

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3965798号
(P3965798)

(45) 発行日 平成19年8月29日(2007.8.29)

(24) 登録日 平成19年6月8日(2007.6.8)

(51) Int. Cl.

G06F 17/30 (2006.01)

F I

G06F 17/30 140
G06F 17/30 419A

請求項の数 9 (全 52 頁)

(21) 出願番号	特願平10-253427	(73) 特許権者	000005496
(22) 出願日	平成10年9月8日(1998.9.8)		富士ゼロックス株式会社
(65) 公開番号	特開2000-90091(P2000-90091A)		東京都港区赤坂九丁目7番3号
(43) 公開日	平成12年3月31日(2000.3.31)	(74) 代理人	100092152
審査請求日	平成16年2月18日(2004.2.18)		弁理士 服部 毅巖
		(72) 発明者	門馬 敦仁
			神奈川県足柄上郡中井町境430 グリー
			ンテクなかい 富士ゼロックス株式会社内
		審査官	辻本 泰隆
		(56) 参考文献	特開平08-030620(JP, A)
			特開平04-281568(JP, A)
			特開平8-190542(JP, A)
最終頁に続く			

(54) 【発明の名称】 データ処理装置、文書処理装置、データ処理プログラムを記録したコンピュータ読み取り可能な記録媒体、文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体、データ処理方

(57) 【特許請求の範囲】

【請求項1】

有向順序木で表現された情報に対する処理を行うデータ処理装置において、
複数のノードと該ノード間の接続関係とを有向順序木で表現した処理対象情報を格納する情報格納手段と、

有向順序木を構成するノードの内容に関する条件を示す述語と、該述語の条件を満たすノード間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、
該検索条件のいずれかのパターン指定において接続関係における下位のノードの置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位ノードの直下の処理対象のノードを、該処理対象のノードの直下のノード列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段と、

を有することを特徴とするデータ処理装置。

【請求項2】

前記マッチング手段は、前記検索条件のいずれかのパターン指定において接続関係における下位のノードの読み飛ばしを許容することが示されている場合は、前記処理対象情報の有向順序木に対し、該パターン指定で示される接続関係の上位ノードの直下の処理対象のノードを頂点とする部分木の一部を読み飛ばした結果得られる均質化有向順序木を処理対象とすることを特徴とする請求項1記載のデータ処理装置。

10

20

【請求項 3】

前記マッチング手段は、ノードの属性の指定を含む前記検索条件が入力された場合には、前記処理対象のノードに定義されている属性を、該処理対象のノードから置換されたノード列の属性とみなして、前記均質化有向順序木内の前記検索条件に適合する有向順序木の出力を行うことを特徴とする請求項 1 記載のデータ処理装置。

【請求項 4】

前記マッチング手段は、ノード属性の変換規則を含む前記検索条件が入力された場合には、前記検索条件に適合する有向順序木の各ノードに定義されている属性を対応する変換規則に従って変換した結果得られる有向順序木を出力することを特徴とする請求項 1 記載のデータ処理装置。

10

【請求項 5】

有向順序木で表現された構造化文書に対する処理を行う文書処理装置において、複数の文書要素と該文書要素間の接続関係とを有向順序木で表現した処理対象構造化文書を格納する情報格納手段と、

有向順序木を構成する文書要素の内容に関する条件を示す述語と、該述語の条件を満たす文書要素間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位の文書要素の置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象構造化文書の有向順序木に対し、該パターン指定で示される接続関係における上位文書要素の直下の処理対象の文書要素を、該処理対象の文書要素の直下の文書要素列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段と、

20

前記マッチング手段が抽出した有向順序木を論理構造とする構造化文書に対して既定の処理を実行する文書処理手段と、

を有することを特徴とする文書処理装置。

【請求項 6】

有向順序木で表現された情報に対する処理を行うデータ処理プログラムを記録したコンピュータ読み取り可能な記録媒体において、

複数のノードと該ノード間の接続関係とを有向順序木で表現した処理対象情報を格納する情報格納手段、

30

有向順序木を構成するノードの内容に関する条件を示す述語と、該述語の条件を満たすノード間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位のノードの置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位ノードの直下の処理対象のノードを、該処理対象のノードの直下のノード列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段、

としてコンピュータを機能させることを特徴とするデータ処理プログラムを記録したコンピュータ読み取り可能な記録媒体。

40

【請求項 7】

有向順序木で表現された構造化文書に対する処理を行う文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体において、

複数の文書要素と該文書要素間の接続関係とを有向順序木で表現した処理対象構造化文書を格納する情報格納手段、

有向順序木を構成する文書要素の内容に関する条件を示す述語と、該述語の条件を満たす文書要素間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位の文書要素の置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象構造化文書の有向順序木に対し、該パターン指定で示される接続関係における上位文書

50

要素の直下の処理対象の文書要素を、該処理対象の文書要素の直下の文書要素列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段、

前記マッチング手段が抽出した有向順序木を論理構造とする構造化文書に対して既定の処理を実行する文書処理手段、

としてコンピュータを機能させることを特徴とする文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項 8】

有向順序木で表現された情報に対する処理を行うためのデータ処理方法において、情報格納手段が、複数のノードと該ノード間の接続関係とを有向順序木で表現した処理対象情報を格納し、

10

マッチング手段が、有向順序木を構成するノードの内容に関する条件を示す述語と、該述語の条件を満たすノード間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位のノードの置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位ノードの直下の処理対象のノードを、該処理対象のノードの直下のノード列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出する、

ことを特徴とするデータ処理方法。

20

【請求項 9】

有向順序木で表現された情報に対する処理を行うための文書処理方法において、情報格納手段が、複数の文書要素と該文書要素間の接続関係とを有向順序木で表現した処理対象構造化文書を格納し、

マッチング手段が、有向順序木を構成する文書要素の内容に関する条件を示す述語と、該述語の条件を満たす文書要素間の接続関係に関する条件を記述した検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位の文書要素の置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象構造化文書の有向順序木に対し、該パターン指定で示される接続関係における上位文書要素の直下の処理対象の文書要素を、該処理対象の文書要素の直下の文書要素列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出し、

30

文書処理手段が、前記マッチング手段が抽出した有向順序木を論理構造とする構造化文書に対して既定の処理を実行する、

ことを特徴とする文書処理方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は有向順序木で表現された情報に対する処理を行うデータ処理装置および構造化文書を対象とする文書処理装置に関し、特に大量の情報に対して処理を行うデータ処理装置および大量の構造化文書から必要な部分構造を検索して処理する文書処理装置に関する。

40

【0002】

【従来の技術】

オフィスなどで利用される文書の再利用に関する利便性を向上させるために、文書を電子的に管理することが広く行われている。さらに、章、節、段落、図表などの文書要素レベルでの再利用を可能とするために、文書要素と要素間の論理的な関係（論理構造）を伴う文書（構造化文書）による文書管理が先進的なユーザによって行われている。なお、以下単に「文書」といった場合には、すべて構造化文書を指すものとする。

【0003】

図 6 8 は、構造化文書の第 1 の例を示す図である。図 6 8 に示す文書 9 1 の矩形は文書要

50

素を表わし、矩形中の文字列は文書要素の名前を表わす。各文書要素には、「001」～「012」の識別子が付与されている。この図に示すように、構造化文書の論理構造は有向順序木として表現される。また、文書要素に接続された楕円は、文書要素に対応する内容を表わす。

【0004】

なお、図68では、文書要素には名前のみが対応付けられているが、名前以外にもさまざまな属性を対応付けることが広く行われている。構造化文書フォーマットの標準としては、SGML(Standard Generalized Markup Language)とXML(eXtensible Markup Language)が広く知られている。SGMLは1986年にISO(International Standard Organization)の標準として規定され、現在に至るまで主に出版分野における電子文書フォーマットとして利用されてきた。いっぽうXMLは、インターネット上の構造化文書フォーマットの標準としてSGMLの一部の機能を継承しつつインターネット文書フォーマットの実質的な標準であるHTML(Hyper Text Markup Language)による運用の知見を取り入れて、1997年の2月にW3C(World Wide Web Consortium)の勧告として制定された。W3Cは、インターネット関連の標準を制定するための国際的な非営利団体である。

【0005】

XML文書は、図68の文書のように有向順序木で表現される論理構造を持つ。HTML文書も有向順序木で表現される論理構造を備えているが、HTMLは文書の表示処理を主眼において設計されているため、表示処理では不要な要素名や要素間の関係を利用できないように、文書構造に対してあらかじめ制限が加えられている。いっぽう、XML文書では、論理構造が有向順序木で表現されていればよく、要素名や要素間の関係に関する制限を文書作成者自身が定義できる。とはいえ、文書作成者が好き勝手に要素名を決めてしまえば、文書を利用するシステムが文書要素を誤った用途に用いることにもなりかねない。このため、W3Cでは、XMLの拡張として、XML文書の要素名の意味を文書処理システムが一意に特定するための機構の標準化活動を進めている。この標準化が達成されれば、インターネット上の構造化文書を文書要素単位で再利用するためのインフラが整うことになる。以上のことから、インターネット上に存在する大量かつ多様な構造をもつ構造化文書から、文書利用者が必要とする文書要素を検索する必要性が増大することは明らかである。

【0006】

ここで、構造化文書を対象とした検索処理について説明する。例えば、検索対象となる構造化文書として、図68に示した文書91以外に以下のような文書を想定する。

【0007】

図69は、構造化文書の第2の例を示す図である。図69の文書92は、図68の文書91と類似した内容を、別の構造で表したものである。この文書92は、根となる識別子「001」の文書要素の子供の要素として、他の全ての文書要素が接続されている。

【0008】

図68の文書91と図69の文書92から、参考文献エントリに対応する文書要素(以下、エントリ要素)を検索する処理を考える。なお、参考文献エントリとは、図68の"1.momma."や"1.numata..."に対応する文書要素、及び図69の"[1]門馬..."や"[2]沼田..."に対応する文書要素を指すものとする。

【0009】

文書91と文書92のエントリ要素を比較すると、さまざまな相違点が見出される。たとえば、文書91のエントリ要素の名前は"ITEM"であるのに対して文書92のエントリ要素の名前は"PARA"である。また、文書91のエントリ要素はリストを意味するLIST要素でまとめられているのに対し文書92のエントリ要素は文書全体を意味するDOC要素の直下に位置している。

【0010】

このようにインターネット上の構造化文書の論理構造では、同一の用途で用いられるべき文書内容が、多様な構造で表現されるようになる。

構造化文書を対象とした検索では、文書要素の名前、属性や文書要素に対応付けられた内

10

20

30

40

50

容に関する条件だけでなく、文書要素間の接続関係も検索条件に利用することにより、再現率を維持したまま適合率を向上させることができる。例えば特開平7-225770号公報記載のデータ検索装置では、検索対象となる構造化文書の文書要素間の祖孫関係に関する索引をあらかじめ作成しておくことにより、文書要素間の接続関係を利用した高速な検索を可能としている。

【0011】

【発明が解決しようとする課題】

しかし、上述の技術は、検索対象の文書構造の類似度が高い場合には大きな効果を発揮するが、前述の文書91と文書92のように、文書構造上の類似度の比較的低い文書を一括して検索する場合には、文書要素間の祖孫関係に関する条件を利用するのが困難となる。たとえば、文書91と文書92のエントリ要素を一括して検索する場合、両者のエントリ要素に共通する祖孫関係は、エントリ要素が文書全体を意味するDOC要素の子孫であるという関係のみである。上記の技術によりこれらの文書を同時に検索するには、DOC要素の子孫として特定のエントリ要素が存在するという検索条件を入力しなければならない。この場合、DOC要素は文書全体を表しているに過ぎないため、文書要素間の接続関係を利用せずに検索した場合と同様の検索結果しか得ることができず、適合率の向上が図れない。

10

【0012】

しかも、過去に一度参照した文書の検索を行う場合、利用者がその文書の構造を詳細に記憶していることは少ない。そのため、文書要素間の接続関係を用いた検索を行う場合においても、曖昧な記憶に基づいて検索条件が入力される。従来の技術では、検索条件において指定された構造に適合した場合にのみ検出されるため、曖昧な記憶に基づいた検索条件では利用者が探している文書が検出されない可能性がある。

20

【0013】

したがって、利用者が文書を閲覧した際に同等の文書と認識される範囲内の一定の処理を施すことで検索条件に適合するような文書であれば、その文書が検出できるような文書検索処理を行えることが望まれている。

【0014】

ところで、有向順序木で表現されるのは構造化文書だけではない。オブジェクト指向プログラミングにおけるプログラムの関係も有向順序木で表現することができる。従って、有向順序木で表現される各種情報においても同様に、一定の処理を行うことで検索条件に適合する論理構造を、処理対象となる情報の中から抽出することが望まれる。

30

【0015】

本発明はこのような点に鑑みてなされたものであり、有向順序木で表現された情報に対してあいまいな検索条件が入力されても、検索条件に合致する部分構造を抽出することができるデータ処理装置を提供することを目的とする。

【0016】

また、本発明の第2の目的は、構造化文書に一定の処理を施すことにより、利用者により指定された論理構造に合致する部分構造が生成される場合には、その合致した部分構造を抽出することができる文書処理装置を提供することである。

【0017】

また、本発明の第3の目的は、有向順序木で表現された情報に一定の処理を施すことにより、利用者により指定された論理構造に合致する部分構造が生成される場合には、その合致した部分構造を抽出するような処理をコンピュータに行わせることができるデータ処理プログラムを記録したコンピュータ読み取り可能な記録媒体を提供することである。

40

【0018】

また、本発明の第4の目的は、構造化文書に一定の処理を施すことにより、利用者により指定された論理構造に合致する部分構造が生成される場合には、その合致した部分構造を抽出するような処理をコンピュータに行わせることができる文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体を提供することである。

【0019】

50

【課題を解決するための手段】

本発明では上記課題を解決するために、有向順序木で表現された情報に対する処理を行うデータ処理装置において、複数のノードと該ノード間の接続関係とを有向順序木で表現した処理対象情報を格納する情報格納手段と、有向順序木を構成するノードの内容に関する条件を示す述語と、該述語の条件を満たすノード間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位のノードの置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位ノードの直下の処理対象のノードを、該処理対象のノードの直下のノード列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段と、を有することを特徴とするデータ処理装置が提供される。

10

【0020】

このようなデータ処理装置によれば、有向順序木を構成するノードの内容に関する条件を示す述語と、該述語の条件を満たすノード間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位のノードの置き換えを許容することが示されている場合は、マッチング手段により、情報格納手段に格納された処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位ノードの直下の処理対象のノードを、該処理対象のノードの直下のノード列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、均質化有向順序木内の検索条件に適合する有向順序木の抽出処理が行われる。

20

【0023】

また、上記課題を解決するために、有向順序木で表現された構造化文書に対する処理を行う文書処理装置において、複数の文書要素と該文書要素間の接続関係とを有向順序木で表現した処理対象構造化文書を格納する情報格納手段と、有向順序木を構成する文書要素の内容に関する条件を示す述語と、該述語の条件を満たす文書要素間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位の文書要素の置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位文書要素の直下の処理対象の文書要素を、該処理対象の文書要素の直下の文書要素列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段と、前記マッチング手段が抽出した有向順序木を論理構造とする構造化文書に対して既定の処理を実行する文書処理手段と、を有することを特徴とする文書処理装置が提供される。

30

【0024】

このような文書処理装置によれば、有向順序木を構成する文書要素の内容に関する条件を示す述語と、該述語の条件を満たす文書要素間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位の文書要素の置き換えを許容することが示されている場合は、マッチング手段により、情報格納手段に格納された処理対象構造化文書の有向順序木に対し、該パターン指定で示される接続関係における上位文書要素の直下の処理対象の文書要素を、該処理対象の文書要素の直下の文書要素列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、均質化有向順序木内の検索条件に適合する有向順序木の抽出処理が行われる。すると、文書処理手段により、マッチング手段が抽出した有向順序木を論理構造とする構造化文書に対して既定の処理が行われる。

40

【0027】

また、上記課題を解決するために、有向順序木で表現された情報に対する処理を行うデータ処理プログラムを記録したコンピュータ読み取り可能な記録媒体において、複数のノードと該ノード間の接続関係とを有向順序木で表現した処理対象情報を格納する情報格納

50

手段、有向順序木を構成するノードの内容に関する条件を示す述語と、該述語の条件を満たすノード間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位のノードの置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位ノードの直下の処理対象のノードを、該処理対象のノードの直下のノード列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段、としてコンピュータを機能させることを特徴とするデータ処理プログラムを記録したコンピュータ読み取り可能な記録媒体が提供される。

10

【0028】

このような記録媒体に記録されたデータ処理プログラムをコンピュータに実行させれば、有向順序木を構成するノードの内容に関する条件を示す述語と、該述語の条件を満たすノード間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位のノードの置き換えを許容することが示されている場合は、情報格納手段に格納された処理対象情報の有向順序木に対し、該パターン指定で示される接続関係における上位ノードの直下の処理対象のノードを、該処理対象のノードの直下のノード列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、均質化有向順序木内の検索条件に適合する有向順序木を抽出するような処理機能がコンピュータ上に構築される。

20

【0029】

また、上記課題を解決するために、有向順序木で表現された構造化文書に対する処理を行う文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体において、複数の文書要素と該文書要素間の接続関係とを有向順序木で表現した処理対象構造化文書を格納する情報格納手段、有向順序木を構成する文書要素の内容に関する条件を示す述語と、該述語の条件を満たす文書要素間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位の文書要素の置き換えを許容することが示されている場合は、前記情報格納手段に格納された前記処理対象構造化文書の有向順序木に対し、該パターン指定で示される接続関係における上位文書要素の直下の処理対象の文書要素を、該処理対象の文書要素の直下の文書要素列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、前記均質化有向順序木内の前記検索条件に適合する有向順序木を抽出するマッチング手段、前記マッチング手段が抽出した有向順序木を論理構造とする構造化文書に対して既定の処理を実行する文書処理手段、としてコンピュータを機能させることを特徴とする文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体が提供される。

30

【0030】

このような記録媒体に記録された文書処理プログラムをコンピュータに実行させれば、有向順序木を構成する文書要素の内容に関する条件を示す述語と、該述語の条件を満たす文書要素間の接続関係に関する条件を記述したパターン指定とを含む検索条件が入力され、該検索条件のいずれかのパターン指定において接続関係における下位の文書要素の置き換えを許容することが示されている場合は、情報格納手段に格納された処理対象構造化文書の有向順序木に対し、該パターン指定で示される接続関係における上位文書要素の直下の処理対象の文書要素を、該処理対象の文書要素の直下の文書要素列に置換する操作を行った結果得られる均質化有向順序木を処理対象として、均質化有向順序木内の検索条件に適合する有向順序木が抽出され、文書処理手段により、マッチング手段が抽出した有向順序木を論理構造とする構造化文書に対して既定の処理を実行するような処理機能がコンピュータ上に構築される。

40

【0031】

【発明の実施の形態】

以下、本発明の実施の形態を図面を参照して説明する。

50

図 1 は、本発明の原理構成図である。本発明のデータ処理装置は、情報格納手段 1 とマッチング手段 2 とからなる。

【 0 0 3 2 】

情報格納手段 1 は、有向順序木で表現された処理対象情報 1 a を格納する。

マッチング手段 2 は、有向順序木を構成するノードの内容及びノード間の接続関係に関する条件を記述した検索条件 3 が入力されると、情報格納手段 1 に格納された処理対象情報 1 a の有向順序木の中間ノードを削除し、中間ノードのあった位置に中間ノード直下のノード列を配置する操作を行った結果得られる均質化有向順序木を生成する。この均質化有向順序木は複数生成される。そして、生成された均質化有向順序木を処理対象として、均質化有向順序木内の検索条件に適合する有向順序木 4 を抽出する。

10

【 0 0 3 3 】

このようなデータ処理装置によれば、検索条件 3 によって指定したノード間の接続関係が曖昧であっても、処理対象情報 1 a の中間ノードを削除し、その下位のノード列を削除した中間ノードのあった位置に配置することで検索条件に適合する構造が生成できれば、その構造を示す有向順序木が抽出される。したがって、曖昧な記憶に基づいて検索条件を定義しても、利用者の意図した構造の情報を処理対象情報 1 a 内から抽出可能となる。

【 0 0 3 4 】

ところで、有向順序木で表現される情報の主なものとして構造化文書がある。構造化文書に対して本発明に係るデータ処理を行えば、文書構造を検索条件に用いて、比較的類似度の低い構造化文書群に対する検索処理を有効に行うことができる。そこで、このような文書検索を行うことができる文書処理装置を、第 1 の実施の形態として以下に説明する。

20

【 0 0 3 5 】

図 2 は、文書処理システムの構成を示すブロック図である。文書処理システムは、文書処理装置 1 0 と入出力装置 2 0 とからなる。文書処理装置 1 0 は、文書保持部 1 1、検索条件保持部 1 2、階層オートマトン生成部 1 3、マッチング部 1 4、及び文書処理部 1 5 で構成される。

【 0 0 3 6 】

文書保持部 1 1 は、処理対象の構造化文書を保持する。検索条件保持部 1 2 は、検索条件を保持する。本実施の形態における検索条件は、有向順序グラフとして表現される。階層オートマトン生成部 1 3 は、入出力装置 2 0 を介した利用者 3 0 からの検索指令によって、検索条件保持部 1 2 中の検索条件を入力として階層オートマトンを生成する。生成した階層オートマトンは、マッチング部 1 4 に入力される。マッチング部 1 4 は、階層オートマトンと文書保持部 1 1 中の構造化文書を入力として、論理構造のマッチング処理を行う。マッチングの結果として得られる文書要素セット（文書要素の集合）の集合は、文書処理部 1 5 に出力される。文書処理部 1 5 は、文書要素セットの集合を入力として、文書の自動生成や表示、印刷などの文書処理を実行し、入出力装置 2 0 を介して処理結果を利用者 3 0 に提供する。

30

【 0 0 3 7 】

入出力装置 2 0 は、キーボードやマウスなどの入力装置と、C R T (Cathode Ray Tube)、L C D (Liquid Crystal Display) などの表示装置、及びプリンタなどの出力装置からなる。利用者 3 0 が入出力装置 2 0 から指示を入力すると、その指示は文書処理装置に入力される。また、文書処理装置の処理結果が表示装置の画面を通じて文書の利用者 3 0 に通知される。あるいは、処理結果が出力装置を介して紙に印刷される。

40

【 0 0 3 8 】

なお、図 2 の構成と図 1 に示した原理図との対応関係は次のとおりである。すなわち、図 1 の情報格納手段 1 は、文書保持部 1 1 に対応する。マッチング手段 2 は、マッチング部 1 4 に対応する。

【 0 0 3 9 】

図 2 のような構成のシステムにおいて以下のような処理が行われる。

図 3 は、文書処理装置の処理手順を示すフローチャートである。以下の処理をステップ番

50

号に沿って説明する（後述する他のフローチャートにおいても同様）。なお、この処理は、利用者30によって検索指令が入力された際に開始される。検索指令には、検索条件保持部12内のどの検索条件によって検索を行うのかが指定されている。

〔S1〕階層オートマトン生成部13が、利用者30によって指定された検索条件を検索条件保持部12から抽出し、その検索条件に基づいて階層オートマトンを生成する。

〔S2〕マッチング部14が、文書保持部11中の構造化文書と、階層オートマトン生成部13が生成した階層オートマトンとのマッチングを行う。

〔S3〕文書処理部15が、マッチング結果の文書要素の集合を対象とした文書処理を行う。

【0040】

10

次に、図3の各ステップの処理内容を、具体例を交えながら詳細に説明する。この例では、文書保持部11内に、図68に示した文書91に加えて、以下のような文書が保持されているものとする。

【0041】

図4は、構造化文書の第3の例を示す図である。この文書93は、図68の文書91と同様の内容を有しているが、論理構造が異なる。文書91では、参考文献エントリの並びを意味するLIST要素（識別子「007」）が参考文献全体を意味するSECT要素（識別子「005」）直下の子供として表現されている。いっぽう、文書93では、参考文献エントリの並びを意味する要素は存在せず、参考文献エントリはエントリ数が1のLISTリスト要素（識別子「007」「008」）として表現されている。

20

【0042】

なお、以下の説明では、文書要素を指示するときに論理構造上の文書要素間の接続関係を用いることがある。たとえば、「図68の文書91の識別子「002」の文書要素の長子」とは、「図68の文書91の識別子が「002」である文書要素の先頭の子供要素」、すなわち識別子が「003」の文書要素を意味するものとする。また、「図68の文書91の識別子が「002」である文書要素の弟要素」とは「図68の文書の識別子「002」の文書要素と親要素を共有する要素で、識別子「002」の文書要素より後に出現する要素のうち最初に出現する要素」すなわち識別子が「005」の文書要素を意味するものとする。

【0043】

また、検索条件保持部12には、以下のような検索条件が保持されているものとする。

30

図5は、検索条件の例を示す図である。この検索条件40は、破線で示された矩形（検索条件ノード）からなる有向順序グラフとして表現される。この検索条件40は、図68の文書91と図4の文書93に含まれている参考文献エントリを一括して得るための検索条件である。

【0044】

図中実線で示した矩形は、文書要素を引数とする述語41, 43, 45, 47である。この例では、文書要素に対応づけられたテキスト内容がマッチすべきパターン、及び論理演算子が矩形中に示されている。たとえば、述語43に示す「参考文献"or"References」という記述は、「参考文献」または「References」という文字列を含むテキスト内容に対応づけられた文書要素が引数に与えられた場合、真を返す」という条件を意味する。また、述語45に示す「[(数字)*]or(数字)*.」という記述は、「[1]、[2]、・・・」または「1.、2.、・・・」という文字列を含むテキスト内容に対応づけられた文書要素が引数に与えられた場合、真を返す」という条件を意味する。

40

【0045】

述語に対応する文書要素の下部構造に関する条件を設定するために、パターン指定42, 44, 46を述語の直下に記述することができる。パターン指定42, 44, 46は、直上の述語に対応する文書要素直下の文書要素列のパターンを指定する。図5では、パターン指定は楕円で表現されている。図5におけるパターン指定42, 44, 46中の文字列はパターンの種類を意味する。たとえば、「SEQ」は、直下の文書要素列が図で示された順序（左側が上位）で出現することを意味する。また、「OPTREP」は、直下の文書要素が0回

50

以上任意の回数出現することを意味する。

【 0 0 4 6 】

図 5 に示したとおり、パターン指定の直下には述語またはパターン指定の列が出現する。検索条件では、パターン指定と述語の関係を階層的に記述できることから、検索条件を階層的なパターン指定と見なすことができる。さらに図 5 では用いられていないが、直下の文書要素列のいずれか 1 つが一度だけ出現することを意味する "CH0" をパターン指定中の文字列 ("SEQ" および "OPTREP") に追加すれば、検索条件によって、前述の述語からなる任意の階層的な正規表現を表現することが可能となる。

【 0 0 4 7 】

また、パターン指定には特別な記号 "+r" を追加することができる。"+r" は、パターン指定直下の文書要素列を決定する処理において、処理対象の文書要素をその文書要素直下の文書要素列と置き換えて処理を進めてもよいことを意味する。たとえば、処理対象の文書要素が図 6 8 の LIST 要素 (識別子「007」) であった場合、LIST 要素と 2 つの ITEM 要素 (識別子「008」、「009」) からなる要素列を置き換えて処理を進めることができる。もちろん、LIST 要素をそのまま用いて処理を進めることもできる。

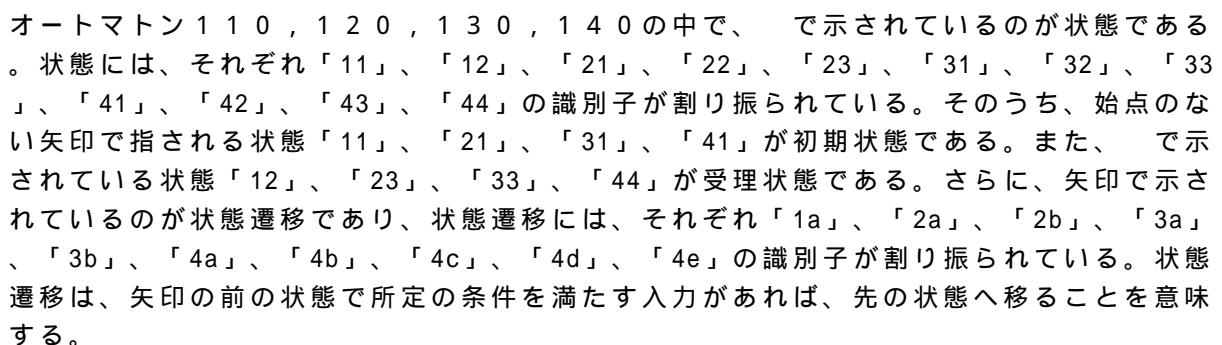
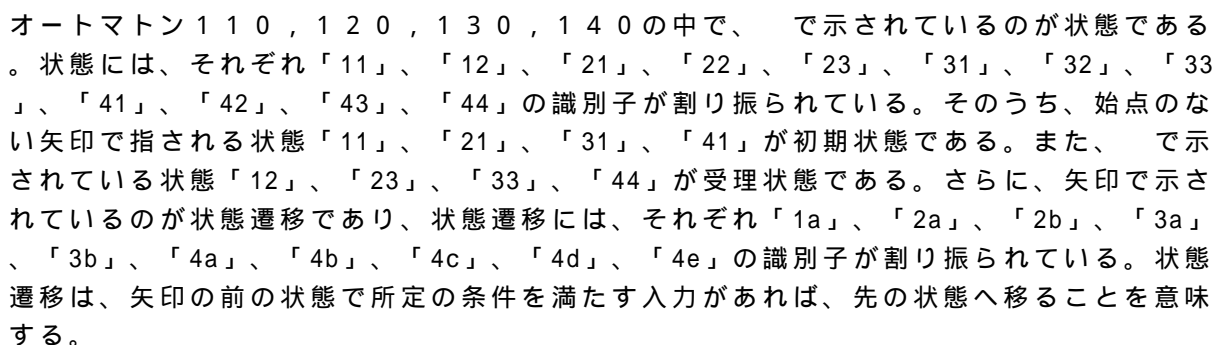
【 0 0 4 8 】

以上のような文書及び検索条件が保持されている状態で、利用者が図 5 に示した検索条件 4 0 に適合する構造を有する文書を検索する場合を考える。その場合、まず階層オートマトン生成部 1 3 により階層オートマトンが生成される。以下に、図 5 の検索条件 4 0 から生成される階層オートマトンを示す。

【 0 0 4 9 】

図 6 は、階層オートマトンの第 1 の例を示す図である。この階層オートマトン 1 0 0 は、4 つのオートマトン 1 1 0, 1 2 0, 1 3 0, 1 4 0 から構成されている。これらのオートマトンの最上位に位置するオートマトン 1 1 0 を、根オートマトンと呼ぶ。

【 0 0 5 0 】

オートマトン 1 1 0, 1 2 0, 1 3 0, 1 4 0 の中で、で示されているのが状態である。状態には、それぞれ「11」、「12」、「21」、「22」、「23」、「31」、「32」、「41」、「42」、「43」、「44」の識別子が割り振られている。そのうち、始点のない矢印で指される状態「11」、「21」、「31」、「41」が初期状態である。また、で示されている状態「12」、「23」、「33」、「44」が受理状態である。さらに、矢印で示されているのが状態遷移であり、状態遷移には、それぞれ「1a」、「2a」、「2b」、「3a」、「3b」、「4a」、「4b」、「4c」、「4d」、「4e」の識別子が割り振られている。状態遷移は、矢印の前の状態で所定の条件を満たす入力があれば、先の状態へ移ることを意味する。

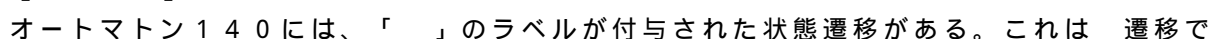
【 0 0 5 1 】

オートマトンの状態遷移には、中央に矩形が記されているものと、矩形の記されていないものとがある。前者は述語を伴う検索条件ノードに対応する状態遷移を意味し、後者は述語を伴わない検索条件ノードに対応する状態遷移を意味する。また、矩形の記述は、状態遷移と文書要素とを対応付けるための条件である。この条件の意味は、検索条件ノードの記述の意味と同一である。図中のオートマトン 1 3 0, 1 4 0 には、右上に "+r" が記されている。この記号は、オートマトン上で文書要素と直下の文書要素列との置き換えが許されることを意味する。

【 0 0 5 2 】

状態遷移からオートマトンへの参照 (オートマトン間の矢印) は、マッチング部 1 4 の処理において状態遷移がマッチング対象の構造化文書の文書要素とマッチしたときに、マッチした文書要素を頂点とする部分木のマッチングに用いるオートマトンへの参照を意味する。部分木のマッチングの結果、参照したオートマトンにおいて受理状態に達すれば、参照した矢印の元に位置する状態遷移の先の状態へ移る。

【 0 0 5 3 】

オートマトン 1 4 0 には、「

10

20

30

40

50

あり、空列を読んで次の状態へ移ることを示している。すなわち、入力が無くても次の状態へ遷移できることを意味する。

【 0 0 5 4 】

以下に、階層オートマトン生成部 1 3 の処理の詳細を説明する。

図 7 は、階層オートマトン生成処理部が行う処理を示すフローチャートである。この処理は、すべて階層オートマトン生成部 1 3 が行う。

[S 1 1] 検索条件の根ノードを引数にして、根オートマトン生成処理を呼び出す。根オートマトン生成処理の結果、変数 T に出力された階層オートマトンが返される。この処理の詳細は、図 8 に示す。

[S 1 2] 変数 T に出力された階層オートマトンをマッチング部 1 4 へ出力し、処理を終了する。 10

【 0 0 5 5 】

図 8 は、根オートマトン生成処理手順を示すフローチャートである。この処理は、全て階層オートマトン生成部 1 3 が行う処理である。この処理の入力は、検索条件の根ノードであり、出力は変数 T に対して出力されたオートマトンである。

[S 2 1] 初期状態と受理状態とを 1 つずつ有するオートマトンを生成し、変数 T に格納する。

[S 2 2] 検索条件の根ノードに保持されている条件を伴う状態遷移を生成する。そして、生成した状態遷移の開始端をステップ S 2 1 で格納したオートマトンの初期状態とし、生成した状態遷移の終了端をステップ S 2 1 で格納したオートマトンの受理状態とする。 20

[S 2 3] 検索条件の根ノードを含む検索条件ノードと状態遷移を引数にして、オートマトン生成処理を呼び出す。この処理の詳細は、図 9 に示す。

【 0 0 5 6 】

図 9 は、オートマトン生成処理手順を示すフローチャートの前半である。この処理は、すべて階層オートマトン生成部 1 3 が行う。この処理の入力は検索条件ノードと状態遷移であり、出力は変数 T へ出力されるオートマトンである。

[S 3 1] 検索条件ノードに対応するオートマトンが、変数 T に含まれるか否かを判断する。含まれる場合は処理をステップ S 3 2 に進め、含まれない場合は処理をステップ S 3 3 に進める。

[S 3 2] 入力の状態遷移から検索条件ノードに対応するオートマトンへの参照を生成し、処理を終了する。 30

[S 3 3] オートマトンを生成し、入力の状態遷移から生成したオートマトンへの参照を生成する。

[S 3 4] 入力の検索条件ノードのパターン指定に "+r" が設定されているか否かを判断する。"+r" が設定されていれば処理をステップ S 3 5 に進め、設定されていなければ処理をステップ S 3 6 に進める。

[S 3 5] 生成したオートマトンに "+r" を設定する。

[S 3 6] 生成したオートマトン上に、初期状態と受理状態を生成する。

[S 3 7] 検索条件ノードのパターン指定の種類を判断する。パターン指定が "SEQ" なら処理をステップ S 3 8 に進め、パターン指定が "CHO" なら処理をステップ S 3 9 に進め、パターン指定が "OPTREP" なら処理をステップ S 4 0 に進める。 40

【 0 0 5 7 】

図 10 は、オートマトン生成処理手順を示すフローチャートの後半である。

[S 3 8] 生成したオートマトンを引数にして SEQ パターン指定生成処理を呼び出す。その後、処理をステップ S 4 1 に進める。SEQ パターン指定生成処理の詳細は、図 11、図 12 において説明する。

[S 3 9] 生成したオートマトンを引数にして CHO パターン指定生成処理を呼び出す。その後、処理をステップ S 4 1 に進める。CHO パターン指定生成処理の詳細は、図 13 において説明する。

[S 4 0] 生成したオートマトンを引数にして OPTREP パターン指定生成処理を呼び出す。 50

その後、処理をステップ S 4 1 に進める。OPTREP パターン指定生成処理の詳細は、図 1 4 において説明する。

[S 4 1] ステップ S 3 8 ~ S 4 0 の各処理の返り値のオートマトンを変数 T に格納する。

【 0 0 5 8 】

図 1 1 は、SEQ パターン生成処理手順を示すフローチャートの前半である。この処理は、全て階層オートマトン生成部 1 3 が行う処理である。この処理の入力はオートマトンと検索条件ノードであり、出力はオートマトンである。

[S 5 1] SEQ パターン指定直下のすべての検索条件ノードを処理したか否かを判断する。すべての検索条件ノードを処理したのであれば、処理をステップ S 6 1 に進め、そうでなければ処理をステップ S 5 2 に進める。

10

[S 5 2] SEQ パターン指定直下の未処理の検索条件ノードのうち、先頭の検索条件ノードを処理対象とする。

[S 5 3] 処理対象の検索条件ノードの位置を判断する。先頭の検索条件ノードであれば処理をステップ S 5 4 に進め、末尾の検索条件ノードであれば処理をステップ S 5 5 に進め、それらのいずれでもない検索条件ノードであれば処理をステップ S 5 6 に進める。

[S 5 4] 状態を 1 つ生成する。そして、初期状態を開始端とし、このステップで生成した状態を終了端とした状態遷移を生成する。その後、処理をステップ S 5 7 に進める。

[S 5 5] 直前に作成された状態を開始端とし、受理状態を終了端とする状態遷移を生成する。その後、処理をステップ S 5 7 に進める。

20

[S 5 6] 状態を 1 つ生成する。そして、直前に生成された状態を開始端として、このステップで生成された状態を終了端とする状態遷移を生成する。

[S 5 7] 処理対象の検索条件ノードに述語が含まれているか否かを判断する。述語が含まれている場合には処理をステップ S 5 8 に進め、そうでない場合には処理をステップ S 5 9 に進める。

[S 5 8] 直前に生成された状態遷移に、処理対象の検索条件ノードの述語を対応付ける。

【 0 0 5 9 】

図 1 2 は、SEQ パターン生成処理手順を示すフローチャートの後半である。

[S 5 9] 処理対象の検索条件ノードにパターン指定が含まれるか否かを判断する。パターン指定が含まれる場合には処理をステップ S 6 0 に進め、そうでない場合には処理をステップ S 5 1 に進める。

30

[S 6 0] 処理対象の検索条件ノードと直前に生成された状態遷移とを引数として、オートマトン生成処理（図 9、図 1 0 に示す）を呼び出す。その後、処理をステップ S 5 1 に進める。

[S 6 1 （図 1 1 に示す）] オートマトンを出力し、処理を終了する。

【 0 0 6 0 】

図 1 3 は、CH0 パターン生成処理手順を示すフローチャートである。この処理は、すべて階層オートマトン生成部 1 3 が行う。この処理の入力はオートマトンと検索条件ノードであり、出力はオートマトンである。

40

[S 7 1] CH0 パターン指定直下のすべての検索条件ノードを処理したか否かを判断する。すべての検索条件ノードを処理していれば処理をステップ S 7 8 に進め、そうでなければ処理をステップ S 7 2 に進める。

[S 7 2] CH0 パターン指定直下の未処理の検索条件ノードのうち、先頭の検索条件ノードを処理対象とする。

[S 7 3] 初期状態を開始端とし、受理状態を終了端とする状態遷移を生成する。

[S 7 4] 処理対象の検索条件ノードに述語が含まれているか否かを判断する。述語が含まれている場合には処理をステップ S 7 5 に進め、述語が含まれていない場合には処理をステップ S 7 6 に進める。

[S 7 5] 直前に生成された状態遷移に処理対象の検索条件ノードの述語を対応付ける。

50

[S 7 6] 処理対象の検索条件ノードにパターン指定が含まれるか否かを判断する。パターン指定が含まれる場合には処理をステップ S 7 7 に進め、パターン指定が含まれない場合には処理をステップ S 7 1 に進める。

[S 7 7] 処理対象の検索条件ノードと、直前に生成された状態遷移とを引数としてオートマトン生成処理（図 9、図 1 0 に示す）を呼び出す。その後、処理をステップ S 7 1 に進める。

[S 7 8] オートマトンを出力し、処理を終了する。

【 0 0 6 1 】

図 1 4 は、OPTREPパターン生成処理手順を示すフローチャートである。この処理は、すべて階層オートマトン生成部 1 3 が行う。この処理の入力はオートマトンと検索条件ノードであり、出力はオートマトンである。

10

[S 8 1] オートマトン上に 2 つの状態（第 1 の状態、第 2 の状態）を生成する。

[S 8 2] 初期状態を開始端として第 1 の状態を終了端とする 遷移、初期状態を開始端として受理状態を終了端とする 遷移、第 2 の状態を開始端として第 1 の状態を終了端とする 遷移、及び第 2 の状態を開始端として受理状態を終了端とする 遷移を生成する。

[S 8 3] 第 1 の状態を開始端として第 2 の状態を終了端とする状態遷移を生成する。

[S 8 4] 処理対象の検索条件ノードに述語が含まれるか否かを判断する。述語が含まれる場合には処理をステップ S 8 5 に進め、述語が含まれない場合には処理をステップ S 8 6 に進める。

[S 8 5] ステップ S 8 3 で生成された状態遷移に、処理対象の検索条件ノードの述語を対応付ける。

20

[S 8 6] 処理対象の検索条件ノードにパターン指定が含まれるか否かを判断する。パターン指定が含まれる場合には処理をステップ S 8 7 に進め、パターン指定が含まれない場合には処理をステップ S 8 8 に進める。

[S 8 7] 処理対象の検索条件ノードと、ステップ S 8 3 で生成された状態遷移とを引数としてオートマトン生成処理（図 9、図 1 0 に示す）を呼び出す。

[S 8 8] オートマトンを出力する。

【 0 0 6 2 】

このような処理を階層オートマトン生成部 1 3 が実行することにより、図 5 に示す検索条件 4 0 から図 6 に示す階層オートマトン 1 0 0 が生成される。生成された階層オートマトン 1 0 0 は、マッチング部 1 4 に渡される。

30

【 0 0 6 3 】

次に、マッチング部 1 4 が行う処理の詳細を説明する。

図 1 5 は、マッチング部の処理手順を示すフローチャートである。この処理は、すべてマッチング部 1 4 によって行われる。また、この処理の入力は階層オートマトンと構造化文書であり、出力は入力 of 構造化文書の文書要素セット（文書要素の集合）の集合である。なお、マッチング部 1 4 では、以下の処理で利用される変数 V を保持する。変数 V には、任意の数の文書要素セットが保持されるものとする。また、以下のステップ S 1 0 3 の処理で利用されているスタックは、文書要素と状態との組を要素とする。

[S 1 0 1] 処理対象の文書要素を構造化文書の根文書要素とする。

40

[S 1 0 2] 階層オートマトンの根オートマトンの初期状態、処置対象の文書要素、空の文書要素セット、及び空のスタックを引数にしてオートマトンマッチング処理を呼び出す。この処理の詳細は、図 1 6 に示す。

[S 1 0 3] 構造化文書の論理構造を前から順に走査した場合に、処理対象の文書要素の次の文書要素が存在するか否かを判断する。次の文書要素が存在する場合には処理をステップ S 1 0 4 に進め、存在しない場合には処理をステップ S 1 0 5 に進める。なお、論理構造上の文書要素の順番とは、親子の文書要素間では親の方が先の順である。兄弟の文書要素間では兄の方が先の順である。さらに、弟の文書要素は、兄の文書要素の子孫にあたる文書要素の最後尾の文書要素の次になる。文書要素の識別子は、論理構造上の順番にしたがって振られている。

50

[S 1 0 4] 処理対象の文書要素の次の文書要素を処理対象とする。

[S 1 0 5] 変数 V の値である文書要素セットを出力する。

【 0 0 6 4 】

図 1 6 は、オートマトンマッチング処理手順を示す第 1 のフローチャートである。この処理は、すべてマッチング部 1 4 が行う。また、この処理の入力は階層オートマトン中のオートマトンの状態、構造化文書の文書要素、文書要素セット、及びスタックであり、出力は変数 V に書き込まれる。なお、入力 of 文書要素として NULL (空の値) を設定することも可能である。

[S 1 1 1] 入力 of 文書要素がスタックの先頭の文書要素であるか否かを判断する。先頭の文書要素であれば処理をステップ S 1 1 2 に進め、そうでない場合は処理をステップ S 1 1 6 に進める。

10

[S 1 1 2] 入力 of 状態とスタック先頭の状態を引数にした到達可能状態特定処理の戻り値である状態の集合にスタック先頭の状態が含まれるか否かを判断する。含まれる場合は処理をステップ S 1 1 3 に進め、そうでない場合は処理を終了する。なお、到達可能状態特定処理の詳細は、図 1 9 に示す。

[S 1 1 3] スタック先頭の文書要素と状態の組を pop し、入力 of 状態を pop した状態と置き換える。

[S 1 1 4] スタックが空であるか否かを判定する。スタックが空の場合は処理をステップ S 1 1 5 へ、そうでない場合は処理をステップ S 1 1 6 へ進める。

[S 1 1 5] 変数 V に入力 of 文書要素セットを書き出す。

20

[S 1 1 6] 入力 of 状態を引数にしてマッチング可能状態遷移特定処理を起動した結果の戻り値の状態遷移を変数 X に格納する。マッチング可能状態遷移特定処理の詳細は、図 2 2 に示す。

[S 1 1 7] 変数 X のすべての状態遷移が処理済みか否かを判断する。すべての状態遷移が処理済の場合は処理を終了し、そうでない場合は処理をステップ S 1 1 8 に進める。

【 0 0 6 5 】

図 1 7 は、オートマトンマッチング処理手順を示す第 2 のフローチャートである。

[S 1 1 8] 変数 X の未処理の状態遷移の 1 つを処理対象とする。

[S 1 1 9] 処理対象の状態遷移と入力 of 文書要素を引数として文書要素マッチング処理を呼び出し、戻り値が真か否かを判断する。戻り値が真なら処理をステップ S 1 2 0 に進め、偽なら処理を S 1 2 7 に進める。なお、文書要素マッチング処理の詳細は、図 2 4 に示す。

30

[S 1 2 0] 入力 of 文書要素セットのコピーを作成し、文書要素セットのコピーに入力 of 文書要素を追加する。

[S 1 2 1] 入力 of 文書要素を引数にして次要素特定処理を起動した結果の戻り値が入力 of 文書要素の長子であるか否かを判定する。長子である場合は処理をステップ S 1 2 2 に進め、そうでない場合は処理をステップ S 1 2 6 に進める。なお、次要素特定処理の詳細は、図 2 5 に示す。

[S 1 2 2] 処理対象の状態遷移がオートマトンを参照しているか否かを判定する。オートマトンを参照している場合は処理をステップ S 1 2 3 へ進め、そうでない場合は処理をステップ S 1 2 6 へ進める。

40

[S 1 2 3] スタックのコピーを生成する。

[S 1 2 4] スタックのコピーに、処理対象の状態遷移の終了端の状態と、入力 of 文書要素を引数とした次要素特定処理の戻り値との組を push する。

【 0 0 6 6 】

図 1 8 は、オートマトンマッチング処理手順を示す第 3 のフローチャートである。

[S 1 2 5] 処理対象の状態遷移が参照するオートマトンの初期状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力 of 文書要素の長子、戻り値の状態、ステップ S 1 2 0 で作成した文書要素セットのコピー、及びステップ S 1 2 4 で作成したスタックのコピーを入力としてオートマトンマッチング処理を呼び出し、処理を

50

ステップ S 1 2 7 に進める。

[S 1 2 6] 処理対象の状態遷移の終了端の状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力 of 文書要素を引数として第要素特定処理を起動した結果の戻り値、戻り値の状態、ステップ S 1 2 0 で作成した文書要素セットのコピー、及び入力 of スタックのコピーを入力としてオートマトンマッチング処理を呼び出す。

[S 1 2 7] 処理対象の状態遷移が属するオートマトンに "+r" が記されているか否か判定する。 "+r" が記されている場合は処理をステップ S 1 2 8 へ進め、そうでない場合は処理をステップ S 1 1 7 に進める。

[S 1 2 8] 入力 of 文書要素を引数にして次要素特定処理を起動した結果の戻り値が入力 of 文書要素の長子であるか否か判定する。入力 of 文書要素の長子である場合は処理をステップ S 1 2 9 に、そうでない場合は処理をステップ S 1 1 7 に進める。

[S 1 2 9] 入力 of 状態、入力 of 文書要素の長子 of 文書要素、入力 of 文書要素セットのコピー、及び入力 of スタックのコピーを入力としてオートマトンマッチング処理を呼び出す。その後、処理をステップ S 1 1 7 に進める。

【 0 0 6 7 】

次に、ステップ S 1 1 2 で行われる到達可能状態特定処理の手順を説明する。

図 1 9 は、到達可能状態特定処理の手順を示すフローチャートである。この処理は、すべてマッチング部 1 4 で行われる。また、この処理の入力は階層オートマトン中のオートマトンの状態の組（第 1 の状態と第 2 の状態）であり、出力は階層オートマトン中のオートマトンの状態の集合である。本処理および本処理で呼び出される処理では、任意の数の状態を保持する共通の変数 Y を保持する。

[S 1 3 1] 入力 of 第 1 の状態と第 2 の状態とを引数にして状態チェック処理を呼び出す。状態チェック処理の詳細は、図 2 0 に示す。

[S 1 3 2] 変数 Y を出力する。

【 0 0 6 8 】

図 2 0 は、状態チェック処理のフローチャートの前半である。この処理は、すべてマッチング部 1 4 で行われる。また、この処理の入力はオートマトンの状態の組（第 1 の状態と第 2 の状態）であり、出力は変数 Y に反映される。

[S 1 4 1] 入力 of 第 1 の状態が変数 Y に含まれているか否か判定する。変数 Y に含まれている場合には処理を終了し、含まれていない場合には処理をステップ S 1 4 2 に進める。

[S 1 4 2] 入力 of 第 1 の状態を変数 Y に追加する。

[S 1 4 3] 入力 of 第 1 の状態と第 2 の状態が同一であるか否か判定する。これらが同一である場合は処理を終了し、そうでない場合は処理をステップ S 1 4 4 に進める。

[S 1 4 4] 入力 of 第 1 の状態を開始端とするすべての状態遷移を処理したか否か判定する。未処理の状態遷移がある場合は処理をステップ S 1 4 5 へ進め、すべての状態遷移が処理済みの場合は処理を終了する。

[S 1 4 5] 未処理の状態遷移を 1 つ選択し、処理対象とする。

[S 1 4 6] 処理対象の状態遷移が 遷移の場合は処理をステップ S 1 4 7 へ進め、述語を伴う状態遷移の場合は処理をステップ S 1 4 4 に進め、いずれでもない場合は処理をステップ S 1 4 8 へ進める。

[S 1 4 7] 処理対象の状態遷移の終了端の状態と第 2 の状態とを引数にして状態チェック処理を呼び出し、処理をステップ S 1 4 9 に進める。

[S 1 4 8] 処理対象の状態遷移が参照するオートマトンの初期状態を引数にして状態チェック処理を呼び出し、処理をステップ S 1 4 9 に進める。

【 0 0 6 9 】

図 2 1 は、状態チェック処理のフローチャートの後半である。

[S 1 4 9] 入力 of 第 1 の状態が受理状態である場合には処理をステップ S 1 5 0 に進め、そうでない場合は処理をステップ S 1 4 4 に進める。

[S 1 5 0] 第 1 の状態を含むオートマトンを参照している状態遷移の終了端の状態と第

10

20

30

40

50

2の状態とを引数にして状態チェック処理を呼び出す。その後、処理をステップS 1 4 4に進める。

【0070】

図22は、マッチング可能状態遷移特定処理の手順を示すフローチャートである。この処理は、すべてマッチング部14が行う。また、この処理の入力はオートマトンの状態の組（第1の状態と第2の状態）であり、出力は状態遷移の集合である。この処理およびこの処理から呼び出される処理では、任意の数の状態を保持する共通の変数Zと、任意の数の状態遷移を保持する共通の状態Wを保持する。

[S 1 6 1] 入力第1の状態と第2の状態を引数にして状態遷移チェック処理を呼び出す。状態遷移チェック処理の詳細は、図23に示す。

[S 1 6 2] 変数Wを出力する。

【0071】

図23は、状態遷移チェック処理の手順を示すフローチャートである。この処理は、すべてマッチング部14が行う。また、この処理の入力はオートマトンの状態であり、出力は変数Zおよび変数Wに反映される。

[S 1 7 1] 入力第1の状態が変数Zに含まれているか否かを判定する。変数Zに含まれている場合は処理を終了し、含まれていない場合には処理をステップS 1 7 2に進める。

[S 1 7 2] 入力第1の状態を変数Zに追加する。

[S 1 7 3] 入力第1の状態と第2の状態が同一であるか否かを判定する。同一である場合は処理を終了し、そうでない場合は処理をステップS 1 7 4に進める。

[S 1 7 4] 入力第1の状態を開始端とするすべての状態遷移を処理したか否かを判定し、処理済みの場合は処理を終了し、そうでない場合は処理をステップS 1 7 5へ進める。

[S 1 7 5] 未処理の状態遷移を1つ選択し、処理対象とする。

[S 1 7 6] 処理対象の状態遷移の種別を判別する。状態遷移が遷移の場合は処理をステップS 1 7 7へ進め、条件を伴う状態遷移の場合は処理をステップS 1 7 9へ進め、いずれでもない場合は処理をステップS 1 7 8へ進める。

[S 1 7 7] 処理対象の状態遷移の終了端の状態と第2の状態を引数にして状態遷移チェック処理を呼び出し、処理をステップS 1 8 0に進める。

[S 1 7 8] 処理対象の状態遷移が参照するオートマトンの初期状態と第1の状態を引数にして状態遷移チェック処理を呼び出し、処理をステップS 1 8 0に進める。

[S 1 7 9] 処理対象の状態遷移を変数Wに追加し、処理をステップS 1 8 0に進める。

[S 1 8 0] 入力第1の状態が受理状態であるか否かを判断し、受理状態の場合には処理をステップS 1 8 1に進め、そうでない場合は処理をステップS 1 7 4に進める。

[S 1 8 1] 第1の状態を含むオートマトンを参照している状態遷移の終了端の状態と第2の状態を引数にして状態遷移チェック処理を呼び出す。その後、処理をステップS 1 7 4に進める。

【0072】

なお、到達可能状態特定処理およびマッチング可能状態遷移特定処理はマッチング部に入力される構造化文書とは無関係に実行できるので、これらの処理をあらかじめ一度だけ実行して処理結果を保持し、オートマトンマッチング処理では保持されている処理結果を適宜利用することにより、同一の引数に対する到達可能状態特定処理およびマッチング可能状態遷移特定処理を繰り返し実行することによるオーバーヘッドを未然に防ぐように構成することも可能である。

【0073】

図24は、文書要素マッチング処理の手順を示すフローチャートである。この処理は、すべてマッチング部14が行う。また、この処理の入力は文書要素と状態遷移であり、出力は真偽値である。

[S 1 9 1] 文書要素に対応付けられた文書内容の文字列が、状態遷移に対応付けられた文字列パターンにマッチするか否かを判断する。文字列パターンにマッチする場合は処理をステップS 1 9 2へ進め、そうでない場合は処理をステップS 1 9 3に進める。

10

20

30

40

50

[S 1 9 2] 真を返し、処理を終了する。

[S 1 9 3] 偽を返し、処理を終了する。

【 0 0 7 4 】

なお、本実施の形態では、状態遷移に対応付けられる条件として文書内容の文字列パターンのみが利用可能と想定したが、状態遷移に対応付けられる条件は、文書要素を入力として真偽値を返す述語であればどのようなものでも構わない。このような条件には、文書要素の名称の完全一致・部分一致・パターンマッチ、文書要素の属性値の完全一致・部分一致・パターンマッチ、属性値の範囲指定、属性値と他の状態遷移にマッチした文書要素の属性値との同値関係・大小関係、および、これらの条件を論理結合子で結合したものが含まれる。

10

【 0 0 7 5 】

図 2 5 は、次要素特定処理の手順を示すフローチャートである。この処理は、すべてマッチング部 1 4 が行う。また、本処理の入力は構造化文書中の文書要素であり、出力は入力の文書要素と同一の構造化文書中の文書要素である。

[S 2 0 1] 入力の文書要素に子供がいるか否か判定する。子供がいる場合は処理をステップ S 2 0 2 へ進め、子供がいない場合は処理を S 2 0 3 へ進める。

[S 2 0 2] 入力の文書要素の長子を出力し、処理を終了する。

[S 2 0 3] 入力の文書要素を引数にして弟要素特定処理を呼び出す。弟要素特定処理の詳細は、図 2 6 に示す。

[S 2 0 4] 弟処理特定処理の返り値を出力し、処理を終了する。

20

【 0 0 7 6 】

図 2 6 は、弟要素特定処理の手順を示すフローチャートである。この処理は、すべてマッチング部 1 4 が行う。また、この処理の入力は構造化文書中の文書要素であり、出力は入力の文書要素と同一の構造化文書中の文書要素である。

[S 2 1 1] 入力の文書要素に弟がいるか否か判定する。弟がいる場合は処理をステップ S 2 1 2 へ進め、弟がいない場合は処理をステップ S 2 1 3 へ進める。

[S 2 1 2] 入力の文書要素の弟要素を出力し、処理を終了する。

[S 2 1 3] 入力の文書要素に親がいるか否か判定する。親がいる場合は処理をステップ S 2 1 4 へ進め、親がいない場合は処理をステップ S 2 1 6 へ進める。

[S 2 1 4] 入力の文書要素の親要素を引数にして弟要素特定処理を呼び出す。

30

[S 2 1 5] ステップ S 2 1 4 で呼び出した弟要素特定処理の返り値を出力し、処理を終了する。

[S 2 1 6] 空 (NULL 値) を出力し、処理を終了する。

【 0 0 7 7 】

以上の処理をマッチング部が行うことにより、文書保持部 1 1 に格納されている文書の中で、階層オートマトン生成部 1 3 が生成した階層オートマトンに適合する文書を抽出することができる。たとえば、図 6 の階層オートマトン 1 0 0 を入力とするオートマトンマッチング処理では、図 6 8 の文書 9 1 と図 4 の文書 9 3 のどちらを入力とする文書としても、参考文献エントリに対応する文書要素を得ることができる。以下に、図 6 8 の文書 9 1 と図 4 の文書 9 3 を対象としたオートマトンマッチング処理の呼び出し関係について説明する。

40

【 0 0 7 8 】

図 6 8 の文書 9 1 に対する処理では、図 6 の識別子「 1 a 」の状態遷移と識別子「 005 」の文書要素がマッチングする場合にのみ参考文献エントリが得られる。識別子「 1 a 」の状態遷移の開始端の状態である識別子「 11 」の状態と識別子「 005 」の文書要素との組を入力とするオートマトンマッチング処理、および、この処理から直接または間接的に呼び出されるオートマトンマッチング処理の呼び出し関係を以下に示す。

【 0 0 7 9 】

図 2 7 は、図 6 8 に示した文書における呼び出し関係を示す第 1 の図である。また、図 2 8 は、図 6 8 に示した文書における呼び出し関係を示す第 2 の図である。

50

【 0 0 8 0 】

図 2 7、図 2 8 のノード 2 0 1 ~ 2 1 5 はオートマトンマッチング処理を意味し、ノード間の実線の矢印は処理の間の直接的な呼び出し関係、ノード間の破線は処理の間の間接的な呼び出し関係を意味する。また、 \square 印は文書要素セットを変数 V に書き込めることを意味し、 \times 印は文書要素セットを変数 V に書き込めないことを意味する。ノード中の 4 つ組は、入力の状態の識別子、入力の文書要素の識別子、入力の文書要素セット、及び入力のスタックを意味する。

【 0 0 8 1 】

オートマトンマッチング処理を実行した場合の呼び出し関係は一般には有向木となるが、オートマトンマッチング処理起動の履歴を内部的に保持し、この履歴を利用して同一の引数に対する処理呼び出しをただ一度に抑えることができる。図 2 7、図 2 8 の有向グラフは、同一の引数に対する処理呼び出しをただ一度に抑えた場合の呼び出し関係を表現している。

10

【 0 0 8 2 】

図 2 7、図 2 8 では、ノード 2 0 4 の処理が実行されるオートマトン（図 6 のオートマトン 1 3 0）に記号 "+r" が設定されているためノード 2 0 5 の処理が呼び出される。この呼び出しにより、識別子「007」の文書要素を読み飛ばし、識別子「008」の文書要素や識別子「009」の文書要素を対象としたオートマトンマッチング処理を呼び出すことを可能としている。この結果、識別子「008」の文書要素と識別子「009」の文書要素を参考文献エントリとして得ることができる。

20

【 0 0 8 3 】

図 4 の文書 9 3 に対する処理でも、図 6 の状態遷移 1 a と識別子「005」の文書要素がマッチングする場合にのみ参考文献が得られる。状態遷移 1 a の開始端の状態である識別子「11」の状態と識別子「005」の文書要素の組を入力とするオートマトンマッチング処理、および、この処理から直接または間接的に呼び出されるオートマトンマッチング処理の呼び出し関係を以下に示す。

【 0 0 8 4 】

図 2 9 は、図 4 に示した文書における呼び出し関係を示す第 1 の図である。また、図 3 0 は、図 4 に示した文書における呼び出し関係を示す第 2 の図である。図 2 9、図 3 0 の表記方法は図 2 7、図 2 8 と同一であり、各ノード 3 0 1 ~ 3 1 9 がオートマトンマッチング処理を示している。

30

【 0 0 8 5 】

図 2 9、図 3 0 では、図 2 7、図 2 8 と同様、ノード 3 0 4 の処理からノード 3 0 5 の処理を呼び出すことによって識別子「007」の文書要素を読み飛ばすことが可能となっている。さらに、ノード 3 0 6 ~ 3 0 9 の処理からノード 3 1 0 ~ ノード 3 1 3 の処理を呼び出すことによって識別子「009」の文書要素を読み飛ばすことを可能としている。この結果、識別子「008」の文書要素と識別子「010」の文書要素を参考文献エントリとして得ることができる。

【 0 0 8 6 】

図 3 1 は、図 6 の階層オートマトンと図 6 8 の文書をマッチング部に入力した結果得られる文書要素セットを示す図である。また、図 3 2 は、図 6 の階層オートマトンと図 4 の文書をマッチング部に入力した結果得られる文書要素セットを示す図である。なお、本実施の形態で示した処理の流れでは、同一の文書要素セットが複数出力されることがあるが、簡単にするため図 3 1、図 3 2 では同一の文書要素セットを 1 つだけ示している。

40

【 0 0 8 7 】

次に、文書処理部 1 5 の処理について説明する。

本実施の形態における文書処理部 1 5 は、マッチング部 1 4 で得られた文書要素セットごとに、処理対象の文書の論理構造を縮退した結果得られる構造化文書から可視化データを生成し、入出力装置 2 0 の C R T ディスプレイに出力するものとする。

【 0 0 8 8 】

50

図 3 3 は、図 6 8 の構造化文書を処理対象としたときに C R T ディスプレイに表示される画面を示す図である。この画面 4 1 0 は、部分構造表示部 4 1 1 と適合内容表示部 4 1 2 とがある。部分構造表示部 4 1 1 には、マッチング部 1 4 で得られた文書要素セットに含まれる文書要素の構造が表示されている。また、適合内容表示部 4 1 2 には、部分構造表示部 4 1 1 に表示された文書要素の内容が表示されている。

【 0 0 8 9 】

図 3 4 は、図 4 の構造化文書を処理対象としたときに C R T ディスプレイに表示される画面を示す図である。なお、図 6 8 の文書 9 1 と図 4 の文書 9 3 とは論理構造が異なるだけで文書の内容は同じであるため、この画面 4 2 0 の部分構造表示部 4 2 1 と適合内容表示部 4 2 2 との表示内容は、図 3 3 の画面 4 1 0 の表示内容と同じである。

10

【 0 0 9 0 】

以上で、第 1 の実施の形態における文書処理装置の処理は終了する。

なお、第 1 の実施の形態で用いた検索条件ノードおよび階層オートマトンでは、パターン指定やオートマトンに記号 "+r" を明示的に表示する表記法を採用したが、他の表記方法を採用してもよい。すなわち、パターン指定直下の文書要素列を決定する処理において、処理対象の文書要素をその文書要素直下の文書要素列と置き換えて処理を進めることができるパターン指定やオートマトンを、他のパターン指定等と区別できればよい。たとえば、上記の例で記号 "+r" を伴わないパターン指定やオートマトンに記号 "-r" を設定し、記号 "+r" を伴うパターン指定やオートマトンには記号を設定しない表記法を採用することもできる。この場合でも、第 1 の実施の形態で示した処理の流れにおいて記号 "+r" を扱う処理を適切に変更することにより、同一の処理の流れを実現することが容易に可能である。

20

【 0 0 9 1 】

また、第 1 の実施の形態のオートマトン生成処理（図 9、図 1 0 に示す）の流れでは、SEQ パターン指定、CHO パターン指定、または OPTREP パターン指定に対応するオートマトンが生成される。しかし、オートマトン生成処理を以下のように実現することで、正規表現に対応するオートマトンを生成することも可能である。その例を、第 1 の実施の形態に関する応用例として以下に説明する。

【 0 0 9 2 】

図 3 5 は、第 1 の実施の形態に関する応用例のオートマトン生成処理手順を示すフローチャートである。この処理は、すべて階層オートマトン生成部 1 3 が行う。

30

[S 3 0 1] 検索条件ノードに対応するオートマトンが変数 T に含まれる場合には処理をステップ S 3 0 2 へ進め、そうでない場合は処理をステップ S 3 0 3 へ進める。

[S 3 0 2] 入力の状態遷移から、検索条件ノードに対応するオートマトンへの参照を生成し、処理を終了する。

[S 3 0 3] オートマトンを生成し、入力の状態遷移からオートマトンへの参照を生成する。

[S 3 0 4] 入力の検索条件ノードのパターン指定に "+r" が設定されているか否か判定し、設定されている場合には処理をステップ S 3 0 5 へ進め、そうでない場合には処理をステップ S 3 0 6 へ進める。

[S 3 0 5] 生成したオートマトンに "+r" を設定する。

40

[S 3 0 6] オートマトンに初期状態と受理状態を 1 つずつ生成する。

[S 3 0 7] 検索条件ノードの正規表現とオートマトンを引数にして属性変換文法評価処理を呼び出す。この処理の詳細は、図 3 6 に示す。

[S 3 0 8] ステップ S 3 0 7 で生成されたオートマトンを変数 T に出力する。

【 0 0 9 3 】

図 3 6 は、属性変換文法評価処理の手順を示すフローチャートである。この処理は、すべて階層オートマトン生成部 1 3 で行われる。また、この処理の入力は正規表現とオートマトンであり、出力はオートマトンである。

[S 3 1 1] 「プログラミング言語処理系」（佐々政孝著、岩波書店）6 5 ページ図 3 . 5 記載の属性変換文法で入力の正規表現を解析し、構文木を構成する。

50

[S 3 1 2] 構文木の根を処理対象とする。

[S 3 1 3] 処理対象に対応する生成規則を上記属性変換文法から特定し、特定した生成規則に対応する操作を行う。

[S 3 1 4] ステップ S 3 1 3 で特定された生成規則が生成規則 (2) であり、かつ、ステップ S 3 1 3 で状態遷移に対応付けられた述語と同一の検索条件ノードに正規表現が含まれる場合には処理をステップ S 3 1 5 へ進め、そうでない場合には処理をステップ S 3 1 6 へ進める。

[S 3 1 5] ステップ S 3 1 3 で生成された状態遷移と、ステップ S 3 1 3 で状態遷移に対応付けられた述語と同一の検索条件ノードに含まれる正規表現を引数にして、図 3 5 のオートマトン生成処理を呼び出す。

[S 3 1 6] 構文木を前順に走査したときに処理対象の次に走査される記号が存在すれば処理をステップ S 3 1 7 へ進め、そうでなければステップ S 3 1 8 に進む。

[S 3 1 7] ステップ S 3 1 6 の記号を処理対象とし、処理をステップ S 3 1 3 に進める。

[S 3 1 8] オートマトンを出力し、処理を終了する。

【 0 0 9 4 】

このようなオートマトン生成処理を用いることで、検索条件ノードのパターン指定として任意の正規表現を用いることができる。ただし、この場合の検索条件ノードには常に述語が含まれる。

【 0 0 9 5 】

さらに、図 3 6 に示した属性変換文法評価処理を以下のように実現することにより、検索条件ノードのパターン指定として任意の正規表現を用いることができるだけでなく、検索条件中に述語を含まない検索条件ノードを用いることが可能となる。

【 0 0 9 6 】

図 3 7 は、属性変換文法評価処理の変形例を示す図である。この処理は、図 3 6 の処理に代えて行われ、すべての処理が階層オートマトン生成部 1 3 で行われる。また、この処理の入力は正規表現とオートマトンであり、出力はオートマトンである。

[S 3 2 1] 「プログラミング言語処理系」(佐々政孝著、岩波書店) 6 5 ページ図 3 . 5 記載の属性変換文法で入力 of 正規表現を解析し、構文木を構成する。

[S 3 2 2] 構文木の根を処理対象とする。

[S 3 2 3] 処理対象に対応する生成規則を上記属性変換文法から特定し、特定した生成規則に対応する操作を行う。

[S 3 2 4] ステップ S 3 2 3 で状態遷移に対応付けられた検索条件ノードに正規表現が含まれる場合には処理をステップ S 3 2 5 へ進め、そうでない場合には処理をステップ S 3 2 6 へ進める。

[S 3 2 5] ステップ S 3 2 3 で生成された状態遷移と、ステップ S 3 2 3 で状態遷移に対応付けられた述語と同一の検索条件ノードに含まれる正規表現を引数にして、図 3 5 のオートマトン生成処理を呼び出す。

[S 3 2 6] 構文木を前順に走査したときに処理対象の次に走査される記号が存在するかどうかを判断する。存在すれば処理をステップ S 3 2 7 へ進め、そうでなければ処理をステップ S 3 2 8 に進める。

[S 3 2 7] ステップ S 3 2 6 の記号を処理対象とし、処理をステップ S 3 2 3 に進める。

[S 3 2 8] オートマトンを出力し、処理を終了する。

【 0 0 9 7 】

次に、第 2 の実施の形態について説明する。第 2 の実施の形態における処理対象の構造化文書は、図 6 9 の構造化文書である。第 1 の実施の形態で示した文書処理装置に図 5 の検索条件と図 6 9 の文書 9 2 を入力しても、識別子「002」の文書要素が検索条件ノード中の述語にマッチしないために参考文献エントリを得ることができない。

【 0 0 9 8 】

10

20

30

40

50

そこで、図 6 9 の文書 9 2 に対しても所望の結果を得るため、本実施の形態の文書処理装置では、検索条件のパターン指定に "+p" を設定することができるようにする。 "+p" は、このパターン指定と文書要素とのマッチング中に、マッチング対象の文書要素を読み飛ばしてもよいことを意味する。

【 0 0 9 9 】

第 2 の実施の形態で用いる検索条件を以下に示す。

図 3 8 は、 "+p" を用いて検索条件の例を示す図である。この検索条件 5 0 は、図 5 の検索条件 4 0 のパターン指定 4 2 に記号 "+p" を追加し、パターン指定 4 4 およびパターン指定 4 6 から記号 "+r" を取り除いたものであり、文書要素を引数とする述語 5 1 , 5 3 , 5 5 , 5 7 と、パターン指定 5 2 , 5 4 , 5 6 とで表されている。

10

【 0 1 0 0 】

" +p " の追加により、述語 5 3 にマッチする文書要素の前に出現する文書要素群、述語 5 3 にマッチする文書要素とパターン指定 5 4 にマッチする文書要素列の間の文書要素群、及びパターン指定 5 4 にマッチする文書要素列の後に出現する文書要素群を読み飛ばしてマッチングを行うことが可能となる。具体的には、図 6 9 の文書に対する処理では識別子が「 002 」、「 003 」、「 007 」、「 008 」の文書要素を読み飛ばして所望の文書要素を得ることができる。

【 0 1 0 1 】

第 2 の実施の形態に係る文書処理装置を実現するための構成要素は、第 1 の実施の形態と同様である。そのため、図 2 に示した構成を用いて第 2 の実施の形態を説明する。第 2 の実施の形態では、 "+p" を含む検索条件を処理するため、本実施の形態の文書処理装置では、第 1 の実施の形態の文書処理装置のオートマトン生成処理とオートマトンマッチング処理に以下の変更を施した処理が実行される。

20

【 0 1 0 2 】

図 3 9 は、第 2 の実施の形態におけるオートマトン生成処理の手順を示すフローチャートである。これは、第 1 の実施の形態の図 9、図 1 0 の処理に代えて階層オートマトン生成部 1 3 で行われる。

[S 4 0 1] 検索条件ノードに対応するオートマトンが変数 T に含まれるか否かを判断する。変数 T に含まれる場合は処理をステップ S 4 0 2 へ進め、そうでない場合は処理をステップ S 4 0 3 へ進める。

30

[S 4 0 2] 入力の状態遷移から、検索条件ノードに対応するオートマトンへの参照を生成し、処理を終了する。

[S 4 0 3] オートマトンを生成し、入力の状態遷移からオートマトンへの参照を生成する。

[S 4 0 4] 入力の検索条件ノードのパターン指定に "+p" が設定されているか否かを判定し、設定されている場合には処理をステップ S 4 0 5 へ進め、そうでない場合には処理をステップ S 4 0 6 へ進める。

[S 4 0 5] 生成したオートマトンに "+p" を設定する。

[S 4 0 6] オートマトン上に初期状態と受理状態を生成する。

[S 4 0 7] 検索条件ノードのパターン指定の種別を判別する。パターン指定が "SEQ" なら処理をステップ S 4 0 8 へ進め、 "CH0" なら処理をステップ S 4 0 9 へ進め、 "OPTREP" なら処理をステップ S 4 1 0 へ進める。

40

[S 4 0 8] オートマトンを引数にして SEQ パターン指定生成処理を呼び出す。その後、処理をステップ S 4 1 1 に進める。なお、SEQ パターン指定生成処理の詳細は、図 1 1、図 1 2 に示した処理と同様である。

[S 4 0 9] オートマトンを引数にして CH0 パターン指定生成処理を呼び出す。その後、処理をステップ S 4 1 1 に進める。なお、CH0 パターン指定生成処理の詳細は、図 1 3 に示した処理と同様である。

[S 4 1 0] オートマトンを引数にして OPTREP パターン指定生成処理を呼び出す。その後、処理をステップ S 4 1 1 に進める。なお、OPTREP パターン指定生成処理の詳細は、図 1

50

4 に示した処理と同様である。

[S 4 1 1] 各生成処理の戻り値のオートマトンを変数 T に格納する。

【 0 1 0 3 】

図 3 8 の検索条件 5 0 を第 2 の実施の形態のオートマトン生成処理に入力した結果出力される階層オートマトンを以下に示す。

図 4 0 は、第 2 の実施の形態の文書処理装置により生成された階層オートマトンを示す図である。階層オートマトン 5 0 0 は、4 つのオートマトン 5 1 0 , 5 2 0 , 5 3 0 , 5 4 0 からなる。各オートマトン 5 1 0 , 5 2 0 , 5 3 0 , 5 4 0 の状態と、状態遷移とには、図 6 の例と同様の識別子が割り振られている。

【 0 1 0 4 】

この階層オートマトン 5 0 0 は、オートマトン 5 2 0 の右上に "P" が設定されている。この記号は、オートマトン上でマッチング対象の文書要素を読み飛ばしてもよいことを意味する。

【 0 1 0 5 】

第 2 の実施の形態におけるオートマトンマッチング処理は、以下のように変更される。

図 4 1 は、第 2 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの前半である。この処理は、第 1 の実施の形態の図 1 6 ~ 図 1 8 の処理に代えて行われる。

[S 4 2 1] 入力 of 文書要素がスタック先頭の文書要素であるか否かを判断する。スタック先頭の文書要素であれば処理をステップ S 4 2 2 へ進め、そうでない場合は処理をステップ S 4 2 6 に進める。

[S 4 2 2] 入力の状態とスタック先頭の状態を引数にした到達可能状態特定処理の戻り値である状態の集合に、スタック先頭の状態が含まれるか否かを判断する。スタック先頭の状態が含まれる場合は処理をステップ S 4 2 3 に進め、そうでない場合は処理を終了する。

[S 4 2 3] スタック先頭の文書要素と状態の組を pop し、入力の状態を pop した状態と置き換える。

[S 4 2 4] スタックが空であるか否かを判定し、空の場合は処理をステップ S 4 2 5 へ進め、そうでない場合は処理をステップ S 4 2 6 へ進める。

[S 4 2 5] 変数 V に入力 of 文書要素セットを書き出す。

[S 4 2 6] 入力 of 状態を引数にしてマッチング可能状態遷移特定処理を起動した結果 of 戻り値 of 状態遷移を変数 X に格納する。マッチング可能状態遷移特定処理の詳細は、図 2 2 に示す通りである。

[S 4 2 7] 変数 X のすべての状態遷移が処理済みの場合は処理を終了し、そうでない場合は処理をステップ S 4 2 8 に進める。

[S 4 2 8] 変数 X の未処理 of 状態遷移の 1 つを処理対象とする。

[S 4 2 9] 処理対象 of 状態遷移と入力 of 文書要素を引数として文書要素マッチング処理を呼び出し、戻り値 of 値を判断する。戻り値が真なら処理をステップ S 4 3 0 に進め、戻り値が偽なら処理をステップ S 4 3 7 に進める。なお、文書要素マッチング処理の詳細は、図 2 4 に示す通りである。

[S 4 3 0] 入力 of 文書要素セット of コピーを作成し、文書要素セット of コピーに入力 of 文書要素を追加する。

【 0 1 0 6 】

図 4 2 は、第 2 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの後半である。

[S 4 3 1] 入力 of 文書要素を引数にして次要素特定処理を起動した結果 of 戻り値が入力 of 文書要素 of 長子であるか否かを判定する。長子である場合は処理をステップ S 4 3 2 に進め、そうでない場合は処理をステップ S 4 3 6 に進める。なお、次要素特定処理の詳細は、図 2 5 に示す通りである。

[S 4 3 2] 処理対象 of 状態遷移がオートマトンを参照しているか否かを判定し、参照して

10

20

30

40

50

いる場合は処理をステップ S 4 3 3 へ進め、そうでない場合は処理をステップ S 4 3 6 へ進める。

〔 S 4 3 3 〕 スタックのコピーを作成する。

〔 S 4 3 4 〕 スタックのコピーに、処理対象の状態遷移の終了端の状態と、入力 of 文書要素を引数とした第要素特定処理の戻り値の組を push する。

〔 S 4 3 5 〕 処理対象の状態遷移が参照するオートマトンの初期状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力 of 文書要素の長子と、戻り値の状態と、ステップ S 4 3 0 で作成した文書要素セットのコピーと、ステップ S 4 3 4 で作成したスタックのコピーを入力としてオートマトンマッチング処理を呼び出す。その後、処理をステップ S 4 3 7 に進める。

10

〔 S 4 3 6 〕 処理対象の状態遷移の終了端の状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力 of 文書要素を引数として第要素特定処理を起動した結果の戻り値と、戻り値の状態と、ステップ S 4 3 0 で作成した文書要素セットのコピーと、入力 of スタックのコピーを入力としてオートマトンマッチング処理を呼び出す。

〔 S 4 3 7 〕 処理対象の状態遷移が属するオートマトンに "+p" が記されているか否か判定し、記されている場合は処理をステップ S 4 3 8 へ進め、そうでない場合は処理をステップ S 4 2 7 に進める。

〔 S 4 3 8 〕 入力 of 状態、入力 of 文書要素を第要素特定処理に入力した場合の戻り値 of 文書要素、入力 of 文書要素セットのコピー、入力 of スタックのコピーを入力としてオートマトンマッチング処理を呼び出す。その後、処理をステップ S 4 2 7 に進める。

20

【 0 1 0 7 】

第 2 の実施の形態のマッチング部 1 4 に図 6 9 の構造化文書と図 4 0 の階層オートマトンを入力したときのオートマトンマッチング処理の呼び出し関係を以下に示す。

【 0 1 0 8 】

図 4 3 は、第 2 の実施の形態におけるオートマッチング処理の呼び出し関係を示す第 1 の図である。また、図 4 4 は、第 2 の実施の形態におけるオートマッチング処理の呼び出し関係を示す第 2 の図である。図 4 3、図 4 4 の表記法は図 2 7、図 2 8 の表記法と同一であり、ノード 6 0 1 ~ 6 1 4 はオートマトンマッチング処理を意味する。なお、簡単のため、図 4 3、図 4 4 では文書要素セット (001,004,005,006) の出力に至る呼び出し関係のみを示し、文書要素セット (001,004) および文書要素セット (001,004,005) の出力に至る呼び出し関係は省略してある。

30

【 0 1 0 9 】

図 4 3、図 4 4 では、ノード 6 0 2 の処理が実行されるオートマトン (図 4 0 のオートマトン 5 2 0) に記号 "+p" が設定されているため、ノード 6 0 3 の処理およびノード 6 0 4 の処理が呼び出される。この呼び出しにより、識別子が「 002 」の文書要素および識別子が「 003 」の文書要素を読み飛ばし、識別子が「 004 」、「 005 」、「 006 」の文書要素を対象としたオートマトンマッチング処理を呼び出すことを可能としている。さらに、ノード 6 1 1 の処理が実行されているオートマトン (図 4 0 のオートマトン 5 2 0) に記号 "+p" が設定されているため、ノード 6 1 2 の処理およびノード 6 1 4 の処理が呼び出される。この呼び出しにより、識別子が「 007 」の文書要素と「 008 」の文書要素を読み飛ばすことができる。この結果、識別子が「 005 」の文書要素と「 006 」の文書要素を参考文献エントリとして得ることができる。

40

【 0 1 1 0 】

図 4 5 は、図 6 9 の構造化文書と図 3 8 の検索条件を第 2 の実施の形態 of 文書処理装置に入力した結果得られるマッチング部の出力を示す図である。このように、第 2 の実施の形態によれば、図 6 9 のような文書の中からも、図 3 8 の検索条件に適合する構造を抽出することができる。

また、第 1 の実施の形態 of 応用例 (図 3 5 に示す) と同様の変更を、第 2 の実施の形態に対して行うこともできる。そのような、第 2 の実施の形態に対する応用例を以下に示す。

【 0 1 1 1 】

50

図46は、第2の実施の形態の応用例に係るオートマトン生成処理のフローチャートである。これは、図39の処理に代えて階層オートマトン生成部13で行われる処理である。
[S501] 検索条件ノードに対応するオートマトンが変数Tに含まれるか否かを判断する。変数Tに含まれる場合には処理をステップS502へ進め、そうでない場合は処理をステップS503へ進める。

[S502] 入力の状態遷移から、検索条件ノードに対応するオートマトンへの参照を生成し、処理を終了する。

[S503] オートマトンを生成し、入力の状態遷移からオートマトンへの参照を生成する。

[S504] 入力の検索条件ノードのパターン指定に"+p"が設定されているか否かを判定する。"+p"が設定されている場合には処理をステップS505へ進め、そうでない場合には処理をステップS506へ進める。 10

[S505] 生成したオートマトンに"+p"を設定する。

[S506] オートマトンに初期状態と受理状態を1つずつ生成する。

[S507] 検索条件ノードの正規表現とオートマトンを引数にして属性変換文法評価処理を呼び出す。

[S508] ステップS507で生成されたオートマトンを変数Tに出力し、処理を終了する。

【0112】

このようにして、第1の実施の形態の応用例と同様の変更を、第2の実施の形態に対して行うことができる。 20

次に、第3の実施の形態について説明する。第3の実施の形態は、第1の実施の形態の機能と、第2の実施の形態の機能とを併せ持った文書処理装置である。

【0113】

第3の実施の形態において処理対象として想定している構造化文書は、次のような文書である。

図47は、構造化文書の第4の例を示す図である。第1の実施の形態で示した文書処理装置に図5の検索条件40と図47の文書94を入力しても、識別子「007」のPARA要素が検索条件ノード中の述語にマッチしないために参考文献エントリを得ることができない。また、第2の実施の形態で示した文書処理装置に図5の検索条件40と図47の文書94 30
を入力しても、識別子「008」のLIST要素が検索条件ノード中の述語にマッチしないために参考文献エントリを得ることができない。

【0114】

そこで、図47の文書94に対しても所望の結果を得るため、第3の実施の形態の文書処理装置では、検索条件のパターン指定に"+r"と"+p"の両方を設定することができるようにする。

【0115】

図48は、"+r"と"+p"を用いた検索条件を示す図である。図48の検索条件60は、図5の検索条件40のパターン指定42に記号"+p"を追加したものであり、文書要素を引数とする述語61、63、65、67と、パターン指定62、64、66とで表されている 40
。この追加により、述語63にマッチする文書要素の前に出現する文書要素群、述語63にマッチする文書要素とパターン指定64にマッチする文書要素列の間の文書要素群、及びパターン指定64にマッチする文書要素列の後に出現する文書要素群を読み飛ばしてマッチングを行うことが可能となる。具体的には、図47の文書94に対する処理では識別子「007」「009」の文書要素を読み飛ばして所望の文書要素を得ることができる。

【0116】

ここで、"+r"と"+p"を含む検索条件を処理するため、第3の実施の形態の文書処理装置では、第1の実施の形態の文書処理装置のオートマトン生成処理とオートマトンマッチング処理に以下の変更を施した処理が実行される。なお、第3の実施の形態に係る文書処理装置に必要な構成要素は、第1の実施の形態と同様であるため、図2に示した構成を用いて 50

第 3 の実施の形態を説明する。

【 0 1 1 7 】

図 4 9 は、第 3 の実施の形態におけるオートマトン生成処理手順を示すフローチャートである。この処理は、第 1 の実施の形態のオートマトン生成処理（図 9、図 1 0 に示す）に代えて行われる処理であり、すべて階層オートマトン生成部 1 3 で行われる。

〔 S 6 0 1 〕 検索条件ノードに対応するオートマトンが変数 T に含まれるか否かを判断し、変数 T に含まれる場合は処理をステップ S 6 0 2 へ進め、そうでない場合は処理をステップ S 6 0 3 へ進める。

〔 S 6 0 2 〕 入力の状態遷移から、検索条件ノードに対応するオートマトンへの参照を生成し、処理を終了する。

〔 S 6 0 3 〕 オートマトンを生成し、入力の状態遷移からオートマトンへの参照を生成する。

〔 S 6 0 4 〕 入力の検索条件ノードのパターン指定に "+r" が設定されているか否かを判定し、設定されている場合には処理をステップ S 6 0 5 へ進め、そうでない場合には処理をステップ S 6 0 6 へ進める。

〔 S 6 0 5 〕 生成したオートマトンに "+r" を設定する。

〔 S 6 0 6 〕 入力の検索条件ノードのパターン指定に "+p" が設定されているか否かを判定し、設定されている場合には処理をステップ S 6 0 7 へ進め、そうでない場合には処理をステップ S 6 0 8 へ進める。

〔 S 6 0 7 〕 生成したオートマトンに "+p" を設定する。

〔 S 6 0 8 〕 オートマトン上に初期状態と受理状態を生成する。

〔 S 6 0 9 〕 検索条件ノードのパターン指定の種別を判別する。パターン指定が "SEQ" なら処理をステップ S 6 1 0 へ進め、"CH0" なら処理をステップ S 6 1 1 へ進め、"OPTREP" なら処理をステップ S 6 1 2 へ進める。

〔 S 6 1 0 〕 オートマトンを引数にして SEQ パターン指定生成処理を呼び出す。その後、処理をステップ S 6 1 3 に進める。なお、SEQ パターン指定生成処理の詳細は、図 1 1、図 1 2 に示した処理と同様である。

〔 S 6 1 1 〕 オートマトンを引数にして CH0 パターン指定生成処理を呼び出す。その後、処理をステップ S 6 1 3 に進める。なお、CH0 パターン指定生成処理の詳細は、図 1 3 に示した処理と同様である。

〔 S 6 1 2 〕 オートマトンを引数にして OPTREP パターン指定生成処理を呼び出す。なお、OPTREP パターン指定生成処理の詳細は、図 1 4 に示した処理と同様である。

〔 S 6 1 3 〕 各生成処理の戻り値のオートマトンを変数 T に格納する。

【 0 1 1 8 】

このようにして、"+r" と "+p" との両方を使用した検索条件に対する階層オートマトンが生成される。

図 5 0 は、第 3 の実施の形態の文書処理装置により生成された階層オートマトンを示す図である。これは、図 4 8 の検索条件を第 3 の実施の形態のオートマトン生成処理に入力した結果出力される階層オートマトン 7 0 0 である。この階層オートマトン 7 0 0 は、4 つのオートマトン 7 1 0、7 2 0、7 3 0、7 4 0 で構成される。オートマトン 7 2 0 には "+p" の記号が付与されており、オートマトン 7 3 0、7 4 0 には "+r" の記号が付与されている。

【 0 1 1 9 】

このような階層オートマトン 7 0 0 を用いてオートマトンマッチング処理が行われる。オートマトンマッチング処理は以下のように変更される。

図 5 1 は、第 3 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの前半である。この処理は、第 1 の実施の形態のオートマトンマッチング処理（図 1 6 ~ 図 1 8）に代えて、マッチング部 1 4 によって行われる処理である。

〔 S 6 2 1 〕 入力の文書要素がスタック先頭の文書要素であるか否かを判断する。スタック先頭の文書要素であれば処理をステップ S 6 2 2 へ進め、そうでない場合は処理をステ

10

20

30

40

50

ップ S 6 2 6 に進める。

[S 6 2 2] 入力の状態とスタック先頭の状態を引数にした到達可能状態特定処理の戻り値である状態の集合に、スタック先頭の状態が含まれるか否かを判断する。含まれる場合は処理をステップ S 6 2 3 に進め、そうでない場合は処理を終了する。なお、到達可能状態特定処理の詳細は、図 1 9 に示す通りである。

[S 6 2 3] スタック先頭の文書要素と状態の組を pop し、入力の状態を pop した状態と置き換える。

[S 6 2 4] スタックが空であるか否かを判定し、空の場合は処理をステップ S 6 2 5 へ進め、そうでない場合は処理をステップ S 6 2 6 へ進める。

[S 6 2 5] 変数 V に入力の文書要素セットを書き出す。

10

[S 6 2 6] 入力の状態を引数にしてマッチング可能状態遷移特定処理を起動した結果の戻り値の状態遷移を変数 X に格納する。マッチング可能状態遷移特定処理の詳細は、図 2 2 に示す通りである。

[S 6 2 7] 変数 X のすべての状態遷移が処理済みであるか否かを判断し、処理済みの場合は処理を終了し、そうでない場合は処理をステップ S 6 2 8 に進める。

[S 8 2 7] 変数 X の未処理の状態遷移の 1 つを処理対象とする。

[S 6 2 9] 処理対象の状態遷移と入力の文書要素を引数として文書要素マッチング処理を呼び出し、戻り値を判断する。戻り値が真なら処理をステップ S 6 3 0 に進め、偽なら処理をステップ S 6 3 7 に進める。なお、文書要素マッチング処理の詳細は、図 2 4 に示す通りである。

20

[S 6 3 0] 入力の文書要素セットのコピーを作成し、文書要素セットのコピーに入力の文書要素を追加する。

【 0 1 2 0 】

図 5 2 は、第 3 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの後半である。

[S 6 3 1] 入力の文書要素を引数にして次要素特定処理を起動した結果の戻り値が入力の文書要素の長子であるか否かを判定する。戻り値が長子である場合は処理をステップ S 6 3 2 に進め、そうでない場合は処理をステップ S 6 3 6 に進める。なお、次要素特定処理の詳細は、図 2 5 に示す通りである。

[S 6 3 2] 処理対象の状態遷移がオートマトンを参照しているか否かを判定し、参照している場合は処理をステップ S 6 3 3 へ進め、そうでない場合は処理をステップ S 6 3 6 へ進める。

30

[S 6 3 3] スタックのコピーを作成する。

[S 6 3 4] スタックのコピーに、処理対象の状態遷移の終了端の状態と、入力の文書要素を引数とした弟要素特定処理の戻り値の組を push する。

[S 6 3 5] 処理対象の状態遷移が参照するオートマトンの初期状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力の文書要素の長子と、戻り値の状態と、ステップ S 6 3 0 で作成した文書要素セットのコピーと、ステップ S 6 3 4 で作成したスタックのコピーとを入力としてオートマトンマッチング処理を呼び出し、処理をステップ S 6 3 7 に進める。

40

[S 6 3 6] 処理対象の状態遷移の終了端の状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力の文書要素を引数として弟要素特定処理を起動した結果の戻り値と、戻り値の状態と、ステップ S 6 3 0 で作成した文書要素セットのコピーと、入力のスタックのコピーを入力としてオートマトンマッチング処理を呼び出す。

[S 6 3 7] 処理対象の状態遷移が属するオートマトンに "+r" が記されているか否かを判定し、記されている場合は処理をステップ S 6 3 8 へ進め、そうでない場合は処理をステップ S 6 4 0 に進める。

[S 6 3 8] 入力の文書要素を引数にして次要素特定処理を起動した結果の戻り値が入力の文書要素の長子であるか否かを判定し、長子である場合は処理をステップ S 6 3 9 に、そうでない場合は処理をステップ S 6 4 0 に進める。

50

〔S 6 3 9〕入力の状態、入力の文書要素の長子の文書要素、入力の文書要素セットのコピー、入力のスタックのコピーを入力としてオートマトンマッチング処理を呼び出す。

〔S 6 4 0〕処理対象の状態遷移が属するオートマトンに"+p"が記されているか否か判定し、記されている場合はステップS 6 4 1へ進め、そうでない場合は処理をS 6 2 7に進める。

〔S 6 4 1〕入力の状態、入力の文書要素を第要素特定処理に入力した場合の返り値の文書要素、入力の文書要素セットのコピー、入力のスタックのコピーを入力としてオートマトンマッチング処理を呼び出す。その後、処理をステップS 6 2 7に進める。

【0 1 2 1】

以下に、第3の実施の形態のマッチング部14に図47の文書94と図50の階層オートマトン700を入力したときの、オートマトンマッチング処理の呼び出し関係を示す。

10

【0 1 2 2】

図53は、第3の実施の形態におけるオートマトンマッチング処理の呼び出し関係を示す第1の図である。また、図54は、第3の実施の形態におけるオートマトンマッチング処理の呼び出し関係を示す第2の図である。図53、図43の表記法は図27、図28の表記法と同一であり、ノード801～813はオートマトンマッチング処理を意味する。なお、簡単のため、図53、図54では、オートマトンマッチング処理に4つ組(11,005,NULL,NUL L)が入力された時点からの呼び出し関係のみを示す。さらに、文書要素セット(005,006,009,010)の出力に至る呼び出し関係のみを示し、文書要素セット(005,006)および文書要素セット(005,006,009)の出力に至る呼び出し関係は省略してある。

20

【0 1 2 3】

図53、図54では、ノード802の処理が実行されるオートマトン(図50のオートマトン720)に記号"+p"が設定されているため、識別子「007」の文書要素を読み飛ばしてノード803の処理を呼び出すことができる。また、ノード805の処理が実行されるオートマトン(図50のオートマトン730)に記号"+r"が設定されているため、識別子「008」の文書要素を読み飛ばしてノード806の処理を呼び出すことができる。さらに、ノード811の処理が実行されるオートマトン(図50のオートマトン720)に記号"+p"が設定されているため、識別子「011」の文書要素を読み飛ばしてノード812の処理を呼び出し、結果的に識別子「009」の文書要素と識別子「010」の文書要素を参考文献エントリとして得ることができる。

30

【0 1 2 4】

図55は、第4の実施の形態におけるマッチング部の出力例を示す図である。これは、図47の構造化文書と図48の検索条件を本実施例の文書処理装置に入力した結果得られるマッチング部の出力結果である。

【0 1 2 5】

このように、第1の実施の形態の機能と第2の実施の形態の機能とを兼ね備えることで、様々な構造化文書の中から、利用者の意図に合った部分構造をマッチング処理により検出することができる。

【0 1 2 6】

ところで、第1の実施の形態において説明した応用例(図35に示す)と同様の応用を、第3の実施の形態に対して行うこともできる。その場合図35に示したオートマトン生成処理は以下のように変更される。

40

【0 1 2 7】

図56は、第3の実施の形態のオートマトン生成処理に関する応用例を示すフローチャートである。この処理は、図49の処理に代えて階層オートマトン生成部13で行われる処理である。

〔S 7 0 1〕検索条件ノードに対応するオートマトンが変数Tに含まれるか否かを判断し、変数Tに含まれる場合には処理をステップS 7 0 2へ進め、そうでない場合は処理をステップS 7 0 3へ進める。

〔S 7 0 2〕入力の状態遷移から、検索条件ノードに対応するオートマトンへの参照を生

50

成し、処理を終了する。

[S 7 0 3] オートマトンを生成し、入力の状態遷移からオートマトンへの参照を生成する。

[S 7 0 4] 入力の検索条件ノードのパターン指定に "+r" が設定されているか否か判定し、設定されている場合には処理をステップ S 7 0 5 へ進め、そうでない場合には処理をステップ S 7 0 6 へ進める。

[S 7 0 5] 生成したオートマトンに "+r" を設定する。

[S 7 0 6] 入力の検索条件ノードのパターン指定に "+p" が設定されているか否か判定し、設定されている場合には処理をステップ S 7 0 7 へ進め、そうでない場合には処理をステップ S 7 0 8 へ進める。

[S 7 0 7] 生成したオートマトンに "+p" を設定する。

[S 7 0 8] オートマトンに初期状態と受理状態を 1 つずつ生成する。

[S 7 0 9] 検索条件ノードの正規表現とオートマトンを引数にして属性変換文法評価処理を呼び出す。

[S 7 1 0] ステップ S 7 0 7 で生成されたオートマトンを、変数 T に出力する。

【 0 1 2 8 】

次に、第 4 の実施の形態について説明する。第 4 の実施の形態は、要素に対して属性が与えられている場合に、その属性を、該当する要素の子孫に対して継承させながらマッチングを行うものである。

【 0 1 2 9 】

第 4 の実施の形態では、以下のような構造化文書进行处理対象として想定している。

図 5 7 は、構造化文書の第 5 の例を示す図である。この文書 9 5 の一部の文書要素には、名前だけでなく属性も設定されている。たとえば識別子「 003 」の文書要素では、要素名として "LIST" が設定されているだけでなく、属性として "ALIGN="C"" が設定されている。本実施例では、文書要素の属性 ALIGN はその文書要素を頂点とする部分木に対応する文書内容を表示するときの配置を規定するものとする。そして、属性値 "C" は、文書内容がセンタリングされて表示されることを意味するものとする。

【 0 1 3 0 】

図 5 7 の構造化文書から参考文献エントリを得るための検索条件を以下に示す。

図 5 8 は、第 4 の実施の形態に用いる検索条件を示す図である。この検索条件 7 0 は、文書要素を引数とする述語 7 1 , 7 3 , 7 5 , 7 7 と、パターン指定 7 2 , 7 4 , 7 6 とで表されている。検索条件 7 0 の述語 7 3 , 7 5 , 7 7 には、"CONTENT" と "ATTR" という 2 種類の条件が示されている。"CONTENT" は文書要素の内容に関する条件を意味し、"ATTR" は文書要素の属性に関する条件を意味する。たとえば図 5 8 の検索条件の述語 7 3 に示されている条件は「 参考文献 または "References" を内容に含み、属性 ALIGN の値が "C" である文書要素」となる。

【 0 1 3 1 】

図 5 9 は、属性を指定した検索条件から生成される階層オートマトンを示す図である。この階層オートマトン 8 0 0 は、4 つのオートマトン 8 1 0 , 8 2 0 , 8 3 0 , 8 4 0 からなる。

【 0 1 3 2 】

なお、第 4 の実施の形態の文書処理装置の装置構成および処理の流れは、基本的には第 3 の実施の形態で示したものと同一である。ただし、図 4 8 に示した検索条件を処理できるようにするため、図 5 1、図 5 2 のオートマトンマッチング処理は次のように変更される。

【 0 1 3 3 】

図 6 0 は、第 4 の実施の形態のオートマトンマッチング処理手順を示す第 1 のフローチャートである。

[S 8 0 1] 入力の文書要素がスタック先頭の文書要素であるか否かを判断し、スタック先頭の文書要素である場合は処理をステップ S 8 0 2 へ進め、そうでない場合は処理をス

10

20

30

40

50

テップ S 8 0 6 に進める。

[S 8 0 2] 入力の状態とスタック先頭の状態を引数にした到達可能状態特定処理の戻り値である状態の集合に、スタック先頭の状態が含まれるか否かを判断する。スタック先頭の状態が含まれる場合は処理をステップ S 8 0 3 に進め、そうでない場合は処理を終了する。

[S 8 0 3] スタック先頭の文書要素と状態の組を pop し、入力の状態を pop した状態と置き換える。

[S 8 0 4] スタックが空であるか否かを判定し、空の場合は処理をステップ S 8 0 5 へ進め、そうでない場合は処理をステップ S 8 0 6 へ進める。

[S 8 0 5] 変数 V に入力の文書要素セットを書き出す。

10

[S 8 0 6] 入力の状態を引数にしてマッチング可能状態遷移特定処理を起動した結果の戻り値の状態遷移を変数 X に格納する。なお、マッチング可能状態遷移特定処理の詳細は、図 2 2 に示す通りである。

[S 8 0 7] 変数 X のすべての状態遷移が処理済みか否かを判断し、処理済みの場合は処理を終了し、そうでない場合は処理をステップ S 8 0 8 に進める。

[S 8 0 8] 変数 X の未処理の状態遷移の 1 つを処理対象とする。

[S 8 0 9] 処理対象の状態遷移と入力の文書要素を引数として文書要素マッチング処理を呼び出し、戻り値が真なら処理をステップ S 8 1 0 に進め、偽なら処理をステップ S 8 1 7 に進める。文書要素マッチング処理の詳細は、図 6 3 に示す。

【 0 1 3 4 】

20

図 6 1 は、第 4 の実施の形態のオートマトンマッチング処理手順を示す第 2 のフローチャートである。

[S 8 1 0] 入力の文書要素セットのコピーを作成し、文書要素セットのコピーに入力の文書要素を追加する。

[S 8 1 1] 入力の文書要素を引数にして次要素特定処理を起動した結果の戻り値が入力の文書要素の長子であるか否かを判定し、長子である場合は処理をステップ S 8 1 2 に進め、そうでない場合は処理をステップ S 8 1 6 に進める。

[S 8 1 2] 処理対象の状態遷移がオートマトンを参照しているか否かを判定し、参照している場合は処理をステップ S 8 1 3 へ進め、そうでない場合は処理をステップ S 8 1 6 へ進める。

30

[S 8 1 3] スタックのコピーを作成する。

[S 8 1 4] スタックのコピーに、処理対象の状態遷移の終了端の状態と、入力の文書要素を引数とした弟要素特定処理の戻り値の組を push する。

[S 8 1 5] 処理対象の状態遷移が参照するオートマトンの初期状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力の文書要素の長子と、戻り値の状態と、ステップ S 8 1 0 で作成した文書要素セットのコピーと、ステップ S 8 1 4 で作成したスタックのコピーを入力としてオートマトンマッチング処理を呼び出し、処理をステップ S 8 1 7 に進める。

[S 8 1 6] 処理対象の状態遷移の終了端の状態を引数として到達可能状態特定処理を呼び出した場合の戻り値の状態ごとに、入力の文書要素を引数として弟要素特定処理を起動した結果の戻り値と、戻り値の状態と、ステップ S 8 1 0 で作成した文書要素セットのコピーと、入力のスタックのコピーを入力としてオートマトンマッチング処理を呼び出す。

40

[S 8 1 7] 処理対象の状態遷移が属するオートマトンに "+r" が記されているか否かを判定し、記されている場合は処理をステップ S 8 1 8 へ進め、そうでない場合は処理を終了する。

[S 8 1 8] 入力の文書要素を引数にして次要素特定処理を起動した結果の戻り値が入力の文書要素の長子であるか否かを判定し、長子である場合は処理をステップ S 8 1 9 に進め、そうでない場合は処理をステップ S 8 2 1 に進める。

[S 8 1 9] 入力の文書要素に設定されている属性を入力の子供の文書要素に設定する。

50

【 0 1 3 5 】

図 6 2 は、第 4 の実施の形態のオートマトンマッチング処理手順を示す第 3 のフローチャートである。

[S 8 2 0] 入力の状態、入力の文書要素の長子の文書要素、入力の文書要素セットのコピー、入力のスタックのコピーを入力としてオートマトンマッチング処理を呼び出す。

[S 8 2 1] 処理対象の状態遷移が属するオートマトンに "+p" が記されているか否か判定し、記されている場合は処理をステップ S 8 2 2 へ進め、そうでない場合は処理をステップ S 8 0 7 に進める。

[S 8 2 2] 入力の状態、入力の文書要素を第要素特定処理に入力した場合の返り値の文書要素、入力の文書要素セットのコピー、入力のスタックのコピーを入力としてオートマトンマッチング処理を呼び出す。その後、処理をステップ S 8 0 7 に進める。

10

【 0 1 3 6 】

図 6 0 ~ 図 6 2 の処理の流れは、図 5 1、図 5 2 の処理の流れにステップ S 8 1 9 の処理が追加されたものである。また、図 2 4 の文書要素マッチング処理は以下のように変更される。

【 0 1 3 7 】

図 6 3 は、第 4 の実施の形態における文書要素マッチング処理の手順を示すフローチャートである。この処理は、すべてマッチング部 1 4 によって行われる。

[S 8 3 1] 文書要素に対応付けられた文書内容の文字列が、状態遷移に対応付けられた述語の CONTENT エントリの文字列パターンにマッチするか否かを判断する。マッチする場合は処理をステップ S 8 3 2 へ進め、そうでない場合は処理をステップ S 8 3 4 に進める。

20

[S 8 3 2] 状態遷移に対応付けられた属性の名前と値の組が、文書要素に対応付けられた述語の ATTR エントリの属性の名前と値の組の中に存在するか否かを判断する。存在する場合は処理をステップ S 8 3 3 へ進め、そうでない場合は処理をステップ S 8 3 4 に進める。

[S 8 3 3] 真を返し、処理を終了する。

[S 8 3 4] 偽を返し、処理を終了する。

【 0 1 3 8 】

ここで、第 4 の実施の形態のマッチング部 1 4 に図 5 7 に示す文書 9 5 と図 5 9 に示す階層オートマトン 8 0 0 とを入力した際の、オートマトンマッチング処理の呼び出し関係を以下に示す。

30

【 0 1 3 9 】

図 6 4 は、第 4 の実施の形態による呼び出し関係を示す図である。図 6 4 の表記法は図 2 7、図 2 8 などの表記法と同一であり、ノード 9 0 1 ~ 9 0 5 はオートマトンマッチング処理を意味する。なお、図 6 4 では、本実施の形態の文書処理装置の特徴的な処理の流れを説明するのに必要な呼び出し関係のみを示している。

【 0 1 4 0 】

図 6 4 では、第 3 の実施の形態と同様、ノード 9 0 2 の処理が実行されるオートマトン (図 5 9 のオートマトン 8 2 0) に記号 "+p" が設定されているため、識別子「 003 」の文書要素から識別子「 009 」の文書要素を読み飛ばし、識別子「 010 」文書要素を処理対象とすることができる (ノード 9 0 3 の処理)。ここで、第 3 の実施の形態の文書処理装置では、識別子「 001 」の文書要素をマッチングの対象にするか、識別子「 010 」の文書要素を読み飛ばして識別子「 011 」の文書要素をマッチング対象にすることしかできなかった。そのため、識別子「 3 a 」の状態遷移の条件を満足する文書要素を得ることはできない。

40

【 0 1 4 1 】

しかし、第 4 の実施の形態の文書処理装置では、識別子「 010 」の文書要素の属性を識別子「 011 」の文書要素に追加して得られる文書要素をマッチングの対象とすることができる (ノード 9 0 4 の処理)。したがって、識別子「 011 」の文書要素を識別子「 3 a 」の

50

状態遷移とマッチさせ、識別子「012」の文書要素以降を対象とした処理を続けることができる。そして、結果的に、識別子「013」の文書要素と識別子「014」の文書要素とを参考文献エントリとして得ることができる。

【0142】

次に、第5の実施の形態について説明する。第5の実施の形態は、検索条件に適合した構造の内容を、別の内容に置き換えて出力するようにしたものである。

構造化文書を対象とした処理では、処理対象の文書の部分構造を検索するだけでなく、検索された部分構造の文書要素の名前や属性を変換することにより、文書作成時に想定された用途以外の用途に利用することが行われている。この処理に対応するため、第5の実施の形態の文書処理装置では、第3の実施の形態の文書処理装置の入力となる検索条件の述語に、述語にマッチした文書要素に設定すべき名前や属性を対応付けておく。このような検索条件の例を以下に示す。

10

【0143】

図65は、第5の実施の形態に用いる検索条件を示す図である。この検索条件80は、文書要素を引数とする述語81, 84, 86, 88と、パターン指定82, 83, 85, 87とで表されている。図65の述語に対応する矩形には、記号" "に続いて文字列が記載されている。これは、述語にマッチした文書要素の名前を" "に続く文字列に置き換えることを意味する。

【0144】

第5の実施の形態における文書処理装置の装置構成および処理の流れは、第3の実施の形態における装置構成および処理の流れと基本的に同一である。唯一の差異は、図51に示したオートマトンマッチング処理のステップS630が、以下のように変更される点である。

20

[S630a] 入力 of 文書要素セットのコピーを作成し、入力の状態遷移に対応する検索条件ノードに設定された名前や属性を入力 of 文書要素のコピーに設定して、このコピーを文書要素セットのコピーに追加する。

【0145】

第5の実施の形態における文書処理装置に図65の検索条件と図69の構造化文書を入力した場合のオートマトンマッチング処理の呼び出し関係を、以下に示す。

【0146】

30

図66は、第5の実施の形態による呼び出し関係を示す図である。図66の表記法は図27、図28などの表記法と同一であり、ノード1001～1008はオートマトンマッチング処理を意味する。図66の呼び出し関係は図43、図44の呼び出し関係と基本的に同一であるが、処理呼び出し時の第3引数の文書要素セットには、文書要素識別子と文書要素に設定された要素名が記号"-"で接続されたものが記されている。たとえば、ノード1002の処理の第3引数である「001-"参考文献"」は、「文書要素セットがマッチング部から出力されるときには、識別子「001」である文書要素(DOC要素)の名前を"参考文献"に変更して出力すること」を意味する。

【0147】

また、本実施の形態における文書処理装置に図65の検索条件と図69の構造化文書を入力した結果得られる文書要素セット(001,004,005,006)についてディスプレイに表示されるイメージの例を以下に示す。

40

【0148】

図67は、第5の実施の形態によって表示される画面を示す図である。この画面430は図33と同様に、部分構造表示部431と適合内容表示部432とで構成されている。そして、部分構造表示部431の表示内容は、要素名が置き換えられた状態で表示されている。

【0149】

なお、第5の実施の形態の文書処理装置の処理の流れは、第3の実施の形態における文書処理装置との差異を示すことによって説明した。しかし、他の実施例における文書処理装

50

置に対しても、上記の説明で示した変更と同様の変更を施すことで、図 6 5 の検索条件を処理することが可能となる。

【 0 1 5 0 】

なお、上記の処理機能は、コンピュータによって実現することができる。その場合、データ処理装置及び文書処理装置が有すべき機能の処理内容は、コンピュータで読み取り可能な記録媒体に記録されたプログラムに記述しておく。そして、このプログラムをコンピュータで実行することにより、上記処理がコンピュータで実現される。コンピュータで読み取り可能な記録媒体としては、磁気記録装置や半導体メモリ等がある。市場に流通させる場合には、C D - R O M (Compact Disk Read Only Memory) やフロッピーディスク等の可搬型記録媒体にプログラムを格納して流通させたり、ネットワークを介して接続されたコンピュータの記憶装置に格納しておき、ネットワークを通じて他のコンピュータに転送することもできる。コンピュータで実行する際には、コンピュータ内のハードディスク装置等にプログラムを格納しておき、メインメモリにロードして実行する。

10

【 0 1 5 1 】

【 発明の効果 】

以上説明したように本発明の第 1 のデータ処理装置では、有向順序木の中間ノードを削除し、中間ノードのあった位置に中間ノード直下のノード列を配置する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うようにしたため、論理構造が異なる多数の情報に対する検索を行う際においても、ノード間の接続関係を利用した検索条件を利用できる。

20

【 0 1 5 2 】

また、本発明の第 2 のデータ処理装置では、有向順序木中のいずれかのノードを頂点とする部分木を削除する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うようにしたため、論理構造が異なる多数の情報に対する検索を行う際においても、ノード間の接続関係を利用した検索条件を利用できる。

【 0 1 5 3 】

また、本発明の第 1 の文書処理装置では、有向順序木の中間ノードを削除し、中間ノードのあった位置に中間ノード直下のノード列を配置する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うようにしたため、論理構造が異なる多数の構造化文書に対する検索を行う際においても、文書要素間の接続関係を利用した検索条件を利用できる。

30

【 0 1 5 4 】

また、本発明の第 2 の文書処理装置では、有向順序木中のいずれかのノードを頂点とする部分木を削除する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うようにしたため、論理構造が異なる多数の構造化文書に対する検索を行う際においても、文書要素間の接続関係を利用した検索条件を利用できる。

【 0 1 5 5 】

また、本発明の第 1 のデータ処理プログラムを記録したコンピュータ読み取り可能な記録媒体では、有向順序木の中間ノードを削除し、中間ノードのあった位置に中間ノード直下のノード列を配置する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うような処理をコンピュータに行わせることができるため、ノード間の接続関係を利用した検索条件を利用して論理構造が異なる多数の情報に対する検索を行う機能を、コンピュータ上に構築することができる。

40

【 0 1 5 6 】

また、本発明の第 2 のデータ処理プログラムを記録したコンピュータ読み取り可能な記録媒体では、有向順序木中のいずれかのノードを頂点とする部分木を削除する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うような処理をコンピュータに行わせることができるため、ノード間の接続関係を利用した検索条件を利用して論理構造が異なる多数の情報に対する検索を行う機能を、コンピュ

50

ータ上に構築することができる。

【 0 1 5 7 】

また、本発明の第 1 の文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体では、有向順序木の中間ノードを削除し、中間ノードのあった位置に中間ノード直下のノード列を配置する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うような処理をコンピュータに行わせることができるため、ノード間の接続関係を利用した検索条件を利用して論理構造が異なる多数の構造化文書に対する検索を行う機能を、コンピュータ上に構築することができる。

【 0 1 5 8 】

また、本発明の第 2 の文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体では、有向順序木中のいずれかのノードを頂点とする部分木を削除する操作を行った結果得られる均質化有向順序木を処理対象として、検索条件に対する適合構造の抽出を行うような処理をコンピュータに行わせることができるため、ノード間の接続関係を利用した検索条件を利用して論理構造が異なる多数の構造化文書に対する検索を行う機能を、コンピュータ上に構築することができる。

【図面の簡単な説明】

【図 1】 本発明の原理構成図である。

【図 2】 文書処理システムの構成を示すブロック図である。

【図 3】 文書処理装置の処理手順を示すフローチャートである。

【図 4】 構造化文書の第 3 の例を示す図である。

【図 5】 検索条件の例を示す図である。

【図 6】 階層オートマトンの第 1 の例を示す図である。

【図 7】 階層オートマトン生成処理部が行う処理を示すフローチャートである。

【図 8】 根オートマトン生成処理手順を示すフローチャートである。

【図 9】 オートマトン生成処理手順を示すフローチャートの前半である。

【図 10】 オートマトン生成処理手順を示すフローチャートの後半である。

【図 11】 SEQ パターン生成処理手順を示すフローチャートの前半である。

【図 12】 SEQ パターン生成処理手順を示すフローチャートの後半である。

【図 13】 CHO パターン生成処理手順を示すフローチャートである。

【図 14】 OPTREP パターン生成処理手順を示すフローチャートである。

【図 15】 マッチング部の処理手順を示すフローチャートである。

【図 16】 オートマトンマッチング処理手順を示す第 1 のフローチャートである。

【図 17】 オートマトンマッチング処理手順を示す第 2 のフローチャートである。

【図 18】 オートマトンマッチング処理手順を示す第 3 のフローチャートである。

【図 19】 到達可能状態特定処理の手順を示すフローチャートである。

【図 20】 状態チェック処理のフローチャートの前半である。

【図 21】 状態チェック処理のフローチャートの後半である。

【図 22】 マッチング可能状態遷移特定処理の手順を示すフローチャートである。

【図 23】 状態遷移チェック処理の手順を示すフローチャートである。

【図 24】 文書要素マッチング処理の手順を示すフローチャートである。

【図 25】 次要素特定処理の手順を示すフローチャートである。

【図 26】 弟要素特定処理の手順を示すフローチャートである。

【図 27】 図 6 8 に示した文書における呼び出し関係を示す第 1 の図である。

【図 28】 図 6 8 に示した文書における呼び出し関係を示す第 2 の図である。

【図 29】 図 4 に示した文書における呼び出し関係を示す第 1 の図である。

【図 30】 図 4 に示した文書における呼び出し関係を示す第 2 の図である。

【図 31】 図 6 の階層オートマトンと図 6 8 の文書をマッチング部に入力した結果得られる文書要素セットを示す図である。

【図 32】 図 6 の階層オートマトンと図 4 の文書をマッチング部に入力した結果得られる文書要素セットを示す図である。

10

20

30

40

50

【図 3 3】 図 6 8 の構造化文書を処理対象としたときに C R T ディスプレイに表示される画面を示す図である。

【図 3 4】 図 4 の構造化文書を処理対象としたときに C R T ディスプレイに表示される画面を示す図である。

【図 3 5】 第 1 の実施の形態に関する応用例のオートマトン生成処理手順を示すフローチャートである。

【図 3 6】 属性変換文法評価処理の手順を示すフローチャートである。

【図 3 7】 属性変換文法評価処理の変形例を示す図である。

【図 3 8】 "+p"を用いて検索条件の例を示す図である。

【図 3 9】 第 2 の実施の形態におけるオートマトン生成処理の手順を示すフローチャートである。 10

【図 4 0】 第 2 の実施の形態の文書処理装置により生成された階層オートマトンを示す図である。

【図 4 1】 第 2 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの前半である。

【図 4 2】 第 2 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの後半である。

【図 4 3】 第 2 の実施の形態におけるオートマッチング処理の呼び出し関係を示す第 1 の図である。

【図 4 4】 第 2 の実施の形態におけるオートマッチング処理の呼び出し関係を示す第 2 の図である。 20

【図 4 5】 図 6 9 の構造化文書と図 3 8 の検索条件を第 2 の実施の形態の文書処理装置に入力した結果得られるマッチング部の出力を示す図である。

【図 4 6】 第 2 の実施の形態の応用例に係るオートマトン生成処理のフローチャートである。

【図 4 7】 構造化文書の第 4 の例を示す図である。

【図 4 8】 "+r"と"+p"とを用いた検索条件を示す図である。

【図 4 9】 第 3 の実施の形態におけるオートマトン生成処理手順を示すフローチャートである。

【図 5 0】 第 3 の実施の形態の文書処理装置により生成された階層オートマトンを示す図である。 30

【図 5 1】 第 3 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの前半である。

【図 5 2】 第 3 の実施の形態におけるオートマトンマッチング処理の手順を示すフローチャートの後半である。

【図 5 3】 第 3 の実施の形態におけるオートマッチング処理の呼び出し関係を示す第 1 の図である。

【図 5 4】 第 3 の実施の形態におけるオートマッチング処理の呼び出し関係を示す第 2 の図である。

【図 5 5】 第 4 の実施の形態におけるマッチング部の出力例を示す図である。 40

【図 5 6】 第 3 の実施の形態のオートマトン生成処理に関する応用例を示すフローチャートである。

【図 5 7】 構造化文書の第 5 の例を示す図である。

【図 5 8】 第 4 の実施の形態に用いる検索条件を示す図である。

【図 5 9】 属性を指定した検索条件から生成される階層オートマトンを示す図である。

【図 6 0】 第 4 の実施の形態のオートマトンマッチング処理手順を示す第 1 のフローチャートである。

【図 6 1】 第 4 の実施の形態のオートマトンマッチング処理手順を示す第 2 のフローチャートである。

【図 6 2】 第 4 の実施の形態のオートマトンマッチング処理手順を示す第 3 のフローチャートである。 50

ャートである。

【図 6 3】 第 4 の実施の形態における文書要素マッチング処理の手順を示すフローチャートである。

【図 6 4】 第 4 の実施の形態による呼び出し関係を示す図である。

【図 6 5】 第 5 の実施の形態に用いる検索条件を示す図である。

【図 6 6】 第 5 の実施の形態による呼び出し関係を示す図である。

【図 6 7】 第 5 の実施の形態によって表示される画面を示す図である。

【図 6 8】 構造化文書の第 1 の例を示す図である。

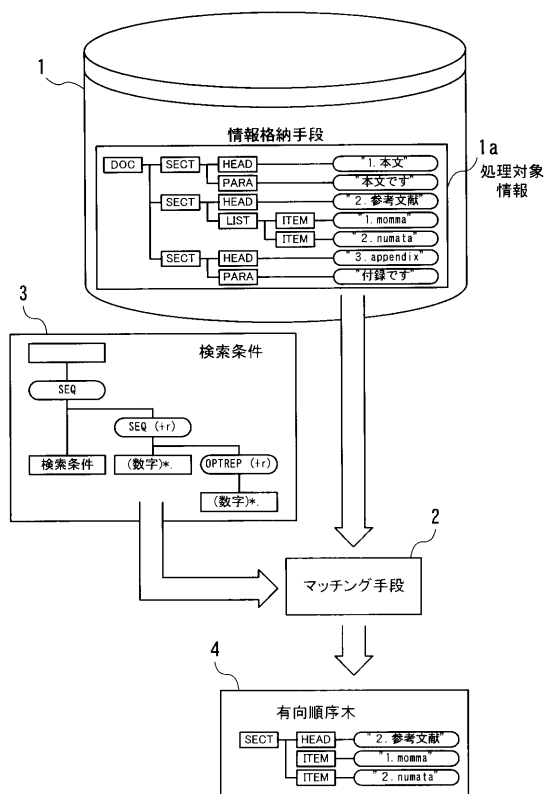
【図 6 9】 構造化文書の第 2 の例を示す図である。

【符号の説明】

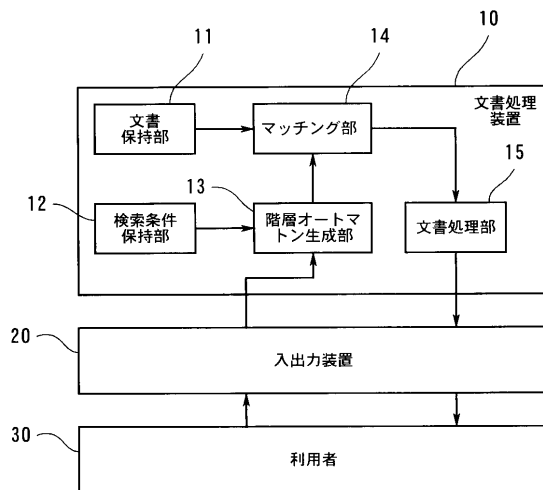
- 1 情報格納手段
- 2 マッチング手段
- 3 検索条件
- 4 有向順序木

10

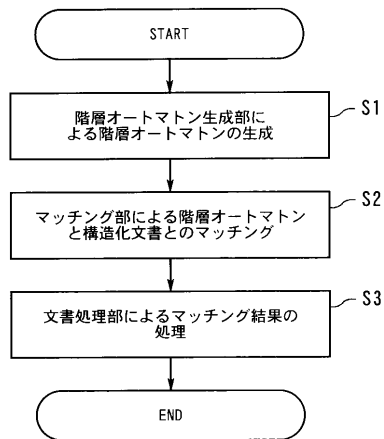
【図 1】



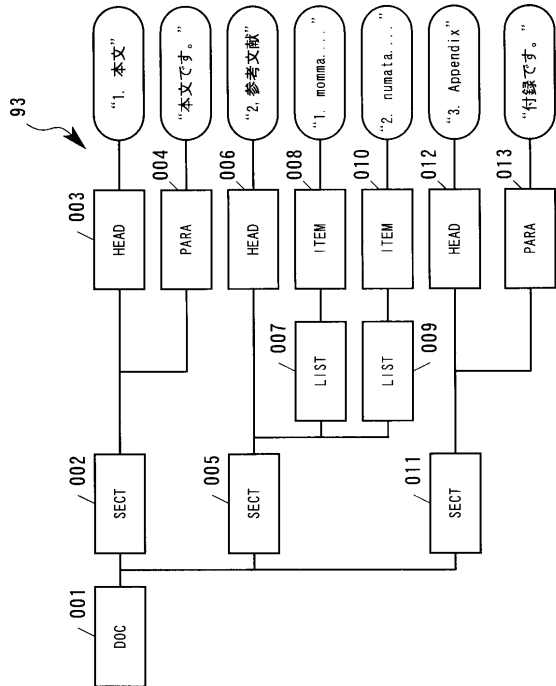
【図 2】



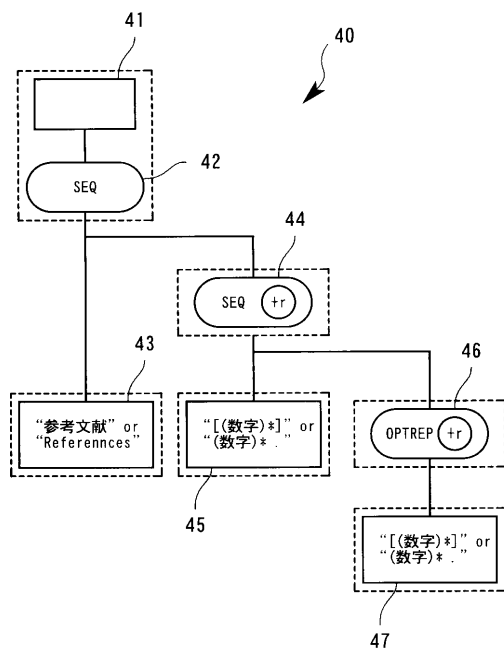
【図 3】



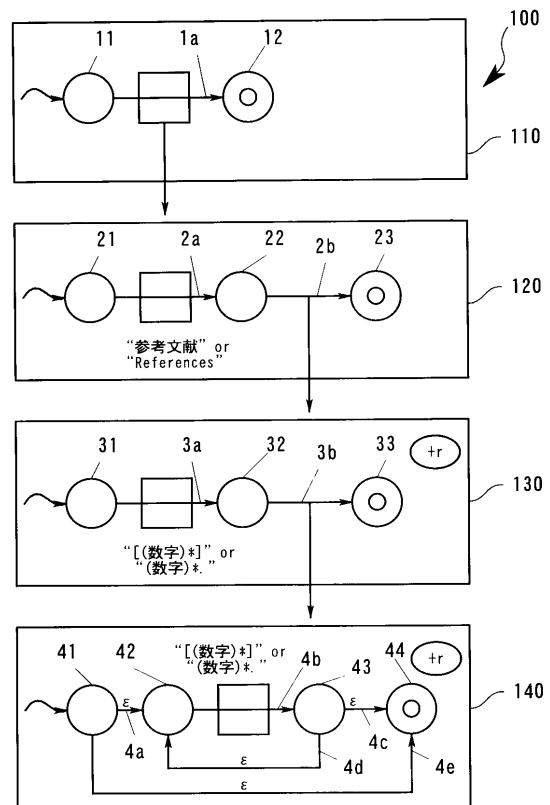
【図 4】



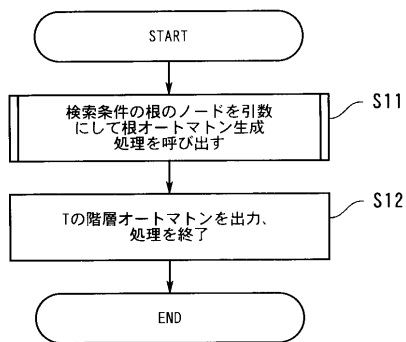
【図 5】



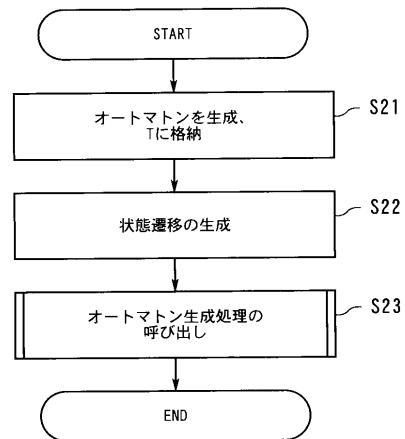
【図 6】



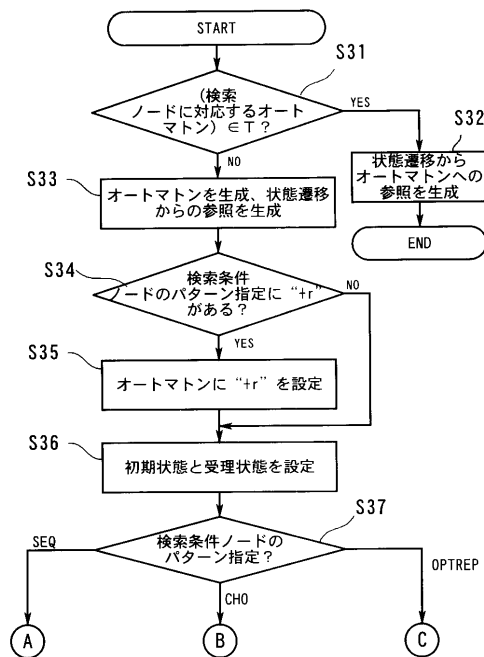
【図 7】



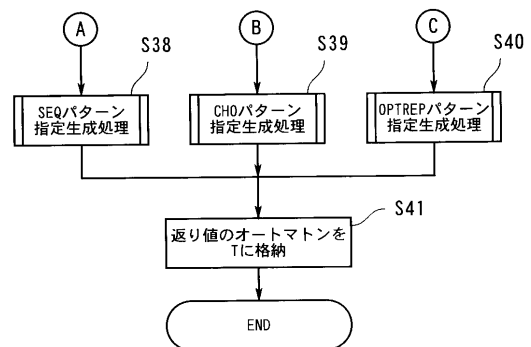
【図 8】



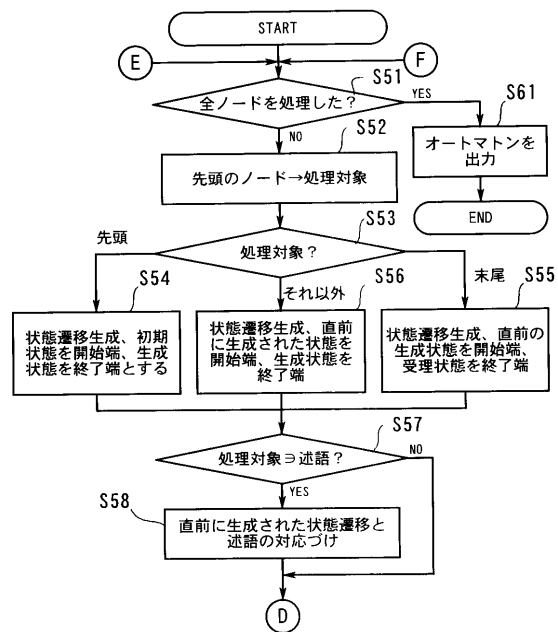
【図 9】



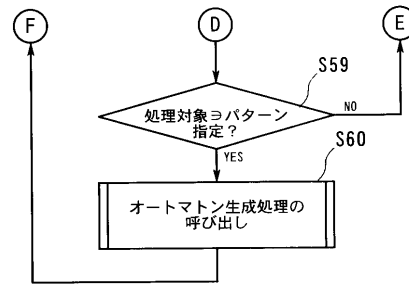
【図 10】



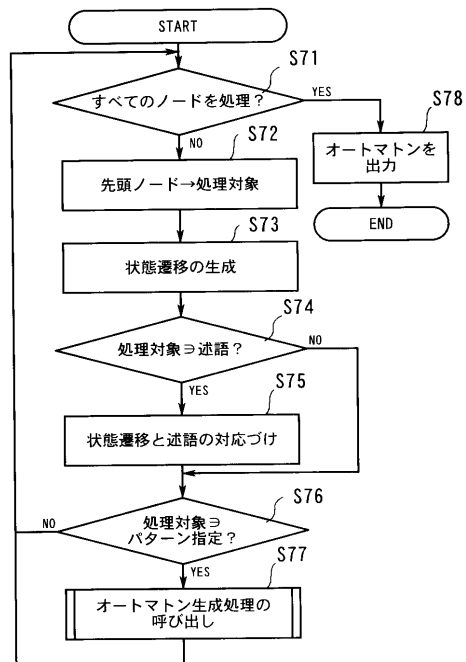
【図 1 1】



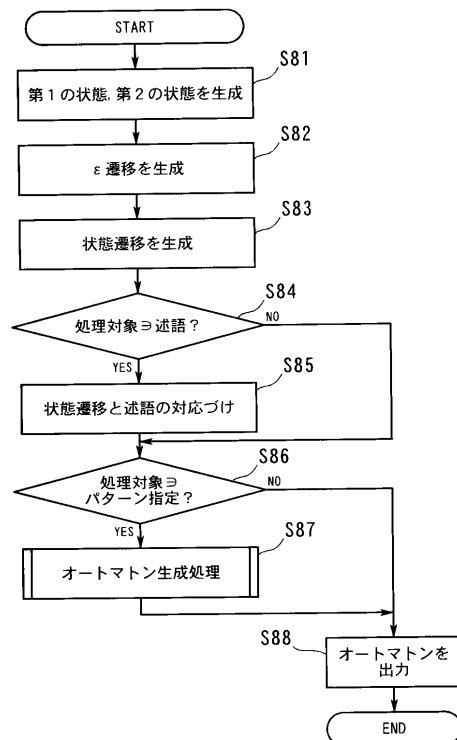
【図 1 2】



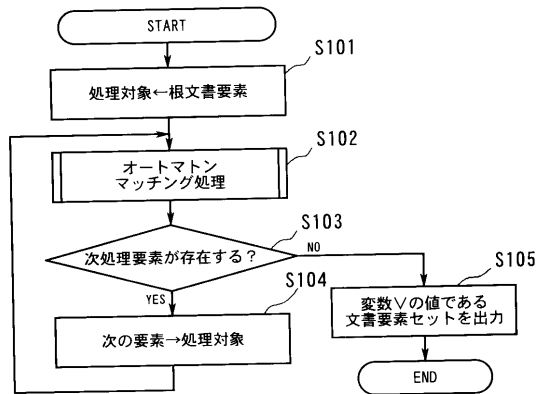
【図 1 3】



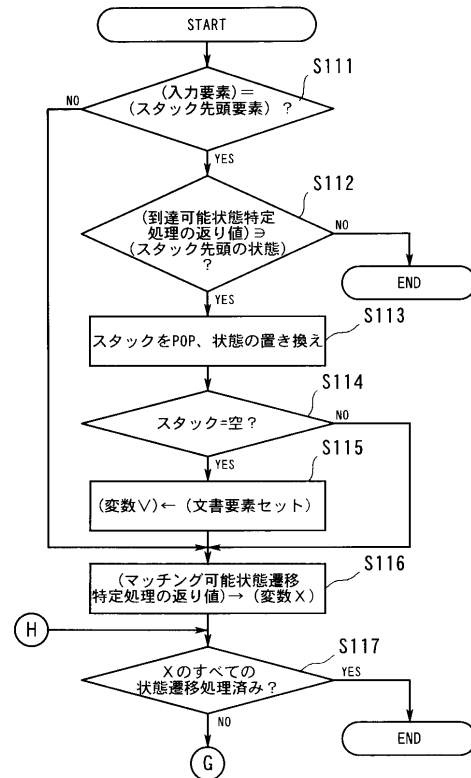
【図 1 4】



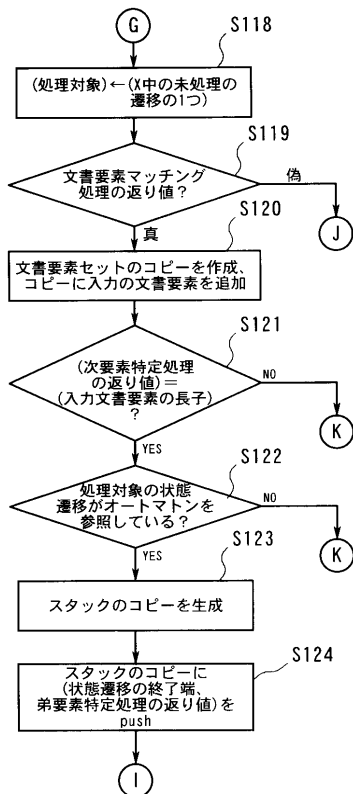
【図 15】



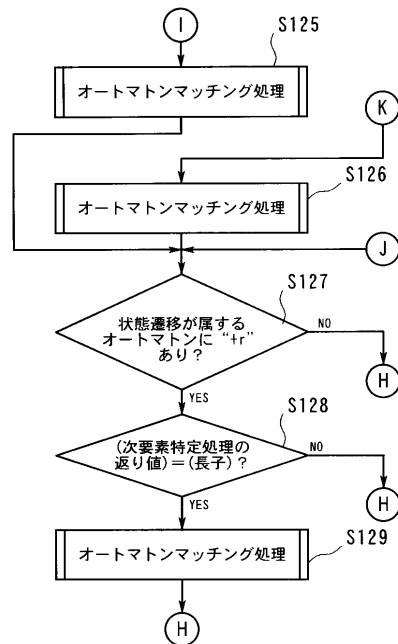
【図 16】



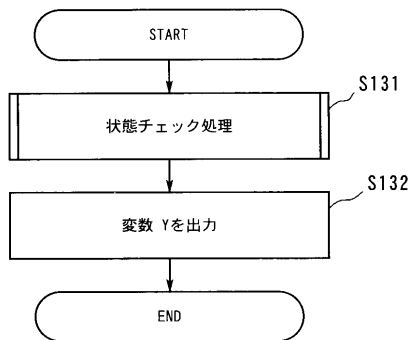
【図 17】



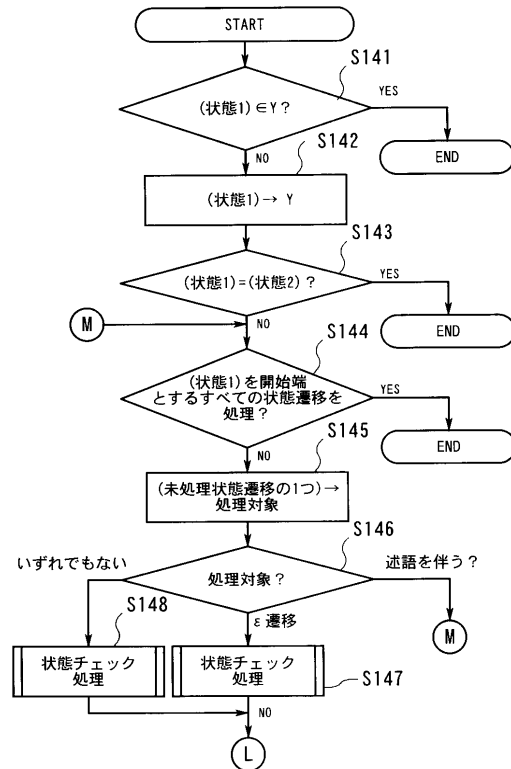
【図 18】



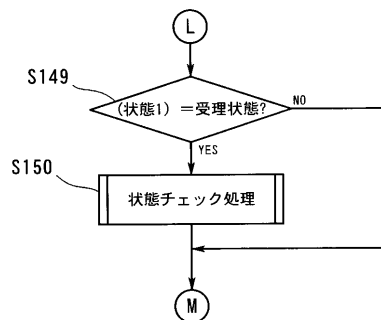
【図 19】



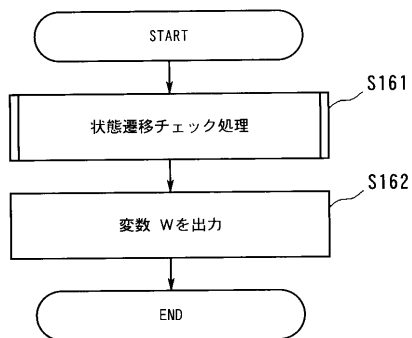
【図 20】



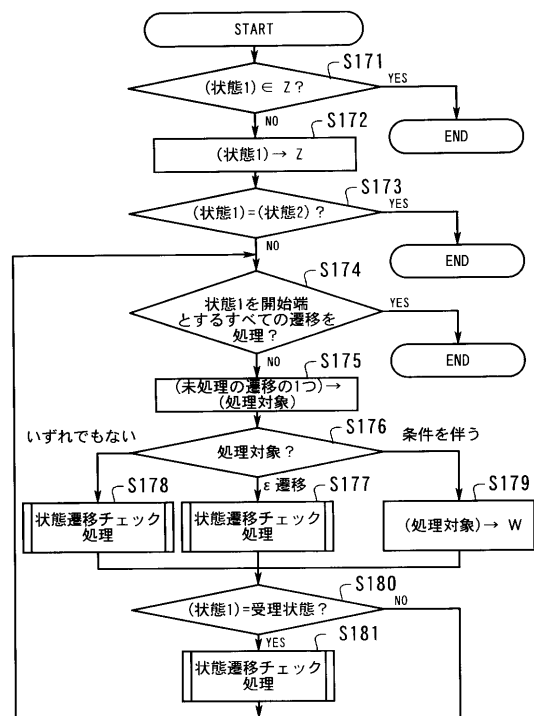
【図 21】



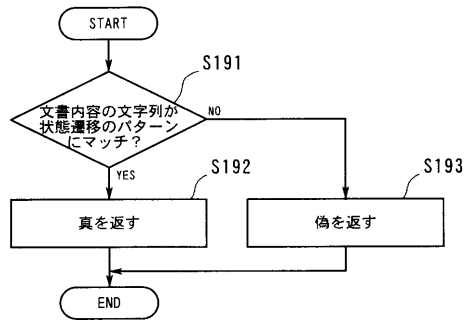
【図 22】



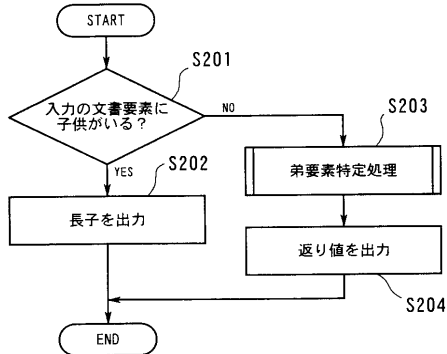
【図 23】



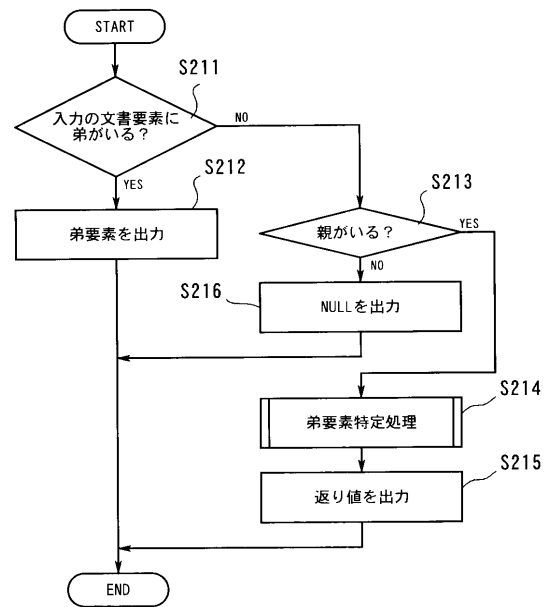
【図 24】



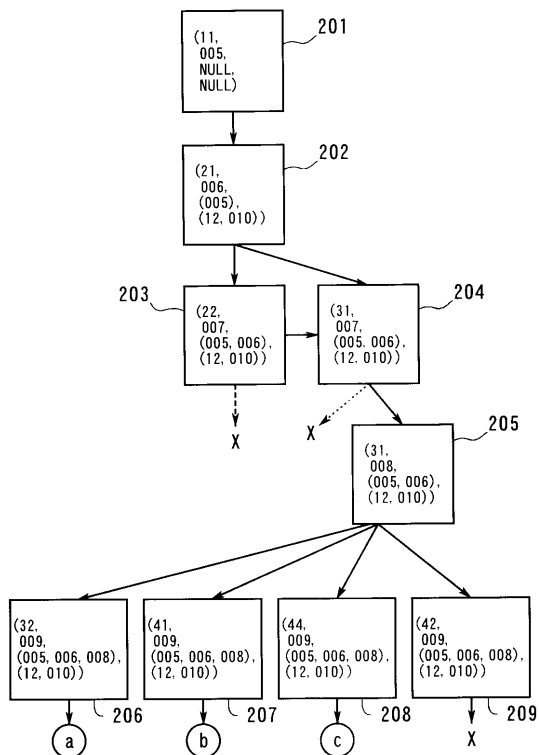
【図 25】



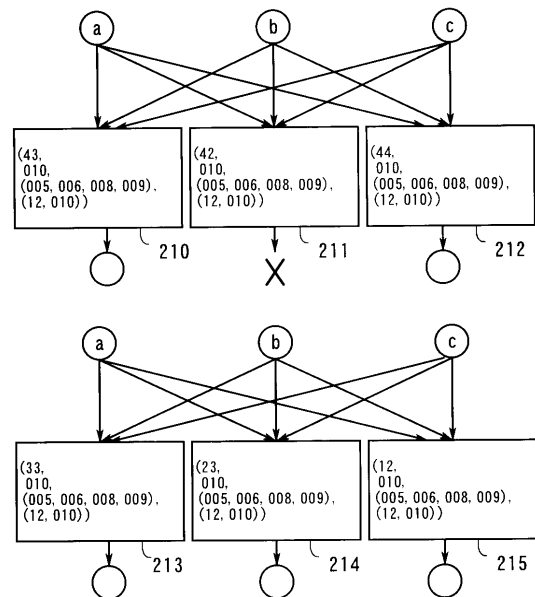
【図 26】



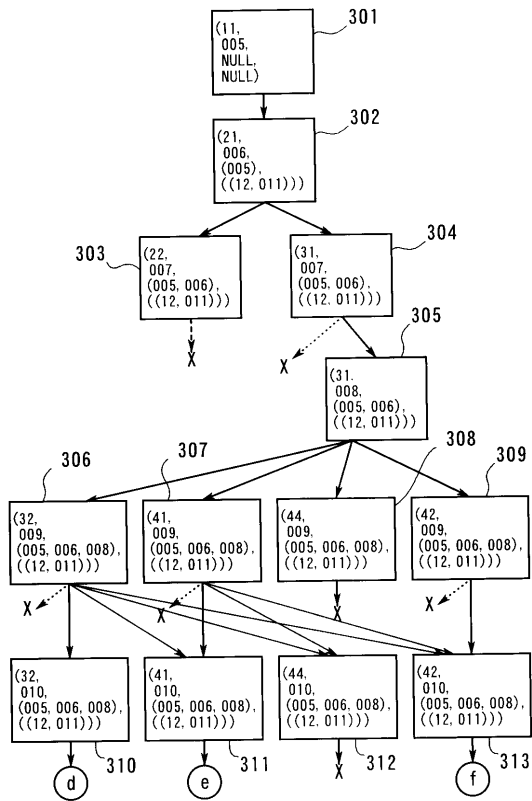
【図 27】



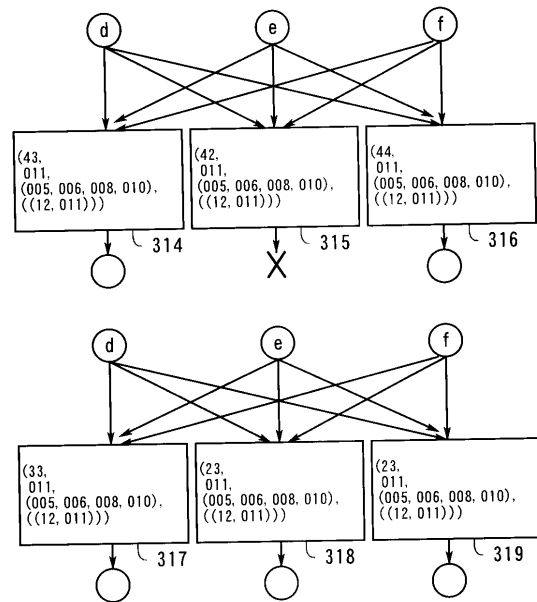
【図 28】



【図 29】



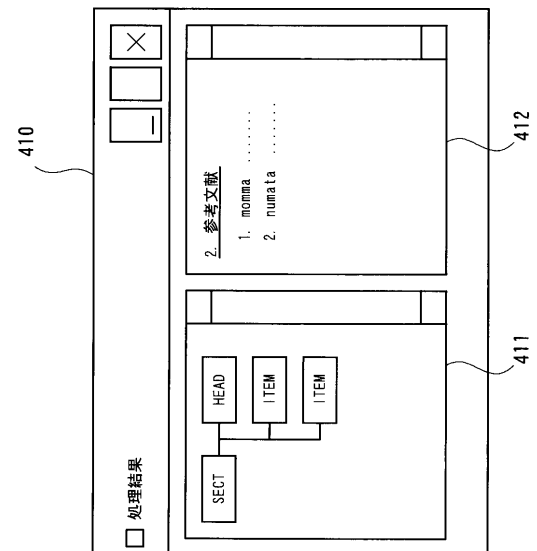
【図 30】



【図 31】

(005, 006, 008)
(005, 006, 008, 009)

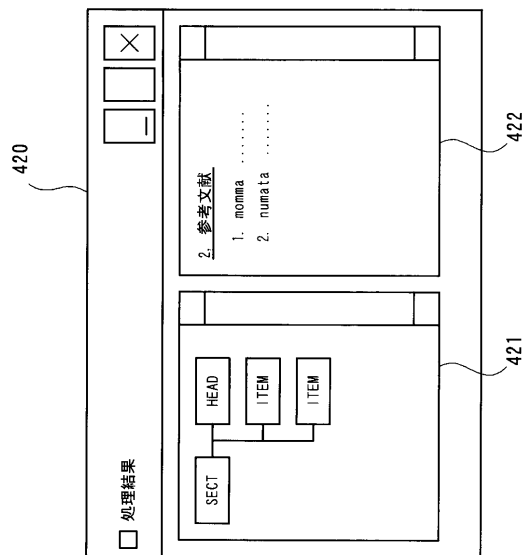
【図 33】



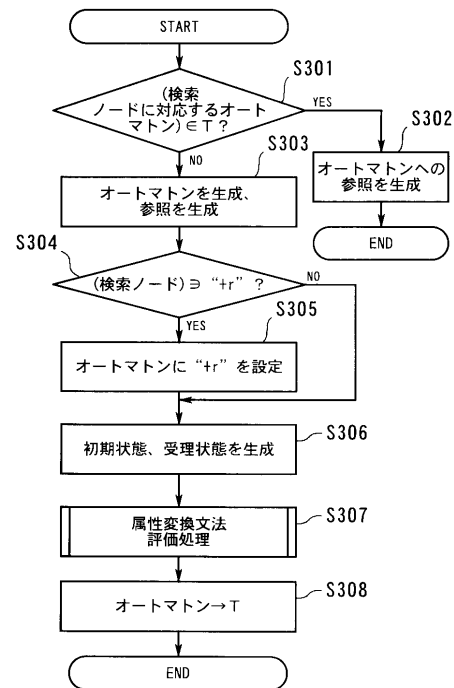
【図 32】

(005, 006, 008)
(005, 006, 008, 010)

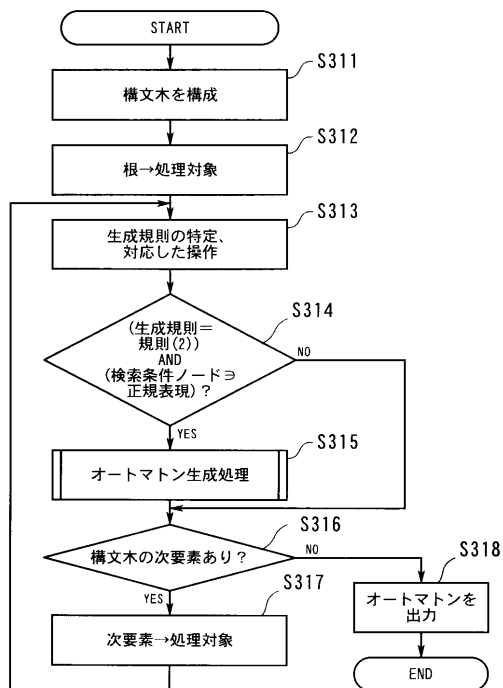
【図 34】



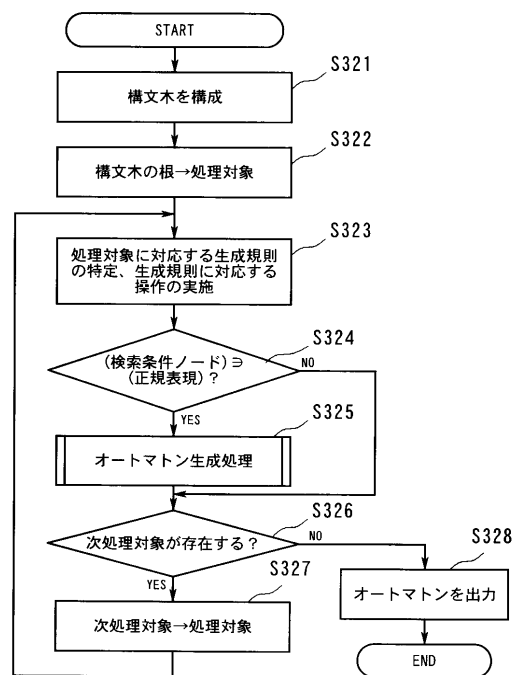
【図 35】



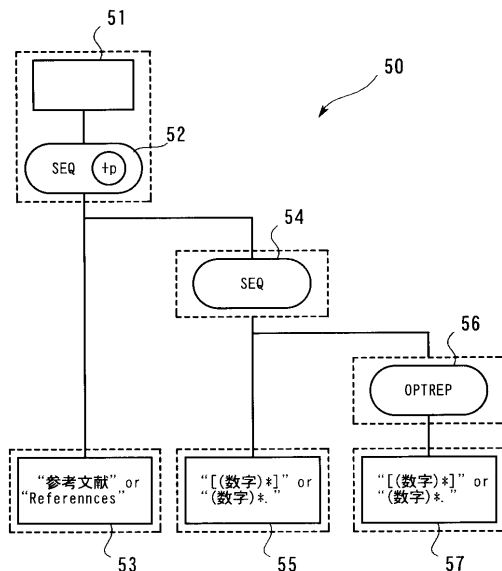
【図 36】



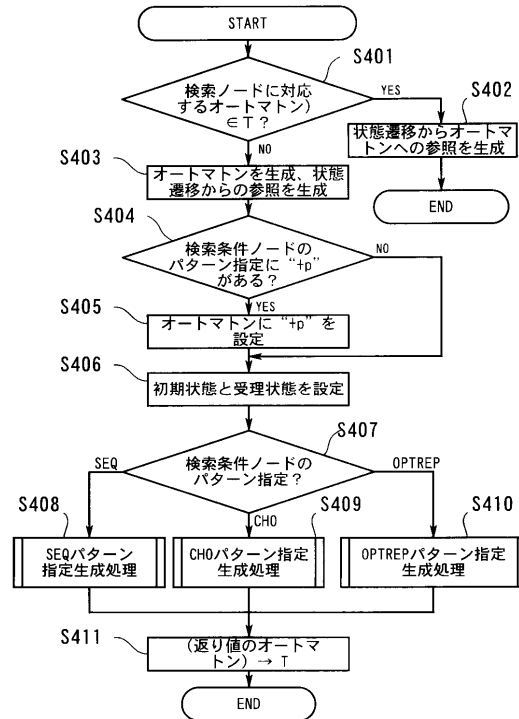
【図 37】



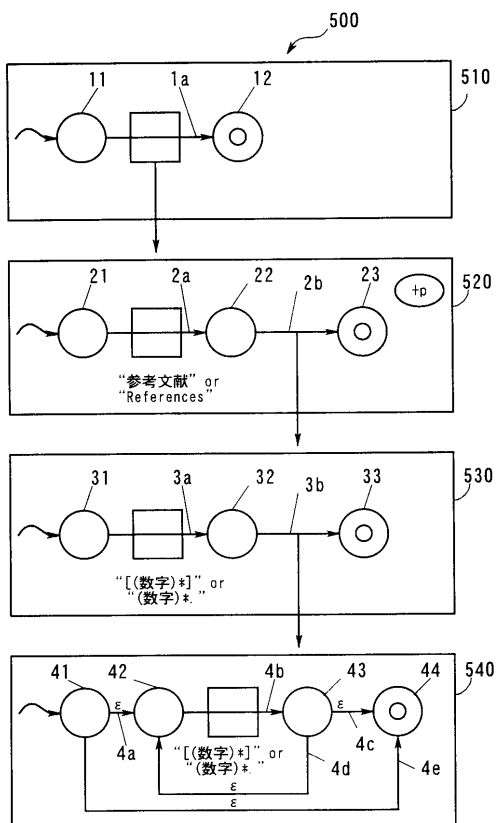
【図 38】



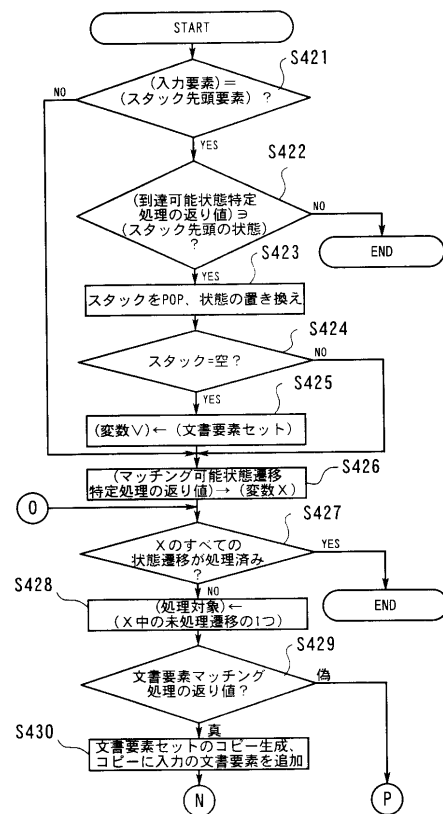
【図 39】



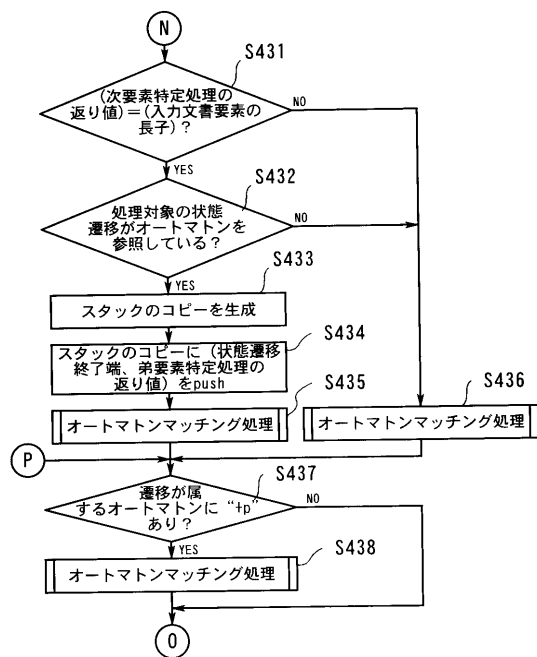
【図 40】



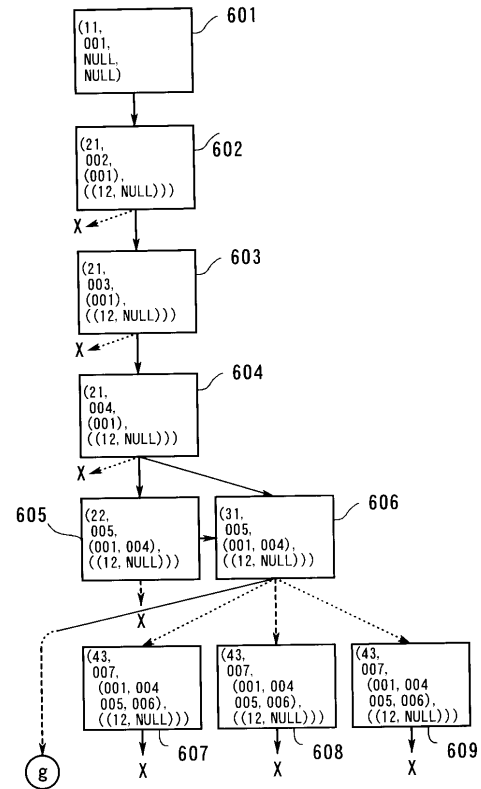
【図 41】



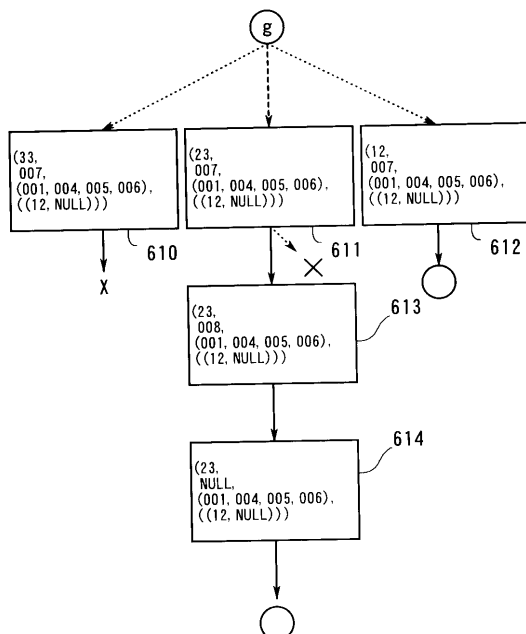
【図 4 2】



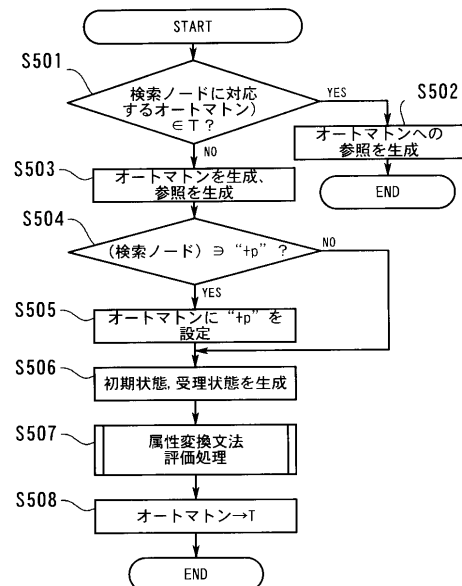
【図 4 3】



【図 4 4】



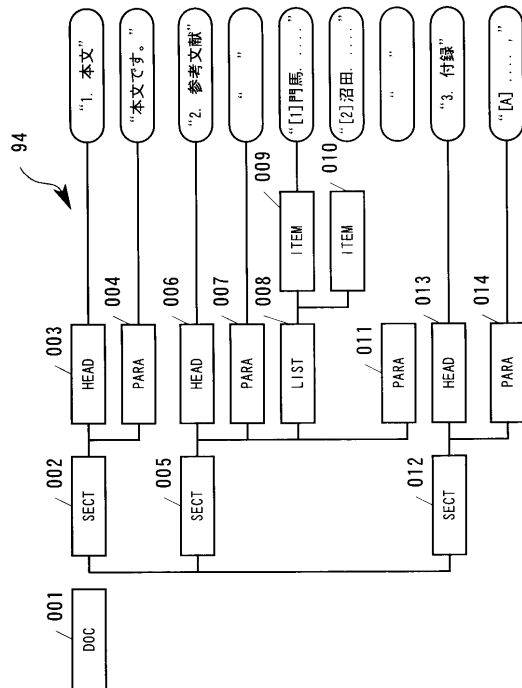
【図 4 6】



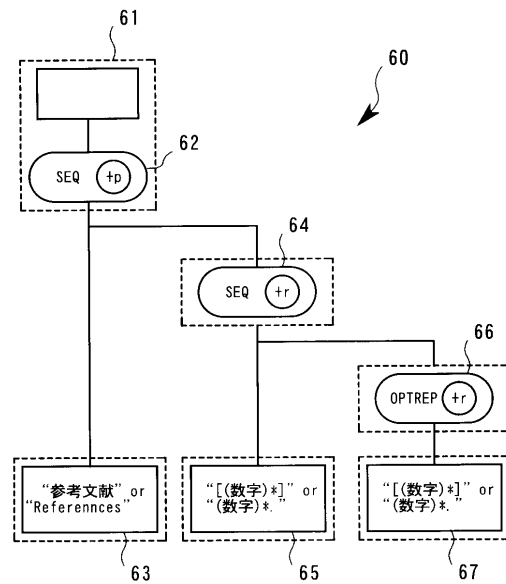
【図 4 5】

(001, 004,)
 (001, 004, 005)
 (001, 004, 005, 006)

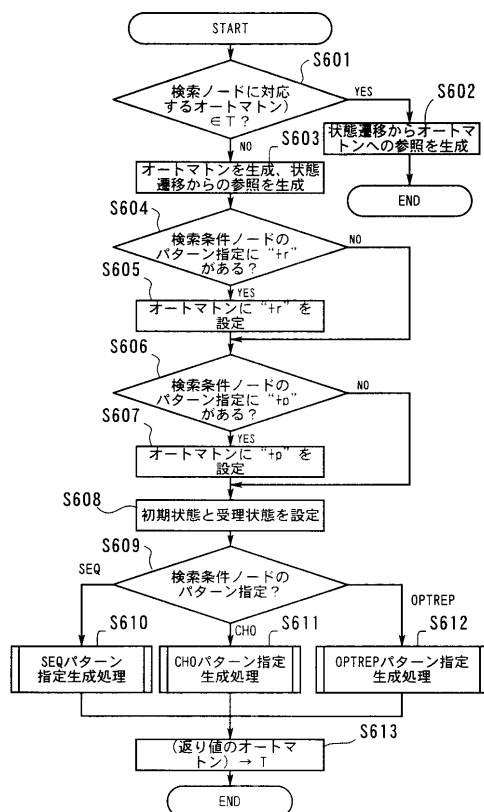
【図 47】



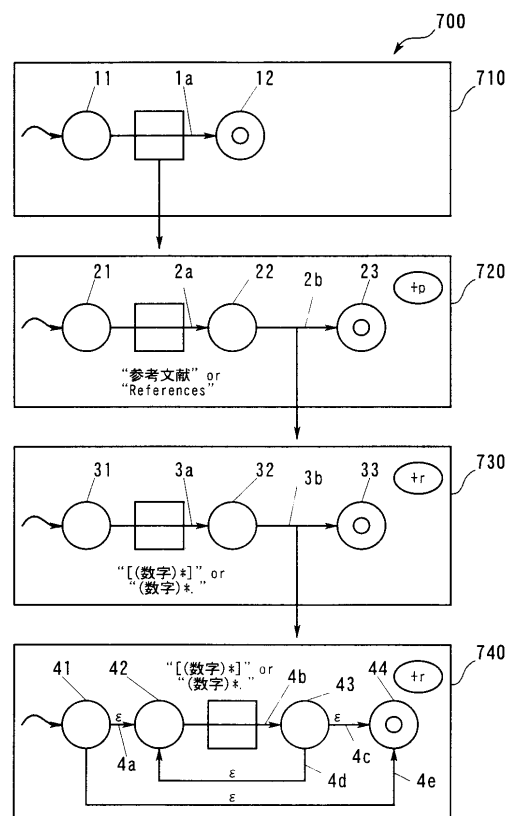
【図 48】



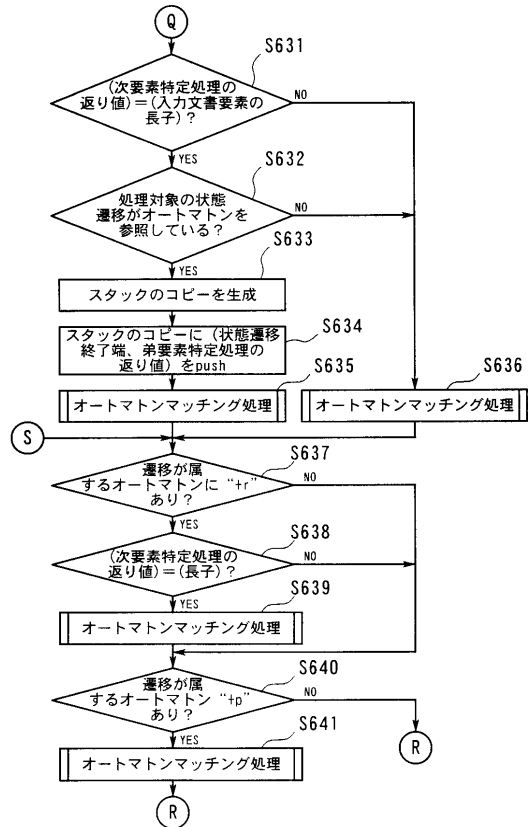
【図 49】



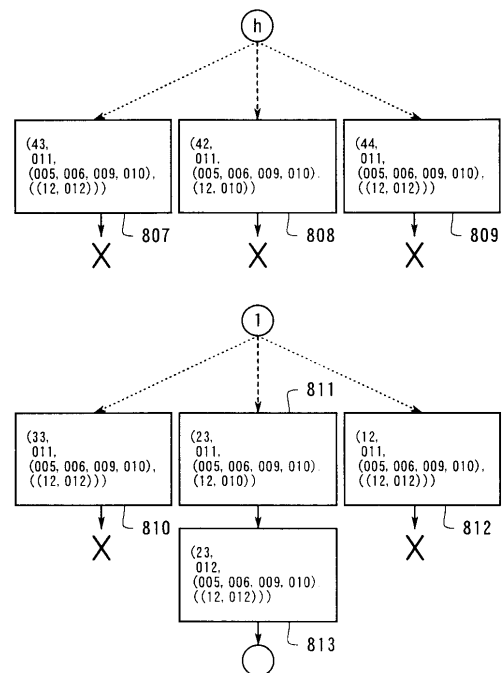
【図 50】



【 図 5 2 】



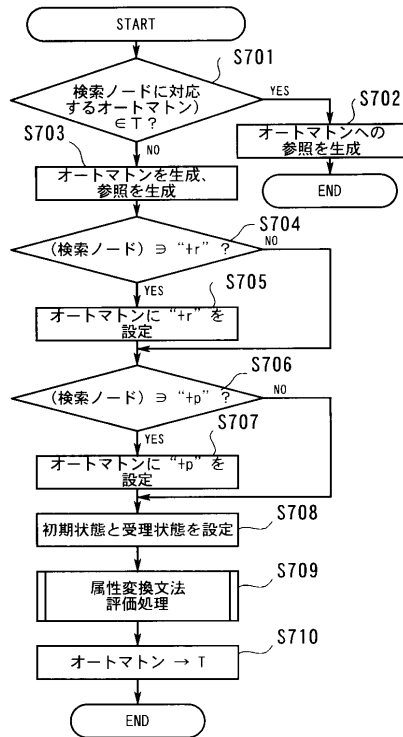
【 図 5 4 】



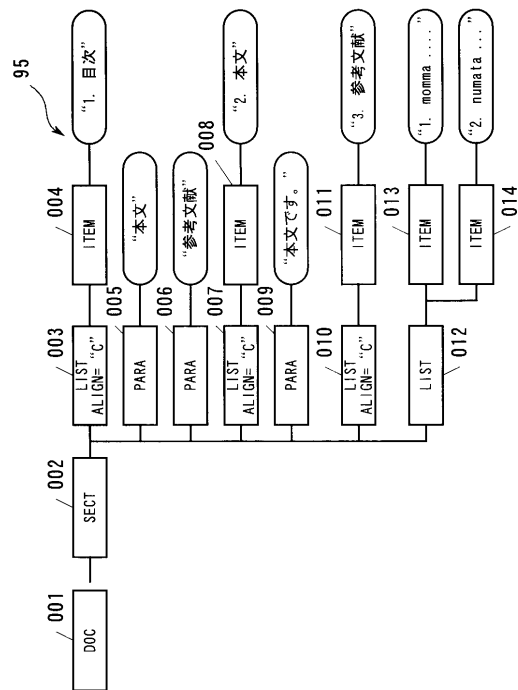
【 図 5 5 】

(005, 006)
(005, 006, 009)
(005, 006, 009, 010)

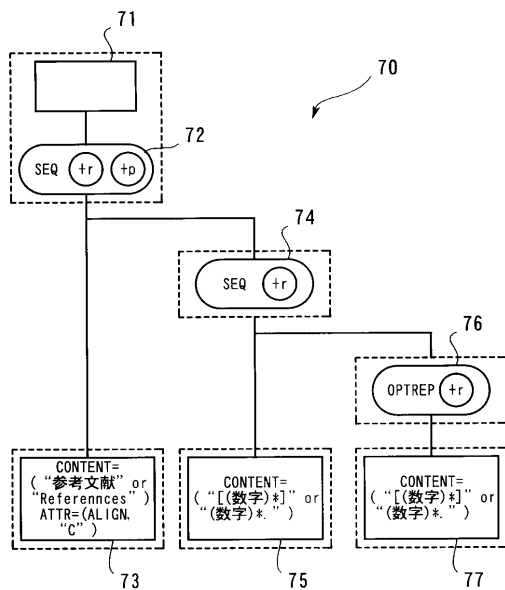
【図 56】



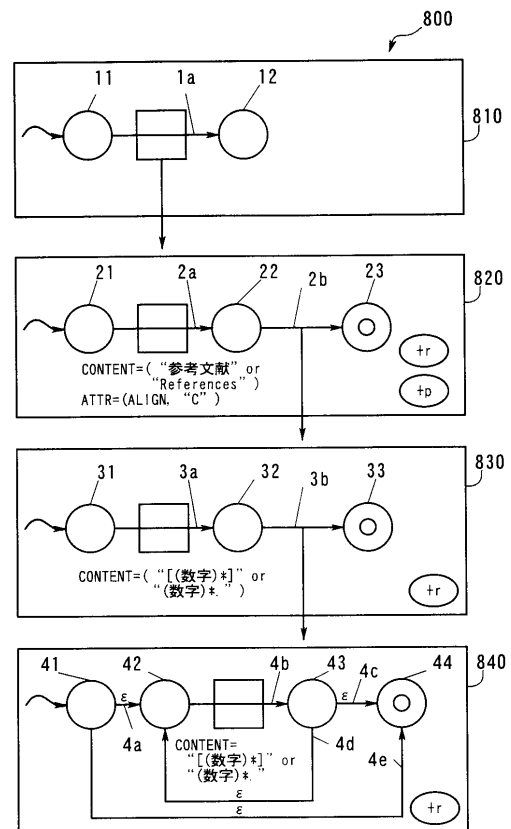
【図 57】



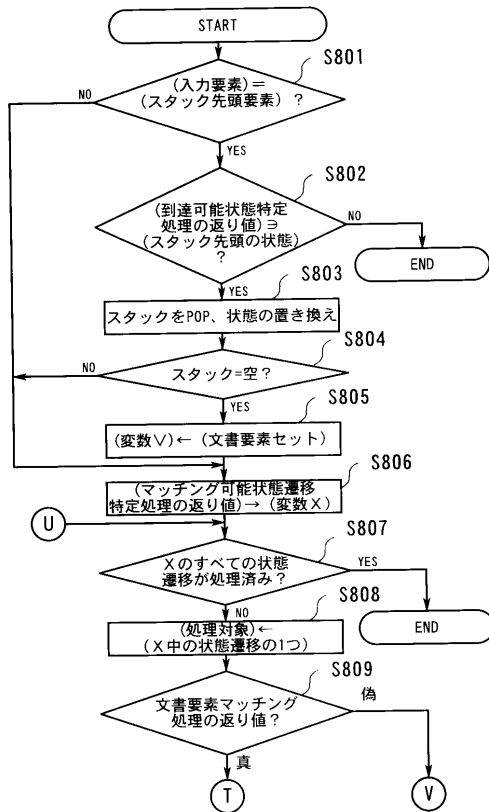
【図 58】



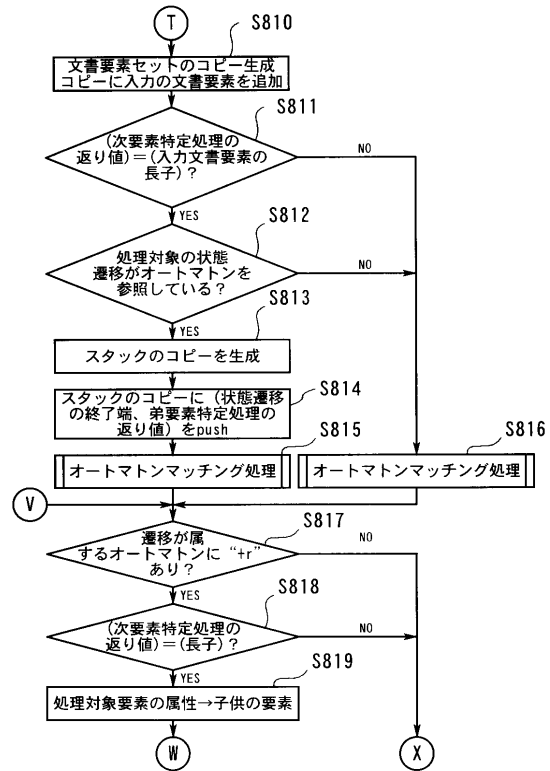
【図 59】



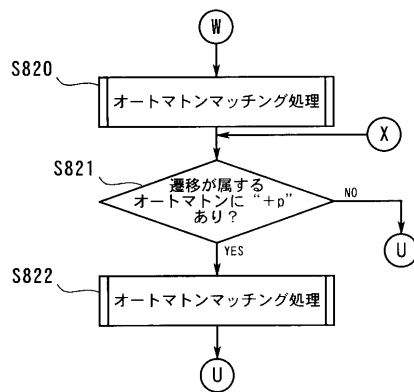
【図 60】



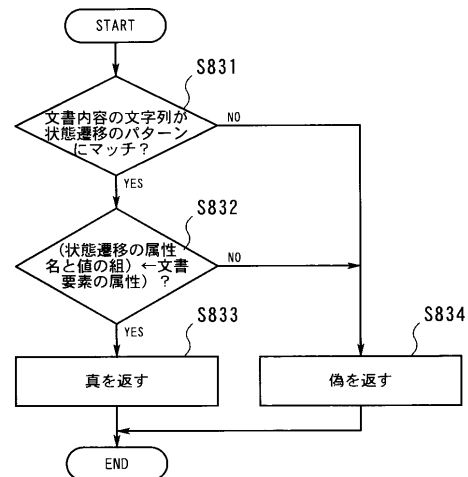
【図 61】



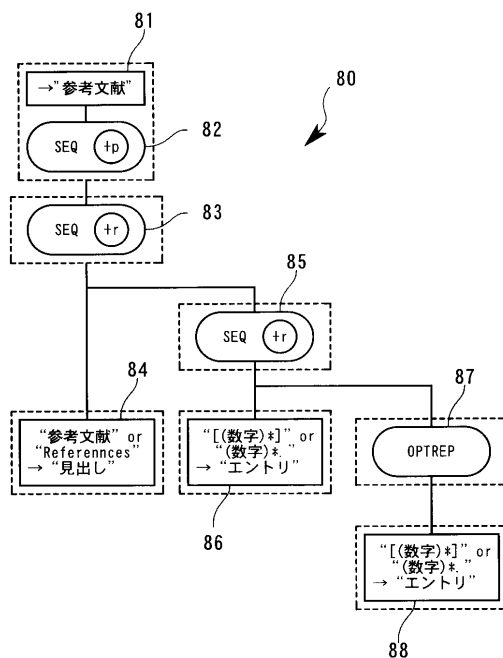
【図 62】



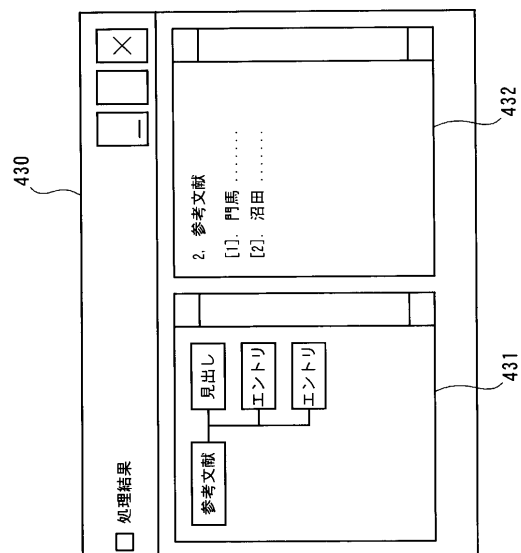
【図 63】



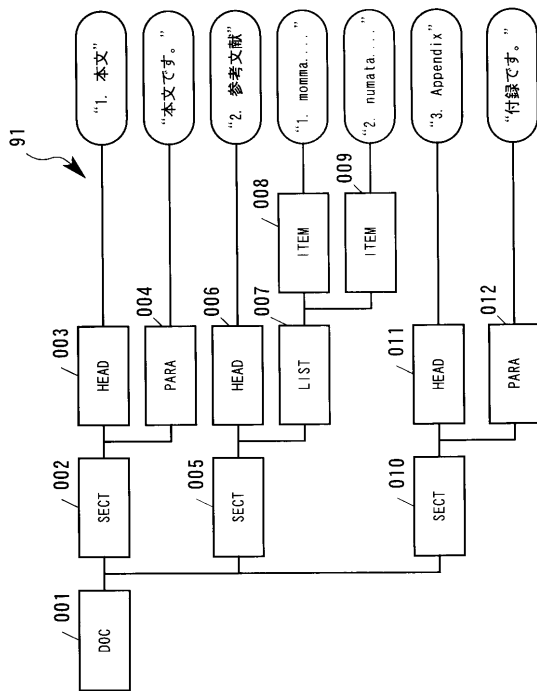
【 図 6 5 】



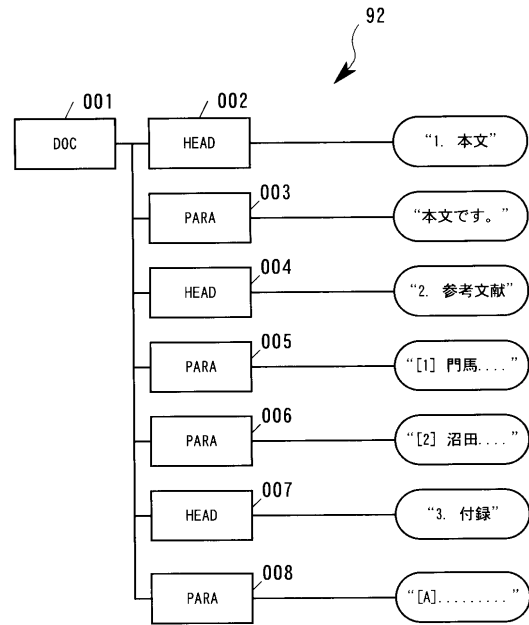
【 図 6 7 】



【図 68】



【図 69】



フロントページの続き

(58)調査した分野(Int.Cl. , D B名)

G06F 17/30

(54)【発明の名称】データ処理装置、文書処理装置、データ処理プログラムを記録したコンピュータ読み取り可能な記録媒体、文書処理プログラムを記録したコンピュータ読み取り可能な記録媒体、データ処理方法、および文書処理方法