

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2020-35115
(P2020-35115A)

(43) 公開日 令和2年3月5日(2020.3.5)

(51) Int.Cl.
G06F 16/00 (2019.01)

F I
G06F 17/30 170F

テーマコード (参考)

審査請求 未請求 請求項の数 11 O L (全 17 頁)

(21) 出願番号 特願2018-159778 (P2018-159778)
(22) 出願日 平成30年8月28日 (2018. 8. 28)

(71) 出願人 504203572
国立大学法人茨城大学
茨城県水戸市文京二丁目1番1号
(74) 代理人 100107766
弁理士 伊東 忠重
(74) 代理人 100070150
弁理士 伊東 忠彦
(72) 発明者 藤芳 明生
茨城県日立市中成沢町四丁目12番1号
国立大学法人茨城大学 工学部内

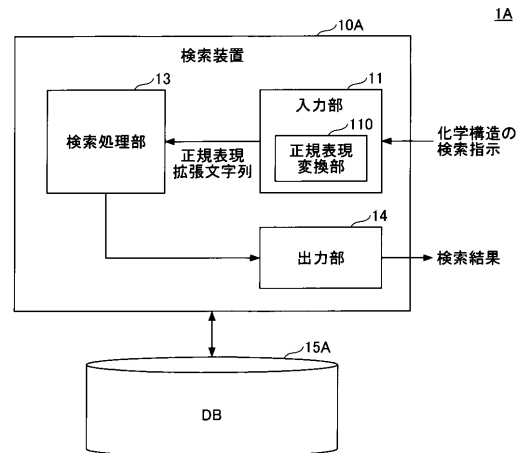
(54) 【発明の名称】 検索システム、検索方法、及び検索プログラム

(57) 【要約】

【課題】簡単かつ柔軟に化学物質を検索することのできる構成と手法を提供する。

【解決手段】検索システムは、化学構造を表わす分子記述言語に正規表現を適用して拡張した表現形式の正規表現拡張文字列を取得する入力部と、前記正規表現拡張文字列をもとに、データベースから該当する化学物質を抽出する検索処理部と、前記検索処理部による検索結果を出力する出力部と、を有する。

【選択図】 図 1



【特許請求の範囲】**【請求項 1】**

化学構造を表わす分子記述言語に正規表現を適用して拡張した表現形式の正規表現拡張文字列を取得する入力部と、

前記正規表現拡張文字列をもとに、データベースから該当する化学物質を抽出する検索処理部と、

前記検索処理部による検索結果を出力する出力部と、

を有することを特徴とする検索システム。

【請求項 2】

前記検索処理部の入力に接続されて、前記正規表現拡張文字列をグラフオートマトンに変換する変換部、

をさらに有し、

前記検索処理部は、前記グラフオートマトンにしたがって、前記グラフオートマトンで定義される状態遷移を満たす化学物質を前記データベースから抽出することを特徴とする請求項 1 に記載の検索システム。

10

【請求項 3】

前記検索処理部は、前記正規表現拡張文字列を先頭から順に読み込み、前記正規表現拡張文字列を、前記化学構造の状態遷移を規定する前記グラフオートマトンに変換することを特徴とする請求項 2 に記載の検索システム。

【請求項 4】

前記グラフオートマトンを保存する記憶部、

をさらに有することを特徴とする請求項 2 または 3 に記載の検索システム。

20

【請求項 5】

前記入力部は、前記正規表現拡張文字列を入力として受け取り、前記正規表現拡張文字列を前記検索処理部へ渡すことを特徴とする請求項 1 に記載の検索システム。

【請求項 6】

前記入力部は、既存の分子記述文字列または化学構造式を入力として受け取り、前記既存の分子記述文字列または前記化学構造式を前記正規表現拡張文字列に変換して、前記検索処理部に渡すことを特徴とする請求項 1 に記載の検索システム。

【請求項 7】

前記検索システムは、ネットワークを介して接続される 1 以上の端末装置を含み、

前記端末装置は、前記入力部の機能の少なくとも一部と、前記出力部の機能の少なくとも一部を有することを特徴とする請求項 1 ~ 6 のいずれか 1 項に記載の検索システム。

30

【請求項 8】

情報処理装置において、化学構造を表わす分子記述言語に正規表現を適用して拡張した表現形式の正規表現拡張文字列を取得し、

前記正規表現拡張文字列をもとに、データベースから該当する化学物質を抽出し、抽出された結果を出力する

工程を含むことを特徴とする検索方法。

【請求項 9】

取得された前記正規表現拡張文字列をグラフオートマトンに変換するステップ、をさらに有し、

前記化学物質の抽出は、前記グラフオートマトンにしたがって、前記グラフオートマトンで定義される状態遷移を満たす化学物質を前記データベースから抽出することを特徴とする請求項 8 に記載の検索方法。

40

【請求項 10】

コンピュータに、

化学構造を表わす分子記述言語に正規表現を適用して拡張した表現形式の正規表現拡張文字列を取得する手順と、

前記正規表現拡張文字列をもとに、データベースから該当する化学物質を抽出する手順

50

と、

抽出された結果を出力する手順と、
を実行させる検索プログラム。

【請求項 11】

取得された前記正規表現拡張文字列をグラフオートマトンに変換する手順、
をさらに有し、

前記化学物質を抽出する手順は、前記グラフオートマトンにしたがって、前記グラフオートマトンで定義される状態遷移を満たす化学物質を前記データベースから抽出することを特徴とする請求項 10 に記載の検索プログラム。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明は、化学構造の検索システム、検索方法、及び検索プログラムに関する。

【背景技術】

【0002】

分子の化学構造を文字列で表現する方法として、SMILES (Simplified Molecular Input Line Entry System)、SMARTS (SMILES Arbitrary Target Specification) などの表現法が用いられている。SMILES は、分子の化学構造を ASCII コードの英数字と記号で文字列化した表記法であり、構造検索やデータ入力などに広く用いられている。

20

【0003】

SMARTS は、SMILES を検索クエリに拡張した表記法である。SMARTS では、結合の種類 (二重結合または三重結合) や基の種類 (芳香族または脂肪族) など、簡単な構造検索の条件を表記することができる (たとえば非特許文献 1 参照)。SMILES や SMARTS の外にも、InChI (Information Chemical Identifier) など、いくつかの表記法がある。

【0004】

一方、文書の検索において、正規表現が利用されている。文書の検索・照合において、文字列の正規表現から有限状態オートマトンを構築し、有限状態オートマトンを用いて文字列のマッチングを行う手法が提案されている (たとえば、特許文献 1 参照)。

30

【先行技術文献】

【特許文献】

【0005】

【特許文献 1】特許第 3852757 号

【非特許文献】

【0006】

【非特許文献 1】Daylight Theory Manual, Daylight version 4.9, Release Date 08/01/11, Daylight Chemical Information Systems, Inc., <http://www.daylight.com/dayhtml/doc/theory/index.html>

【発明の概要】

40

【発明が解決しようとする課題】

【0007】

既存の分子記述言語による表現法では、検索したい化学物質の任意の集合を一つの文字列で表現することが困難である。任意の集合には、化学物質の部分構造、特定の性質を持つ分子の集合、反応により生成される生成物の集合などが含まれる。SMILES で化学構造の集合を表現する場合、その集合に含まれる化学物質の文字列をすべて列挙しなければならない。列挙するには大きなスペースが必要であり、そもそも無限集合は表現できない。類似の構造を持つ化合物の集合を示す表記法としてマルクーシュ構造 (Markush structure) が存在するが、化学構造の表現の制限が大きく、選択肢に含まれる置換基を列挙する自然語の説明文が、別途必要となる。

50

【0008】

本発明は、簡単かつ柔軟に化学物質を検索することのできる構成と手法を提供することを目的する。

【課題を解決するための手段】

【0009】

簡単かつ柔軟な化学構造の検索を実現するために、化学構造を表わす分子記述言語に正規表現を適用して拡張した正規表現拡張文字列を用いる。正規表現拡張文字列をもとに、データベースで該当する化学物質を検索する。

【0010】

本発明の一態様では、検索システムは、

10

化学構造を表わす分子記述言語に正規表現を適用して拡張した表現形式の正規表現拡張文字列を取得する入力部と、

前記正規表現拡張文字列をもとに、データベースから該当する化学物質を抽出する検索処理部と、

前記検索処理部による検索結果を出力する出力部と、
を有する。

【発明の効果】

【0011】

上記の構成により、簡単かつ柔軟に化学構造を検索することができ、検索範囲を拡張することができる。

20

【図面の簡単な説明】

【0012】

【図1】実施形態の検索システムの一例を示す模式図である。

【図2】検索システムの別の例を示す模式図である。

【図3】検索システムの実現に適したハードウェア構成図である。

【図4】実施形態の検索システムの変形例を示す図である。

【図5A】化学構造の正規表現拡張文字列を例示する図である。

【図5B】化学構造の正規表現拡張文字列を例示する図である。

【図6】グラフオートマトンへの変換を説明する図である。

【図7】グラフオートマトンに基づいて検索された検索結果の一例を示す図である。

30

【図8】検索画面の一例を示す図である。

【図9】「最小一致」の検索条件が選択されたときの検索結果の表示例を示す図である。

【図10】グラフオートマトンへの変換の別の例を示す図である。

【図11】データベースに格納されるテーブルの一例を示す図である。

【図12】検索方法のフローチャートである。

【発明を実施するための形態】

【0013】

実施形態では、化学構造を表わす分子記述言語に正規表現を適用して、多様な化学構造とその検索条件を一つの文字列で表現する。この明細書では、化学構造を表わす分子記述言語に正規表現を適用して拡張した表現形式の文字列を、「正規表現拡張文字列」と呼ぶ。分子構造の記述に正規表現を適用することで、化学構造における任意の集合（所定の構造の有無、任意の回数繰り返しの繰返し、置換基どうしの結合/分離など）や検索条件を、1つの文字列で記述することができる。

40

【0014】

正規表現拡張文字列に基づいて化学物質を検索する際に、正規表現拡張文字列をグラフオートマトンに変換して、グラフオートマトンにしたがって検索してもよい。グラフオートマトンを用いることで、完全一致の他に、部分構造検索や、任意の集合の検索が容易になる。

【0015】

グラフオートマトンは、グラフ理論を適用したオートマトンである。オートマトンは、

50

入力に対して、状態に応じた処理を行って次の状態に遷移させる仮想マシンであるが、ここでは、オートマトンにより規定された状態遷移の集合も「オートマトン」に含めることとする。

【0016】

グラフ理論における「グラフ」は、点（頂点又はノード）の集合と、点から延びる辺（枝またはエッジ）の集合で表される構造である。化学物質データベースに保存されている個々の化学物質の構造は、グラフ理論におけるグラフとみなすことができる。データベースに保存されている既知の化学物質が、グラフオートマトンで規定された状態遷移に一致する構造を有するか否かを判定することで、検索が容易になる。

【0017】

図1は、実施形態の検索システム1Aの模式図である。検索システム1Aは、検索装置10Aと、データベース(DB)15Aを含む。検索装置10は、入力部11、検索処理部13、及び出力部14を有する。

【0018】

入力部11は、化学構造の検索指示を入力として受け付ける。検索対象となる化学構造の入力形態は任意である。ユーザによって、直接、正規表現拡張文字列が入力される場合は、入力された文字列はそのまま検索処理部13に入力される。検索対象が、既存の分子記述言語、構造式、マルクーシュ構造などで入力された場合は、正規表現変換部110で正規表現拡張文字列に変換されてから、検索処理部13に入力される。

【0019】

上述したように、正規表現拡張文字列は、分子記述言語に正規表現を適用して拡張された表現形式の文字列である。以下では、分子記述言語としてSMILESを例にとって説明するが、SMARTやその他の分子記述言語に正規表現を適用してもよい。

【0020】

SMILESの主要な表記規則として、次のようなものがある。元素は元素記号で表示され、C, N, O, P, S, Cl, Br, Iに結合する水素は記載されない。二重結合は「=」、三重結合は「#」、分岐は「()」（小括弧）で表される。環の開始と終点となる原子に番号を付けるなどである。たとえば、プロパン(C₃H₈)は「CCC」、イソブタン(C₄H₁₂)は「C(C)CC」、シクロヘキサン(C₆H₁₂)は「C1CCCCC1」と表記される、等である。

【0021】

正規表現は、文字列の集合を単一の文字列で表現することができる。たとえば、「太郎または花子」は「太郎|花子」と記述される。ここで縦棒「|」は、選択肢を区切る表現である。大文字で始まり小文字が続く英単語は、[A-Z][a-z]+と記述される。プラス符号「+」は、直前の表現が1個以上あることを示す。アスタリスク「*」は、直前の表現が0個以上あることを示し、疑問符「?」は直前の表現が0個または1個あることを示す。携帯電話の電話番号は0[7-9]0-[0-9]{4}-[0-9]{4}と表現される。波括弧「{ }」内の数字は、直前の要素の繰り返しの回数を示す。

【0022】

SMILESに正規表現を適用した文字列を「正規表現拡張SMILES」と呼ぶ。正規表現拡張SMILESの定義の例をいくつか挙げる。

- (1) すべてのSMILESの文字列は正規表現拡張SMILESとすることができる。
- (2) 文字列 $w_1 = a b_1 c$, $w_2 = a b_2 c$, ..., $w_n = a b_n c$ を正規表現拡張SMILESの文字列とする。この場合、 $a \{ b_1 | b_2 | \dots | b_n \} c$ は、集合 $\{ w_1, w_2, \dots, w_n \}$ を表わす正規表現拡張SMILESである。
- (3) 文字列 $w = a b c$ を正規表現拡張SMILESの文字列とする。この場合、 $a \{ b \}^* c$ は、集合 $\{ a c, a b c, a b b c, a b b b c, \dots \}$ を表わす正規表現拡張SMILESである。
- (4) 文字列 $w = a b c$ を正規表現拡張SMILESの文字列とする。この場合、 $a \{ b \}^+ c$ は、集合 $\{ a b c, a b b c, a b b b c, \dots \}$ を表わす正規表現拡張SMILESであ

10

20

30

40

50

る。

(5) 文字列 $w = a b c$ を正規表現拡張 S M I L E S の文字列、 i を整数値とする。この場合、 $a\{b\}_i c$ は、集合 $\{ a b^i c \}$ を表わす正規表現拡張 S M I L E S である。

(6) 文字列 $w = a b c$ を正規表現拡張 S M I L E S の文字列、 $i < j$ を整数値とする。この場合、 $a\{b\}_{i,j} c$ は、集合 $\{ a b^i c , a b^{i+1} c , \dots , a b^j c \}$ を表わす正規表現拡張 S M I L E S である。

【0023】

S M A R T S 等の他の分子記述表現を正規表現に拡張する場合も、同様に定義される。正規表現変換部 110 は、入力された分子記述表現に上記の規則を適用して、化学構造の正規表現拡張文字列を生成する。正規表現を適用することで、化学構造の中の特定の部分の繰り返し、特定の部位における選択肢、特定の構成要素の有無、置換基の結合または離脱などを、1つの文字列で記述することができる。さらに、繰り返し回数や選択肢の範囲を、無限、有限を含めて表現することができる。正規表現拡張文字列は、選択肢の範囲や、集合に含まれる要素数が多いほど効果的である。

10

【0024】

標準 S M I L E S などの既存の分子記述言語に替えて、化学構造式が検索対象として入力されたときは、正規表現変換部 110 は、化学構造を文字列化し、この文字列に正規表現を適用する。化学構造の文字列化は、たとえば、化学構造のあるひとつの頂点を選び、その頂点と辺で連結される隣接する頂点を順に選択して符号(原子記号を含む)を与えて文字列化する。環を形成しているところは、切り開いて、グラフ理論の「スパニングツリー」に変換する。このとき、環を切り開いたところにラベル付けをして、連結されていたもの同士を明示することで、文字列にすることができる。

20

【0025】

検索処理部 13 は、外部の化学物質データベースを参照し、また、必要に応じてデータベース 15 A を参照して、該当する化学物質を抽出する。化学物質は、たとえばパターンマッチング等によって抽出される。パターンマッチングには最短一致、最長一致などが含まれてもよい。

【0026】

出力部 14 は、検索処理部 13 による検索結果を出力する。検索結果は、正規表現拡張文字列を関連付けてデータベース 15 A に保存されて、次回以降の検索に利用されてもよい。

30

【0027】

図 2 は、実施形態の検索システム 1 B の模式図である。検索システム 1 B は、検索装置 10 B と、データベース (DB) 15 B を含む。検索装置 10 B は、入力部 11、変換部 12、検索処理部 13、及び出力部 14 を有する。

【0028】

入力部 11 は、化学構造の検索指示を入力として受け付ける。検索対象となる化学構造の入力形態は任意である。ユーザによって、直接、正規表現拡張文字列が入力される場合は、入力された文字列はそのまま変換部 12 に入力される。検索対象が、既存の分子記述言語、構造式、マルクーシュ構造などで入力された場合は、正規表現変換部 110 で正規表現拡張文字列に変換されてから、変換部 12 に入力される。

40

【0029】

変換部 12 は、入力された正規表現拡張 S M I L E S の文字列を、グラフオートマトンに変換する。入力文字列からグラフオートマトンへの変換方法は、後述する。変換されたグラフオートマトンは、正規表現拡張文字列及び化学構造式と対応付けて、データベース 15 B に保存されてもよい。データベース 15 B に保存された情報は、次回以降の変換処理や検索処理に利用されてもよい。

【0030】

検索処理部 13 は、外部の化学物質データベースを参照し、また、必要に応じてデータベース 15 B を参照して、グラフオートマトンで定義される状態遷移を満たす化学物質を

50

抽出する。グラフオートマトンで定義される遷移状態を満たす化学物質は、正規表現拡張 S M I L E S で特定された集合に含まれる物質である。

【 0 0 3 1 】

検索は、たとえばマッチング判定によって行われ、外部のデータベースに格納されている化学物質の任意の頂点（ノード）から順に、グラフオートマトンで定義される状態遷移が満たされるどうかを判定していく。したがって、完全一致だけではなく、部分構造の一致も検索することができる。また、検索指示された集合の中の最小サイズの要素を検索する最小一致や、集合中の最大サイズの要素を検索する最大一致なども指定することができる。

【 0 0 3 2 】

出力部 1 4 は、検索処理部 1 3 による検索結果を出力する。

【 0 0 3 3 】

図 3 は、図 1 の検索システム 1 A、及び/または図 2 の検索システム 1 B の実現に適したハードウェア構成図である。検索システムは、たとえば、ネットワークに接続されたパーソナルコンピュータ（P C）1 0 0 によって実現可能である。P C 1 0 0 は、C P U（Central Processing Unit）1 0 1、主記憶装置 1 0 2、補助記憶装置 1 0 3、入力装置 1 0 4、表示装置 1 0 5、通信インターフェース（I / F）1 0 7、及びドライブ装置 1 0 8 を有し、これらの要素はバス B によって相互に接続されている。

【 0 0 3 4 】

C P U 1 0 1 は、主記憶装置 1 0 2 に格納されたプログラムに従って検索装置 1 0 の動作を制御する。検索装置 1 0 A 及び 1 0 B の検索処理部 1 3 と、検索装置 1 0 B の変換部 1 2 は、C P U 1 0 1 によって実現可能である。

【 0 0 3 5 】

主記憶装置 1 0 2 には、R A M（Random Access Memory）、R O M（Read Only Memory）等が用いられ、C P U 1 0 1 で実行されるプログラム、C P U 1 0 1 による処理に必要なデータ、C P U 1 0 1 の処理で得られたデータ等を記憶又は一時保存する。

【 0 0 3 6 】

補助記憶装置 1 0 3 には、S S D（Solid State Drive）、H D D（Hard Disk Drive）等が用いられ、各種の処理を実行するためのプログラム等のデータが格納される。補助記憶装置 1 0 3 に格納されているプログラムの一部を主記憶装置 1 0 2 にロードし、ロードされたプログラムを C P U 1 0 1 が実行することで、各種の処理が実現される。図 2 の変換部 1 2 によって生成されたグラフオートマトンを保存するデータベース 1 5 B は、補助記憶装置 1 0 3 によって実現されてもよいし、外部のメモリを利用してよい。

【 0 0 3 7 】

入力装置 1 0 4 は、マウス、キーボード等を有し、ユーザが検索装置 1 0 で検索を行うときに必要な情報を入力する。表示装置 1 0 5 は、C P U 1 0 1 の制御のもとに、入力画面、検索結果の出力画面などを含む各種の情報を表示する。P C 1 0 0 がタブレット P C の場合、入力装置 1 0 4 と表示装置 1 0 5 が一体化されたタッチパネル式のディスプレイ（ユーザインターフェース）であってもよい。

【 0 0 3 8 】

通信 I / F 1 0 7 は、ケーブル配線又は無線により、ネットワークを通じて通信を行う。たとえば P C 1 0 0 からネットワークを介して外部の化学物質データベースにアクセスして検索する場合、通信 I / F 1 0 7 によって通信が行われる。

【 0 0 3 9 】

ドライブ装置 1 0 8 は、ドライブ装置 1 0 8 にセットされた C D - R O M（Compact Disc Read-Only Memory）等の記憶媒体 1 9 と、P C 1 0 0 との間のインターフェースをとる。

【 0 0 4 0 】

P C 1 0 0 を検索装置 1 0 として動作させるために、検索プログラムが用いられてもよい。検索プログラムは、C D - R O M 等の記憶媒体 1 0 9 によって P C 1 0 0 に提供され

10

20

30

40

50

てもよいし、通信 I / F 1 0 7 を介してダウンロードされてもよい。P C 1 0 0 にインストールされた検索プログラムは、C P U 1 0 1 によって実行される。

【 0 0 4 1 】

プログラムを保存する記憶媒体 1 0 9 は C D - R O M に限定されず、コンピュータで読み取り可能なデータの構造を有する一時的でない (non-transitory) 有形の (tangible) 媒体であればよい。C D - R O M の他に、D V D (Digital Versatile Disk)、U S B メモリ等の可搬の記録媒体であってもよいし、フラッシュメモリ等の半導体メモリであってもよい。

【 0 0 4 2 】

図 4 は、実施形態の変形例としての検索システム 1 C の模式図である。検索システム 1 C は、サーバ装置 2 0 0 と、ネットワーク 2 を介してサーバ装置 2 0 0 に接続される一つ以上の端末装置 3 A ~ 3 N (適宜、「端末装置 3」と総称する)を含む。サーバ装置 2 0 0 は、検索処理部 2 0 と、データベース 2 5 を有する。サーバ装置 2 0 0 は、ネットワーク 2 を介して外部の化学物質データベース (D B) 2 6 と接続されていてもよい。

10

【 0 0 4 3 】

各端末装置 3 は、図 1 の検索装置 1 0 A (または図 2 の検索装置 1 0 B) の入力部 1 1 の少なくとも一部の機能と、出力部 1 4 の少なくとも一部の機能を果たす。ユーザは、たとえば端末装置 3 のタッチパネルを操作して検索画面を開き、検索対象として所望の化学物質を入力する。検索対象は、正規表現拡張文字列で特定されてもよいし、標準 S M I L E S、S M A R T S、化学構造式、マルクーシュ構造等の、他の表現形式で特定されてもよい。

20

【 0 0 4 4 】

化学物質の検索要求は、化学構造の指定とともに、ネットワーク 2 を介してサーバ装置 2 0 0 に送信される。

【 0 0 4 5 】

サーバ装置 2 0 0 の検索処理部 2 0 は、図 1 の検索装置 1 0 A または図 2 の検索装置 1 0 B と同様の機能を果たす。検索処理部 2 0 は、受信した検索要求から化学物質の正規表現拡張文字列を取り出す。検索対象の化学物質が、その他の表現形式で指定されている場合は、正規表現拡張文字列に変換する。正規表現拡張文字列に基づいて、外部の化学物質データベース 2 6 を検索して、該当する化学物質を抽出する。図 2 のように正規表現拡張文字列をグラフオートマトンに変換する場合は、グラフオートマトンにしたがって、外部の化学物質データベース 2 6 を検索する。グラフオートマトンで定義される遷移状態とマッチングが得られた化学物質が、検索条件に一致する化学物質として抽出される。

30

【 0 0 4 6 】

データベース 2 5 は、図 1 のデータベース 1 5 A、または図 2 のデータベース 1 5 B と同じ機能を果たし、検索の過程で生成された正規表現拡張文字列やグラフオートマトンを、検索により抽出された化学構造式と対応付けて記録する。データベース 2 5 に格納されるデータは、以降の処理で、グラフオートマトンへの変換、マッチング処理等に利用されてもよい。

【 0 0 4 7 】

検索処理部 2 0 による検索結果は、ネットワーク 2 を介して端末装置 3 に送信され、端末装置 3 に表示される。この検索システム 1 C で用いられるサーバ装置 2 0 0 も、図 3 のハードウェア構成で実現可能である。

40

【 0 0 4 8 】

図 5 A と図 5 B は、化学構造の正規表現拡張 S M I L E S を例示する図である。比較例として、標準 S M I L E S の文字列を記載する。

【 0 0 4 9 】

例 1 で、トルエンまたはフェノールを含む物質を検索する場合、ベンゼンの水素原子の一つをメチル基、またはヒドロキシル基で置換したものが検索対象となる。標準 S M I L E S による検索では、トルエンの文字列と、フェノールの文字列を個別に入力するが、正

50

規表現拡張SMILESでは、一つの文字列の中で、選択肢の置換基を「{C|O}」と表現すればよい。

【0050】

例2で、ビシクロヘキシルまたはシクロヘキシルシクロペンタンまたはビシクロペンチルを含む物質を検索する場合、標準SMILESによる検索では、シクロヘキシルとシクロペンタンの結合の方向を含めて4通りの文字列を用いる。これに対し、正規表現拡張SMILESでは、シクロヘキシルとシクロペンタンの4通りの組み合わせを表わすのに、2つの集合を含む文字列を用いるだけでよい。

【0051】

例3で、シクロアルカンの集合を含む物質を検索する場合、標準SMILESでは、メチレン基の数（構造式中のnは0以上の整数）に応じて、すべての構造を表わす文字列を入力するので、無限範囲の表記は不可能である。これに対し、正規表現拡張SMILESでは、0回以上の繰り返しを表わす記号「*」を用いて、「C1C{C}*C1」と表記するだけでよい。ここで、2つの「1」の文字は、分子の同じ位置で連結して炭素の環を形成していることを示すラベルである。

10

【0052】

例4で、ベンゼンとアセン類の集合を含む物質を検索する場合、標準SMILESでは直線状に縮合するベンゼン環の数によって、すべての構造を文字列で特定するので、環の数が多くなるほど入力が高くなる。これに対し、正規表現拡張文字列では、0回以上の繰り返しを表わす「*」を用いて、「c1ccc{c(c1c1)c}*cc1」と表記すればよい。ベンゼン環を含む芳香族の場合、炭素を小文字の「c」で表記している。

20

【0053】

上記の例以外にも、直前の集合の1回以上の繰り返しを含む構造を検索したい場合は「+」を用いて表記すればよいし、所定範囲（i回以上、j回以下）の繰り返しを特定することも可能である。

【0054】

図6は、正規表現拡張SMILESの文字列からグラフオートマトンへの変換を説明する図である。たとえば、例4のベンゼンとアセン類の集合を含む物質を検索する場合、変換部12には、正規表現拡張SMILESの文字列「c1ccc{c(c1c1)c}*cc1」が入力される。この正規表現拡張SMILESの文字列は、図6の上段の化学構造式に対応する。

30

【0055】

変換部12は、入力された文字を先頭から順に読み込んで、開始状態からの状態遷移を規定する。入力された正規表現拡張SMILESの文字列の最初の「c」が開始状態q0となる。この最初の「c」は、構造式の繰り返し部分（角括弧の中）を除く頂点に対応する炭素原子である。一例として、左側のベンゼン環の上側の頂点をq0とする。

【0056】

入力文字列の最初の「c」は直後に数字の「1」を伴い、開始点で環が閉じられることが示されている。入力文字列の「c1」は、グラフオートマトンの状態「q0(c(:,@))」に変換される。「c」は芳香族炭素を表わし、コロン「:」は、1つの芳香族炭素が結合していることを表わす。「@」はq0への戻りパスがあることを示している。

40

【0057】

この例では、説明を簡単にするために結合の種類（単結合、多重結合）を特定していないが、コロン「:」に替えて、単結合を表わす記号「-」や、二重結合を表わす記号「=」を用いて結合の種類を表わしてもよい。

【0058】

入力文字列で、「c1」に続いて「c」が記述されている。開始状態q0には、次の頂点「c」への遷移パスと、反対方向からのq0への戻りパスがある。q0からq1への遷移と、q0への戻りパスの存在は、グラフオートマトンで

q0(c(:,@)) > c(q1(:),p1(@))

50

と規定される。

【0059】

入力文字列が4番目の「c」まで読み込まれると、状態は、q1、q2、q3と順に遷移する。ここまでが、グラフオートマトンの1行目から3行目に規定される遷移である。

【0060】

入力文字列の4番目の「c」の後に、0回以上繰り返される集合「 $\{c(c1c1)c\}^*$ 」が記述されている。したがって、状態q3には、繰り返し回数0のときの遷移(集合の後ろの「c」へのパス)と、繰り返しがあるときの遷移(集合内の「c」へのパス)の2通りの遷移が存在する。グラフオートマトンでは、この2通りの遷移が以下のように規定される。

10

【0061】

$$q3(c(:)) > c(q4(:))$$

$$q3(c(:)) > c(q8(:))$$

状態q3から状態q4への遷移は、構造式の角括弧の中の繰り返しへの遷移である。状態q3から状態q8への遷移は、角括弧の外の頂点への遷移である。

【0062】

状態q4では、分岐が行われる。集合 $\{c(c1c1)c\}$ の中で、分岐を表わす記号「()」が記述されており、括弧内の最初の「c1」への分岐と、括弧の後ろの「c」への分岐である。この分岐は、状態q4から状態q5への遷移パスと、状態q4から状態q7への遷移パスとして、グラフオートマトンで、

20

$$q4(c(:, :)) > c(q5(:), q7(:))$$

と規定される。

【0063】

括弧内の最初の「c1」の後に、もう一つ「c1」が続く。すなわち、状態q5には、入力文字列の最初の「c1」に対応する状態q0に戻るパスp1(:)と、次の「c1」に対応する状態q6への遷移パスがある。

【0064】

これは、グラフオートマトンで

$$q5(c(:, :)) > c(p1(:), p6(:))$$

と規定される。状態q6では、括弧内の2つめの芳香族炭素「c1」への戻りパスを待ち受ける。

30

【0065】

一方、状態q7で、繰り返しが続く場合は、状態q4(集合内の最初の「c」)に戻る。繰り返しが有限回数で指定されている場合は、繰り返しの終了により状態q8(集合の後ろの最初の「c」)に遷移する。これらの遷移は、グラフオートマトンの9行目と10行目に規定される。

【0066】

入力文字列の最後の「c1」は、状態q9として定義され、状態q6への戻りパスp1をたどる(グラフオートマトンの11行目と12行目)。

【0067】

別の例として、図5Bの例3のシクロアルカンの集合を検索する場合は、変換部12は入力文字列「 $C1C\{C\}^*C1$ 」を先頭から読み込んで、以下のようなグラフオートマトンを生成する。

40

【0068】

$$q0(C(-, @)) > C(q1(-), p1(@))$$

$$q1(C(-)) > C(q2(-))$$

$$q1(C(-)) > C(q3(-))$$

$$q2(C(-)) > C(q3(-))$$

$$q2(C(-)) > C(q2(-))$$

$$q3(C(-)) > C(p1(-))$$

50

このように、グラフオートマトンへの変換は、変換部 1 2 に入力された正規表現拡張文字列を先頭から順に読み込み、グラフに見立てた化学構造式の頂点間の遷移状態を規定する処理である。

【0069】

図 7 は、検索処理部 1 3 による検索結果の一例を示す図である。この例では、検索対象として「ベンゼンとアセン類の集合」が指定されており、1 つ以上のベンゼン環が直線状に縮合した構造を含むすべての化学物質が抽出される。検索番号 2 9 3 では、3 つのベンゼン環が直線状に縮合したアセトラセンを含む物質がリストされ、検索番号 3 0 0 では、2 つのベンゼン環が直線状に縮合したナフタレンを含む物質がリストされている。

【0070】

マッチング処理にグラフオートマトンを用いる場合、グラフオートマトンで規定された状態遷移を満たす化学物質が特定される。任意のデータベースに格納されている既存の化学物質のノードを順にたどって、グラフオートマトンで規定された状態遷移が満たされるどうかを判断するので、化学物質に含まれる一部分が、指定された化学物質と一致する場合も、正確に抽出することができる。

【0071】

検索処理部 1 3 による検索結果は、出力部 1 4 によって出力され、たとえば表示装置の表示画面に表示される。図 7 の出力例で、灰色の枠内の番号は検索結果の通し番号、白枠の番号は、使用された化学物質データベースでの登録番号である。

【0072】

図 7 では、「ベンゼンとアセン類の集合」に含まれる物質を有するすべての化学物質が抽出され、表示されている。しかし、最小一致の物質や、最大一致の物質を検索したい場合がある。たとえば、ベンゼン環を 1 つでも含む物質をすべて検索したい場合や、最大数のベンゼン環を含むアセン類を検索したい場合などである。

【0073】

図 8 は、入力部 1 1 のインターフェースである検索画面 1 1 1 の一例を示す。検索対象入力ボックス 1 1 2 の他に、検索条件選択ボックス 1 1 3 が表示されている。ユーザは、検索対象入力ボックス 1 1 2 に、「ベンゼンとアセン類の集合」を表わす正規表現拡張 SMILES 「c1ccc{c(c1c1)c}*cc1」を入力し、検索条件として「最小一致」を選択する。入力部 1 1 への検索対象の入力は、必ずしも正規表現拡張 SMILES でなくてもよいが、正規表現拡張 SMILES で入力する場合は、入力パターンが短く、入力作業が簡単になる。

【0074】

この入力文字列は、変換部 1 2 によってグラフオートマトンに変換されてもよい。検索処理部 1 3 は、グラフオートマトンにしたがって、データベース中の化学物質をひとつずつ調べる。検索条件として「最小一致」が選択されているので、化学物質中に 1 つでもベンゼンが含まれていれば、検索条件に合致すると判定される。たとえば、図 6 のオートマトンで「 $q_5(c(:,:)) \rightarrow C(p_1(:), q_6(:))$ 」まで遷移できたところで、「一致」と判断して、次の化学物質の検索に進んでもよい。

【0075】

検索条件は、図 8 の例に限定されず、「完全一致」、「類似構造検索」などの検索条件を選択可能にしてもよい。「類似構造検索」とは、たとえば、入力文字列で記述される化学構造と、構成元素や置換基の種類、位置等が異なっても構造が類似する化学物質の検索である。化学物質において、グラフオートマトンで規定される状態と元素の種類が異なっても同じ遷移をたどる場合は、類似物質として抽出される。

【0076】

図 9 は、「最小一致」の検索条件が選択されたときの、「ベンゼンまたはアセン類の集合」を含む化学物質の検索結果の表示例である。「最小一致」が検索されたとき、たとえば直線状に縮合する複数のベンゼン環のうち、1 つだけを実線で表示し、他のベンゼン環を破線で表示してもよい。あるいは、1 つのベンゼン環の色を変える等、任意のハイライ

10

20

30

40

50

ト表示が可能である。

【0077】

図10は、グラフオートマトンの別の変換例を示す図である。ここでは「トルエンまたはフェノール」を含む物質(図5Aの例1)を検索する。検索装置10の変換部12は、入力された正規表現拡張SMILESの文字列「{C|O}c1ccccc1」を先頭から順に読み込む。文字列の先頭に、置換基の選択肢を表わす集合{C|O}が記述されているので、グラフオートマトンで開始状態q0と開始状態q1が生成される。開始状態q0は、炭素原子に対応するノードを表わしている。もうひとつの開始状態q1は、酸素原子に対応するノードを表わしている。

【0078】

集合を表す文字列「{C|O}」の直後に、「c」が記述されている。グラフオートマトンで、開始状態q0から状態q2への遷移と、開始状態q1から状態q2への遷移が規定される。状態q2は、芳香族炭素に対応するノードである。

【0079】

入力文字列において、集合の直後の「c」は、後ろに「1」を伴うので、この「c」でベンゼン環が閉じられ、逆方向からのパスがあることが示されている。文字列では、「1」に続いて「c」が記述されている。すなわち、状態q2は、状態q3に遷移可能であるとともに、逆方向からのパスp1が待ち受け可能(「@」)である。グラフオートマトンで、

$$q2(c(:,@)) \rightarrow c(q3(:),p1(@))$$

の遷移が規定される。ここでは、図6のグラフオートマトンと整合をとるために、芳香族炭素cの結合を「:」で示しているが、単結合か二重結合かに応じて「-」と「=」を使い分けてもよい。

【0080】

以下、状態q3から状態q7へと順番に遷移し、入力文字列の最後の「c1」で、状態q7から戻りパスp1でq2に戻る。

【0081】

このようにして生成されたグラフオートマトンは、正規表現拡張文字列及び構造式と対応付けて、データベース15に保存されてもよい。

【0082】

検索処理部13は、データベース15及び/または外部の化学物質データベースを参照して、グラフオートマトンに従って化学物質をひとつずつ調べる。グラフオートマトンで定義される状態遷移を満たす化学物質が抽出され、出力される。

【0083】

図11は、データベース15B(または25)に保存されるグラフオートマトンの記録例を示す。生成されたグラフオートマトンを、正規表現拡張文字列と、化学構造式とに関連付けて保存する。これ以外にも、「ベンゼン」、「シクロアルカン」などという名称や、分子式等に関連付けて保存してもよい。

【0084】

変換部12は、正規表現拡張文字列が入力されたときに、データベース15Bを参照して、すでに対応するグラフオートマトンが保存されている場合は、保存されたグラフオートマトンを読み出して検索処理部13に渡せばよい。これによって検索時間を短縮することができる。

【0085】

図12は、化学構造の検索方法のフローチャートである。この処理フローは、検索システム1Aの検索装置10A、または検索システム1Bの検索装置10Bで実行されてもよいし、検索システム1Cのサーバ装置200によって実行されてもよい。

【0086】

まず、化学構造の正規表現拡張文字列を取得する(S11)。検索装置10A,10B、またはサーバ装置200に、直接、正規表現拡張文字列が入力されてもよいし、その他

10

20

30

40

50

の表現形式で化学構造が特定されている場合は、正規表現拡張文字列に変換する。

【0087】

図1の検索装置10Aを用いる場合は、ステップ13に飛んで検索処理を行う。図2の検索装置10Bを用いる場合は、入力された正規表現拡張文字列を、グラフオートマトンに変換する(S12)。グラフオートマトンへの変換処理は、上述した通り、入力された文字列を最初から順に読み込み、グラフに見たてた化学構造のノードからノードへの遷移を規定する。

【0088】

入力された正規表現拡張文字列、または生成されたグラフオートマトンにしたがって、検索処理を行う(S13)。検索処理では、任意の化学物質データベースに記録されている化学物質について、正規表現拡張文字列に一致するか否か、またはグラフオートマトンで規定される遷移を満たすか否かが判定される。グラフオートマトンで規定される遷移を満たす物質は、検索条件に一致すると判断される。

10

【0089】

最後に、検索結果を出力する(S14)。図1の検索システム1Aまたは図2の検索システム1Bの場合は、検索装置10Aまたは10Bの表示装置等の出力部14に検索結果を表示する。図4の検索システム1Cの場合は、サーバ装置200から端末装置3に、検索結果を送信し、端末装置3の表示画面に検索結果が表示される。

【0090】

化学構造を正規表現拡張文字列で表現することで、置換、集合、繰り返し等を含む多様な構造を一つの文字列で簡単に表現することができる。

20

【0091】

正規表現拡張文字列をグラフオートマトンに変換する場合は、既存のデータベースに保存されている化学物質とのマッチング判定が容易になり、検索速度が速くなる。

【0092】

検索の過程で、生成されたグラフオートマトンを入力文字列と化学式に関連付けて保存することで、グラフオートマトンのデータベースを構築することができる。グラフオートマトンのデータベースを、化学物質の検索に利用することもできる。

【0093】

実施形態の検索をプログラムで実現する場合は、プログラムに記述された以下の手順をコンピュータによって実行する。

30

(a) 化学構造を表わす分子記述言語に正規表現を適用して拡張した正規表現拡張文字列を取得する手順；及び

(b) 正規表現拡張文字列をもとに、データベースから該当する化学物質を抽出する手順。

【0094】

グラフオートマトンを利用する場合は、上記の手順に加えて、正規表現拡張文字列をグラフオートマトンに変換する手順をコンピュータに実行させてもよい。この場合は、グラフオートマトンにしたがってデータベースから該当する化学物質を抽出する。

【0095】

これによって、検索装置10A、検索装置10B、またはサーバ装置200を実現することができる。

40

【符号の説明】

【0096】

1A、1B、1C 検索システム

2 ネットワーク

3A～3N 端末装置

10A、10B 検索装置

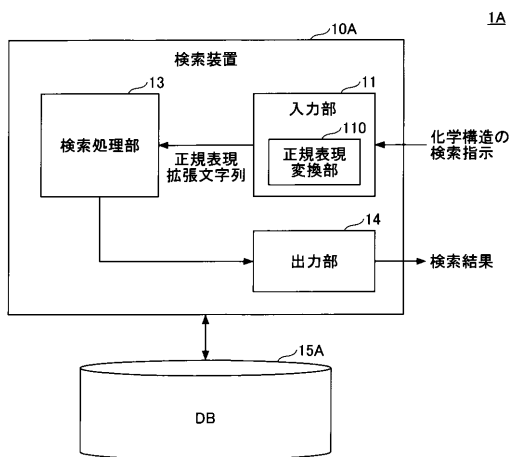
11 入力部

12 変換部

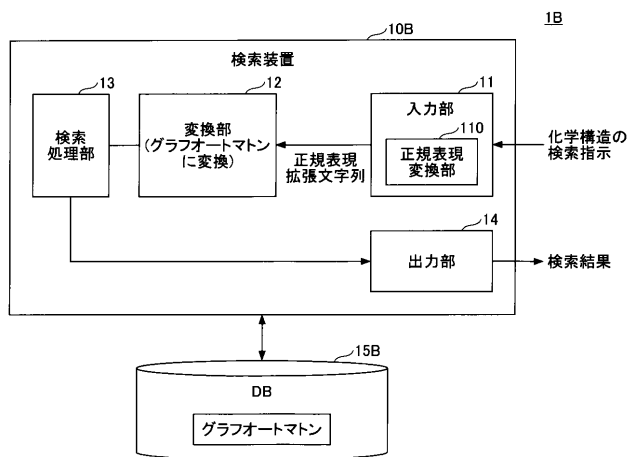
50

- 1 3 検索処理部
- 1 4 出力部
- 1 5 A、1 5 B、2 5 データベース（記憶部）
- 2 0 検索処理部
- 1 1 1 検索画面
- 1 1 2 検索対象入力ボックス
- 1 1 3 検索条件選択ボックス
- 2 0 0 サーバ装置

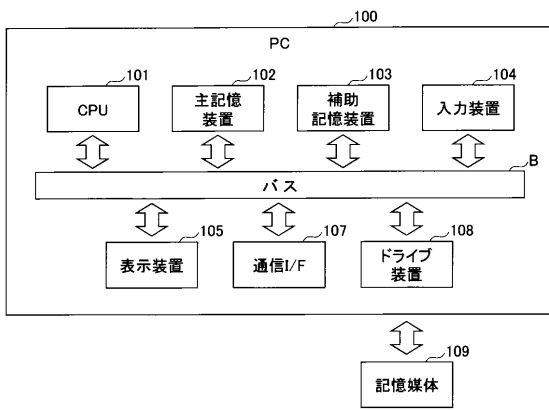
【 図 1 】



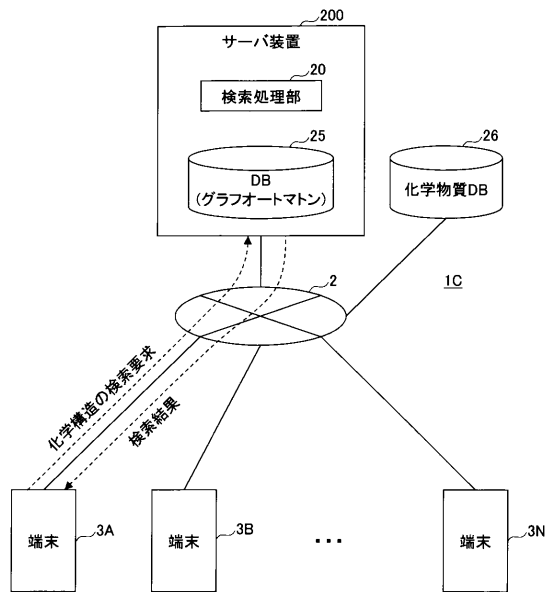
【 図 2 】



【 図 3 】



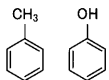
【 図 4 】



【 図 5 A 】

化学構造の正規表現拡張文字列の例

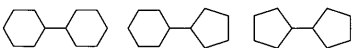
1. トルエンまたはフェノール



実施形態の入力文字列
: {C(O)c1ccccc1}

c.f. 既存の入力文字列
: {Cc1ccccc1, Oc1ccccc1}

2. ビシクロヘキシルまたはシクロヘキシルシクロペンタンまたはビシクロペンチル



実施形態の入力文字列
: {C1CCCC1|C1CCCC1}{C1CCCC1|C1CCCC1}

c.f. 既存の入力文字列
: {C1CCCC1C1CCCC1, C1CCCC1C1CCCC1, C1CCCC1C1CCCC1, C1CCCC1C1CCCC1}

【 図 5 B 】

化学構造の正規表現拡張文字列の例

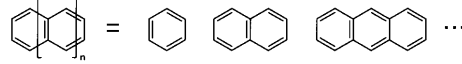
3. シクロアルカンの集合



実施形態の入力文字列
: C1C{C}*C1

c.f. 既存の入力文字列
: {C1CC1, C1CCC1, C1CCCC1, C1CCCC1, ...}

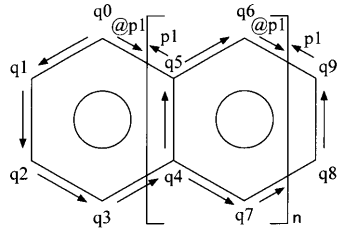
4. ベンゼンとアセン類の集合



実施形態の入力文字列
: c1ccc{c(c1)c}*cc1

c.f. 既存の入力文字列
: {c1ccccc1, c1cccc(c1)c1ccc1, c1cccc(c1)c1cc(c1)c1ccc1, ...}

【 図 6 】

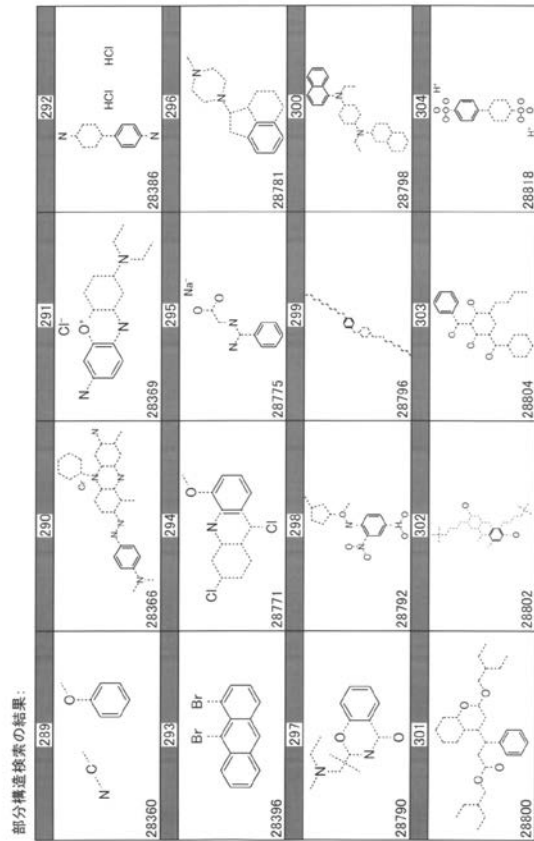


入力された正規表現拡張文字列: c1ccc{c(c1c1)c}*cc1

↓ グラフオートマトンへの自動変換

- q0(c(:,@)) → c(q1(:,p1(@)))
- q1(c(:)) → c(q2(:))
- q2(c(:)) → c(q3(:))
- q3(c(:)) → c(q4(:))
- q3(c(:)) → c(q8(:))
- q4(c(:,:)) → c(q5(:,q7(:)))
- q5(c(:,:)) → c(p1(:,q6(:)))
- q6(c(@)) → c(p1(@))
- q7(c(:)) → c(q4(:))
- q7(c(:)) → c(q8(:))
- q8(c(:)) → c(q9(:))
- q9(c(:)) → c(p1(:))

【 図 7 】



【 図 8 】

● 検索対象入力 111

c1ccc{c(c1c1)c}*cc1 112

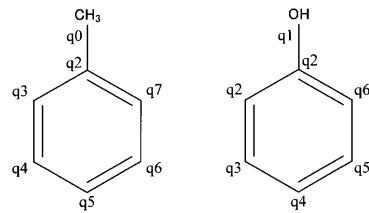
● 検索条件

すべて

最大一致

最小一致 113

【 図 10 】

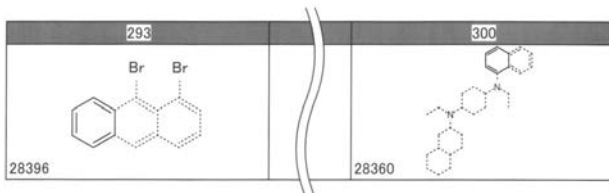


入力文字列: {C|O}c1ccccc1

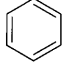
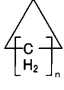
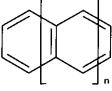
↓ グラフオートマトンへの自動変換

- q0(C(-)) → C(q2(-))
- q1(O(-)) → O(q2(-))
- q2(c(:,@)) → c(q3(:,p1(@)))
- q3(c(:)) → c(q4(:))
- q4(c(:)) → c(q5(:))
- q5(c(:)) → c(q6(:))
- q6(c(:)) → c(q7(:))
- q7(c(:)) → c(p1(:))

【 図 9 】



【 図 1 1 】

正規表現拡張文字列	化学構造式	グラフオートマトン
<chem>c1ccccc1</chem>		$q0(c(:,@)) \rightarrow c(q1(:),p1(@))$ $q1(c(:)) \rightarrow c(q2(:))$ $q2(c(:)) \rightarrow c(q3(:))$ $q3(c(:)) \rightarrow c(q4(:))$ $q4(c(:)) \rightarrow c(q5(:))$ $q5(c(:)) \rightarrow c(q6(:))$
<chem>C1C{C}*C1</chem>		$q0(C(-,@)) \rightarrow C(q1(-),p1(@))$ $q1(C(-)) \rightarrow C(q2(-))$ $q1(C(-)) \rightarrow C(q3(-))$ $q2(C(-)) \rightarrow C(q3(-))$ $q2(C(-)) \rightarrow C(q2(-))$ $q3(C(-)) \rightarrow C(p1(-))$
<chem>c1ccc{c(c1c1)c}*cc1</chem>		$q0(c(:,@)) \rightarrow c(q1(:),p1(@))$ $q1(c(:)) \rightarrow c(q2(:))$ $q2(c(:)) \rightarrow c(q3(:))$ $q3(c(:)) \rightarrow c(q4(:))$ $q3(c(:)) \rightarrow c(q8(:))$ $q4(c(:)) \rightarrow c(q5(:),q7(:))$ $q5(c(:)) \rightarrow c(p1(:),q6(:))$ $q6(c(@)) \rightarrow c(p1(@))$ $q7(c(:)) \rightarrow c(q4(:))$ $q7(c(:)) \rightarrow c(q8(:))$ $q8(c(:)) \rightarrow c(q9(:))$ $q9(c(:)) \rightarrow c(p1(:))$

【 図 1 2 】

