

□ データベース関連

① データベースの基本

② データベースの内容／形態による分類

・文献データベース／ファクトデータベース

これらの呼称はデータベース内のデータの内容に従ったものである。データベースの内容が文献に関するものである時、文献データベースと称される。また、データベース中に保存されているデータの種類の種類がファクトデータ（物性データ、薬理データ、スペクトルデータ等）である時、このデータベースはファクトデータベースと呼ばれる。化学分野ではこのファクトデータベースが多数存在する。

以下に化学やバイオテクノロジー関連分野で利用されているデータベースの内容を列挙する。

文献／特許／数値データ（種々物性）／スペクトル／薬理データ／X線結晶解析／プロテイン／DNA／その他

・パブリックデータベース／インハウスデータベース

データベース自体が一般に公開されているのか、されていないかで2種類のデータベースに分類される。一般に公開されるものをパブリックデータベース、公開されずに企業内、あるいはグループ内で専用利用されるものをインハウスデータベースという。

③ データベースモデルの種類と概要

データベースはそのデータを管理する方法で様々なモデルがあり、この管理（データ間のネットワーク化）方法の違いにより幾つかのモデルに分類する事が可能である。実用化されている典型的なモデルとしては3種類ある。

- 階層型モデル
- ネットワーク型モデル
- リレーショナル（関係）型モデル

これら3種類のモデルのうち階層型及びネットワーク型モデルはデータを互いにネットワーク化するもので歴史が古く、実用システムとして大規模なものが構築されている。最後のリレーショナルモデルの歴史は3モデルの中では最も新しいものである。このモデルは化学（自然科学一般ともいえる）分野では文献検索等を除けば、物性／毒性／スペクトル等様々なデータを扱うデータベースとして最も利用されているモデルである。このリレーショナルモデルの基本概念（データ構造）は、科学者が日常的に用いているデータ整理手法とさほどかわりがなく、科学者には受け入れやすいモデルである。

以下前記3種類のデータモデルについて順にのべる。

□ 階層型モデル及びネットワーク型モデル

階層型モデル及びネットワーク型モデルはいずれもデータ同志の連携をさせる為にネットワークを形成させるもので、2モデルの差はネットワークの形式の差による。

・階層型モデル

階層型モデルではデータの相互関係が階層的になり、且つデータ間の横のつながりが比較的弱いようなデータ構造を持つものに的したデータベースである。このモデルのデータ間のつながりの形は基本的に木（TREE）構造をなしている。この木構造にも様々

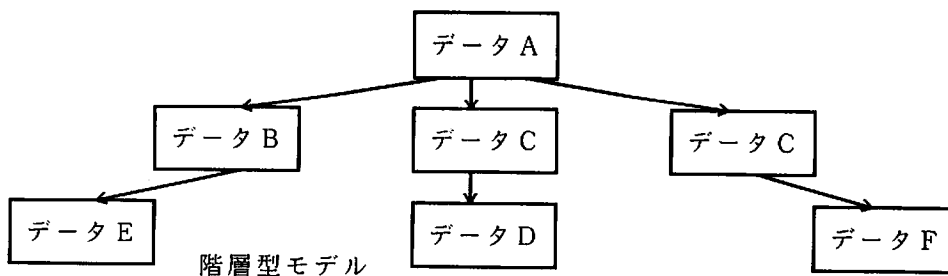


図 . 階層型モデルのデータ構造

なバリエーションが存在し、個々のデータベースに適した形態をとるようにする必要がある。

階層構造を持つデータは一般的な分野では多数存在する。最も典型的な例としては家系図がある。また、さまざまな構造体とその部品との関係もツリーとして構成することが可能である。

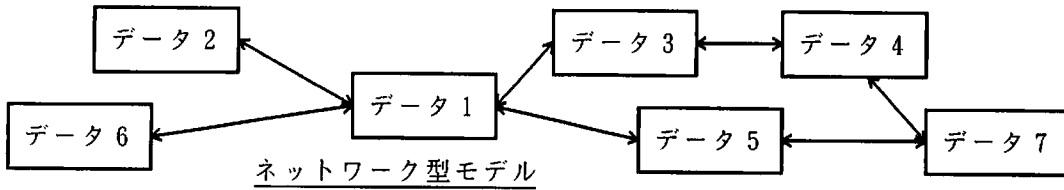


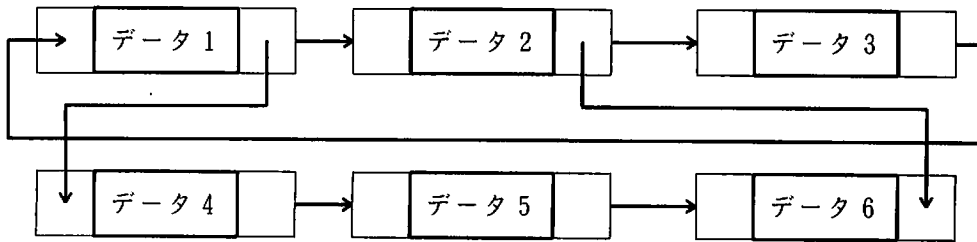
図 . ネットワーク型モデルのデータ構造

ネットワーク型モデルはより複雑な相互関係を持つデータを扱う時に利用されるモデルである。個々のデータ間に高度な階層関係や相互関係が存在し、これらが互いに交錯しているようなデータを扱う時に利用される。

□ 多重つなぎ編成によるデータベース構築

この多重つなぎ編成はネットワーク型モデルの一種である。このアプローチではデータ間のネットワーク関係（数学的にはリスト：LISTと呼ばれ、データベース関連では鎖：CHAINと呼ばれる）をそのままデータベース構造に反映させてデータベースを構築するものである。従ってこのデータベースは高速化を狙うよりも、データ間の結合関係を重視したデータ検索を目的としている。

データは本来のデータ自身に対する情報と、そのデータと関連するデータへの結合情報部分とから構成される。



*太線枠内は保存用データ、細線枠内が結合情報を持っている部分

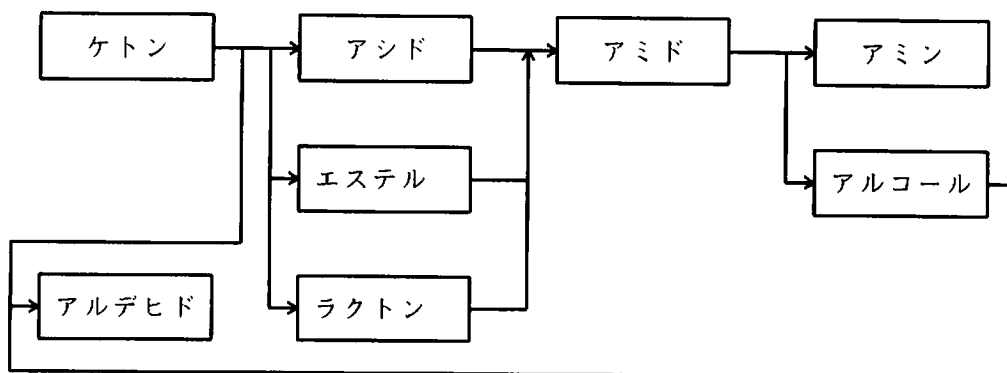
図 . 多重つなぎ編成によるデータ連携関係例

図からもわかるがこの多重つなぎ編成によるデータベースではデータ間の“リンク情報”が情報として大きな部分を占め、この情報がデータ間の複雑な相互関係の反映に役立っている。従って、データ間に複雑な階層関係や結合関係がある時にはこの方法がとられる事が多い。この多重つなぎ編成の特徴を以下にまとめる。

データが相互にリンクされているので、この結合関係を用いて連続的なデータの呼出が可能である
複雑な関係を有する情報も取扱可能である

このアプローチの欠点としては、結合関係を重視している為にデータの独立性が相対的に低くなる事。また、データの結合関係が複雑になるとデータベースの操作やメンテナンスが複雑になる事などである。

例) ケトンを基点とした官能基群を考慮し、化合物データベースを構築する



□ リレーショナルモデル (関係モデル: RELATIONAL MODEL)

リレーショナルモデルは自然科学分野で頻繁に採用されている。これは、リレーショナルモデルで扱うデータ形式が自然化学分野におけるデータ構造の表現に都合のよい形式であることが大きな理由である。すなわち、リレーショナルモデルは1サンプル対多項目というマトリクスを基本としたデータ構造を持つものを扱うのに適している。例えば一つの化合物(サンプル)に関する様々な物性(項目)といった関係を持つデータであり、このような関係を持つデータは化学に限らず自然化学分野では頻繁に現れる。

このようなデータは、先にのべた階層型やネットワーク型のモデルにしてデータベース化することは困難である。むしろ研究者が昔から行っていたデータの管理形式、すなわちサンプルを行とし、項目を列としたマトリクスで容易に管理する事が可能である。リレーショナルモデルはこのマトリクスを扱うためのデータベースといえる。

このリレーショナルモデルは以下に示す関係式で表すことが可能である。

$$R(X_1, X_2, \dots, X_n) = \{ (x_1, x_2, \dots, x_n) \mid x_i \in X_i \wedge \text{命題 } R(x_1, x_2, \dots, x_n) \text{ が真} \} \subseteq X_1 \times X_2 \times \dots \times X_n$$

X_i は関係の定義域 (DOMAIN)、定義域の数を関係の度数 (DEGREE)、 x_i は関係の要素、 (x_1, x_2, \dots, x_n) は n 組 (n -TUPLE) と呼ぶ。

第1マトリクス 関係名	定義域				
	1	2	3	4	5
実験者	実験者名	実験番号	実験日付	天気	実験結果
試薬	試薬名	購入日付	製造メーカ	残量	
化合物	化合物ID	製造者名	製造日付	機器分析	薬理データ
.....
第2マトリクス 化合物	1	2	3	4	5
	化合物ID	製造者名	製造日付	機器分析	薬理データ
化合物1	621A	KY	3.01.22	IR, NMR	11C
化合物2	372X	KK	3.02.15	IR, NMR	11C
化合物3	277C	TO	3.02.05	IR, MASS	12C
.....

前例ではマトリクスが2階層になって管理されている。第1マトリクスは種々のデータの関係に関する情報を一つのマトリクスとした管理用(またはインデックス用)のマトリクスである。第2マトリクスは第1マトリクスにより定義された個々の関係毎に設けられ、これら関係に関する実際のデータが入っているマトリクスである。

スプレッドシート型モデル

最近パソコンやワークステーション上で頻繁に利用されるデータ管理形式として、スプレッドシートというものが利用されている。このスプレッドシートは名前のおり、概念的にはデータが書かれた一枚のシートをあるまとまったデータ単位として扱うものである。前記リレーショナルデータベースを基本として考えるならば、一つのスプレッドシートは一つのマトリクスに該当することになる。

このマトリクス1個を一つのシートとみなし、一枚あるいは複数枚のシートを用いて計算機上で利用される。このスプレッドシートはLOTUSやEXCEL等の表管理プログラムの基本となっている。

化学の分野でもこのスプレッドシートが利用されはじめている。他の分野のスプレッドシートとの違いの大きな点は、第1列目の項目として化合物IDや化合物の構造式を用いることである。それ以外の列に個々の化合物に関する情報が配置されている。このような形式のスプレッドシートを特に“ケミカルスプレッドシート”と呼んでいる。ケミカルスプレッドシートの一例を表に示す。

表 . ケミカルスプレッドシート例

化合物構造式	融点	沸点	薬理データ	LOGP
化合物 1	11.4	120.8	134	4.3
化合物 2	-5.8	186.4	155	2.7
化合物 3	123.3	256.1	176	3.3
.....
.....

オブジェクト思考型モデル

④ 実際のデータベース構築時に利用される様々な技術

データベースの構築では検索対象となるデータの量がかなり大きなものとなる為、いくつかの点で注意が必要である。

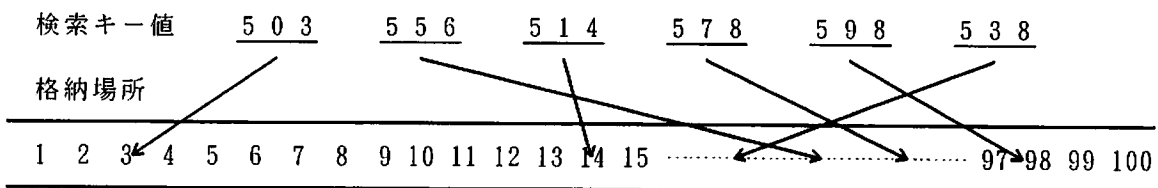
1. 検索スピードがデータベース中のデータ量と関係なく一定で速い事。
2. データの量が多い事から記憶領域の効率的な利用が必要。
3. データベース内のデータの特徴を反映する構造である事。

これらの条件を満たす為の技術としていくつかの技術が確立されている。ここではこのような技術の中で計算機上での一般的な技術について簡単にまとめる。

□ HASHING FUNCTION (散らし関数) による高速化

散らし関数により検索キーの数値データを変換し、その変換された値をデータベース上の格納場所、あるいは被検索体のID番号等にする手法である。

いま格納場所/IDを示す値が1~100の整数で示されているとして、被検索体の検索キーの値が500~600の値であるならば、この値を100に近い素数97で割り、残った余りの値をデータが格納されている場所の番号やIDに対応させることで直接被検体を捜し出す事が出来る。(割り算法)



つまり全く関係の無い2つの数値をハッシング関数を用いる事で関係付け、1対1の高速検索を実現させるものである。この関係は検索用の数値データをx、被検索体が収容されている場所を示す値をn、ハッシング関数をHとすると以下の式で示される。

$$n = H(x) \quad \text{—————} \quad ()$$

この検索の特徴は検索キーの値による1体1の高速検索が可能となる事である。このハッシング(散らし)により、被検体が格納場所上に均等に分散されているならば、データベースの大きさに関係無く常に一定の速さで検索する事が可能となる。

ハッシングの特徴を生かした高速検索を保証する為にはこのハッシング(散らし)がうまく行われるか否か(データ格納アドレスがメモリー上で均等に割り振られるか否か)に依存する。一般的にこのハッシングをうまく行う為には以下に示す2つの要素を検討する事が必要である。

① 実際のデータ(被検索体)の数よりも大きな格納場所を用意する事が必要

データベースの大きさが大きくなるとハッシング関数を用いても同じ番号に被検索体が複数アサインされる事が生じてくる。この同義語(SYNONYM)の発生が多くなると検索効率が著しくダウンする。この同義語の発生を少なくするにはハッシング関数を再考する事と、被検体の格納領域を大きくする事が考えられる。この格納領域は計算機の主記憶の容量や記憶領域の有効利用等からベストと思われる大きさを判断する事が必要である。(一般的にはデータ数の2~5割増しといわれている)

② 数値データの内容に応じた適切な散らし関数を採用する事が必要

散らし関数の良否が検索効率を大きく左右する。散らしの目的は被検索体を格納場所に均一に割り当てる事にある。この散らしが均一に行われないと、一つの場所に複数の被検索体がアサインされたり、かたよってアサインされる事になり検索スピードのダウンやメモリー利用効率の低下等の問題が発生してくる。

この散らし関数もさまざまなものがある。最初の例に示した割り算法の他、混合法(キー値x定数A/定数Bの余りを用いる)、中央二乗法(キー値を二乗し、中央部分のある一定の桁数を取り出し、その値を用いる)等目的に合わせて選択する事が必要である。

例え理想的な設計をしたとしても同義語の発生を0とする事は困難である。この同義

語が発生した時の回避方法として幾つか実用化されている手法がある。

- 一つの番地に複数の被検体の番号が格納できるようにする
 - 同義語が発生した番地の前後を開放し、そこに同義語の被検体を収納する
 - 同義語だけをまとめて収納する場所を予め設定しておく
- (このような同義語の格納場所はバケツ (BUCKET) と呼ばれる)

これらのアプローチのどれを採用するかは、メモリーの容量や検索スピードの問題等の関係で、個々のシステム単位で最適のものを採用する事が必要である。

利用例)

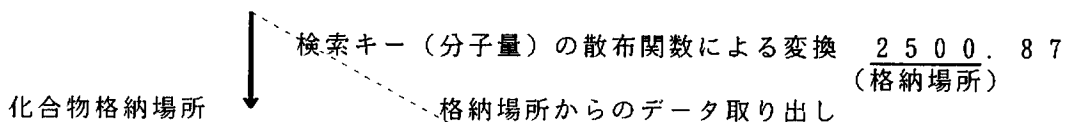
化合物の原子量をキー値として化合物 (4000個) の検索をする。
 化合物の分子量Mは200~600の間の値を取り、化合物の格納場所Dは1~5000 (25%増) とする。この時用いる散らし関数としては以下のようなものを用いる。

$$D = \text{int} \left[\left(\left[\frac{M}{5001} \right] - 0.03997 \right) \times \frac{5}{8} \times 10^5 \right] \text{ --- ()}$$

上式でintは切り捨てを意味し、基本的な手法は割り算を採用するとする。従って、分子量を格納場所の数5000に最も近い素数5001で割り、値の基準を合わせる為にその値から0.03998を引く。この値を格納場所の1~5000にあわせる為に500000/8をかけ、この値を切り捨てることで格納場所のID番号を得る。

実際の手続きとしては効率良い検索を実現する為にもデータの格納場所は、被検索体4000の2割から5割程度 (ここでは格納場所として5000を用いた) 大きめにするのが良い。さらに、散らし関数を用いても同義語 (同一アドレス) をさける事は困難であること。また、今回用いた分子量は分子式が同じならば異なった化合物であっても同じ分子量となり、全く同一の検索キーを持ち同義語の発生原因となる。これらの点からも同義語が少なからず発生する事は明白である。このような同義語を取扱う為、同義語の格納領域を予め設定しておくことが望ましい。

化合物 : 検索キー値 400 (分子量)



1 2 3 4 5101 102 1032500 2501 2502..... 3001 3002 3003 4999 5000

同義語格納場所 同義語存在の確認

1<2 3 4 5 98 99 100

図 . ハッシング関数を用いた化合物の保存形式

- インバーテッドファイル (INVERTED FILE : (転置ファイル)) の作成による高速化

文字どおりファイルの中身 (行列) を反転させる手法を意味している。形式上は行列を反転するのであるが、各行あるいは列単位にある特定の条件を満足するデータ群を取り出して新たなファイルをつくる。このファイルの利用目的は検索時のインデックスに利用するというものである。行を単位にするのか列を単位として展開するかは個々の目的に応じて採用される。

一般的にデータはマトリクス (行列) の形式で表現される。検索という観点から考える時、行の要素としては検索目的となる項目を取り、列の要素としては検索を行うに必要なデータ項目を指定しておくことが多い。

被検体 I D	検 索 項 目						
	1	2	3	4	n - 1	n
1	0	0	1	1	0	1
2	1	0	1	0	1	1
3	0	1	0	1	1	0
.....
N	1	0	1	0	0	0



インバーテッドファイル化（データ中1のものだけ抜き出してくる）

検索項目	被 検 索 I D						
	1	2	3	4	D - 1	D
1	2	5	10	12
2	3	4	7	18
3	1	2	8	11
.....
n	1	2	6	21

* N : 被検索対象サンプル数 (行)

n : 検索項目数 (列)

D : インバーテッドファイルで新たに取り出されたサンプルの数

総てのサンプルNがインバーテッドファイル中に採用された時 $N = D$ となる

通常は $N > D$ である

このようなマトリクスを利用する時、計算機のファイルにはそのままの形式でストアさせるだけでなく、個々の項目の高速検索を目的として特別のファイルを設ける事がある。これがインバーテッドファイル（転置ファイル）と呼ばれるものである。

このインバーテッドファイルは基本となるファイルの行と列とを入れ換え、新たに検索用のファイルとして作成されたものをいう。通常このファイルは列となっている項目（例えば融点/沸点/分子屈折率等）毎に設けられる。従って、項目数がn個の時、インバーテッドファイルの数は最大n個となる。

このインバーテッドファイルを用いない時、検索項目のパターンに従って検索するならば検索は総ての被検体サンプル(N)についてチェックを行う事が必要になり、被検体の数に比例して検索時間が増大してゆく。しかし予め検索項目を行に設定し、その検索項目を満たすパターンについてのID番号を列に設定したインバーテッドファイルを用意しておけば検索スピードが向上し、検索上での様々な操作(AND/OR検索等)を簡単に実行する事が可能となる等の利点が生じる。

例えばある一つの項目で検索を行いたい時、このインバーテッドファイルがなければこの検索項目に従って対象となる総ての被検体をチェックすることが必要である。計算機ではこのようなデータ検索はディスクへのアクセスを要求する事がおおく、且つこのディスクアクセスは計算機泣かせの最も時間のかかる作業である。データベースの設計ではこのディスクアクセスをいかにして少なくさせるかという事がデータベース成否の鍵をにぎっている。先のインバーテッドファイルに関するアクセスを考えた時、検索キーに従って被検体が分類わけされている為、検索項目が一つであるならばその項目に該当する行の部分に存在する被検体を取り出す事で検索は完了する。この検索項目が複数ある時はインバーテッドファイルが無い時と比べるとさらに効果が出てくる。

このインバーテッドファイル作成による検索時の効果をまとめると以下のようなになる。

検索スピードが速くなる

検索に要する時間は被検索体の数に影響されることが少なく、ほぼ一定である

複数項目を組み合わせた複合検索が簡単で、スピードも速い

注意すべき点としては本来のデータ格納領域の他にインバーテッドファイル用の格納領

域も必要となりメモリーを多く必要とすること、及びデータベースの更新が発生するとインバーテッドファイルの更新も必要になるという事である。従って、頻繁にデータの更新が発生するデータベースではメンテナンスという観点からはインバーテッドファイルの長所をいかし切れなくなる事がある。

*インバーテッドファイルは大きなデータベースとなると項目毎に設ける事も多い。項目毎にインバーテッドファイルを用意し、且つ項目の内容を細かに分類する事でより細かな検索を高速に達成する事が可能となる。

例) ケトンの有無に関するインバーテッドファイル

ケトンの有無だけの検索ではこのファイルの総ての被検索体を出力する事が必要であるが、このケトンをも更に細かく分類する(ケトン、アルデヒド、エステル、アミド、ラク톤等)事でより細かな情報をより高速で検索する事が可能となる。

⑤ 日本及び外国のデータベース

- ・化学文献データベース (CAS, DARC)
- ・化合物物性データベース (CIS)
- ・化合物スペクトルデータベース (Tool-IR, SDBS)
- ・X線結晶データベース
- ・バイオ/プロテイン関連データベース
- ・総合データベース (DIALOG等)

⑥ 日本の団体

- ・科情協
- ・JICST
- ・BIDEC

⑦ データベース関連テーブル

第3章 化合物の検索 (化合物データベース)

1. 化合物検索概論

化合物の検索技術は計算機化学における基本である。化合物の検索技術は計算機のハード関連技術の向上と共に変化して来ている。

化合物の検索では以下の手続きをとることが必要となる。

- ① 化合物構造式 (検索対象化合物) の計算機への入力
- ② データベース中の化合物から化合物を取り出してくる

この2段階の手続きの中に様々な技術が必要となるが、これらの技術はその基本を計算機におくのか化学におくかで2種類に分類できる。これら様々な技術を駆使して化合物の検索が行われることになる。

計算機に基本をおく技術としては例えば“インバーテッドファイルの利用”や“ハッシング (散らし) 関数の利用”等がある。その他、ファイルのストア形式として階層型データベースやリレーショナルデータベース等様々な形式のものがある。これらの技術に関しては他の計算機に関する著書に詳しいのでここでは取り上げない。化学に近い、あるいは化学と計算機を結びつける為の技術を中心としてのべる。

□ 検索技術の歴史的変遷

当初計算機のハードが貧弱であった時には、計算速度も遅く、しかもメモリの単価が高いため、このメモリを節約しながらデータベースを作成すること及び検索を高速に行うための技術が最優先技術として要求された。

化合物検索の初期では、化合物の構造式を1次元の文字及び数値からなる文字コードへと変換し、そのコードを検索する事が一般的に行われていた。化合物構造式の1次元文字列へのコード化により、メモリを節約した化合物データベースの実現や検索が可能となる。さらに、この文字コードを用いた化合物検索には計算機分野で一般的に利用されてきた文字列検索技術をそのまま適用することが可能となる。このため、化合物を対象とした新たな検索技術を開発する必要もないので初期の化合物検索システムにはよく利用されたアプローチである。

その後計算機のハードの向上と共にメモリー単価も低くなり、メモリーを十分に用いた化合物表現 (結合表等) 及び検索技術が幅広く展開されるようになった。またグラフィックディスプレイの充実とともに化合物構造式のコード化を人間が行う必要はなくなり、化合物構造式をディスプレイ上に直接書き込む事で、あとは計算機が自動的にコード化する事が一般的アプローチとしてとられるようになった。この結果、当初計算機と化学者との間に存在した大きなコミュニケーションギャップは大幅なMMI (マンマシンインタフェース) の向上とともに可能な限り小さなものとなりつつある。

以上の特徴を簡単に表1にまとめる。

表1. 化合物検索の推移

化合物検索	計算機	検索手法	長所○/欠点X
当初	演算速度遅い メモリー不足 貧弱なグラフィック ディスプレイ	数値/文字 データの検索	○ メモリーが少なくて済む ○ 検索アルゴリズムが簡単である X 構造式のコード化は学習が必要 X 構造式からのコード化は人が行う X 総ての化合物に対応すると、 コード化のルールが複雑になる
現在	演算速度速い 大きなメモリー 強力なグラフィック ディスプレイ	結合表を用 いた検索	○ どのような化合物も対応可能 ○ 研究者が親しんだ構造式の利用 ○ 構造式のコード化は計算機が行う X メモリーを多く必要とする (但し、現在の計算機では問題ない)

2. 化合物の検索様式

化合物の検索様式は大きく2種類に分類される。即ち、①完全一致検索と②部分構造検索とである。完全一致検索は検索キーとして用いた化合物と被検索化合物とで構造式が完全に一致した化合物のみを取り出してくるものである。部分構造検索とは化合物の部分構造を検索キーとし、被検索化合物中検索キーとして用いた部分構造をその構造式中に持つ化合物を取り出してくる検索手法である。当然、部分構造検索は検索部分構造と全く同じ化合物を取り出してくる時は完全一致検索となる。

利用される検索手法によるが、一般的には完全一致検索の方が部分構造検索よりも高速に検索出来る事が多い。

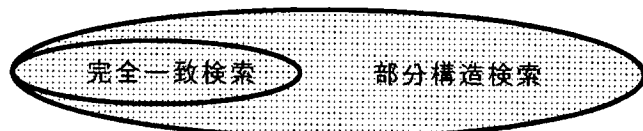


図1. 完全一致検索と部分構造検索との包含関係

① 完全一致検索

完全一致検索は化合物検索における最も基本的なものである。この完全一致検索にたいするアプローチは幾つか存在するが、これらの手法は大きく2種類に分類する事が可能である。一つは化合物を一旦数値/文字データによる線型データに変換し、このデータを用いて検索を行うものである。残る手法は化合物の結合表を直接用いて検索を行うもので、最近の検索はこの手法を取っている事が多い。ここでは、これらの2手法について順を追って説明する。

(1) 変換コードを用いた検索手法

変換コードを用いた検索は、化合物を一旦数値データや文字データに変換し、その数値/文字データを検索するものである。従って、一旦数値及び文字データに変換されればその後の検索は単なる数値/文字データの検索となり、通常用いられている検索技術の適用が可能となる。

一般的にはこのコード変換に工夫があり、線型表記の章でも述べたように一元一項対応が取られるように化合物はコード化されなければならない。ここでは化合物の検索を最優先させた時に用いられるコード化(線型表記)を中心として述べる。従って、説明される手法は化合物を一元一項対応で数値データに変換する事が主目的であり、一旦数値データに変換されたものを元の化合物に戻す事は考慮されていない(殆ど不可能)。

現在、このような目的で利用されている化合物のコード化手法として幾つかの手法が存在する。ここではMORGAN名とその拡張型であるSEMA名及び最近展開されてきたSMILESについて簡単に紹介する。SMILESは化合物を数値データでなく、文字データへと変換するものであり、先にのべた様々な線型表記手法に似ているが、コード化の為にルールが極めて簡単で元の構造式への変換及び結合表への変換も容易である等様々な特徴を有している手法である。

1. MORGAN NAME (立体情報を含まない化合物の時)

1. 1. MORGAN NAME 概要

MORGAN NAME (以下MORGAN名と略す)は化合物のユニークナンバリング*1(化合物に対応した1個だけの番号付け)を行う為のアルゴリズムとしてMORGANにより提唱された手法である。このMORGANアルゴリズムにより、化合物中の個々の原子について唯一種の番号付けがなされる。この番号に基づき、原子の種類に関する情報と結合の種類に関する情報とをまとめて一つの数値データとしたものがMORGAN名である。

MORGAN名とは化合物に一元一項対応で付けられた化合物名であるが、その名前から元の構造式を再現する事は困難である。先に述べたようにMORGAN名の目的はユ

ニークナンバリングと検索であり、もともと構造式の再現性を目的にはしていない。

MORGAN名 = ユニークナンバリング + 原子/結合情報

* 1) ユニークナンバリングについて：

ユニークナンバリングとは化合物を構成する原子につけられる番号を、1化合物1通り(一元一項対応)に決定する事を意味する。

通常、化合物検索を行う時は検索キー化合物と被検索化合物との原子対応をとる事が必要となる。人間は複数の化合物を比較する時、無意識のうちに互いに対応する原子を認識するが、計算機ではこのような事は不可能である。従って、計算機を用いて原子に番号をつける時には、同一化合物には全く同じ番号付けがされるようにしておく必要がある。この一元一項対応した番号付けをユニークナンバリングという。

このユニークナンバリングに関する問題は、化合物検索を行う時に重要な問題となってくる。つまり、完全一致検索を行う時に検索化合物と被検索化合物を構成する原子それぞれに原子番号が付けられるが、この原子番号が検索及び被検索化合物とで互いに対応がとれている事が必要となる。この原子番号の対応が取れていない時には、例え同じ化合物であっても異なる化合物と認識されてしまう事になる。(図2参照)

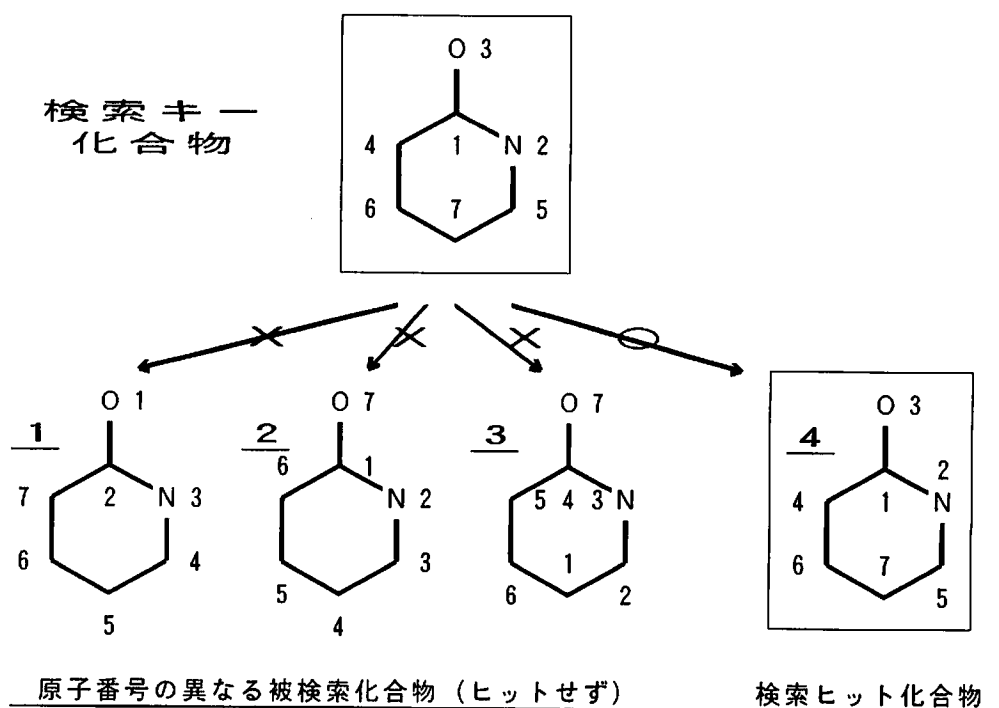


図2. 原子番号の差による化合物検索結果

化合物に付けられる原子番号が異なると、上図でも示されたように計算機の内部では全く同じ化合物であっても異なる化合物として認識される。従って、化合物の原子に付けられる番号はユニーク(一元一項対応)でなければならない。(図中検索キー化合物と4番の化合物が同じ) 但し、MORGANアルゴリズムによるユニーク番号付けでも限界がある。例えば対象性の高い化合物では、このアルゴリズムで一義的に番号を付ける事は不可能であり、複数の番号付けの可能性が発生するが検索上特に大きな問題ではない。

1. 2. ユニークネーミングの為のMORGANアルゴリズム

MORGAN名を発生するには以下に示される4段階の手続きを踏まねばならない。

1. EC (Extended Connectivity) の決定。
2. 各ノード (NODE) に対する番号付け。
3. 以下のリストの作成。
 - a. From Attachment List
 - b. Ring Closure List
 - c. Atom Type (Node Value) List
 - d. Bond Type (Line Value) List
 - e. Modification List
4. MORGAN名の作成。

以下にはWipkeらの報文を基本として具体的にMORGAN名の決定について述べる。用いた化合物は図1に示されている。

1. 3. MORGAN名決定事例

対象化合物構造式

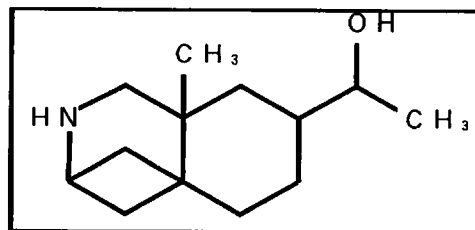


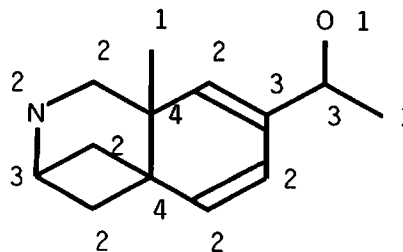
図1. MORGAN名算出サンプル化合物

手順1: EC (Extended Connectivity) の決定

① 初期EC値の決定。

EC値は化合物中の各原子をNODEとした時、各NODE (原子) が結合しているNODEの数 (次数) のことである。

* 水素は除く



② NECV値の決定。

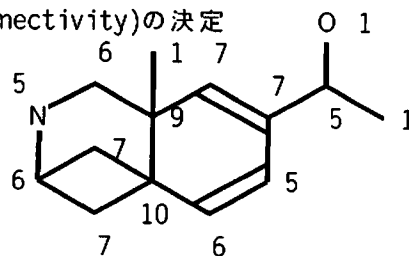
EC値の種類数をNECVとする。①ではEC値は1、2、3、4の4種類である。従って、①のNECV値は4となる。

$$\underline{NECV = 4}$$

③ 試行拡張分岐数 (TEC: Trial Extended Connectivity)の決定

各NODEに隣接しているNODEのEC値を加算し、各NODEに対して新たにTECを算出する。

*末端原子は常に1とする。



④ NTECV値の決定。

TEC値の種類数をNECVとする。③ではTEC値は1、5、6、7、9、10の6種類である。従って、①のNTECV値は6となる。

$$\underline{NTECV = 6}$$

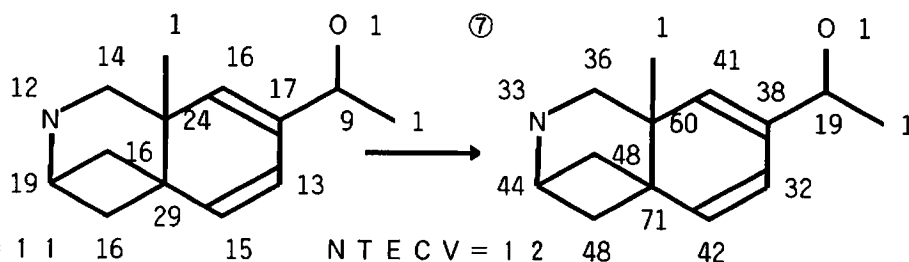
⑤ NECVとNTECVとの比較。

NECVとNTECVの値を比較し、

- NECV \geq NTECV の時は最終ステージに行く。
- NECV < NTECV の時は各NODEのNTECV値を新たにNECV値とみなし、ステップ③に戻る。

これ以降はaの状態になるまで③~⑤のステップを繰り返す。

⑥



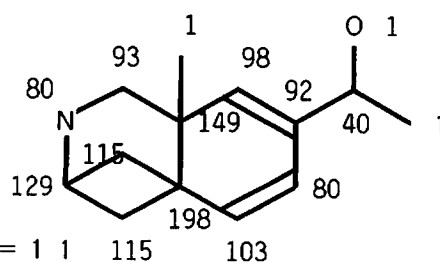
$$NTECV = 11 \quad 16$$

$$NTECV = 12 \quad 48 \quad 42$$

⑧ EC値の決定。

NECV (12) \geq NTECV (11) となるので作業を終了する。

この一つ前の時点でのEC値を各NODEの拡張分岐数 (EC) として採用する。



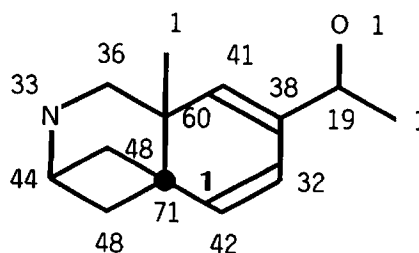
$$NTECV = 11 \quad 115 \quad 103$$

手順2: NODEの番号付け

① 1番目のノードの決定。

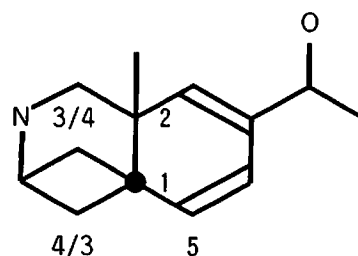
手順1の⑦で得られたEC値の内、最大値を有するNODEを選択し、1番とする。

例ではEC値71が最大なので、このEC値を持つNODEを1番とする。



② 1番のNODEのまわりのNODEの番号決定

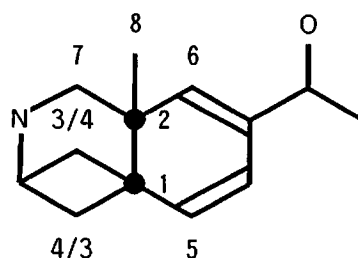
先に決定された1番のNODEに隣接するNODE総てについて、EC値の大きい順に2番目以降の番号を付ける。



* EC値が同じ時、総ての可能性について番号付けが行われる。この時点では2種類の番号付けが可能であり、最終的には4種類の番号付けが可能となる (NODE番号1に結合するNODE中、EC値が48のものが2個あり、番号3と4が順位的には同値である。また、同じNODEに結合している測鎖のOとCはEC値が共に1であるので番号付けの順位は同じである)。化合物の対象性が大きくなると番号付けの可能性がおおきくなり、多数の番号付けが出てくることになる。これらの番号付けのなかで、MORGAN番号が最も小さくなるものが最終的なMORGAN番号として採用される。

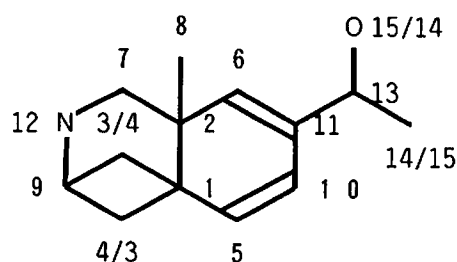
③ 1番から2番のNODEまわりへの拡張。

1番のNODEに隣接する総てのNODEに番号がふられたならば、次に2番のNODEに隣接する残りのNODEに対し、EC値の順に番号を付ける。



④ 総てのNODEに対する番号付け。

③の手続きを化合物上の総てのNODEに対し適用する。総てのNODEに対し番号を割付、番号付けを完了する。



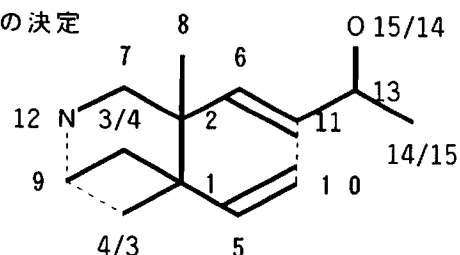
この手順2を完了した時点で与えられた化合物を構成するすべての原子について一元一項対応の番号 (ユニークナンバリング) がつけられた事になる。

図からもわかるように対象性の強い化合物では番号付けの可能性が多数発生し、一元一項対応を守る事は不可能となる。これがこのアルゴリズムの限界である。

手順3: 種々リストの作成及びMORGAN名候補の作成

① From Attachment Listの決定

構造式中の各NODEの結合関係を示すリストである。各NODEに番号が付けられた時の基準となったNODEの番号が各NODEについて与えられる。結果としてFrom Listは化合物中の結合を定義するものとなる。



Atom Number	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
From List	-	1	1	1	1	2	2	2	3	5	6	7	11	13	13

*From List は右図の実戦で描かれた結合を定義している。点線で描かれた結合に関してはFrom Listでは定義されない。

② Ring Closure Listの決定

①の構造式で、From Attachment Listでは定義されなかった点線で示される結合をこのRing Closure Listで定義する。

この時は(4, 9)、(9, 12)、(10, 11)で示される。この数値の並びをそれぞれのペアの番号を昇順に並べたものをRing Closure Listとして採用する。従ってこの場合は4, 9, 9, 12, 10, 11となる。

③ Atom Type Listの決定

原子の種類を定義するリストである。このリストは原子に付けられた番号順に並べられる。

Atom Number	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Atom Type List	C	C	C	C	C	C	C	C	C	C	C	N	C	C	O

④ Bond Type Listの決定

結合に関する情報を定義するリストである。単結合、2重結合、3重結合、芳香族結合にそれぞれ1, 2, 3, 4の番号を割り当てる。このリストはFrom ListとRing Closure Listにより示される結合について定義される。

From List	-	1	1	1	1	2	2	2	3	5	6	7	11	13	13
Ring Closure List															4,9; 9,12; 10,11
Bond Type List	1	1	1	1	1	1	1	1	2	2	1	1	1	1	1

⑤ Modification Listの決定

このリストは電荷、同位体、異常原子価等の記述に使用するものである。以上のリストをまとめて一つのリストにする。

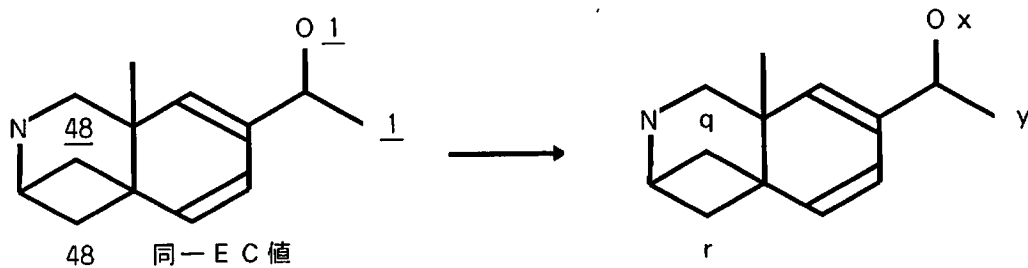
Atom Number	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
From List	-	1	1	1	1	2	2	2	3	5	6	7	11	13	13
Ring Closure List															4,9; 9,12; 10,11
Atom Type List	C	C	C	C	C	C	C	C	C	C	C	N	C	C	O
Bond Type List	1	1	1	1	1	1	1	1	1	2	2	1	1	1	1

ここでC, N, Oとにそれぞれ1, 2, 3の番号をアサインすると、前記Listを総て結合する事によりMORGAN名が生成される。従って、今回用いた対象化合物に対するMORGAN名の一つは以下のような数値データとなる。但し、この時点では番号付けは絶対的なものでなく、複数のMORGAN名が発生している。以下にこれら複数(4種類)発生されるMORGAN名の一例を示す。このMORGAN名では酸素原子が15番に、末端メチルが14番にアサインされている。

候補MORGAN名:

01010101020203050607111313 FROM LIST	040909121011 RING CLOSURE
0101010101010101010102010103 ATOM TYPE LIST	01010101010101010202010101010101 BOND TYPE LIST

手順4: Better Nameの取り出し及びMORGAN名の最終決定



手順3の時点でMORGAN名の発生が可能となる。この時点で唯一個だけのMORGAN名が発生されていれば、その名をそのまま対象化合物のMORGAN名と採用する事が可能である。しかし、一般的には例で示されるように複数のMORGAN名が発生されている。つまり対象化合物には2ヵ所の同一EC値を持つ所があり、4ヵ所のNODEが関与している(上図中番号が3、4及び14、15)。従って番号付けの可能性はそれぞれの場所で2通りの可能性を持ち、総計で4種類の番号付けの可能性を含んでいる。

従って、一元一項対応の目的から、これら複数発生したMORGAN名を可能な限り少なくする事が必要である。即ち、複数MORGAN名が発生した時は、その中でBETTER NAMEと呼ばれるものを選択する。

* BETTER NAMEとは、MORGAN名を数値とした時、より小さな値を持つ名前を指している。

① 番号付け可能ケースに関する表の作成

対象化合物中の同位NODE総てにアルファベットによりID符号をアサインする。この場合はq、rとx、yとラベルする。

先ず、各NODEのq、r、x、yとが取りうる総ての番号についてまとめ、簡単な表を作成する。この表を作成する時、NODE番号は昇順に並べておく事が必要である。

番号付けのケース	NODE番号			
	3	4	14	15
ケース1	q	r	y	x
ケース2	q	r	x	y
ケース3	r	q	y	x
ケース4	r	q	x	y

② 原子の種類に関するIDコードを用いた作表

続いてこの表の各NODEについて先にアサインされている数字(C=1、N=2、O=3)を割当てて新たな表を作成する。

番号付けのケース	NODE番号			
	3	4	14	15
ケース1	1	1	1	3
ケース2	1	1	3	1
ケース3	1	1	1	3
ケース4	1	1	3	1

③ Better Nameの選択及びMORGAN名の作成

4通りのケースについて3、4、14、15とをまとめ、4桁の数字と考えた時、その値が最も小さいものをBETTER NAMEとして採用する。

例に用いた化合物はケース1と3とが1113でケース2と4の1131よりも小さいので採用される。ケース1と3とでは値が同じなのでこれ以上の選択は不可能であり、最終的な番号付けは2種類採用される事になる。つまり、yが14でxが15となり、qとrは3及び4の両方採用される。従って、最終MORGAN名は2種類となる。

この効果はAtom Type Listの番号で該当する炭素と酸素の順番を入れ換える事で小さな値の方を採用する事を意味している。

Atom Number	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15		
Atom Type List	C	C	C	C	C	C	C	C	C	C	C	N	C	C	O		
○Atom Type List (数値)	1	1	1	1	1	1	2	1	1	1	1	1	1	1	3	(小)	
Atom Type List (数値)	1	1	1	1	1	1	2	1	1	1	1	1	1	1	3	1	(大)

ここで用いた化合物においては、先に示した候補MORGAN名がそのまま最終MORGAN名として採用される。

候補MORGAN名：
 0101010102020203050607111313040909121011010101010101010101010102010103
 0101010101010101020201010101010101

2. SEMA (立体情報を含む化合物に適用可能な名)

2. 1. SEMA概要

SEMA (STEREOCHEMICALLY EXTENDED MORGAN ALGORITHM)名はその名の通り、MORGAN名が立体情報を含む化合物の識別が不可能であったものを、立体を含む化合物まで識別可能となるようにMORGANアルゴリズムを拡張したものである。従ってSEMA名とは、立体情報を含む化合物の識別を目的とするユニークネーミングの算出法を意味する。このSEMA名はMORGAN名の後半部分に立体に関する情報を付加させた形式を取っている。

$$\text{SEMA名} = \text{MORGAN名} + \text{立体情報に関するコード}$$

* W.Todd Wipke and Thomas M. Dyott, "Stereochemically Unique Naming Algorithm", J.A.C.S., 96, 4834(1974).

2. 1. 1. MORGANアルゴリズムの立体化学への拡張

MORGANアルゴリズムを基本とし、立体や幾何情報を含む化合物に適用可能なように拡張する時、検索で問題となる一元一項対応を守る為に化合物中に存在する立体中心*1 (Stereo Center) に関する配座情報を総て含ませる事が必要となる。

- * 1. ここで定義される立体中心とは、反転 (Inversion)する事で異なる立体異性体を形成する構造的特徴を意味し、化学で一般的にいわれている立体中心とは異なるものである。ここでいう立体中心に該当するものは、化学では不斉炭素及び、炭素-炭素2重結合 (炭素以外にも拡張可能) が該当する。
- * 2. 文献中ではヘテロ原子を含む立体に関する情報は議論の対象としていないが、容易に適用拡張可能である。

MORGANアルゴリズムの拡張によりSEMA (Stereocally Extended Morgan Algorithms) 名へと導く事ができる。このSEMA名によりMORGAN名では識別不可能であった立体異性体や幾何異性体等の識別が可能となる。しかし、このSEMA名にも限界があり、回転異性体等の識別は不可能である。

2. 1. 2. 立体異性体や幾何異性体を識別する為のルール

立体情報を識別し、記号情報へと変換する目的で2種類のパリティを定義する。このパリティはEVEN及びODD (以下適宜 (E) 及び (O) と略す) で表現され、化合物中の立体中心に、このどちらかのパリティが設定される。

以下にこのEVEN及びODDの定義について不斉炭素と2重結合の場合を例に取って説明する。

A. 不斉炭素の時

① 水素を含まない不斉炭素の時

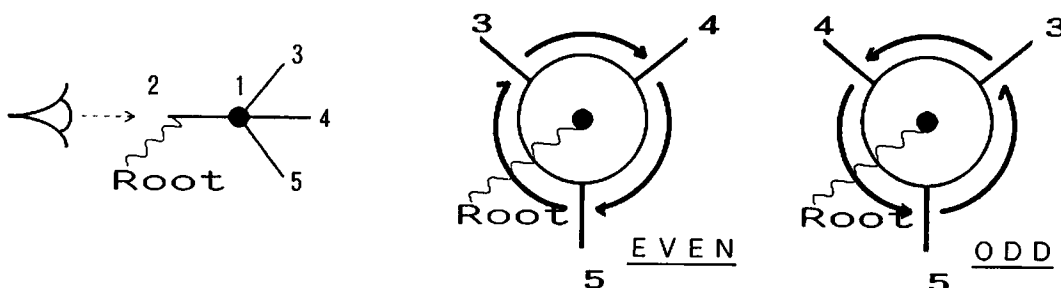


図1. 水素を含まない不斉炭素に関するパリティ (EVEN, ODD) の決定

②水素を含む不斉炭素の時

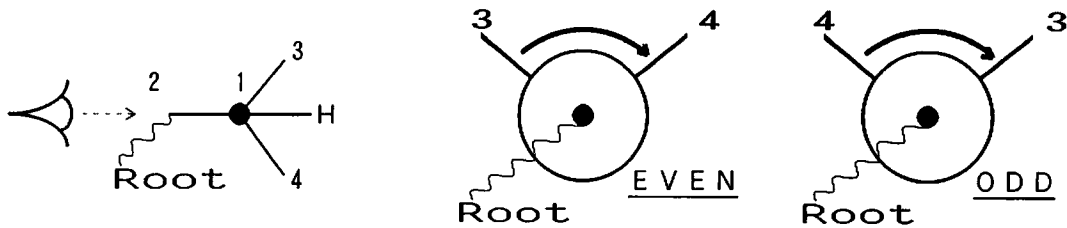


図 2. 水素を含む不斉炭素に関するパリテイ (EVEN, ODD) の決定

図中の原子に付けられた番号はMORGANアルゴリズムにより決定された番号である。パリテイ決定の為の手続きは、先ず目的とする立体中心に結合した原子の番号中、最も小さい番号の原子を視点とする。この視点から立体中心原子を見た時、数字の並びが時計回りの時EVEN、反時計回りの時ODDと定義する。

B. 2重結合の時

2重結合に関しては、その2重結合を形成する2個の原子について個別にパリテイがアサインされる。つまり不斉炭素の時と同様、一個の2重結合原子を中心としてみた時、隣接原子の番号の並びが時計回りの時はEVEN、反時計回りの時はODDがアサインされる。但し、2重結合原子の隣接2原子が互いに非等価である事が前提である。

このパリテイのアサインに先立ち、2重結合を固定する事が必要となる。この固定は2重結合を形成する2個の原子に付けられた番号に注目し、その原子番号が小さい原子を左に置き、大きい方の原子を右に置くように定義する。このようにして固定された状態の2重結合に関しパリテイを決定する。以下にその事例を示す。

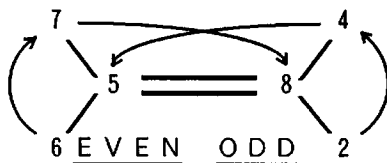


図3中、原子5に結合している隣接原子の番号は小さい方から順にみた時、時計回りである。従って、原子5のパリテイはEVENである。また、原子8に関しては2、4、5であるので反時計回りとなりODDがアサインされる。

図 3. 2重結合のパリテイ定義

2. 1. 3. 不斉炭素及び2重結合に関するパリテイの決定事例

① 不斉炭素のパリテイ

不斉炭素及びその不斉炭素に直結する4原子にMORGANアルゴリズムにより左図のような番号が付けられていたとすると、不斉炭素5の廻りの原子のうち最も小さな番号を持つ原子の方向に視点を置く。この視点から7番の原子を通し、不斉炭素を眺めた時、残る3原子に付けられた番号の並びをチェックする。この時12~14~15という原子の並びは反時計回りであるので、この不斉炭素の並びはODDとなる。

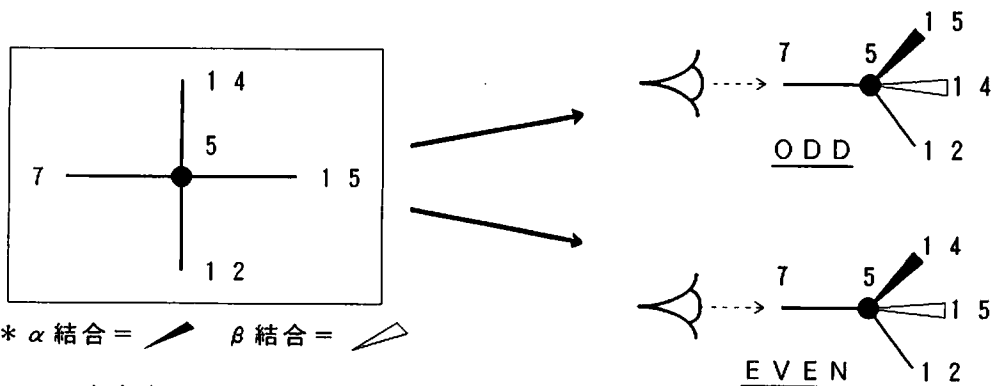


図 4. 不斉炭素のパリテイ決定事例

2. 2. SEMA名によるMORGAN名限界の克服

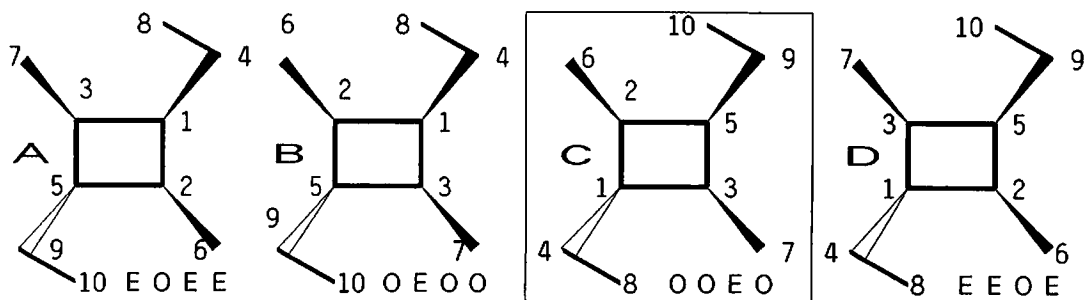


図5. MORGANアルゴリズムにより創出された4種の原子番号

上図で示された化合物に対し、MORGANアルゴリズムでは4種類の番号付けが可能であり、ユニークな番号をつける事は不可能である。この問題が、先のパリティを利用する事で解決可能となる事を示す。

先ず前記4化合物中の立体中心(1, 2, 3, 5)についてパリティを調べる。決定されたパリティに対し、EVEN=2, ODD=1を与え、その数値の大小を比べる。

原子番号	1	2	3	5	数値列	原子番号	1	2	3	5	数値列
A	E	O	E	E	= 2 1 2 2	B	O	E	O	O	= 1 2 1 1
C	O	O	E	O	= 1 1 2 1	D	E	E	O	E	= 2 2 1 2

これら4個の数値を比べた時、最小の数値となるC(1 1 2 1)がこの化合物のユニーク名として採用する。従って、Cに付けられた番号が最終番号となる。

2. 3. MORGAN名からSEMA名の作成

2. 2節でユニーク番号付けに関し、立体に関するパリティを導入する事でMORGANアルゴリズムの限界をクリアする事が出来る事をしめた。一旦このユニーク番号が決定されれば、この番号を基準としてSEMA名を発生する事が可能となる。

このSEMA名は先にも示した通り、MORGAN名に立体に関する情報を付加する事で得られる。この立体に関する情報は各原子毎のパリティを数値データに変換して得られた数値列としてえられる。2. 2節の説明で用いた化合物を例に取るならば以下の様な表が得られる。表1中、0は立体中心とは関係の無い原子を、1及び2はODDとEVENとを示している。

表1. 最終原子配置リスト

原子番号:	1	2	3	4	5	6	7	8	9	10
パリティ:	1	1	2	0	1	0	0	0	0	0

この数値データ列は2重結合に関する情報であれば2重結合配置リスト(Double Bond Configuration List)として、不斉炭素であれば原子配置リスト(Atom Configuration List)と称される。これらのリストをMORGAN名の最後に付加する事でSEMA名となる。

2. 4. 立体化学とSEMA名との関係

SEMAは単に立体情報を含む化合物のユニークネーミングの発生による検索だけに利用されるわけではない。SEMAの特徴をいかす事で立体化学情報に関する様々な利用が可能である。以下にはその例の幾つかについて述べる。

① ENANTIOMERに対するSEMA名発生

化合物が与えられた時、その化合物のパリティを検討する事でENANTIOMERのパリティを決定する事が可能である。



図6. ENANTIOMERを有する化合物

化合物Aの立体パリティの検討により、ENANTIOMERである化合物Bの立体パリティの決定が可能である。 先ず、化合物Aに対しMORGANアルゴリズムによる番号付けは以下のa~dで示された4種類が発生される。

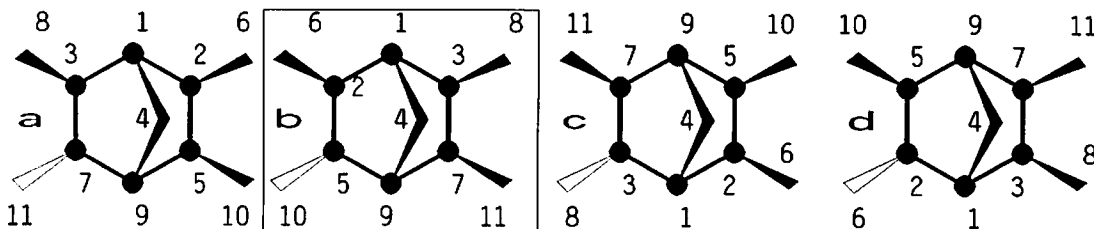


図7. MORGANアルゴリズムにより創出された4種類の原子番号

この4種類の番号付けに従った時の立体パリティをチェックし、その順番に並べる。

表2. 図7の化合物に対する立体パリティ

原子番号 :	1	2	3	5	7	9	数値列	序列
a	E	O	E	O	O	O	= 2 1 2 1 1 1	
b	O	E	O	O	O	E	= 1 2 1 1 1 2	最低序列
c	O	E	E	E	O	E	= 1 2 2 2 1 2	
d	E	E	E	O	E	O	= 2 2 2 1 2 1	最高序列

前記4種の番号付けに対する立体パリティは表2の様になる。 従って、このパリティを1 (ODD) と2 (EVEN) とで置換し、相当する数値列へと変換する。 ここで得られた数値列のうち、立体パリティ決定に関する2. 3節で示した定義より、化合物Aに対する最終的な番号付けは値の最も小さなbに決定される。

AのENANTIOMERであるBのパリティは、Aの最低及び最高序列のパリティを比較し、両方に共通なパリティを最低序列のものに関し反転させる事により得られる。

原子番号 :	1	2	3	5	7	9		1	2	3	5	7	9
b	O	E	O	O	O	E (最低)	→	O	O	O	O	O	E
d	E	E	E	O	E	O (最高)	入れ換え						

この解析で得られた共通パリティを有する原子2と原子5は次に述べるキラリティの判定時に利用される重要な情報となるので、Reduced Set of Chiral Centers (Src) として、Src = { 2, 5 } として表現する。

このSrcは他の立体中心原子と区別する為、SEMA名の中でも特に明記して記述するようになっている。 つまり、これらSrcの原子の立体パリティの値に3を加えたもの (ODD = 1 + 3 = 4, EVEN = 2 + 3 = 5) が最終原子配置リスト (Atom Configuration List) として採用される。 表3には化合物A及びBの原子配置リストが示され

ている。

表 3. 化合物 A 及び B の原子配置リスト及び S r c 原子

原子番号 :	1	2	3	4	5	6	7	8	9	10	11
化合物 A	1	5	1	0	4	0	1	0	2	0	0
化合物 B	1	4	1	0	5	0	1	0	2	0	0
Src 情報	Src		Src			Src					

② キラリティ有無の判定

化学において化合物のキラリティ有無を判定することが必要となる場合がある。 S E M A 名では S r c の存在を判定する事で、化合物のキラリティの判定を行える。

S r c = 0 の時、化合物はアキラル (A C H I R A L) である。
S r c ≠ 0 の時、化合物はキラル (C H I R A L) である。

図 6 の化合物は S r c = 2 である為キラルである。一方、図 5 の化合物は最低序列 (O O E O) と最高序列 (E E O E) との間で、共通なパリティが存在しないので S r c は 0 であり、アキラルとなり、 E N A N T I O M E R は存在しない。

③ 配座異性体の識別

S E M A 名を拡張した『拡張 S E M A』により配座異性体の識別が可能となる。

この拡張 S E M A では 2 面体角を形成する 4 個の原子の決定とその角度情報を決定し、 S E M A 名の後に結合する事で得られる。この時、選択される 2 面体角は S E M A ルールに適合して選択する (最も小さい値を有する S E M A ナンバリングによる値を示す 2 面体)。以下に trans-Decalin を例に取り、手順に従って順番に説明する。

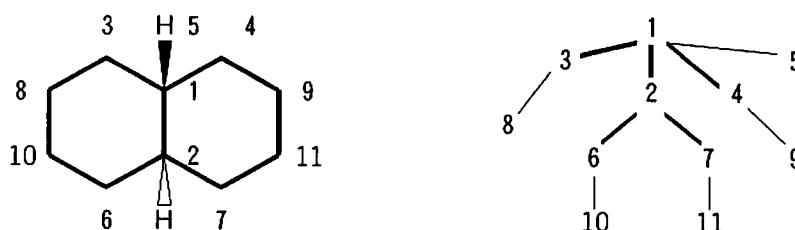


図 8. trans-Decalin とスパンニングツリー

(1) スパンニングツリーの作成

trans-Decalin の配座を指定する為には原子 1 (O D D) と 2 (O D D) (S E M A の番号付けにより決定された原子番号利用) に関する原子配置を特定する事で可能である。この 1 と 2 の原子に関する結合環境を明記する為にスパンニングツリーを作成する。このスパンニングツリーは目的とする原子 (1 と 2) を出発原子とし、 3 番目の原子迄の結合関係を示している。

(2) S E M A 名に組み込む 2 面体角の選択

このスパンニングツリーを基本とし、原子 1 及び 2 が関係する 2 面体角を抽出する。個々の原子について複数の 2 面体角候補があるが、ここでは S E M A の定義と同様に 2 面体角を形成する 4 原子の原子番号の並びが最も小さい番号を有する 2 面体角を採用する。これにより情報の重複を最小限にとどめる事が可能である。

例では 1 及び 2 番目の原子に関係するすべての結合 ((1 , 2) (1 , 3) (1 , 4) (2 , 6) (2 , 7)) に関し 2 面体角の設定が行われる。例えば結合 (1 , 3) に関

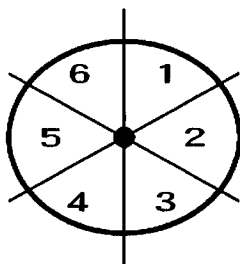
しては2面体角を形成する原子の組み合わせは(2 1 3 8)、(4 1 3 8)、(5 1 3 8)の3種類あるが、この3種のうち値が最も小さいものは(2 1 3 8)の組み合わせである。従って、結合(1, 3)の設定は原子(2 1 3 8)による2面体角の設定で充分となる。同様に残る結合について2面体角を決定する。最終的な2面体角をMORGAN名のFROMリスト(結合が定義される)と対照して示したものが表4である。

表4. 最終決定2面体角構成原子IDリスト

原子番号	From List	2面体角構成原子	結合角情報
1			
2	1	3 - 1 - 2 - 6	
3	1	2 - 1 - 3 - 8	
4	1	2 - 1 - 4 - 9	
5	1		
6	2	1 - 2 - 6 - 10	
7	2	1 - 2 - 7 - 11	
8	3		
9	4		
10	6		
11	7		

以上の手続きにより多数存在する2面体角の中で、SEMA名に情報として与えるべき2面体角を選択する事が出来た。続いて、これらの指定された2面体角を実際にデータとして与える事が必要となる。

(3) 2面体角情報の取り出し



2面体角の情報は角度で与える事も可能であるが、ユニークネーミングという目的から判断すれば必ずしも細かな角度情報である必要はない。例えば、図9で示される様に角度を6等分し、この1~6の数字を角度情報として与える事で充分であり、同時に計算機による取り扱いも簡単になる。

図9. 2面体角の6分割

(4) 配座リスト(Conformation List)の作成

2面体角の情報と対応する原子番号が決定されたならば、これらの情報を原子番号順に並べる。2面体角の情報が無い原子は0を代入する。

原子番号:	1	2	3	4	5	6	7	8	9	10	11
2面体角:	0	2	3	4	0	1	1	0	0	0	0

この配座リストをSEMA名の最後に付加する事で『拡張SEMA名』となり、配座異性体の識別が可能となる。

2. 5. MORGAN名、SEMA名及び拡張SEMA名との関係

化合物のユニークナンバリング及びユニークネーミング手法として3種類の手法がある事を示した。これらは互いに独立するものでなく、MORGANアルゴリズムを基本とし、その機能拡大により対照化合物範囲の拡大を行っている。図10にこれら3手法の関係を示す。



図10. MORGAN名、SEMA名及び拡張SEMA名との関係

2. 6. SEMA名の限界

SEMA及び拡張SEMA名にも扱えない化合物が存在する。これらの手法の限界として明記する必要がある。これらの化合物はSEMA等の運用により解決が可能なものと、アルゴリズム上絶対的に不可能なものとの分けられる。

SEMA/拡張SEMAにより表現不可能な化合物

- 絶対立体不斉の決定 (相対立体不斉のみ識別可能)
- 回転異性体の識別
- イオンの形になった化合物
- 2、3量体等の複数分子で形成される化合物
- 錯体化合物

これらの化合物中、c、d、eは簡単な運用によりSEMAの適用が可能である。aに関しては、SEMAでは立体パリティの決定がMORGANアルゴリズムにより付けられた原子番号を基準とし、原子量及び置換基の大小による比較ではないので絶対立体不斉の決定は出来ない。しかし、立体パリティの決定アルゴリズムを利用し、原子量及び置換基による比較が出来るようにする事で可能となる。

bは現アルゴリズムでは解決不可能な問題である。