

□ 化合物構造情報の計算機上での表現

① 化合物構造情報表記の一般的問題点

② 線型表記法

(1) 構造式の再現不可能な構造表記法

- ・ CAS 番号
- ・ MORGAN 名
- ・ SEMA 名

(2) 構造式の再現可能な構造表記法

- ・ 一般的命名法
- ・ WLN, SKOLNIK, GREMAS, SMILES

③ 結合表による化合物の表現

④ 3次元構造式の表現

- ・ 結合表による表現
- ・ Z-マトリクスによる表現

⑤ 化合物の立体化学関連表示

⑥ 関連計算機技術

- ・ グラフィックからの直接入力
- ・ CHEMICAL タイプライター
- ・ WLN → CT 変換
- ・ 化合物名 → CT 変換

1. 化合物構造情報の計算機上での表現 (概論)

1. 1 計算機の情報体系と化合物の情報体系の差と融合

□ 計算機における情報の特徴

計算機が理解できるデータは0/1からなるバイナリーデータである。人間がこのバイナリーデータをもて必要な情報を取り出すことは困難である。計算機では、このバイナリーコードを基本とし、より高度なコード体系として数値データや文字データを扱うことが可能である。この数値/文字データは人間でも理解可能である。計算機と人間の両方が理解できる数値/文字データをもちいることで、基本的な計算機と人間とのコミュニケーションがなりたつ。通常はこの数値および文字データを計算機の入/出力として用いる。この数値/文字データの他に重要な情報として画像や音声データがある。この画像/音声データは0/1情報の集合体として扱うことで計算機上での処理が可能となる。しかし、このアプローチはメモリーを大量に必要とし、処理等にも時間がかかることから最近になり実用が可能となってきたもので、初期の計算機では実現されていなかった。

□ 化合物構造式における情報の特徴

化合物構造式を情報という観点から考えてみると、①2次元/3次元の形を持つ絵としての情報、②原子同志の結合関係の情報の二種類存在することがわかる。

化合物の構造式を絵として表現することは大事であり、化学者間のコミュニケーション手段として重要である。しかし、化合物構造式は絵であっても、化合物としての本質的な情報は原子の結合関係(トポロジカル/トポグラフィカル)であり、化学上の様々な処理はこの結合関係を変化させることで行われる。従って、画像情報は化学の分野では特殊な目的(コミュニケーション、イメージ把握等)の時にのみ利用され、化合物情報の実質的な受け渡しは原子の結合関係を中心として行われる。

□ 化合物構造式情報の計算機への取り込み

先に、計算機と人間の接点は数値/文字データと画像情報の扱いにあることをのべた。このうち画像情報の扱いは初期の計算機にとり実行困難であり、且つ画像情報は化学の分野では必ずしも本質的な情報伝達手段ではないことものべた。従って、初期の計算機による化合物構造式の取り扱いには例外を除き、必然的に数値/文字データによる扱いから出発した。

化合物構造式 ⇄ 数字および文字データ ⇄ 計算機

1. 2 化合物構造式の数値/文字データへの変換

化合物構造式を数値/文字データへと変換する手法として様々な手法が存在する。これらの手法は化学の分野で古くから利用されてきたものから、計算機の進歩に伴って最近展開されてきたものまで多種多様存在する。現在利用されている化合物表記法の代表的なものを表1にまとめる。

表1. 化合物の数値/文字データへの変換手法

	線形表記	マトリクス形式
2次元構造データ	CAS番号* 化合物名(IUPAC) WLN, GREMAS SMILES MORGAN* SEMA*	結合表(2次元)
3次元構造データ		結合表(3次元) Zマトリクス

* 構造式の再現不可能

これらの表記法のうち、計算機と全く関係のない目的で開発されたものとして、CAS

番号と化合物命名法（IUPAC名）がある。CAS番号は化合物に人為的に与えられる登録あるいはID番号であり、このID番号を発行する作業は米国のケミカルアブストラクツサービス（CAS）が行っている。CASでは世界中の文献やパテントに発表された化合物をチェックし、新奇化合物についてこの番号を発行している。従って、化合物構造式情報とCAS番号との間には何の相関もなく、単なる化合物のIDである。この番号は現在一千万を超えている。IUPAC命名法は化学者が必要とする共通言語としての化合物命名法であり、世界の化学者の合意のもとに定められたものである。残りの表記法は計算機での利用を前提として開発された。

□計算機の進歩と変換様式の変化

計算機による化合物構造式の変換様式は、計算機のハード上での進歩と密接に関係している。計算機の当初はメモリーの容量（コア及び外部メモリー）が少なく、計算速度も遅いため、大量の数値／文字データを計算機上で扱うことは極めて困難であった。この結果、当初の計算機による化合物構造式の扱いは、なるべく少ない数値／文字情報で扱うことを最優先としたアプローチに限定された。この要求に答えるものとして、線型表記法が展開された。

しかし、化合物構造式そのものを文字データだけで表すこと自体が困難である事に加え、さらに可能なかぎり少ない数字／文字であらわすことは、抽象化が進み、ルールが複雑となり化学者にとって扱いにくくなることを意味している。計算機からの一方的な要求を満たすために考案された線型表記法は、最後まで化学者とのギャップを埋めることは出来なかった。現在では一部のアプローチで利用されているのみである。

計算機のハード技術の進歩（特にグラフィック関連技術）にともない、化合物構造式そのものをグラフィックディスプレイ上で表現することが現実のものとなった。計算機上で化合物構造式を絵として扱うには従来の構造式変換技術とは異なる扱いが必要である。この目的で結合表に代表されるいくつかの変換技術が新たに開発された。現在ではこの結合表形式が主流となっており、結合表自体も国際的に統一されつつある。

□計算機の進歩と利用形態の変化

変換様式のみならず、利用形態そのものも計算機のハード技術の進歩と密接に関係して変化してきた。当初の利用は化合物データベースとしての利用が主であった。分子軌道法での利用等も行われていたが、計算速度も遅く、メモリーの制限等から経験的及び半経験的分子軌道法の適用が限界で、実用に供するレベルとは程遠かった。しかも、化合物は数字／文字で表現されるだけで、化合物構造式そのものを表示することが出来ないために化学者が計算機を利用することは困難であった。

グラフィックディスプレイの進歩とともに、化合物構造式をディスプレイ上に表示することが可能となった。この結果、一般の化学者でも構造式をコミュニケーション手段とし、計算機と直接対話することが可能となってきた。グラフィックのみならず、計算速度の高速化やメモリーの増大といったハードの進歩は化学分野における計算機利用を急速に拡大した。

構造式を直接扱うことでプログラム自体も構造式を主体とするマンマシンインタフェースを中心とした設計がなされるようになり、化学者と計算機とのギャップは従来と比べてかなり小さくなってきた。

1. 3 化合物構造式変換法の分類

化合物構造式の数値／文字データへの変換は、その変換対象化合物が2次元構造式か3次元構造式かで大きく2分類可能である。ここで2次元構造式は紙上に描かれる化合物構造式であり、3次元構造式というのは分子モデルのように3次元座標情報を持つものを意味する。

さらに、変換された数値データあるいは文字データを単に線條に並べた一次元の表現（線形表記）で行うのか、あるいは2次元のマトリクス形式で表現するかでさらに2種類に分類される。マトリクス形式で表現する変換式では構造式を再現出来るが、線型表記法では構造式を再現できるものと出来ないものとの二種類存在する。これらの分類は既に表1に示されている。ここに示された手法以外にも化合物の変換形式は様々なものが存在し、ケースバイケースまたはローカルに利用されている。

□個々の化合物表記法の特徴

CAS以外の表記法は、基本となる化合物構造式に対し個々の手法により定められた変換アルゴリズムを適用することで数字／文字データへと変換される。これらの変換法のうち化合物名（IUPAC命名法）については先にのべた。WLN、GREMAS、SMILES等の線型表記および結合表形式の総ての表記法は、変換データから化合物構造

式を再現することが可能である。これら以外のMORGANやSEMAといった表記法は化合物の一元／一項対応による化合物検索を主目的とし、変換コードからの化合物構造式の再現は不可能である。

□化合物の表記における一元一項対応の問題について

化合物構造式を何らかの形で表現する時に問題となることとして、化合物構造式と化合物名との一元一項対応という問題が存在する。ここでは簡単にこの一元一項対応についてのべる。

この一元一項対応が特に問題となるのは化合物の命名法を定める時である。化合物名というものは厳密に1化合物に1つの名前というように1対1で対応していることが理想的な形態である。しかし、3次元的に複雑な形態を持つものを1次元の文字集合で完全に表すのは困難である。さらに、化合物命名法が定まる前から使われていた慣用名等が存在し、このような事実が化合物と化合物名の一元一項対応を困難にしている。

一元一項対応の元という言葉はこのばあい化合物に対応している。また、項という言葉は化合物名に対応している。従って、一元一項対応とは化合物とその化合物名とが1対1で厳密に対応していることを意味する。

多元一項対応とは一つの化合物名が複数の化合物に対応する時を、一元多項対応とは一つの化合物が複数の名前をもっていることになる。このように、一つの名前に複数の化合物が対応したり、一つの名前に複数の化合物が対応することは、化合物検索を行う時に化合物名をもちいて検索することは能率が悪いことを意味している。

一元一項対応 ・CAS番号	化合物(1)	化合物名(1)
多元一項対応 ・	化合物(N)	化合物名(1)
一元多項対応 ・	化合物(1)	化合物名(N)
多元多項対応 ・	化合物(N)	化合物名(N)

1.4 線型表記法 (Linear Line Notation) による化合物構造式の表現

線型表記法は化合物構造式を一次元の数字／文字列として表すものであり、化合物命名法の基本である。本来、二次元及び三次元である化合物を一次元で表現するため、情報を完全に保ちながら一次元に落とすには様々な変換ルールが必要となる。使用にあたっては、その変換ルールを十分に理解することが必要であり、化学者にとっては大きな負担となる。特に化合物の複雑な3次元情報を完全な形で1次元情報へと変換することには種々の点で限界が存在し、その限界をカバーする為にルールがさらに複雑になるといった悪循環に陥ることになる。

化合物名 (IUPAC命名法) は化学者の常識、コミュニケーション手段として必須であり、この変換ルールを知っていることは化学者である必須条件である。従って、このIUPAC名を用いて計算機を利用することができれば化学者の負担は少ないのだが、不幸にしてこのIUPAC化合物名をそのまま計算機上で扱うことは困難である。この主な原因として以下のようなことがある。

①キーボードには無い特殊な記号を用いている

②化合物名の順序づけに自由度があり過ぎる

③構造が複雑になると表現が長くなる

①の問題は化合物構造式を計算機に入力出来ないという決定的な問題となる。②の問題は化合物検索時に、該当する化合物の検索が出来ない、検索効率が悪いといった問題が発生する。これはIUPAC命名法が化合物の命名を目的とするものであり、効率的な検索を目的とするものではないことによる。③ではメモリー等の制限に関する問題が発生する。

これらの欠点をカバーし、特に計算機での取扱を目的として開発されたのがWLNやGREMAS等で代表される線型表記法であり、後にのべる結合表とよばれる表記法である。

これらの表記法を用いることで、先にのべた様々な欠点が解消される。

以下の表には化合物名とWLN、GREMAS等の線型表記法との特徴的な差がまとめられている。

特徴	化合物名 (IUPAC)	線型表記法 (WLN, GREMASその他)
Unique/ Unambiguous	同一化合物名の表現に種々の名前が存在し、化合物と名前が1対1に対応しない。	化合物と変換コードは良く対応する。厳密な一元/一項対応を目的とした表記法もある。
使用文字	特殊記号を用いる。 『添字、大文字/小文字の混在』	キーボード上にある文字のみを使用
表現文字数	一般的に、化合物の表現に必要な文字数は多くなる。 (慣用語は除く)	化合物表現の為の文字は少なくなるように設計されている。
例)	BENZENE ETHANE 1-METHYL-2-METHOXY-4-AMINO BENZENE	—→ R —→ 2 —→

上記線型表記法の開発により、化合物構造式の計算機への入力が可能となった。つまり、2次元および3次元の複雑なトポロジカル情報が1次元情報へと変換され、この結果構造式が単なる文字列へと変換されたためである。この結果、化合物検索は単なる文字検索の問題へとすりかえられ、計算機による従来からの検索技術を利用する事で化合物検索が可能となった。

しかし、前記WLNやGREMAS等を用いるためには特殊な変換ルールをマスターすることが必要である。しかし、主なユーザである化学者はこれらの変換ルールを知らない事が多い。従って、化合物検索を行おうとするならば化学者は普段利用している化合物命名法(IUPAC命名法)の他にWLN、GREMAS等の線型表記法を改めて学ぶ事が必要となる。この事実は多くの化学者にとって大きな負担となる。

化学者の思考体系は化合物構造式を起源とし、総てのコミュニケーションも化合物構造式が基本である。従って、WLN等の数字/文字データを基本として思考することはありえず、必ず構造式に変換する事が必要である。従って、計算機を離れた途端WLNはその存在意義を失い、化合物構造式への変換が必要となる。このように、日常的に構造式→WLN→構造式といった変換業務が発生し、特に複雑な化合物の時は、この業務が大きなネックとなってくる。

1. 5 個々の表記法による化合物構造式の変換

□WLN(WISWESSER LINE-FORMULA CHEMICAL NOTATION)の概要

WLN表記法は計算機による使用を目的とし、W. J. WISWESSERにより提唱された化合物の線型表記法である。この表記法は計算機での使用を前提とするため、使用する文字は数字、26個の英大文字、3個の区切り記号&、-、/およびブランクだけである。

これらの数字/文字/記号は一字単位で特有の情報を表現することを原則とし、例外(出現頻度の低い原子で通常二字で表現される(例、Fe→FE、Na→NA等)ものやUUで示される3重結合等)を除き二文字以上で一つの単位となることはない。表2にWLNで利用される代表的な文字/記号とその情報内容をまとめて示す。

表2. WLNで利用される英文字/特殊記号とその情報内容

Q	-OH	E	Br	N	3級窒素
V	ケトン	F	F	K	4級窒素
W	-NO ₂ 、-SO ₂ -	G	Cl	X	4級炭素
M	imino、imido	I	I	Y	3級炭素(メチレン可)
Z	amino、amido	J	ハロゲンの総称的表現、環情報の区切り		
H	水素(通常は省略される)	U	2重結合	UU	3重結合
R	ベンゼン環	L	炭化水素環	T	ヘテロ環及び飽和環
&	切断を示す	-	繰り返す	/	区切り

化合物構造式はこれらの文字/記号を組み合わせることで線型な文字列へと変換される。

て化合物検索等を行うならば、このSMILES名を検索に都合の良いユニーク性をもたせることが必要となる。

ここではSMILESをユニークSMILESとするための手続きを簡単にのべる。

[手順1]

最初に個々の原子について以下の条件を見出す。

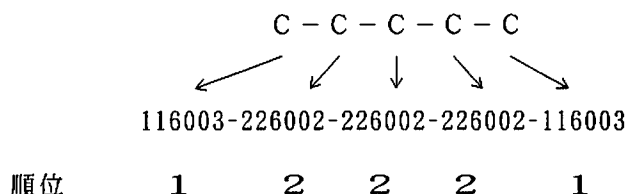
- 1) 結合数
- 2) 非水素原子の結合数
- 3) 原子番号
- 4) チャージの符号 (0, -, +)
- 5) チャージの絶対数 (0-10)
- 6) 結合水素数 (0-9)

例) 末端メチル原子 : 1 1 6 0 0 3

メチレン : 2 2 6 0 0 2

[手順2]

各原子にアサインされた前記数値データを比べて、各原子の順位付けを行う。



[手順3]

先に付けられた順位を基準とし、その隣接順位の数を加えて新たな番号を付ける。この番号を比較し、新たに順位をつけなおす。

順位	1	2	2	2	1
	+2	1+2	2+2	2+1	+2
番号のつけなおし	2	3	4	3	2
新順位	1	2	3	2	1

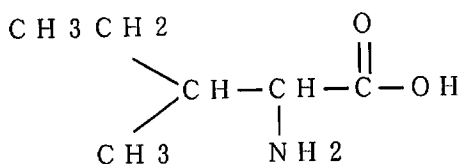
[手順4]

手順3を繰り返す、順位の変化がなくなる、または順位の数が増えた時はその一つ前の順位で停止する。

新順位	1	2	3	2	1
	+2	1+3	2+2	1+3	+2
	2	4	4	4	2
新順位	1	2	2	2	1

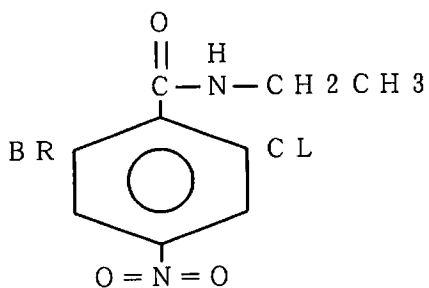
ここで付けられた順位は1と2だけで2種類しか存在しない。従って、3種類の順位が存在していた一つ手前の順位付けを最終順位として採用する。

最終順位	1	2	3	2	1
	C	-	C	-	C
		-	C	-	C
			-	C	-
				-	C



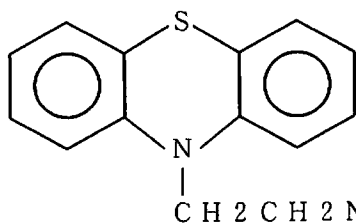
QVYZY2&1

CCC(C)C(N)C(=O)O



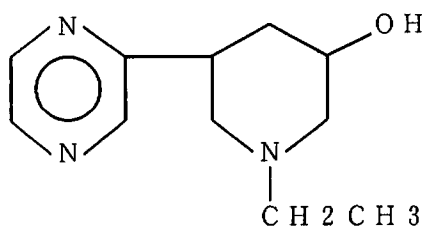
WNR CG EE DVM2

CCNC(=O)c1c(CL)cc(N(=O)=O)c1(BR)



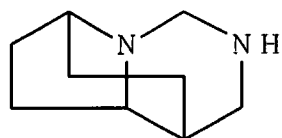
T C666 BN ISJ B2N1&1

C(C)NCCn1c2ccccc2sc3ccccc13



T6N DNJ B- CT6NTJ A2 EQ

CCN1CC(c2cnccn2)CC(O)C1



T566/BH 2AB K AN JMTJ

C1CC2CCC3C1CNCN23

図 . WLN及びSMILESによる化合物構造式の線型表記