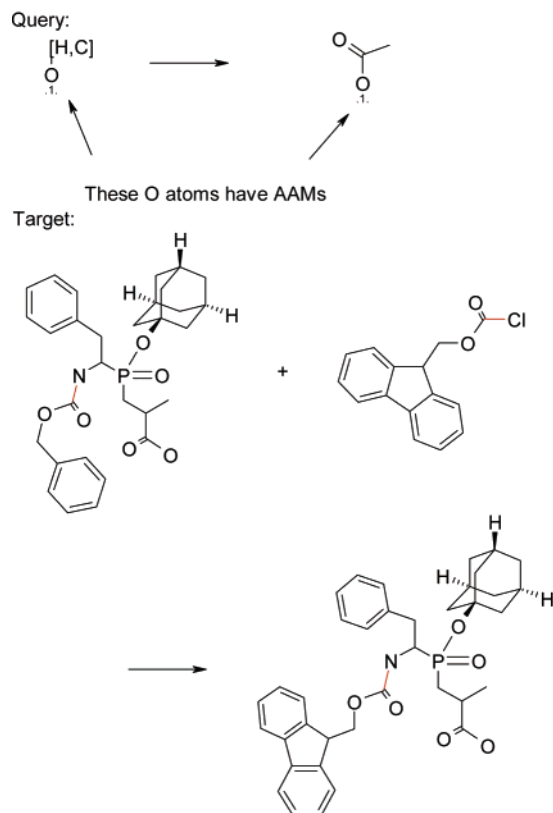
The third generation of the RSS algorithm has been implemented in MDL Relational Chemistry Server, a new MDL product based on Oracle data cartridge technology. This new algorithm is based on a depth-first search approach and is able to fully and prospectively use reaction specific information, such as reacting center and AAM information. It allows the user to precisely find reaction instances in

Query:



**Figure 22.** The query and the target are identical except that the query has AAM but target does not.





**Figure 23.** The query has AAM but no reacting center information. The target has reacting center but no AAM. (Note: '[H,C]' is an atom list, implying either hydrogen or a carbon atom.)

databases while minimizing unrelated hits. The performance of RSS3 is enhanced by the screening technique based on the combination of the molecular keys and an improved version of the reaction keys. The detailed discussion on the key screening methods and the evolution of MDL keys will be given elsewhere.<sup>10,11</sup>

The following list summarizes the changes in RSS3 that can cause differences in search results between RSS3 and RSS2. RSS3:

- (1) Requires that all reacting bonds for atoms with ExactChange flag must be explicitly drawn and marked, including bonds to hydrogens.
- (2) Guarantees correct atom-atom mapping relationships between query and target.
- (3) Guarantees correct matching of queries with explicit hydrogens.
- (4) More flexibly handle queries involving ether and ester reactions.
- (5) Might not match a reaction to itself if the reaction contains errors in reacting center status values. Therefore, RSS can also be a useful tool for the detection of database errors.
- (6) Does not match a query to any target reactions if the query contains conflicting reaction center specifications.
- (7) Supports both conventions 2 and 3 but allows them to be turned off.
- (8) The RSS2 routines can only deal with reactions with a maximal number of eight component molecules (reactants

plus products) and each molecule can consist of up to 255 atoms. The RSS3 program does not have these limitations.

In summary, RSS3 is superior to RSS2 in several aspects: it supports several important new features, improves error-checking, gives fewer extraneous hits, and yields more accurate search results.

#### ACKNOWLEDGMENT

We would like to thank Dr. B. Snyder and Dr. T. Wright for offering the updated information about MDL reaction databases as well as Dr. J. Durant and Dr. D. Henry for their valuable comments on the manuscript.

#### REFERENCES AND NOTES

- (1) Wipke, W. T.; Dill, J. D.; Peacock, S.; Hounshell, D. *Search and Retrieval Using an Automated Molecular Access System*; Presented at the 182nd National Meeting of the American Chemical Society, New York, August, 1981.
- (2) Molecular Connection (The MDL Newsletter for Communicating with Customers), The 20th Anniversary Issue, Vol. 17, No. 1, January 1998; p 19.
- (3) The most updated information can be found at MDL's Web site: <http://www.mdl.com>.
- (4) Chen, L. Substructure and Maximal Common Substructure Searching. In *Computational Medicinal Chemistry and Drug Discovery*; Tollenaere, J., Bultinck, P., Winter, H. D., Langenaeker, W., Eds.; Marcel Dekker: New York, in press.
- (5) Moock, T. E.; Nourse, J. G.; Grier, D.; Hounshell, W. D. The Implementation of AAM and Related Reaction Features in the Reaction Access System (REACCS). In *Chemical Structures*; Warr, W. A., Ed.; Springer-Verlag: Berlin, 1988; pp 303-313.
- (6) At the time of this writing (December 2001), there are a total of 1 011 178 reactions in all the public databases available from MDL (excluding MDL Beilstein Database):

database	no. of reactions	database	no. of reactions
JSM	72 244	ORGSYN	5 487
CHC	42 376	SPORE	13 253
REFLIB	171 111	CIRX (1991-2001)	877 818
		total	1 011 178

Note 1: JSM: Journal of Synthetic Methodology; CHC: MDL Comprehensive Heterocyclic Chemistry; REFLIB: MDL Reference Library of Synthetic Methodology; ORGSYN: Organic Syntheses; SPORE: Solid-Phase Organic Reactions; CIRX: ChemInform Reaction Library. Note 2: Theilheimer (containing 46 000 reactions) is a subset of REFLIB, and CSM (containing 68 482 reactions) is a subset of CIRX. Therefore, they are not counted in the above calculation. Note 3: The MDL Beilstein Database alone contains over 9 million reactions (See MDL's CrossFire Gmelin Database Press Release - 12/19/01).

- (7) Golomb, S. W.; Baumert, L. D. Backtrack Programming, *J. Assoc. Comput. Mach.* **1965**, *12*, 516-524.
- (8) Twelve MDL reaction databases are as follows: seven ChemInform Reaction Library databases (CIRX98, CIRX97, CIRX96, CIRX95, CIRX94, CIRX93, CIRX92), Journal of Synthetic Methods (RX-JSM), MDL Reference Library of Synthetic Methodology, Theilheimer, Organic Syntheses, and MDL Comprehensive Heterocyclic Chemistry.
- (9) 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*, 244-255.
- (10) Durant, J. L.; Leland, B. A.; Henry, D. R.; Nourse, J. G. Reoptimization of MDL Keys for Use in Drug Discovery. *J. Chem. Inf. Comput. Sci.*, in press.
- (11) Chen, L.; et al. Over 20 Years of Chemical Structure Access Systems from MDL: Evolution of MDL Keys and Their Applications. Manuscript in preparation.

CI020023S