

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5265378号  
(P5265378)

(45) 発行日 平成25年8月14日(2013.8.14)

(24) 登録日 平成25年5月10日(2013.5.10)

(51) Int.Cl.

F I

G O 6 F 17/30 (2006.01)

G O 6 F 17/30 4 1 5

G O 6 F 17/30 1 7 O A

請求項の数 38 (全 32 頁)

(21) 出願番号	特願2008-543412 (P2008-543412)	(73) 特許権者	508137947
(86) (22) 出願日	平成18年11月29日 (2006.11.29)		エクセジエ・インコーポレイテッド
(65) 公表番号	特表2009-517782 (P2009-517782A)		アメリカ合衆国、ミズーリ・63127、
(43) 公表日	平成21年4月30日 (2009.4.30)		セント・ルイス、サウス・ガイヤー・ロー
(86) 国際出願番号	PCT/US2006/045653		ド・3668、スイート・300
(87) 国際公開番号	W02007/064685	(74) 代理人	100103920
(87) 国際公開日	平成19年6月7日 (2007.6.7)		弁理士 大崎 勝真
審査請求日	平成21年10月23日 (2009.10.23)	(74) 代理人	100114188
(31) 優先権主張番号	11/293,619		弁理士 小野 誠
(32) 優先日	平成17年12月2日 (2005.12.2)	(74) 代理人	100140523
(33) 優先権主張国	米国 (US)		弁理士 渡邊 千尋
		(74) 代理人	100119253
			弁理士 金山 賢教
		(74) 代理人	100124855
			弁理士 坪倉 道明

最終頁に続く

(54) 【発明の名称】 高性能正規表現パターンマッチングのための方法および装置

(57) 【特許請求の範囲】

【請求項 1】

決定性有限オートマトン (DFA) により入力文字列をパターンにマッチさせるための装置であって、DFAが、現在状態を含む複数の状態を備え、入力文字列が、複数の入力シンボルを含み、装置が、

DFAを実施するパターンマッチング回路を備え、パターンマッチング回路が、m個の入力シンボルからなるグループとして入力文字列の入力シンボルを受け取り、連続して処理するように構成され、ここで、mは1以上の整数であり、パターンマッチング回路が、(1) 複数の保存された遷移を含む遷移テーブルメモリであって、各保存された遷移が、少なくとも1つの入力シンボルに対応するデータによってインデックス付けされ、各遷移が、次状態識別子を含む、遷移テーブルメモリと、(2) 受け取られた入力シンボルグループの少なくとも1つの入力シンボルに対応するデータによってインデックス付けされる各遷移を、現在状態を考慮することなく、遷移テーブルメモリから取り出すように構成される遷移取り出しロジックと、(3) 各取り出された遷移を受け取り、取り出された遷移のどれが現在状態に対応するかを決定し、決定された遷移の次状態識別子に基づいてDFAについての次状態を決定し、入力文字列とパターンの間にマッチが存在するかどうかを決定するように構成される状態選択ロジックと、を備える、装置。

【請求項 2】

入力文字列の各入力シンボルが、アルファベットの要素である複数の入力シンボルの1つであり、パターンマッチング回路が、各受け取られた入力シンボルグループを同値類識

別子（ECI）に変換するように構成されるアルファベット符号化ロジックをさらに備え、遷移テーブル内で各遷移がそれによってインデックス付けされる入力シンボルグループに対応するデータが、その入力シンボルグループに対するECIに対応するデータを含む、請求項1に記載の装置。

【請求項3】

パターンマッチング回路と通信する正規表現コンパイラであって、パターンマッチング回路によって実施されるDFAを決定するためにユーザから正規表現を受け取り、処理し、処理された正規表現に基づいて遷移テーブルにデータ投入し、処理された正規表現に基づいてECIを定義するように構成される正規表現コンパイラをさらに備え、正規表現は入力文字列がそれに対してマッチさせられるパターンを識別する、請求項2に記載の装置。

10

【請求項4】

遷移テーブルが、メモリ内で実施され、メモリが、メモリアドレスによってインデックス付けされる複数のワードを含み、各ワードが、ランレングス符号化された複数の遷移を含む、請求項3に記載の装置。

【請求項5】

パターンマッチング回路が、インダイレクションテーブルと、インダイレクションテーブルロジックとをさらに備え、インダイレクションテーブルが、複数のエントリを含み、各エントリが、ECIによってインデックス付けされ、各エントリが、そのECIによってインデックスされる少なくとも1つの遷移を含むメモリ内でワードが保存されるメモリアドレスを指すポインタを含み、インダイレクションテーブルロジックが、受け取られた入力シンボルグループによってインデックス付けされるインダイレクションテーブルからポインタを取り出し、取り出されたポインタを遷移取り出しロジックに提供するように構成され、各遷移がそれによってインデックス付けされるその入力シンボルグループに対するECIに対応するデータが、取り出されたポインタによって識別されるメモリアドレスを含む、請求項4に記載の装置。

20

【請求項6】

各インダイレクションテーブルエントリが、遷移インデックスと、遷移カウントとをさらに含み、遷移インデックスが、入力シンボルグループに対するECIによってインデックス付けされる第1の遷移が、そのエントリのポインタによって識別されるワード内のどこに配置されているかを識別し、遷移カウントが、いくつの遷移が、入力シンボルグループに対するECIによってインデックス付けされるかを識別し、インダイレクションテーブルロジックが、受け取られた入力シンボルグループによってインデックス付けされるインダイレクションテーブルエントリから、遷移インデックスおよび遷移カウントを取り出し、取り出された遷移インデックスおよび遷移カウントを、遷移取り出しロジックに提供するように構成される、請求項5に記載の装置。

30

【請求項7】

遷移取り出しロジックが、取り出されたポインタおよび取り出された遷移カウントによって識別される各ワードを遷移テーブルから取り出すように構成され、状態選択ロジックが、遷移に対応する各取り出されたワードの、入力シンボルグループに対するECIによってインデックス付けされないどの部分もマスクアウトするようにさらに構成される、請求項6に記載の装置。

40

【請求項8】

各インダイレクションテーブルエントリが、先頭遷移インデックスと、末尾遷移インデックスと、ワードカウントとをさらに含み、先頭遷移インデックスが、入力シンボルグループに対するECIによってインデックス付けされる第1の遷移が、そのエントリのポインタによって識別されるワード内のどこに配置されているかを識別し、ワードカウントが、受け取られた入力シンボルグループに対するECIによってインデックス付けされるすべての遷移がその中に連続的に保存されるワードの数を識別し、末尾遷移インデックスが、入力シンボルグループに対するECIによってインデックス付けされる最後の遷移が、

50

そのエントリのワードカウントによって識別される最後のワード内のどこに配置されているかを識別し、インダイレクションテーブルロジックが、受け取られた入力シンボルグループによってインデックス付けされるインダイレクションテーブルエントリから、先頭遷移インデックス、末尾遷移インデックス、およびワードカウントを取り出し、取り出された先頭遷移インデックス、末尾遷移インデックス、およびワードカウントを、遷移取り出しロジックに提供するように構成される、請求項 5 に記載の装置。

【請求項 9】

遷移取り出しロジックが、取り出されたポインタおよび取り出されたワードカウントによって識別される各ワードを、遷移テーブルから取り出すように構成され、状態選択ロジックが、遷移に対応する各取り出されたワードの、入力シンボルグループに対する E C I によってインデックス付けされないどの部分もマスクアウトするようにさらに構成される、請求項 8 に記載の装置。

10

【請求項 10】

各遷移が、事前計算ランレングスプレフィックス合計をさらに含む、請求項 5 に記載の装置。

【請求項 11】

m が、1 よりも大きい、請求項 5 に記載の装置。

【請求項 12】

正規表現コンパイラが更に、処理された正規表現に基づいてインダイレクションテーブルにデータ投入するように構成される、請求項 5 に記載の装置。

20

【請求項 13】

パターンマッチング回路が、FPGA 上で実施される、請求項 5 に記載の装置。

【請求項 14】

パターンマッチング回路が、ASIC 上で実施される、請求項 5 に記載の装置。

【請求項 15】

各遷移が、マッチフラグをさらに含み、状態選択ロジックが、決定された遷移のマッチフラグに基づいて、入力文字列とパターンの間にマッチが存在するかどうかを決定するようにさらに構成される、請求項 3 に記載の装置。

【請求項 16】

DFA が、位置独立の DFA を含む、請求項 15 に記載の装置。

30

【請求項 17】

各遷移が、マッチリスタートフラグをさらに含み、状態選択ロジックが、決定された遷移のマッチリスタートフラグに基づいて、パターンマッチングプロセスが DFA においてリスタートされたかどうかを決定するようにさらに構成される、請求項 16 に記載の装置。

【請求項 18】

m が、1 よりも大きい、請求項 17 に記載の装置。

【請求項 19】

パターンマッチング回路と通信する正規表現コンパイラであって、パターンマッチング回路によって実施される DFA を決定するためにユーザから正規表現を受け取り、処理し、処理された正規表現に基づいて遷移テーブルにデータ投入するように構成される正規表現コンパイラをさらに備え、正規表現は入力文字列がそれに対してマッチさせられるパターンを識別する、請求項 1 に記載の装置。

40

【請求項 20】

パターンマッチング回路が、FPGA 上で実施される、請求項 19 に記載の装置。

【請求項 21】

パターンマッチング回路が、ASIC 上で実施される、請求項 19 に記載の装置。

【請求項 22】

パターンマッチング回路が、FPGA 上で実施される、請求項 1 に記載の装置。

【請求項 23】

50

パターンマッチング回路が、A S I C上で実施される、請求項 1 に記載の装置。

【請求項 2 4】

パターンマッチング回路が、マルチプロセッサシステム上で実施される、請求項 1 に記載の装置。

【請求項 2 5】

遷移取り出しロジックと状態選択ロジックが、並列に動作するパイプラインステージとして構成される、請求項 1 から 2 4 のいずれか一項に記載の装置。

【請求項 2 6】

状態機械を使用して入力文字列をパターンにマッチさせる方法であって、入力文字列が、複数の入力シンボルを含み、各入力シンボルが、複数のビットを含み、方法が、

状態機械についての複数の状態を定義するステップであって、状態機械が、現在状態を有し、入力文字列の入力シンボルが受け取られ、処理されるのにつれて、複数の状態を経ながら進行し、それによって、入力文字列の複数の入力シンボルがパターンとマッチするかどうかを検出するように構成される、ステップと、

状態機械の状態間の複数の遷移を定義するステップであって、各遷移が、状態機械の現在状態に対応するデータと、少なくとも 1 つの入力シンボルに対応するデータとによってインデックス付けされ、各遷移が、状態機械の次状態についての識別子と、複数の入力シンボルとパターンの間にマッチが存在するかどうかを示すマッチフラグとを含む、ステップと、

入力文字列の少なくとも 1 つの入力シンボルを受け取るステップと、

少なくとも 1 つの受け取られた入力シンボルに基づいて及び現在状態を考慮することなく、少なくとも 1 つの定義された遷移を取り出すステップと、

少なくとも 1 つの取り出された遷移のどれが現在状態によってインデックス付けされているかを決定するステップと、

状態機械の現在状態を、決定された取り出された遷移の次状態識別子によって識別される状態として更新するステップと、

取り出された遷移のマッチフラグに基づいて、複数の入力シンボルとパターンの間にマッチが存在するかどうかを決定するステップと、

入力文字列の追加の入力シンボルが受け取られた場合に、取り出し、遷移決定、更新、およびマッチ決定のステップを繰り返すステップと、を備える、方法。

【請求項 2 7】

各遷移が、受け取られた少なくとも 1 つの入力シンボルがパターンに対するマッチングプロセスのリスタートをもたらすかどうかを示すマッチリスタートフラグをさらに含み、方法が、

取り出された遷移のマッチリスタートフラグに基づいて、マッチリスタートが存在するかどうかを決定するステップをさらに備え、

繰り返すステップが、入力文字列の追加の入力シンボルが受け取られた場合に、取り出し、遷移決定、更新、マッチ決定、および部分的マッチ開始決定のステップを繰り返すステップを備える、請求項 2 6 に記載の方法。

【請求項 2 8】

受け取るステップが、各クロックサイクルについて複数の m 個の入力シンボルを受け取るステップを備え、取り出すステップが、m 個の受け取られた入力シンボルに基づいて、少なくとも 1 つの定義された遷移を取り出すステップを備える、請求項 2 7 に記載の方法。

【請求項 2 9】

状態機械によって処理され得る入力シンボルのアルファベットに対して、各入力シンボルが同値類識別子にマッピングされるように、複数の同値類識別子を定義するステップであって、各遷移がそれによってインデックス付けされる少なくとも 1 つの入力シンボルに対応するデータが、その少なくとも 1 つの入力シンボルに対する同値類識別子を含む、ステップと、

受け取られた少なくとも1つの入力シンボルをその同値類識別子にマッピングするステップと、をさらに備え、

取り出すステップが、少なくとも1つの受け取られた入力シンボルに対する同値類識別子に基づいて、少なくとも1つの定義された遷移を取り出すステップを備え、

繰り返すステップが、入力文字列の追加の入力シンボルが受け取られた場合に、マッピング、取り出し、遷移決定、更新、マッチ決定、および部分的マッチ開始決定のステップを繰り返すステップを備える、請求項28に記載の方法。

【請求項30】

状態機械によって処理され得る入力シンボルのアルファベットに対して、各入力シンボルが同値類識別子にマッピングされるように、複数の同値類識別子を定義するステップであって、各遷移がそれによってインデックス付けされる少なくとも1つの入力シンボルに対応するデータが、その少なくとも1つの入力シンボルに対する同値類識別子に対応するデータを含む、ステップと、

受け取られた少なくとも1つの入力シンボルをその同値類識別子にマッピングするステップと、をさらに備え、

取り出すステップが、少なくとも1つの受け取られた入力シンボルに対する同値類識別子に基づいて、少なくとも1つの定義された遷移を取り出すステップを備え、

繰り返すステップが、入力文字列の追加の入力シンボルが受け取られた場合に、マッピング、取り出し、遷移決定、更新、マッチ決定、および部分的マッチ開始決定のステップを繰り返すステップを備える、請求項26に記載の方法。

【請求項31】

ランレングス符号化によってメモリ内で遷移を圧縮するステップをさらに備える、請求項30に記載の方法。

【請求項32】

圧縮するステップが、事前計算プレフィックス合計を有するランレングス符号化によってメモリ内で遷移を圧縮するステップをさらに備える、請求項31に記載の方法。

【請求項33】

メモリが、複数のワードを含み、各ワードが、メモリアドレスによってインデックス付けされ、各ワードが、複数のx個の圧縮された遷移を含み、方法が、

少なくとも1つの受け取られた入力シンボルに対する同値類識別子に基づいて、圧縮された遷移にアクセスするためのインダイレクションテーブルを定義するステップであって、インダイレクションテーブルが、複数のエントリを含み、各エントリが、同値類識別子によってインデックス付けされ、各エントリが、少なくとも1つの受け取られた入力シンボルに対する同値類識別子によってインデックス付けされる少なくとも1つの圧縮された遷移を含むワードが保存されるメモリアドレスを指すポインタを含む、ステップをさらに備え、

各遷移がそれによってインデックス付けされる同値類識別子に対応するデータが、その同値類識別子に対するインダイレクションテーブルポインタを含み、方法が、

同値類識別子によってインデックス付けされるポインタを決定するために、少なくとも1つの受け取られた入力シンボルに対する同値類識別子に基づいてインダイレクションテーブルにアクセスするステップをさらに備え

取り出すステップが、少なくとも1つの受け取られた入力シンボルの同値類識別子に対する決定されたポインタに基づいて、少なくとも1つの定義された遷移を取り出すステップを備え、

繰り返すステップが、入力文字列の追加の入力シンボルが受け取られた場合に、マッピング、アクセス、取り出し、遷移決定、更新、マッチ決定、および部分的マッチ開始決定のステップを繰り返すステップを備える、請求項30に記載の方法。

【請求項34】

各インダイレクションテーブルエントリが、遷移インデックスと、遷移カウントとをさらに含み、遷移インデックスが、少なくとも1つの受け取られたシンボルの同値類識別子

10

20

30

40

50

に対する第 1 の圧縮された遷移が、ポインタによって識別されるワード内のどこに配置されているかを識別し、遷移カウントが、少なくとも 1 つの受け取られたシンボルの同値類識別子に対していくつの圧縮された遷移が、メモリ内に連続的に保存されているかを識別し、アクセスするステップが、同値類識別子によってインデックス付けされるポインタ、遷移インデックス、および遷移カウントを決定するために、少なくとも 1 つの受け取られた入力シンボルに対する同値類識別子に基づいて、インダイレクションテーブルにアクセスするステップを備え、取り出すステップが、少なくとも 1 つの受け取られた入力シンボルの同値類識別子に対する決定されたポインタ、遷移インデックス、および遷移カウントに基づいて、各定義された遷移を取り出すステップをさらに備える、請求項 3 3 に記載の方法。

10

#### 【請求項 3 5】

各インダイレクションテーブルエントリが、先頭遷移インデックスと、末尾遷移インデックスと、ワードカウントとをさらに含み、先頭遷移インデックスが、少なくとも 1 つの受け取られたシンボルの同値類識別子に対する第 1 の圧縮された遷移が、ポインタによって識別されるワード内のどこに配置されているかを識別し、ワードカウントが、少なくとも 1 つの受け取られたシンボルの同値類識別子に対する少なくとも 1 つの圧縮された遷移がメモリ内で連続的に保存されるワードの数を示し、末尾遷移インデックスが、最後の圧縮された遷移が、ワードカウントによって識別される最後のワード内のどこに配置されているかを識別し、アクセスするステップが、同値類識別子によってインデックス付けされるポインタ、先頭遷移インデックス、末尾遷移インデックス、およびワードカウントを決定するために、少なくとも 1 つの受け取られた入力シンボルに対する同値類識別子に基づいて、インダイレクションテーブルにアクセスするステップを備える、請求項 3 3 に記載の方法。

20

#### 【請求項 3 6】

取り出すステップが、決定されたポインタおよび決定されたワードカウントによって識別される各ワードを、メモリから取り出すステップをさらに備え、遷移決定ステップが、決定された先頭遷移インデックス、決定された末尾遷移インデックス、および決定されたワードカウントによって識別される範囲内にはない各取り出されたワード内に配置されるどの圧縮された遷移もマスクアウトするステップを備える、請求項 3 5 に記載の方法。

#### 【請求項 3 7】

状態定義ステップが、正規表現の集合から状態機械についての複数の状態を定義するステップを備える、請求項 3 6 に記載の方法。

30

#### 【請求項 3 8】

取り出すステップと少なくとも 1 つの取り出された遷移のどれが現在状態によってインデックス付けされているかを決定するステップが、パイプライン化された方法で並列に実行される、請求項 2 6 から 3 7 のいずれか一項に記載の方法。

#### 【発明の詳細な説明】

#### 【技術分野】

#### 【0 0 0 1】

本発明は、一般に、データシンボルストリームに属するいずれかの文字列がパターンとマッチするかどうかを決定するための、データシンボルからなるストリームの処理の分野に関する。

40

#### 【背景技術】

#### 【0 0 0 2】

ネットワークおよびストレージサブシステム設計における進歩が、コンピュータシステム間およびコンピュータシステム内でデータストリームが処理されなければならない速度を押し上げ続けている。一方で、すべてのレベルのコンポーネントが、時間的制約のある動作をトリガし得るパターンに関してストリームを調べるので、そのようなデータストリームの内容は、ますます増加する精査に委ねられる。パターンは、文字列定数（例えば「dog」および「cat」）ばかりでなく、いくつかを挙げれば、クレジットカード番号

50

、通貨価値、または電話番号などを表す、明細事項も含むことができる。広く使用されているパターン指定言語は、正規表現言語である。正規表現と、決定性有限オートマトン (DFA) によるそれらの実施は、よく発達した分野である。Hopcroft および Ullman、*「Introduction to Automata Theory, Languages, and Computation」*、Addison Wesley、1979 年を参照されたく、その開示の全体は、参照により本明細書に組み込まれる。以下で説明されるように、DFA は、状態機械 (state machine) の動作を定義する論理表現である。しかし、高性能パターンマッチングとの関連において、正規表現の使用を改善する必要性が当技術分野に存在すると、発明者らは本明細書において確信する。

10

#### 【0003】

パケットヘッダフィルタリングなど、いくつかのアプリケーションでは、与えられたパターンの位置は、固定されることができ、位置の固定は、パターンがデータストリーム内の 1 組の定められた位置で開始または終了する場合にのみ、マッチが発生する状況を示す。より一般的に、多くのアプリケーションでは、パターンは、データストリーム (例えば、非構造化データストリーム、パケットペイロードなど) 内の任意の場所で開始または終了することができる。いくつかのアプリケーションは、データストリームのあらゆるバイトにおいて、数千のパターンの同時組み付けを必要とする。そのようなアプリケーションの例は、

(約 10000 個のパターンからなるルールベースを使用して一般に動作する) ネットワーク侵入検出 / 防止システム (Roesch, M., *「Snort - lightweight intrusion detection for networks」*、LISA '99: 13<sup>th</sup> Systems Administration Conference、229 - 238 ページ、1999 年を参照されたく、その開示の全体は、参照により本明細書に組み込まれる)、

20

不適切または不法な内容を探し出すために発信電子メールをスキャンする電子メール監視システム、

ユーザ指定パターンを課して着信電子メールをフィルタリングするスパムフィルタ、

有害であることが知られているプログラムのシグニチャについてフィルタリングを行うウイルススキャナ、および

30

海賊版コンテンツを探し出すためにメディアファイルまたはソケットストリームをスキャンする著作権実施プログラム

を含むがこれらに限定されない。上記のようなアプリケーションでは、データストリーム内で探されるパターンの集合は、日々変化し得る。

#### 【0004】

今日の従来型ハイエンドワークステーションは、高速ネットワークおよびストレージサブシステムから発信されるデータストリームの速度を与えられるパターンマッチングアプリケーションと同じ速度を維持することができない。この性能ギャップに対処するため、発明者らは本明細書において、パイプラインアーキテクチャでの DFA の定式化および実現 (例えば、ハードウェアロジック、ネットワーク化プロセッサ、またはその他のパイプライン化処理システム) におけるアーキテクチャ革新に解決策を求める。

40

#### 【0005】

正規表現  $r$  は、正規言語  $L(r)$  を表し、その場合、言語は、(有限) 文字列からなる (おそらく無限) 集合である。各文字列は、アルファベット から取り出されたシンボルで構成される。正規表現の構文は、以下の基本表現を用いて、帰納的に定義される：

シンボル は、 $\{ \}$  を表し、

シンボル は、空 (ゼロ幅) 文字列を含む単集合を表し、

シンボル は、空集合を表す。

上記の各々は、正規言語である。より複雑な正規表現は、当技術分野でよく知られているように、合併演算子、連接演算子、およびクリーネ閉包演算子を使用して、構成されるこ

50

とができる。シンボル範囲指示子および節反復係数が、統語上の便宜のため、一般に提供される。よく知られた正規表現表記および拡張のいずれもが、本発明を実施する際の使用に適しているが、perlが普及しているため、本明細書の説明および本発明の好ましい実施形態は、正規表現用にperl表記および拡張を支持する。

#### 【0006】

上で述べられたように、正規表現は、ファイル検索システムおよびネットワーク侵入検出システムを含むが、これらに限定されない、きわめて多くの検索アプリケーションにおいて、実際の用途を見出す。ほとんどのテキストエディタおよび検索ユーティリティは、何らかの形式の正規表現構文を使用して、検索ターゲットを指定する。perl構文を使用する例示的な一例として、図1に示されるパターンは、米国の通貨価値を表す文字列とマッチするように意図されている。

10

#### 【0007】

バックスラッシュ「\」は、文字通りに解釈されるべき任意の特別なシンボルの前に置かれる、

文字範囲の下限値および上限値は、ダッシュ「-」を使用して指定され、

「+」符号は、直前の表現が1回以上繰り返され得ることを表し、

中括弧内の単一の数は、直前の表現が指定された回数だけ正確に繰り返され得ることを表し、

中括弧内の対をなす数は、直前の表現についての繰り返しの範囲を表す。

したがって、上記の表現とマッチする文字列は、シンボル「\$」で始まり、それに10進数字からなる何らかの正の数が続き、その文字列には、任意選択で、10進小数点「.」と、さらに正確に2つの10進数字が続く。実際には、そのようなマッチについてのパターンは、マッチが(空白など)何らかの区切り符号によって囲われることも指定することができ、その結果、文字列「\$347.12」は、4つのマッチ(すなわち、「\$3」、「\$34」、「\$347」、「\$347.12」)ではなく、1つのマッチをもたらす。

20

#### 【0008】

正規表現を使用して対象とするパターンを指定するアプリケーションは、一般に以下のように動作する。正規表現rと、ターゲット文字列t(一般にファイルなど何らかの入力ストリームの内容)を所与として、L(r)内でtの部分文字列をすべて見出す。部分文字列は、t内のその位置によって一般に報告される。したがって、別途指摘がない限り、パターンrがターゲット内のあらゆる位置で適用され、すべてのマッチが報告されることが、一般に意図される。

30

#### 【0009】

正規表現を使用して指定されるパターンを認識するための最も簡単で最も実用的なメカニズムは、DFAであり、DFAは、形式的に5つ組

$(Q, \Sigma, q_0, A)$

として記述され、ここで、

Qは、状態の有限集合であり、

$\Sigma$ は、入力シンボルのアルファベットであり、

$q_0 \in Q$ は、DFAの初期状態であり、

Aは、遷移関数：

40

#### 【数1】

$$Q \times \Sigma \mapsto Q$$

であり、

A  $\subseteq$  Qは、アクセプト状態の集合である。

#### 【0010】

DFAは、以下のように動作する。DFAは、状態 $q_0$ で始まる。DFAが状態qにある場合、次の入力シンボルaは、 $(q, a)$ によって決定される遷移を引き起こす。DFAが状態qからAへの遷移をおこす場合、その時点までに処理された文字列はアクセプト

50



され、D F Aによって認識される言語内にある。例示的な一例として、図1の正規表現は、マッチの開始位置を任意（非位置固定）にするステップを含む一連のよく知られたステップを使用して、図2に示される正準D F Aに翻訳されることができる（上で引用されたH o p c r o f tおよびU l l m a nの参考文献を参照されたい）。便宜上、図2のD F Aは、集合{ 0 , 1 , 2 , 3 , 4 , 5 , 6 , 7 , 8 , 9 }を表すために、用語「[ 0 - 9 ]」を使用し、集合{ 0 , 1 , 2 , 3 , 4 , 5 , 6 , 7 , 8 , 9 , \$ , . }内にはい のすべてのシンボルを表すために、シンボル「~」を使用する。

#### 【 0 0 1 1 】

D F Aの構成は、非決定性有限オートマトン（N F A）が構成される中間ステップを一般に伴う。D F Aが各入力シンボルおよび状態についてたかだか1つの遷移を許容する有限状態機械であるのに対して、N F Aは各入力シンボルおよび状態について2つ以上の遷移を許容する有限状態機械である点で、N F AはD F Aとは異なる。また、あらゆる正規言語は、その言語を認識するのに必要とされる状態の数を最小化することによって取得される、正準（c a n o n i c a l）D F Aを有する。本明細書において別途指摘がない限り、すべてのオートマトンは正準（決定性）形式にあることが仮定されるべきである。

#### 【 0 0 1 2 】

しかし、パターンマッチングの目的にとって、図2に示されるD F Aは以下の点で不十分であると、発明者らは本明細書において確信する。

#### 【 0 0 1 3 】

正規表現のアルファベット内にはいシンボルは、D F Aがブロックされる原因となる。パターンマッチングの場合、D F Aがマッチの探索を継続し得るように、そのようなシンボルは、D F Aによって無視されるべきである。この不完全性は、以下のようにD F Aを完成することによって、克服されることができる。

ターゲット文字列内で発生し得るどのようなシンボルも含むように、アルファベットが広げられる。本明細書では、は、256個のシンボルを含むA S C I Iキャラクタセットであることが仮定される。

新しい状態Uを有するように、D F Aが拡大される。 $Q \cup \{U\}$ 。

それまで が定義されていなかった、すべての $q \in Q$ 、 $a \in \Sigma$ に対して、 $(q, a) \rightarrow U$ と定義することによって、遷移関数 が完成される。

#### 【 0 0 1 4 】

マッチは、ターゲット文字列の先頭文字からマッチが始まる場合にのみ見出される。パターンマッチングアプリケーションは、ターゲット内の任意の位置における示されたパターンの発生をすべて見出すことに関心がある。この不完全性は、D F Aがあらゆる位置でリスタートすることを可能にすることによって、克服されることができる。形式的には、あらゆる $q \in Q$ から $q_0$ への 遷移が挿入される。

#### 【 0 0 1 5 】

上述の拡大の結果は、D F Aを取得するための知られた技法を介して、正準D F Aに変換され得るN F Aである。図3は、正準D F Aの例示的な一例を提供している。

#### 【 0 0 1 6 】

D F Aは、その遷移 をテーブルとして実現することによって、一般に解釈的に実施され、各行は、D F Aの状態に対応し、各列は、入力シンボルに対応する。図3のD F Aについての遷移テーブルが、図4に示されている。D F Aのためのアルファベット が、（多くのアプリケーションでしばしばそうであるように）A S C I Iキャラクタセットである場合、図4の遷移テーブルは、256個の列を有する。図4の遷移テーブル内の各エントリは、次状態識別子を含む。したがって、図4の遷移テーブルは、以下のように機能する。D F Aの現在状態がBであり、次の入力シンボルが2である場合、「D」が現在状態Bと入力シンボル2によってインデックス付けされた次状態識別子であるので、遷移テーブルは、状態Dへの遷移を呼び出す。本明細書の説明では、シンボル符号化との混同を避けるため、状態は文字によってラベル付けされる。しかし、実際には、状態は遷移テーブル内で整数インデックスによって一般に表されることに留意されたい。

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【0017】

パイプライン化アーキテクチャでDFAを実施するためのパターンマッチング技法は、本明細書で開示される新規なパターンマッチングアーキテクチャによって大きく改善され得ると、発明者らは本明細書において確信する。

## 【課題を解決するための手段】

## 【0018】

本発明の一態様によれば、すべての状態依存（反復、フィードバック依存）動作をパイプラインの最終ステージに委ねるパイプライン化戦略が開示される。好ましくは、遷移テーブル検索は、DFAによって処理される入力シンボルに対応するすべての遷移テーブルエントリを取り出すように動作する。遷移テーブルメモリからの遷移エントリの取り出しは、DFAの現在状態に基づいていない。代わりに、遷移テーブルメモリからの取り出しは、処理される入力シンボルに対応するデータに基づいて、1組の保存された遷移エントリを取り出すように動作する。

10

## 【0019】

入力データストリームの入力シンボルを同値類（equivalence class）識別子（ECI）にマッピングするためにアルファベット符号化が使用される好ましい実施形態では、これらの遷移テーブルエントリは、ECIに対応するデータによって、1つまたは複数の入力シンボルに間接的にインデックス付けされる。この改良は、単一サイクル状態遷移決定の実行を可能にし、より複雑な圧縮および符号化技法の使用を可能にし、アーキテクチャのスループットおよびスケーラビリティを向上させる。

20

## 【0020】

本発明の別の態様によれば、遷移テーブルの遷移は、好ましくは、遷移を引き起こした入力シンボルの受け取り時に、入力シンボル文字列のパターンとのマッチが発生したかどうかを表すマッチフラグを含む。同様に、遷移テーブルの遷移は、好ましくは、遷移を引き起こした入力シンボルの受け取り時に、マッチングプロセスがリスタートしたかどうかを表すマッチリスタートフラグを含む。アクセプト状態は取り除かれ、遷移のマッチフラグ内にまとめられ得るので、各遷移におけるマッチフラグの存在は、DFA内の状態の数が、従来のDFAに比べて減少されることを可能にする。マッチリスタートフラグの存在は、DFAが、非位置固定パターンとマッチする入力ストリームの部分文字列を識別することを可能にする。併せて、遷移におけるこれらのフラグの存在は、本発明の別の態様に貢献し、その態様では、好ましいDFAは、サイクル当たりに処理されるバイトの数をスケールアップする能力を有するように構成される。状態遷移は、（クロックサイクル当たり単一の入力シンボルのみを処理するように制限されるのではなく）m個の入力シンボルからなる系列によってトリガされることができ、mは1以上である。本明細書で開示されるように、遷移がマッチフラグおよびマッチリスタートフラグを含む含み方のため、DFAは、DFAによってグループとして処理されるm個の入力シンボルからなる系列に属する先行または中間入力シンボルの結果として、マッチが入力ストリーム内のどこでいつ発生したかを依然として検出することができる。

30

40

## 【0021】

本発明のまた別の態様によれば、高スループットDFAを実現するのに必要とされる資源を著しく低減するために、増加スケリング技法、圧縮技法、および文字符号化技法が使用される。例えば、DFA遷移テーブルによって消費されるメモリの量を低減（すなわち圧縮）するために、ランレングス符号化が使用されることができ。さらに、その後、DFAについての次状態を決定するために、ランレングス符号化された遷移に基づいて、状態選択ロジックが動作する。遷移テーブルメモリの検討対象部分から、処理される入力シンボルのECIに対応する遷移を含まないワードを除去するために、マスキングが状態選択ロジックで使用されることができ。

## 【0022】

50

また、本発明のまた別の態様によれば、E C Iを遷移テーブルメモリ内の遷移にマッピングするために、インダイレクションレイヤが使用されることができる。このインダイレクションレイヤは、遷移テーブルメモリ内の遷移エントリについてのランレングス符号化プロセスを最適化するのに、また遷移テーブルメモリに対して必要なアクセスの回数を最小化し得るように、ランレングス符号化された遷移エントリを遷移テーブルメモリのワードに効果的にパッキングするプロセスを最適化するのに有効な、様々な最適化技法の使用を可能にする。インダイレクションを使用して、インダイレクションテーブルメモリ内のインダイレクションエントリは、遷移テーブルメモリ内の遷移エントリへのE C Iのマッピングを、遷移テーブルメモリ内の遷移エントリに対して実行されたいずれの最適化プロセスも考慮に入れて構成するように、データ投入されることができる。

10

**【0023】**

さらに、本発明の別の態様によれば、遷移テーブル内でD F A状態を順序よく配置し、それによって、ランレングス符号化された遷移の効率を高めてD F Aのメモリ要件を改善するための最適化アルゴリズムが、本明細書で開示される。

**【0024】**

さらにまた、複数のメモリワードにまたがる共通の対応する入力シンボル（またはE C Iなどのその派生物）を共有する遷移テーブルエントリの数最小化されるように、遷移テーブルエントリをメモリワード内に効率的にパッキングするための最適化アルゴリズムが、本明細書で開示される。メモリの効率的なパッキングは、1つまたは複数の入力シンボルを処理するときに必要なメモリアクセスの回数を低減することができるので、このメモリパッキングプロセスは、D F Aのスループットを改善するように動作する。

20

**【0025】**

本発明の別の態様によれば、検索中に適用されるパターンは、パイプラインアーキテクチャ自体のロジックを変更することなく、動的に変更されることができる。正規表現コンパイラは、新しい正規表現に対してパターンマッチングパイプラインを再プログラムするために、遷移テーブルメモリ、インダイレクションテーブル、E C Iマッピングテーブル、および関連レジスタにデータ投入しさえすればよい。

**【0026】**

本明細書で提示されるD F A設計に対する改良に基づくことで、本発明の好ましい実施形態によって達成されるスループットおよび密度は、その他の知られたパターンマッチングソリューションを大きく上回ると、発明者らは本明細書において確信する。

30

**【0027】**

本発明の上記およびその他の特徴が、これ以降で説明され、それらは、以下の明細書および図面を検討することで、当業者には明らかとなる。

**【発明を実施するための最良の形態】****【0028】**

図5は、本発明の好ましい実施形態の概要を示している。好ましい実施形態のアーキテクチャは、正規表現回路502内に示されており、正規表現回路502は、複数の連続的な入力シンボルを含む入力データストリームに基づいて動作するパターンマッチング回路として機能する。好ましくは、正規表現回路502は、ハードウェアロジックで実装される（例えば、F P G Aなどの再構成可能なハードウェアロジック、またはA S I Cなどの再構成不可能なハードウェアロジック）。本発明の実施者によって望まれるのであれば、正規表現回路502の1つまたは複数が同じデバイス上に実装され得ることに留意されたく、図24には、そのことも反映されている。また、正規表現回路は、マルチプロセッサシステムなど、その他のパイプライン化アーキテクチャで実施されることもでき、その場合、各プロセッサは、正規表現回路のパイプラインステージとして機能する。そのような例においては、パイプラインに属する異なるプロセッサは、互いにネットワーク接続されることができる。

40

**【0029】**

正規表現回路のデータテーブルおよび関連レジスタは、好ましくは、正規表現コンパイ

50

ラ 5 0 0 の出力によってデータ投入される。正規表現コンパイラ 5 0 0 は、本明細書で説明されるような正規表現回路 5 0 2 によって実現される D F A を生成するために、指定された（好ましくはユーザ指定された）正規表現を処理するように動作する。好ましくは、正規表現コンパイラ 5 0 0 は、パーソナルコンピュータ、ワークステーション、またはサーバの C P U などの汎用プロセッサによって実行されるソフトウェアで実施される。

#### 【 0 0 3 0 】

正規表現コンパイラ 5 0 0 は、ネットワークデータ通信、直接インタフェース、およびシステムバスを含むが、これらに限定されない、任意の適切なデータ通信技法によって、正規表現回路 5 0 2 と通信する。

#### 【 0 0 3 1 】

正規表現回路 5 0 2 は、好ましくは、複数のパイプラインステージを介して、1つまたは複数の指定された正規表現によって定義される D F A を実現する。第 1 のパイプラインステージは、好ましくは、m 個の入力シンボルからなる入力から E C I 出力を生成するアルファベット符号化ステージ 5 0 4 であり、m は、1 以上の整数とすることができる。第 2 のパイプラインステージは、好ましくは、インダイレクションテーブルメモリステージ 5 0 6 である。インダイレクションテーブルメモリ状態 5 0 6 は、様々な方法でアドレッシングされることができる。好ましくは、インダイレクションテーブルメモリステージ 5 0 6 は、アルファベット符号化ステージ 5 0 4 の E C I 出力によって直接的にアドレッシングされる。第 3 のパイプラインステージは、インダイレクションテーブルメモリステージ 5 0 6 の出力からインダイレクションテーブルエントリを受け取り、受け取られたインダイレクションエントリを、遷移テーブルメモリステージ 5 1 0 内の 1 つまたは複数のアドレスに変換するように動作する、遷移テーブルロジックステージ 5 0 8 である。遷移テーブルロジックステージ 5 0 8 は、好ましくは、（マスキング動作に関連して以下で説明されるように）受け取られたインダイレクションテーブルエントリを、状態選択ロジックステージ 5 1 2 によって使用されるデータにも変換する。

#### 【 0 0 3 2 】

遷移テーブルメモリステージ 5 1 0 は、D F A の次状態を決定し、またマッチが見出されたかどうかを決定するために、D F A によって使用される遷移を保存する。状態選択ロジックステージ 5 1 2 は、遷移テーブルメモリステージ 5 1 0 から出力された遷移エントリの 1 つまたは複数を受け取り、また D F A の現在状態と受け取られた遷移とに基づいて D F A の次状態を決定するように動作する。任意選択で、以下で説明される状態選択ロジックステージ 5 1 2 内のマスキング動作 5 1 4 および 5 1 6 は、独立したマスキングパイプラインステージ、または 2 つの独立したマスキングパイプラインステージ（先頭マスキングパイプラインステージおよび終端マスキングパイプラインステージ）に分割されることができる。これらのステージの各々についてのさらなる詳細が、本明細書で提示される。

#### 【 0 0 3 3 】

##### 高スループット D F A

従来の D F A は、一度に 1 つの入力シンボル（バイト）を処理し、次状態を決定するために各バイトに基づいてテーブル検索を実行する。しかし、現代の通信インタフェースおよび相互接続は、サイクル当たり複数のバイトをしばしば転送し、そのことが、より高いスループットの達成に関して、従来の D F A を「ボトルネック」にする。スループットとは、データストリームが処理され得る速度、すなわち、設計およびその実施によって受け入れられ得る 1 秒当たりのバイト数のことである。

#### 【 0 0 3 4 】

従来の D F A の拡張は、m 個のシンボルからなる文字列に基づいた単一の遷移の実行を可能にする D F A である。C l a r k および S c h i m m e l、*「Scalable pattern matching for high speed networks」*、I E E E S y m p o s i u m o n F i e l d - P r o g r a m m a b l e C u

10

20

30

40

50

stom Computing Machines、2004年4月を参照されたく、その開示の全体は、参照により本明細書に組み込まれる。すなわち、DFAは、 $m$ 個の入力シンボルをグループとして、入力ストリームを処理する。形式的には、この適応は、アルファベット<sup>m</sup>に基づいたDFAをもたらし、対応する遷移テーブルは、 $|Q| \times |Q|^m$ のサイズをもつ。この明らかに劇的な資源要件の増大は、本明細書で説明される圧縮技法によって緩和される。便宜上、<sup>m</sup>は、長さ $m$ の系列に基づいて動作する遷移関数を表すとし、 $\delta^m = \delta^1$ である。

【0035】

例示的な一例として、一度に2バイト( $m=2$ )を処理することによって、図3に示されたDFAの実効スループットを2倍にすることについて検討する。図4のテーブルに基づいて、現在状態がEである場合、入力系列「2\$」は、状態Bへの遷移をもたらし、すなわち、 $\delta^2(E, 2\$) = (\delta(E, 2), \$) = B$ である。そのような2文字系列をすべて考慮することによって、以下で説明されるように、完全な遷移テーブルが、このスループットがより高いDFAについて計算されることができる。

【0036】

一般に、与えられたDFAについて<sup>m</sup>を構成するアルゴリズムは単純である。状態の集合は、変化せず、遷移関数(テーブル)は、各状態からの進行を長さ $m$ のあらゆる可能な系列についてシミュレーションすることによって計算される。そのアルゴリズムは、サイズが $(|Q| \times |Q|^m)$ のテーブルを計算するのに、 $(|Q| \times |Q|^m m)$ の時間を要する。より高速なアルゴリズムは、以下の形式の動的プログラミングによって取得されることができる。

【数2】

$$m=2^i, i > 0 \text{ かつ文字列 } x_i x_r, |x_i| = |x_r| = \frac{m}{2}$$

について検討する。その場合、

【数3】

$$\forall q \delta^m(q, x) = \delta^{\frac{m}{2}}(\delta^{\frac{m}{2}}(q, x_i), x_r)$$

となる。上述の説に基づいたアルゴリズムが、図6に示されている。

【0037】

スループットがより高いDFAを取得するため、図6のアルゴリズムは、繰り返して呼び出されることができ、各回のたび、サイクル当たり処理されるシンボルの有効数は2倍になる。これは、<sup>m</sup>を $O(|Q| \times |Q|^m \log m)$ まで計算する複雑さを低減する。さらに、冗長な列を識別することによって、各テーブルの有効アルファベットは、以下の「アルファベット符号化」セクションで説明されるような符号化を用いて、実際には著しく減少されることができる。

【0038】

高スループットDFA：アクセプト

スループットがより高いDFAは、サイクル当たり複数の遷移を実行するので、遷移中に元のDFAのアクセプト状態を通過することがあり得る。したがって、一連の遷移(trace)中にアクセプト状態が通過されたかどうかを含むように、遷移関数を拡大する。

$$\delta^m : Q \times Q^m \rightarrow Q \times \{0, 1\}$$

値域の第2の成分は、遷移を引き起こしたシンボルの系列が、元のDFAにアクセプト状態を通過させる非空のプレフィックスを含むかどうかを示す。

【0039】

「アクセプト」(マッチ)情報は、スループットがより高いDFAの辺に関連付けられるので、この形式の遷移関数は、アクセプト状態Aの集合の必要性を回避する。これは、

以下の「ストライドと符号化の相乗的組み合わせ」セクションで定義される修正 D F A によって形式化される。

#### 【 0 0 4 0 】

好ましい D F A は、どの中間シンボルが実際にアクセプト（マッチ）を発生させたかを追跡していないので、 $m > 1$  の場合、アクセプトは不明確である。速度を優先するため、高スループット D F A は、不明確なアクセプトを許容して、アクセプト点の正確な決定をソフトウェア後処理に委ねるように構成されることができる。

#### 【 0 0 4 1 】

高スループット D F A : リスタート

先に説明されたように、正規表現のためのパターンマッチング D F A は、好ましくは、ターゲット文字列のいたる所でマッチが発生することを可能にする遷移を用いて拡大される。マッチはターゲット内のどこを開始位置としても発生し得るので、アクセプトは、ターゲット内におけるマッチの開始位置を報告すべきである。開始点がいつ設定されるかは、図 3 のオートマトンでは明らかでない。例えば、状態 A へのすべての遷移は、開始点を設定するが、「\$」に基づいた状態 B への遷移もそうである。E から A への遷移を考えると、「~」または「.」はリスタートであるが、数字 ( d i g i t ) はアクセプトである。

#### 【 0 0 4 2 】

マッチの位置独立を達成するために導入される 遷移は、通常の構成を介して D F A に変換され得る N F A をもたらす。以下の「ストライドと符号化の相乗的組み合わせ」セクションは、オートマトンのマッチングプロセスをリスタートするためにのみ役立つ遷移を識別するために、その構成をどのように修正すべきかを説明する。

#### 【 0 0 4 3 】

形式的には、遷移関数は、もう一度、今度はいつリスタートが発生するかを示すために拡大される。

$$^1 : Q \times ^m Q \times \{ 0 , 1 \} \times \{ 0 , 1 \}$$

第 1 のフラグは、リスタート遷移（「マッチリスタート」フラグ）を表し、第 2 のフラグは、アクセプト遷移（「マッチ」フラグ）を表す。したがって、D F A 図は、これ以降、緑色の辺を用いてリスタート遷移を、赤色の辺を用いてアクセプト辺を示す。例えば、図 7 は、一度に 1 つのシンボルを処理し、図 1 によって表される言語を認識するオートマトンについての図の例示的な一例を示している。任意選択で、状態 A への遷移に関する辺は、黒色として符号化され、しかるべくフラグ設定されることができ、状態 B への遷移に関する辺だけが、緑色として符号化され、しかるべくフラグ設定される。色付けされた辺を有する D F A の動作は、以下になる。オートマトンは、マッチの開始および終了を記録するための文脈変数 ( c o n t e x t v a r i a b l e ) b および e を含み、最初、 $b = e = 0$  であり、ターゲットの第 1 のシンボルのインデックスは 1 である。これらの変数は、マッチング文字列の位置が、図 2 4 に示されるような文脈バッファ ( c o n t e x t b u f f e r ) 内で見出されることを可能にする。その後、以下のように遷移が実行される。

黒色：e が m だけ、すなわち、遷移を引き起こした入力文字列の長さである、オートマトンのストライド ( s t r i d e ) だけ進められる。図 6 では、 $m = 1$  である。マッチは進行中であり、ここまで関与しているターゲット文字列の部分は、両端を含んで、位置 b で開始し、位置 e で終了する。

赤色のみ：e が m だけ増やされ、マッチが宣言される。マッチを引き起こしたターゲット部分文字列は、位置 b で開始し、位置 e で終了する。

緑色のみ：b が m だけ増やされ、e は b と同じに設定される。オートマトンは、マッチプロセスをリスタートする。

赤色および緑色：緑色動作の前に、赤色動作が実行される。

#### 【 0 0 4 4 】

図 9 は、図 7 の D F A に関連する遷移テーブルを示しており、遷移テーブルの各遷移エ

10

20

30

40

50

ントリは、次状態識別子に加えて、マッチフラグおよびマッチリスタートフラグを含む。アクセプトについての情報は、辺（関数）に関連付けられるので、色付けされた辺を有するDFAは、正準DFAよりも少ない状態を有することができる。

【0045】

マッチフラグおよびマッチリスタートフラグの使用は、DFAを拡大して、サイクル当たり複数の入力シンボルを処理する場合に特に有用である。図10(b)は、図1の正規表現のための遷移テーブルを示しており、mは2に等しい。図10(a)は、この例についてのシンボル符号化を示している。したがって、サイクル当たり複数の入力シンボルを処理する場合であっても、フラグは状態よりもむしろマッチステータス情報を保つので、DFAは、m個の入力シンボルからなるグループの先行または中間シンボル上でマッチおよびリスタートがいつ発生したかを検出することができる。

10

【0046】

アルファベット符号化

上で説明されたように、遷移テーブル（ ）のサイズは、各サイクルで消費される入力系列の長さに伴って指数関数的に増大する。本セクションでは、シンボルアルファベットを符号化するための技法が提示され、その目標は、遷移テーブルのサイズを低減し、ひいてはサイクル当たりに処理されるシンボルの数を最大化することである。

【0047】

しばしば、与えられた正規表現で使用するシンボルの集合は、検索ターゲットのアルファベットと比べて小さい。ターゲット内に存在するが、パターン内には示されないシンボルは、正規表現に関してDFAにおいて必然的に同じ扱いを与えられる。より一般的には、シンボルのある集合上でのDFAの挙動は、その集合内のすべてのシンボルについて同一となる場合があり得る。例示的な一例として、図1の正規表現は、ASCIIキャラクタセットのわずかな部分のみを使用する。数字についての遷移は、図3に示されたDFAでは同じであり、シンボル「~」は、集合{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, \$, .}内にないASCIIセットのすべてのシンボルを表す。

20

【0048】

正規表現は、「[0-9]」のように、文字クラスを明示的に表現することができるが、より汎用的な手法は、等価の状態-遷移挙動(state-transition behavior)についてDFAを分析することによって達成される。形式的には、

30

$$(a)(b)(q \rightarrow Q)(q, a) = (q, b)$$

である場合、aとbは「遷移同値(transition equivalent)」であると言われることができる。

【0049】

遷移テーブル

【数4】

$$\delta : Q \times \Sigma \mapsto Q$$

を所与として、を同値類に区分するための $O(|\Sigma|^2|Q|)$ のアルゴリズムが、図8に示されている。図1の例を使用すると、アルゴリズムは、図1で暗示されたシンボルの同値類を生成して、図9(b)に示される4つの列を形成する。区分は、整数からなる集合Kと、からKへのマッピングkによって表される。アルファベットシンボルの同値性は状態同値によって決定されるので、最良の結果は、状態同値がすでに決定された正準DFAに、が対応する場合に得られる。

40

【0050】

関連する正規表現を検査する代わりに、DFAを使用して同値類を計算することは、以下の理由から好ましい。

【0051】

正規表現は、同値類を暗示し得る範囲またはその他の表記法(gesture)を使用せずに指定されることができる。例示的な一例として、aとbは、正規表現「(ac)|

50

(bc)」に関して、DFAにおいて同値とされ得るが、それらの同値性は、正規表現では明白ではない。

#### 【0052】

同値類は、正規表現の局所的検査によってはしばしば決定することができない。例示的な一例として、正規表現「 $[0-9]a|5c$ 」は、句「 $[0-9]$ 」を含むが、その範囲の1つの要素(5)は、その他の句のために、その他のシンボルと同じ遷移を共有しない。この場合の適切な区分は、 $\{\{0, 1, 2, 3, 4, 6, 7, 8, 9\}, \{5\}\}$ である。

#### 【0053】

形式的には、図8の関数は、区分、すなわち集合を生成し、その各要素は、DFAにおいて同じに扱われる内のシンボルからなる集合である。関数は、その区分のより便利な形式をkにより生成し、それは、各シンボルをその同値類を表す整数にマッピングする(整数は「同値類識別子」(ECI)として機能する)。その後、シンボル符号化に基づいたDFAは、以下の2つのステージで動作する。第1に、次の入力シンボルが、図8のアルゴリズムによって決定されるその同値類を表すECIに、kによってマッピングされる。第2に、DFAは、その現在状態とECIを使用して、DFAの次状態を決定することができる。

#### 【0054】

これまでに提示されたアイデアに基づいて、図7は、一度に単一のシンボルを処理するDFAの例示的な一例を示しており、赤色の辺は「アクセプト」を表し、緑色の辺は「リスタート」を表す。図8のアルゴリズムを使用する分析は、図9(a)に示されるシンボル符号化と、図9(b)に示される遷移テーブルをもたらす。図9(b)に示されるように、遷移テーブルの各エントリは、DFA状態に対応するデータと、m個の入力シンボルからなるグループについてのECIに対応するデータによってインデックス付けされる。図10(a)および図10(b)は、2つの入力シンボルが一度に処理される、図1の正規表現についての、同等のECIテーブルおよび遷移テーブルを示している。

#### 【0055】

遷移テーブル内の各遷移は、次状態識別子と、マッチリスタートフラグと、マッチフラグとを含む3つ組である。例えば、状態DとECI 0とによってインデックス付けされた遷移は、(B, 1, 0)であり、Bが次状態識別子、1がマッチリスタートフラグ、0がマッチフラグである。したがって、ECI 0によって引き起こされる状態Dから状態Bへの遷移は、ECI 0はマッチを発生させないが、マッチングプロセスをリスタートさせるというように、解釈されることができる。

#### 【0056】

ストライドと符号化の相乗的組み合わせ

上で説明されたスループット向上とアルファベット符号化のアイデアは、今から組み合わせられて、(テーブルを構成する際の)時間および(実行時にテーブルを実現する際の)空間を節約するために、サイクル当たり複数のバイトを消費し、その入力を符号化するアルゴリズムをもたらす。

#### 【0057】

そのような新しい高スループットDFA<sup>m</sup>は、形式的に6つ組( $Q$ , ,  $q_0$ ,  $K$ ,  $k$ , )として今は記述されることができ、ここで、

$Q$ は、状態の有限集合であり、

は、ターゲットの入力シンボルのアルファベットであり、

$q_0$ は、初期状態であり、

$K$ は、サイズが $Q \times |K|^m$ である(しかし実際にはより小さいことが予想される)整数からなる集合であり、

$k$ は、一度にm個の入力シンボルをそれらの符号化にマッピングする、関数

10

20

30

40



【数 5】

$$\Sigma^m \mapsto K$$

であり、

は、現在状態およびm個のシンボルからなる次の部分文字列を、次状態、リスタートの可能性、およびアクセプトの可能性にマッピングする、関数

【数 6】

$$Q \times K \mapsto Q \times \{0, 1\} \times \{0, 1\}$$

である。

10

【0058】

一連の変換は、正規表現 r で開始し、DFA<sup>m</sup> を取得するために以下のステップを実行する。

1. 1つまたは複数の正規表現 r について通常の方法（上で引用されたHopcroftおよびUllmanの参考文献を参照されたい）でDFA<sup>d</sup> が構成される。例えば、上で説明されたように、図1の正規表現は、図3に示されたDFAをもたらし。

2. ターゲット内の任意の位置での開始に基づいて、オートマトンが受け入れることが可能な遷移 p の集合が計算される。これは、各状態について q<sub>0</sub> への遷移をシミュレーションすることによって達成される。具体的には、は以下のように計算される。

【数 7】

20

$$\rho \leftarrow \phi$$

foreach  $\phi \in Q$  do

$$\text{foreach } a \in \Sigma \text{ do } \rho \leftarrow \rho \cup \{(p, a) \mapsto \delta(go, a)\}$$

3. dおよびpから、基本DFA<sup>1</sup> が、図11のアルゴリズムによって構成される。

4. 状態最小化が、標準アルゴリズム（上で引用されたHopcroftおよびUllmanの参考文献を参照されたい）によって、高スループットDFAに対して実行されるが、ただし、状態は、従来のDFAのように最終状態かそれとも非最終状態かによって分割される代わりに、入来する辺の色（黒、緑、赤）によって最初に分割される。

30

5. DFA<sup>k</sup> を所与として、アルファベット符号化を伴う、スループットがより高いDFA<sup>2<sup>k</sup></sup> が、図6に示されたアルゴリズムによって構成される。

【0059】

ランレングス符号化による遷移テーブル圧縮

好ましい高スループットDFAのための遷移テーブルは、 $|K| \times |Q|$  個のエントリを含むことができる。状態最小化は、 $|Q|$  を最小化することを試み、より高いスループットとアルファベット符号化の組み合わせに関する先の説明は、 $|K|$  を最小化することを試みる。それにも関わらず、ストレージ資源は一般に制限され、したがって、できるだけ多くのテーブルを収容するための技法が、検討されるべきである。以下では、テーブル自体をどのように圧縮するかを説明することによって、この問題を検討する。

40

【0060】

上述の説明に基づいて、遷移テーブルセルは、3つ組（次状態、スタートフラグ、アクセプトフラグ）を含む。ランレングス符号化は、十分な冗長性を見せるシンボル系列についてストレージ要件を低減することができる単純な技法である。そのアイデアは、文字列 a<sup>n</sup> を、ランレングス n およびシンボル a として符号化するというものであり、表記 n(a) が使用されることができる。したがって、文字列 a a a a b b b c b b a a a は、4(a) 3(b) 1(c) 2(b) 3(a) としてランレングス符号化される。各シンボルおよび各ランレングスが1バイトの記憶域を必要とする場合、ランレングス符号化は、この例では、ストレージ要件を3バイトだけ低減する（13バイトから10バイト）。

【0061】

50

図 9 ( b ) の例を調べると、遷移テーブルの列にはランレングス符号化のための好条件が豊富に存在する。符号化 E C I シンボル 0 および 3 の場合、テーブルは、直前の状態といずれも同じ 3 つ組を指定しており、そのため、E C I 0 および 3 によってインデックス付けされるテーブル内の遷移に対するランレングス符号化プレフィックスは、ともに「5」である。一般に、遷移テーブルにおいて期待されることは、共通の「次状態」挙動である。したがって、遷移テーブルの各列内にただ 1 つしかないエントリの数は、D F A における状態の数よりも一般に少ない。図 1 2 は、図 9 の遷移テーブルのランレングス符号化バージョンを含んでいる。

#### 【 0 0 6 2 】

列圧縮は記憶域を節約するが、所望のエントリを取得するために遷移テーブルにアクセスする際のコストを増大させるように思える。圧縮に先立ち、行は、現在状態によってインデックス付けされ、列は、E C I によってインデックス付けされる。図 1 2 に示されるように、列がランレングス符号化されると、各列の圧縮された内容は、図 1 3 に示されるように、隣接して保存される。この例では、状態自体も、整数として符号化されている ( A は 0、B は 1 によって表されるなど )。D F A の次状態および動作を決定するために、今は 2 つのステップが存在する。

1 . 入力から符号化された次の E C I に基づいて、図 1 3 に示される記憶域レイアウト内で、エントリの関連範囲が見出される。この検索は、E C I から記憶域レイアウトにおけるオフセットへのマッピングを使用して実行されることができ。エントリの範囲は、図 1 2 からの圧縮された列である。図 1 3 では、下線が引かれたエントリが、図 1 2 から

2 . 現在状態に基づいて、次状態および動作フラグが、エントリの関連範囲内で見出されなければならない。状態選択と呼ばれるこのロジックは、どのエントリが現在状態のインデックスに対応するかを見出すために、エントリを必然的に伸張する必要がある。

#### 【 0 0 6 3 】

圧縮列全体が利用可能である場合、図 1 4 に示される回路は、適切な組エントリを決定するために圧縮形式および現在状態を解釈することによって、状態選択がどのように実現され得るかの一例を示している。ランレングス符号化列内の各エントリについて、そのランレングスプレフィックスと先行するエントリのランレングスプレフィックスの合計が計算される。その後、現在状態インデックスが、各合計と比較され、その合計が現在状態インデックスよりも大きい最初 ( 図 1 4 の例では最も左 ) のエントリが、プライオリティ符号器によって選択される。その後、プライオリティ符号器は、D F A のための次状態を決定することができる。

#### 【 0 0 6 4 】

図 1 4 は、図 1 2 の E C I 1 に対して進行中の状態選択の例示的な一例を示している。各「<」ノードは、現在状態インデックス ( この例では 3、すなわち状態 D ) を、圧縮列内のランレングスプレフィックスの合計と比較する。状態インデックス 0 ( 状態 A ) が現在状態として供給された場合、3 つの比較器すべてが、「1」を出力し、プライオリティ符号器は、最も左のものを取って、内容として ( A , 1 , 0 ) を選択する。図 1 4 では、現在状態は 3 であり、それは、第 2 の合計 ( 3 + 1 ) および第 3 の合計 ( 3 + 1 + 1 ) より小さいので、右 2 つの比較器が、「1」を出力する。プライオリティ符号器は、最も左のものを取り、その結果、E C I 1 および状態 3 の検索として、( C , 0 , 0 ) が選択される。

#### 【 0 0 6 5 】

メモリにおける可変長列のサポート

図 1 3 に示された記憶域レイアウトは、物理メモリにマッピングされなければならない、物理メモリにおいて、テーブル全体は、1 回のメモリアクセスでフェッチされるには大き過ぎる。フィールドプログラマブルゲートアレイ ( F P G A ) および同様のデバイスは、サイズおよびワード長に関して構成可能なメモリバンクをしばしばサポートする。さらに、与えられたエントリのサイズは、各フィールド ( ランレングス、次状態識別子、マッチ

リスタートフラグ、マッチフラグ)のために割り当てられたビットの数に依存する。以下の分析は、サイクル当たり $x$ 個の遷移テーブルエントリが取り出され得るという全般的な仮定に基づいている。シングルポートメモリでは、これは、 $x$ がワード当たりの遷移エントリ数に一致することを意味する。マルチポートメモリでは、これは、 $x$ がポート数にワード当たりの遷移エントリ数を乗じた値に一致することを意味する。一例として、サイクル当たり5回のアクセスをサポートし、ワード当たり3つのエントリを保有する物理メモリは、好ましい実施形態では、 $x = 5 \times 3 = 15$ と設定することによって適合される。しかし、サイクル当たり1回のみのアクセスをサポートし、ワード当たり3つのエントリを保有する物理メモリは、好ましい実施形態では、 $x = 3$ と設定することによって適合される。

10

#### 【0066】

いくつかのアーキテクチャは、 $x$ の選択可能性に関して、その他のアーキテクチャよりも大きな柔軟性を提供する。例えば、FPGAブロックRAMのビットは、時には、ワードの数および各ワードの長さに関して構成可能である。以下の考察は、 $x$ の最良の選択に対して一般に適用される。

#### 【0067】

メモリアクセスは、 $x$ をできるだけ高く押し上げる(d r i v i n g)ことによって、一般に低減される。

#### 【0068】

図14の論理回路は、一度に処理されなければならないエントリの数に伴って成長する。回路の全体的サイズに及ぼすその成長の影響は、ターゲットプラットフォームおよび実装に依存する。図14の論理回路を実現するためには、相当なFPGA資源が必要とされる。

20

#### 【0069】

メモリにおける可変長列のサポート：遷移テーブルメモリ

$x$ が選択されると、圧縮列は、できるだけコンパクトに物理メモリ内に配置される。図15(b)は、遷移テーブルの列が、 $x = 3$ のメモリ内にパッキングされた一例を示している。メモリ内の各ワードは、メモリアドレスによってインデックス付けされる。例えば、メモリアドレス0によってインデックス付けされたワードは、以下の遷移、すなわち、5(B, 1, 0)、3(A, 1, 0)、および1(C, 0, 0)を含む。各列の変化する長さのため、与えられた列は、行内のいずれのエントリからでも開始することがある。

30

#### 【0070】

遷移テーブルにインダイレクションレイヤを導入することによって、ランレングス符号化と、遷移テーブルメモリ(TTM)内へのエントリのコンパクト配置とによって提供される、メモリ効率性を利用することが可能である。図15(a)は、各ECIについて1つのエントリを含む、そのようなインダイレクションテーブルの一例を示している。ECIは連続的に割り当てられ得るので、インダイレクションテーブルは、与えられた入力文字列についてのECI値を使用して、直接的にアドレッシングされることができる。各インダイレクションテーブルエントリは、以下のものから構成されることができる。

ポインタ：ランレングス符号化された遷移テーブル列の第1のエントリを含む、遷移テーブルメモリ内のメモリワードのアドレス。

40

遷移インデックス：ランレングス符号化された遷移テーブル列に対する第1のメモリワードにおける、その列の第1のエントリのインデックス。

遷移カウント：(または略して「カウント」)ランレングス符号化された遷移テーブル列内のエントリの数。

#### 【0071】

入力シンボルのECIを使用してインダイレクションテーブルエントリが取り出されると、TTMから第1のメモリワードを読み取るために、取り出されたエントリ内のポインタが使用される。 $x$ がTTM内のメモリワード当たりのエントリ数であることを思い出されたい。アドレス遷移インデックスから開始して、TTMから $w$ 個の連続するワードを

50

読み取ることによって、列全体がアクセスされ、ここで、 $w$ は、

【数 8】

$$w = \left\lceil \frac{\text{transition.count} + \text{transition.index}}{x} \right\rceil \quad (1)$$

によって与えられる。遷移インデックス値および遷移カウント値は、最初のメモリワードのどのエントリから最後のメモリワードのどのエントリまでが列に関与するかを決定する。図 15 の例では、各 T T M ワードは、3 つのエントリを保存することが可能であり、エントリは、ランレングス符号化された遷移組である。図 15 の場合、網掛け部分から分かるように、E C I 1 の列をフェッチするには、遷移テーブルメモリを 2 回読み取ることが必要であることが分かる。E C I 1 の遷移インデックスおよび遷移カウントの具体的な値は、列が、フェッチされた最初のワードの第 2 のエントリから始まり、フェッチされた最後のワードの第 1 のエントリまで続くことを示している。各列がインデックス 0 から始まるように配置することによって、T T M エントリが浪費される場合、アクセスの回数は、

【数 9】

$$\left\lceil \frac{\text{transition.count}}{x} \right\rceil$$

まで減少させることができる。0 遷移インデックス  $< x$  であるので、T T M におけるコンパクト記憶は、アクセス回数をたかだか 1 だけ増大させるに過ぎない。

【0072】

以下で説明され、また図 25 に示されるように、インダイレクションテーブルおよび T T M へのアクセスは、互いにパイプライン化されることができ、また本発明の設計のその他のコンポーネントともパイプライン化されることができ、マルチポートメモリが利用可能である場合、両テーブルは、性能を悪化させることなく、同じ物理メモリ内に保存されることができる。しかし、インダイレクションレイヤは、T T M におけるランレングス符号化列の内容の配置に応じて、状態遷移当たり様々な回数の遷移テーブルメモリアクセスをもたらす。特定の E C I について、メモリから列を取り出すためのメモリアクセスの回数は、直接アドレッシング手法で必要とされる回数を超えることはできない。平均で、状態遷移当たりのメモリアクセスの回数は、かなり小さい。ランレングス符号化およびインダイレクションによって達成されるメモリ効率性は、インダイレクションテーブルおよび追加のランレングスフィールドを保存することによって背負い込むオーバーヘッドを補償して余りあると、発明者らは一般に確信する。

【0073】

さらに、各エントリは同じサイズであり、エントリ数は入力シンボル同値類の最大数に等しいので、インダイレクションテーブルのためのメモリの割り当ては、相対的に単純である。

【0074】

メモリにおける可変長列のサポート：状態選択

好ましくは T T M の効果的な記憶域レイアウトを考慮し、その他の最適化も提供する、状態選択論理回路の実施が、今から説明される。T T M は圧縮列のコンパクトな記憶を提供するが、状態選択ロジックは、より複雑になる。図 14 に示されるロジックは、圧縮列が、無関係なエントリなしに同時にロジックに提示され得ることを仮定する。そのロジックは、T T M を使用して状態選択を実行するためには、以下の理由で、準最適である。

圧縮列が T T M の複数のワードにまたがることがある。

圧縮列の先頭が T T M ワードの中間から始まることがある。したがって、先頭前のエントリは状態選択のために抑制されなければならない。

圧縮列の末尾が T T M ワードの末尾前に生じることがある。したがって、末尾後のエントリは状態選択のために抑制されなければならない。

## 【 0 0 7 5 】

図 1 4 に示されるロジックは、各エントリの前のすべてのランレングスの合計を累算するために、加算器を使用する。ランレングスは各エントリ内で固定であるので、プレフィックス合計を事前計算し、ランレングスそのものの代わりに、それらを各組の「係数」として保存することによって、加算器は除去されることができる。合計を事前計算することによって、図 1 5 に示されるテーブルは、図 1 6 に示されるテーブルに変換される。

## 【 0 0 7 6 】

圧縮列の先頭および末尾を決定するために使用されるロジックの量も、削減されることができる。各列の先頭は、第 1 のエントリを含む T T M ワードと、第 1 のエントリのそのワード内でのインデックスとを提供する、ポインタフィールドおよび遷移インデックスフィールドを使用して、インダイレクションテーブルにおいて指定される。その後、圧縮列によって占められるワード数  $w$  が、式 ( 1 ) によって与えられる。完全に占められた各ワードは、圧縮列の  $x$  個のエントリを含む。最後のワードでは、圧縮列によって占められる最大インデックスは、

$$( \text{カウント} + \text{インデックス} - 1 ) \bmod x \quad ( 2 )$$

によって与えられる。

## 【 0 0 7 7 】

式 2 を計算するために、状態選択回路内にロジックが配備されることができる。しかし、 $x$  は、設計時パラメータである。ハードウェア定義言語 ( H D L ) コードの適切なパラメータ化によって、式 2 は、インダイレクションテーブルおよび T T M テーブルが生成されるときに、計算されることができる。

## 【 0 0 7 8 】

したがって、計算ロジックの量は、各エントリについて以下の変数をインダイレクションテーブルに保存することによって、削減されることができる。

ポインタ：遷移テーブル列の第 1 のエントリを含む T T M ワードのアドレス。

先頭遷移インデックス：遷移テーブル列がまたがる第 1 の T T M ワードにおける、( 遷移テーブル列の ) 第 1 のエントリのインデックス。

末尾遷移インデックス：遷移テーブル列がまたがる最後の T T M ワードにおける、( 遷移テーブル列の ) 最後のエントリのインデックス。

[ 付加的な ] ワードカウント：  $w - 1$ 、ここで、 $w$  は式 1 によって計算される。

この例を続けると、図 1 7 は、図 1 6 に示される遷移テーブルに関する、インダイレクションテーブルおよび T T M エントリを示している。基本的に、遷移カウント値は、ワードカウント値と、末尾遷移インデックス値とに変換される。この変換は、T T M の内容に影響を与えないが、これらのテーブルを処理するために必要とされるロジックを削減する。

## 【 0 0 7 9 】

図 1 8 に示される状態選択ブロックロジックは、1 つの T T M ワードにまたがる、2 つのエントリを含む遷移テーブル列に基づいて動作する。各ワードは、4 つのエントリを保存する。図 1 8 に示されるワードカウントの出力は、現在の遷移テーブル行について検査されたメモリワードの数を反映する。現在のメモリワードが、その行がまたがる第 1 のワードである場合、先頭遷移インデックスマスク ROM から 1 組のマスクビットを取り出すために、先頭遷移インデックスが使用される。これらのビットは、遷移テーブル列の部分ではない、メモリワード内の先行エントリをマスクオフするために使用される。マスキングは、ランレングス合計を強制的にゼロにすることによって達成される。現在のメモリワードが、その列がまたがる第 1 のワードでない場合、このステージでは、どのエントリもマスクされない。

## 【 0 0 8 0 】

次のステージは、遷移テーブル列の部分ではない、最後のメモリワード内のエントリを「マスク」する。遷移テーブル列の部分ではないエントリについてのランレングス合計は、強制的に最大ランレングスレジスタの値にされる。この値は、遷移テーブル列内のエントリの最大数 ( すなわち、符号化されていない遷移テーブルの列数であり、各符号化遷移

10

20

30

40

50

テーブル列内の最終エントリのランレングス合計の値でもある)を記録する。現在のメモリワードが、遷移テーブル列がまたがる最後のメモリワードである(ワードカウンタの値がワードカウントに等しい)場合、末尾遷移インデックスが、末尾遷移インデックスマスクROMに対するアドレスとして使用される。最後のメモリワードでない場合、このステージ中、どのエントリもマスクされない。後続エントリのランレングス合計を強制的に最大ランレングス合計値にすることは、次状態を選択するマルチプレクサのための選択ビットを生成するプライオリティ符号器を簡略化する。このマスキングプロセスは、小なり(less-than)比較から以下の特性を有する出力ベクトルを生成し、その特性とは、最も左の「1」ビットのインデックスが次状態エントリのインデックスであり、このビットより右のビットはすべて「1」に設定される、というものである。先に言及されたように、スループットを向上させるために、マスキングステージはパイプライン化されることができる。一代替実施形態では、小なり比較、プライオリティ符号器、および次状態選択ロジックのみが、最終パイプラインステージで発生する必要がある。

10

#### 【0081】

##### 最適化

高スループットの正規表現パターンマッチングエンジンを達成することは、本明細書で開示される、高スループットDFA、文字符号化、および遷移テーブル圧縮技法を開発する主要な動機である。以下では、何らかのメモリ効率性を犠牲にしてシステムのスループットを最適化する技法が検討され、したがって、以下の技法の各々は、TTMによって制約される。具体的には、TTMは、以下の制約を課す。

20

ワード当たりのエントリの数。

テーブル内のメモリワードの数。

テーブル内のエントリの総数。

#### 【0082】

本セクションで説明される最適化問題は、ビンパッキング問題(bin packing problem)またはナップザック問題(knapsack problem)の類に含まれる。Cormen他、「Introduction to Algorithms」、Cambridge, MA, The MIT Press、1990年を参照されたく、その開示の全体は、参照により本明細書に組み込まれる。ワード当たりのエントリ数は、パッキング問題のビン(またはナップザック)サイズを定義する。符号化遷移テーブルの構造は、エントリの総数またはワード(ビンまたはナップザック)の総数が遷移テーブルメモリによって課される限界を超えない限り、テーブルを表すのに必要とされるエントリおよび/またはワードの総数を増やすことによってメモリアクセスの回数を最小化するために、変更されることができる。

30

#### 【0083】

##### 最適化：スループットの最適化

検索のために必要とされるメモリアクセスの回数は、TTM内における圧縮列の配置と、それらの列がアクセスされるパターンとによって決定される。パターンは、エンジン内における1組の正規表現と、エンジンによって処理される具体的な入力データとに依存する。DFA<sup>m</sup>の場合、m個の入力シンボルは、1列の検索を引き起こすECIに変換される。メモリアクセスの回数は、符号化遷移テーブル内の列の長さと、列アクセスパターンとに依存する。列アクセスパターンは、エンジン内における正規表現(または1組の正規表現)と、入力データとに依存する。与えられた検索についてのメモリアクセスの総数は、

40

#### 【数10】

$$W = N \sum_{i=1}^{|K|} f_i w_i \quad (3)$$

のように表現されることができ、ここで、 $w_i$  は、遷移テーブルメモリにおいて行  $i$  がま

50

たがるワードの数、 $f_i$  は、行  $i$  がアクセスされる相対頻度、 $N$  は、入力データによって生成される同値類識別子の数である。

【 0 0 8 4 】

入力データについての事前知識は存在しないが、符号化遷移テーブルの構造を変更する能力が存在する。直接アドレス遷移テーブル内の行を並べ替えることによって、符号化遷移テーブル内の列の長さに影響を与えることができる。最適化問題は、メモリアクセス  $W$  の総数を最小化する列配列を選択することである。遷移テーブル列アクセスパターンは一様分布  $f_i = N / |K|$  に従うと仮定する。この場合、

【 数 1 1 】

$$W = \frac{N^2}{|K|} \sum_{i=1}^{|K|} w_i \quad (4)$$

10

となる。

【 0 0 8 5 】

これらの条件において、最適化問題は、量

【 数 1 2 】

$$v = \sum_{i=1}^{|K|} w_i = \sum_{i=1}^{|K|} \left\lceil \frac{\text{count}_i}{x} \right\rceil \quad (5)$$

20

を最小化することである。 $\text{count}_i$  は、遷移カウント  $i$ 、またはランレングス符号化遷移テーブルの行  $i$  内のエントリ数であり、 $x$  は、TTMにおけるワード当たりのエントリ数である。

【 0 0 8 6 】

最適化問題を簡略化するため、 $x = 1$  と仮定することができ、その結果、今最小化される必要がある量は、

【 数 1 3 】

$$v = \sum_{i=1}^{|K|} \text{count}_i \quad (6)$$

30

となる。これは、任意の  $x$  を有する関数を最小化すると同様の結果を一般にもたらす。

【 0 0 8 7 】

図 19 は、本明細書で提示された実行例に対する状態並べ替えの利点を示している。

【 0 0 8 8 】

状態並べ替えには多くの手法が存在する。1つの手法は、行内のエントリのソートされた順序に従って、直接アドレス遷移テーブルの行を整列することによって、符号化遷移テーブルの単一の列の長さを最小化することである。これは、その1つの列に対するランレングス符号化の効率を最大化する。しかし、並べ替えは、その他の列に対するランレングス符号化の効率を低下させることもある。

40

【 0 0 8 9 】

好ましい手法は、欲張りな手法であり、好ましくは、大多数の列についてランの長さを最大化し、それによって、各符号化列の長さを最小化することが望まれる。

【 0 0 9 0 】

2つの状態を所与として、相違しており、そのためランを継続しないECIの数を示す、相違行列 ( $\text{difference matrix}$ ) を生成することによって開始することができる。このアルゴリズムが、図 20 に示されている。

【 0 0 9 1 】

次に、相違行列内のエントリに基づいて、何らかの開始点から、状態を整列する。好ま

50

しくは、次のラベルを取得するのに大多数のランレングスを保つ状態を選択する。選択される開始状態は、好ましくは、相違行列において最大の列合計を有する状態である。その状態を最初に選んだ主意は、その他のすべてから最も異なるのがその状態だからである。その状態を（中間ではなく）末尾まで移動することによって、最長のランを保つことができる。このアルゴリズムの概要が、図 2 1 に示されている。図 2 0 および図 2 1 のアルゴリズムは、一緒になって、「遷移テーブル状態並べ替え」アルゴリズムとして機能する。

【 0 0 9 2 】

最適化：メモリパッキング

インダイレクションレイヤは、符号化遷移テーブルの列が T T M 内の任意の位置で開始および終了することを可能にすることを思い出されたい。符号化テーブル列の物理メモリへの愚直なパッキングは、各テーブル列について余計なメモリアクセスを招くことによって、上述の最適化を妨害することがあり得る。図 1 5 において、入力シンボル「 . 」 ( E C I 1 ) に関連付けられたランレングス符号化遷移テーブル列は、3つのエントリを含むが、T T M において2つのメモリワードにまたがることに留意されたい。列を単一のメモリワードに保存することは可能であるが、この列へのアクセスは、図 1 5 においてレイアウトされているように、2回のメモリアクセスを必要とする。符号化遷移テーブル行がたかだか w 個のワードにまたがるに過ぎないことを保証することによって、インダイレクションレイヤによって提供される柔軟性を利用することができ、ここで、

【数 1 4】

$$\omega \leq \left\lceil \frac{\text{count}}{x} \right\rceil \quad (7)$$

である。図 2 0 は、符号化遷移テーブルを T T M にパッキングするためにこの制約を使用した一例を示している。E C I 1 に関連付けられた列の取り出しは、ただ1回のメモリアクセスを必要とすることが分かる。この例では、メモリ効率性における損失はないが、常にそうであるとは限らない。例えば、メモリワードが3つのエントリを保有し、どの符号化遷移テーブル列も2つのエントリを含む場合を考えられたい。

【 0 0 9 3 】

このメモリパッキング問題は、w が制約関数または目的関数である、古典的な分数ナップザック問題 ( fractional knapsack problem ) の変形である。Black, P. E., 「Dictionary of Algorithms and Data Structures」、NIST、2004 年を参照されたく、その開示の全体は、参照により本明細書に組み込まれる。ここでの好ましい実施形態における相違は、符号化遷移テーブル列の連続記憶域を必要とすることである。これは、古典の問題において対象 ( 符号化遷移テーブル列 ) を複数のナップザック ( メモリワード ) に区分するとき、追加の制約を課す。

【 0 0 9 4 】

この問題に対する1つの解決策は、部分集合和に基づいている。これは、一般的な場合における NP 完全問題であるが ( Garey および Johnson, 「Computer and Intractability: A Guide to the Theory of NP-Completeness」、W. H. Freeman and Co., 1979 年を参照されたく、その開示の全体は、参照により本明細書に組み込まれる )、それが多項式時間で実行されるある条件が存在し、すなわち、それは、和を生成するためにそこから要素が選択される要素の数よりも和がはるかに小さいという条件である。好ましい実施形態の和は、常にメモリワードの幅であり、そのため、好ましいアルゴリズムも、多項式時間で実行される。

【 0 0 9 5 】

基本的なアイデアは、最長のランレングス符号化列を見出し、それを最初に選択することである。その後、可能な限り最良のパッキングを達成することを保証しながら、それをメモリワード内にパッキングする。その後、直前の列で残されたエントリの数を取り、残

10

20

30

40

50



りのランレングス符号化列を用いてそれに部分集合和を適用することができる。これは、追加のメモリアクセスを引き起こすことなく、できるだけ満杯になるようメモリをパッキングする。このプロセスは、残りの符号化列がなくなるまで繰り返される。本明細書で「最適メモリワードパッキング」アルゴリズムと呼ばれる、このアルゴリズムの概要が、図 21 に示されており、図において、R は、ランレングス符号化列の集合であり、w は、メモリワードの幅である。

【0096】

実装アーキテクチャ

このセクションでは、好ましい高スループット D F A およびパイプライン化遷移テーブル技法に基づいた、高性能正規表現検索システムの実施が説明される。この実施の主眼は、高速ディスクのアレイへの高帯域幅相互接続を有する複数のスーパースカラマイクロプロセッサおよび再構成可能なハードウェア装置を含む、ハイブリッド処理プラットフォームである。図 24 は、システムレベルアーキテクチャの一例を示しており、ユーザインタフェースと、正規表現コンパイラと、ファイル I / O コントローラと、正規表現ファームウェアモジュールと、結果プロセッサとを含む。

【0097】

アーキテクチャの目的は、正規表現エンジンのアレイを再構成可能なハードウェア装置（例えば F P G A）に埋め込むことによって、スループットを最大化することである。制御ロジックおよび文脈バッファをサポートするエンジンのアレイは、論理的に単一のファームウェアモジュールと見なされることができる。正規表現回路 / エンジンが F P G A 上に実現される一実施形態では、これらのエンジンは、ハードウェア定義言語（H D L）表現に統合され、知られた技法を使用して F P G A にロードされることができる。

【0098】

各正規表現エンジンの主要なタスクは、高速ディスクからの入力ファイルストリーム内の正規表現を認識することである。正規表現の集合は、好ましくは、ユーザインタフェースを介してユーザによって指定され、高スループット D F A にコンパイルされ、ソフトウェアコンポーネントの集まりによって 1 組のテーブルおよびレジスタ値に変換される。1 組のテーブルおよびレジスタ値は、検索を開始する前に、ファームウェアモジュールに書き込まれる。正規表現エンジンは、パターンを認識した場合、文脈（パターンを含むファイルの部分）と、ファイル内のパターンの開始および終了インデックスと、アクセプト状態ラベルとを含むメッセージを、結果プロセッサに送る。動作環境および統合レベルに応じて、ユーザインタフェースは、簡単なコマンドラインインタフェース、グラフィカルユーザインタフェース（G U I）、または特定言語向け A P I とすることができる。以下のサブセクションは、残りのコンポーネントについての詳細な説明を提供する。

【0099】

実装アーキテクチャ：正規表現コンパイラ

図 5 に詳しく示されるように、正規表現コンパイラは、ユーザによって指定された正規表現の集合をユーザインタフェースから受け取る。その後、当技術分野の標準アルゴリズムが、指定された正規表現を解析して、正規表現から N F A を生成し、N F A を位置独立な D F A に変換し、位置独立な D F A を最小 D F A に縮小する。正規表現コンパイラによって実行される次のステップは、従来の最小 D F A を、本発明の好ましい高スループット D F A に変換することである。このステップは、（遷移テーブルについての適切なマッチフラグおよびマッチリスタートフラグの決定と、入力シンボルの E C I への符号化とを含む）サイクル当たり m 個の入力シンボルからなるストライドに対応するための D F A のスケールアップと、遷移テーブルのランレングス符号化とに関して上で説明されたプロセスを含む。次に、図 20 および図 21 のアルゴリズムが、遷移テーブルメモリへのランレングス符号化遷移の保存を最適化するために実行されることができ、図 23 のアルゴリズムが、遷移テーブルメモリのワードへのランレングス符号化遷移エントリのパッキングを最適化するために実行されることができ、その後、インダイレクションテーブル内のエントリが決定されることができ、コンパイラは、入力シンボル対 E C I テーブル、インダイレ

クションメモリテーブル、および遷移テーブルのための回路メモリにデータ投入するように動作する正規表現回路502に、コマンドを発行する準備が整う。

#### 【0100】

実装アーキテクチャ：結果プロセッサ

結果プロセッサを介して検索結果を報告するために、様々な技法のいずれもが使用され得ることが予想される。好ましい結果プロセッサは、正規表現回路（エンジン）によって生成された結果を使用して、各マッチに対する正しい表現および入力文字列セグメントを決定するように構成されることができる。図24に示されるもののような好ましい実施形態では、正規表現回路（エンジン）によって生成される結果は、一意のエンジン識別子（ID）と、開始状態と、アクセプト状態と、アクセプト遷移をトリガするm個の入力シンボルに対するECIを含む。結果は、パターンにマッチした文字列を含む入力ストリーム的一部分である、マッチについての文脈も含む。結果プロセッサは、マッチング文字列および関連する正規表現（パターン）をユーザインタフェースに報告する。

10

#### 【0101】

実装アーキテクチャ：ファイルI/Oコントローラ

ファイルI/Oコントローラは、入力ストリームを制御する、システムのコンポーネントである。図24の例示的なシステムでは、ファイルI/Oコントローラは、データストアから正規表現回路に流れる、ファイルのストリームを制御する。当技術分野で知られているように、入力ストリームも、ネットワークインタフェース（またはその他のデータインタフェース）によって供給され得ることに留意されたい。

20

#### 【0102】

実装アーキテクチャ：正規表現ファームウェア

正規表現ファームウェアモジュールは、図24に示されるシステムアーキテクチャにおける主要なデータパスコンポーネントである。正規表現ファームウェアモジュールは、好ましくは、正規表現エンジン（またはパターンマッチング回路）のアレイと、少量の制御ロジックとを含む。アレイ内のエンジンの数は、システム内の再構成可能ハードウェア装置の能力に依存する。制御ロジックは、入力ファイルストリームを各正規表現エンジンにブロードキャストし、したがって、エンジンは、同じ入力データに基づいて並列に動作する。制御ロジックはまた、入力のコピーを文脈バッファに送る。文脈バッファのサイズは、エンジンがパターンを認識したときに、結果プロセッサに送られる文脈の量に依存する。このパラメータは、ファームウェア/ソフトウェアインタフェースに過剰負荷をかけない間は、返され得る文脈の量を最大化するように調整されることができる。

30

#### 【0103】

先に言及されたように、正規表現エンジンのスループットは、正規表現エンジンが決定性有限オートマトンについての状態遷移を計算することができる速度によって基本的に制限される。現在状態および入力シンボルに基づいた次状態の決定は、本来的には直列動作である。好ましい実施プラットフォーム上で利用可能な再構成可能ロジック資源を利用するために、並列性を最大化することが望まれる。パイプライン化は、直列計算において並列動作の数を増加させるための一般的な技法であるが、処理パイプラインにフィードバックループがないことが必要とされる。パイプラインの与えられたステージにおける動作の出力は、パイプライン内の後続ステージの結果に依存することができない。図25に示されるように、正規表現エンジンは、一連のパイプライン化された計算ブロックである。ブロック間にフィードバックループがなく、各ブロックは、先行ブロックの結果にのみ基づいて動作することに留意されたい。これは、本発明の好ましい実施形態の著しい特徴である。アルファベット符号化は、1組の入力シンボルにのみ基づいて動作する状態独立な動作である。インダイレクションテーブル検索は、エントリを見つけるために、結果の入力シンボルECIを使用する。遷移テーブル検索は、インダイレクションテーブルエントリ内に含まれるポインタおよびインデックスにのみ依存する。現在状態に依存する唯一の動作は、状態選択ブロック内の最後の計算である。この唯一のフィードバックループをパイプラインの最終ステージに効果的に「押し込む」ことによって、好ましい実施形態は、並

40

50

列性を、ひいては正規表現エンジンのスループットを最大化する。以下のサブセクションは、正規表現エンジン内の各ブロックの設計について説明する。

#### 【0104】

正規表現ファームウェア：アルファベット符号化

アルファベット符号化ブロックは、 $m$ 個の入力シンボルからなる組に同値類識別子（ECI）を割り当てる。各入力シンボルが  $i$  個のビットを使用して指定され、ECI が  $p$  個のビットを使用して指定される場合、アルファベット符号化ブロックは、入力を基本的に  $m \cdot i$  ビットから  $p$  ビットに短縮する。この動作を実行するための単純な方法は、直接アドレステーブルを使用して対組み合わせを実行することである。図26に示されるように、テーブルの第1の組は、1つの  $i$  ビット入力シンボルを  $j$  ビットの ECI に変換する。このステップは、単一の入力シンボルを同値類にマッピングする。テーブルの次のステージは、単一の入力シンボルに対応する  $j$  ビット ECI を2つ単に連結し、ECI テーブルを直接アドレスすることによって、2つの入力シンボルに対して  $k$  ビット ECI を生成する。第2のステージにおいて ECI テーブルを実施するために使用されるメモリのアドレス可能性によって、 $j$  の上限が定められることに留意されたい。具体的には、 $2 \cdot j$  は、メモリによってサポートされるアドレスビットの数より小さいか、または等しくなければならない。同様に、次のステージは、2つのシンボルに対応する  $k$  ビット ECI を2つ連結し、ECI テーブルを直接アドレスすることによって、4つの入力シンボルに対して  $p$  ビット ECI を生成する。同様に、ECI テーブルによってサポートされるアドレス空間が、 $k$  に対して上限を定める。理論的には、この技法は、任意の数の入力シンボルに ECI を割り当てるために使用されることができ、実際には、最終の同値類によってカバーされるシンボルの数が増加するにつれて、メモリ効率は著しく悪化する。

#### 【0105】

図25のパイプラインの後続ステージを実施するために必要とされるロジックは、すでに上で説明された。

#### 【0106】

正規表現ファームウェア：バッファ

各正規表現エンジンは、好ましくは、小規模な入力バッファおよび出力バッファを含む。入力バッファは、単一のエンジンが、TTM内の複数のワードにまたがる遷移テーブルを取り出さなければならない場合に、アレイ内のすべてのエンジンを機能停止にさせることを防止する。いずれかのエンジンの入力バッファが満杯である場合、アレイ全体は機能停止しなければならないが、入力バッファは、ファイル入力速度の瞬間的変動を分離するのに役立つ。出力バッファは、正規表現エンジンがマッチを見出した後、そのマッチが結果プロセッサに送られる以前に、正規表現エンジンが処理を継続することを可能にする。文脈バッファは、好ましくは、ラウンドロビン方式で正規表現エンジンの出力バッファにサービスを行う。いずれかのエンジンの出力バッファが満杯である場合、別の結果を出力バッファに送る以前に、エンジンは機能停止しなければならない。エンジンの入力バッファが満杯である場合、好ましくは、アレイが機能停止しなければならない。

#### 【0107】

本発明がその好ましい実施形態に関連して上で説明されたが、本発明の範囲内に依然として収まる様々な修正が、本発明に施されることができる。本発明に対するそのような修正は、本明細書の教示を検討することで認識可能であろう。例えば、遷移テーブルは、行が状態に対応し、列が ECI に対応するものとして本明細書では説明されたが、本明細書で説明されたどのテーブルの行と列も反対にされ得ることは容易に理解されよう。そのため、本発明の完全な範囲は、添付の特許請求の範囲とその法的均等物によってのみ確定されるものとする。

#### 【図面の簡単な説明】

#### 【0108】

【図1】例示的な正規表現を示す図である。

【図2】図1の正規表現についての DFA を示す図である。

【図 3】図 1 の正規表現についての改良 D F A を示す図である。

【図 4】図 3 の D F A についての従来の遷移テーブルを示す図である。

【図 5】本発明の好ましい一実施形態の概略ブロック図である。

【図 6】D F A のストライドを 2 倍にするための好ましいアルゴリズムを示す図である。

【図 7】状態の数が削減され、マッチおよびマッチリスタート用のフラグを遷移内に有する、本発明の好ましい実施形態による D F A を示す図である。

【図 8】入力シンボルを同値類識別子 ( E C I ) に符号化するための好ましいアルゴリズムを示す図である。

【図 9】図 7 の D F A のための好ましい一遷移テーブルを示す図である。

【図 10】D F A のストライドがサイクル当たり 2 つの入力シンボルに等しい、図 1 の正規表現のための遷移テーブルを示す図である。

【図 11】d および p から基本 D F A <sup>1</sup> を構成するための好ましいアルゴリズムを示す図である。

【図 12】例示的なランレングス符号化遷移テーブルを示す図である。

【図 13】図 12 のランレングス符号化遷移テーブルの隣接保存バージョンを示す図である。

【図 14】与えられた E C I に対応する取り出された 1 組のランレングス符号化遷移に基づいて次状態を決定するための例示的な状態選択論理回路を示す図である。

【図 15】図 12 および図 13 のランレングス符号化遷移テーブルが保存される ( a ) インダイレクションテーブルおよび ( b ) メモリを示す図である。

【図 16】図 15 のランレングス符号化遷移テーブルが事前計算ランレングスプレフィックス合計を含む ( a ) インダイレクションテーブルおよび ( b ) メモリを示す図である。

【図 17】図 14 のランレングス符号化遷移テーブルが事前計算ランレングスプレフィックス合計を含む ( a ) インダイレクションテーブルおよび ( b ) T T M の代替定式化を示す図である。

【図 18】例示的な状態選択論理回路を示す図である。

【図 19】( a ) 状態並べ替えによって最適化された例示的な遷移テーブルおよび ( b ) 状態並べ替え遷移テーブルのランレングス符号化バージョンを示す図である。

【図 20】T T M のためのランレングス符号化を最適化するのに使用される相違行列を計算するための好ましいアルゴリズムを示す図である。

【図 21】T T M におけるランレングス符号化の最適化のための好ましいアルゴリズムを示す図である。

【図 22】T T M にパッキングされた符号化遷移テーブルの一例を示す図である。

【図 23】メモリワードの最適パッキングのための好ましいアルゴリズムを示す図である。

【図 24】好ましい正規表現システムアーキテクチャを示す図である。

【図 25】一連のパイプライン化された計算ブロックとしての好ましい実施形態の正規表現エンジンを示す図である。

【図 26】直接アドレステーブルおよび対組み合わせを使用して E C I を入力シンボルに割り当てるための好ましいプロセスを示す図である。

10

20

30

40

【図 1】

$$\backslash \$[0-9] + (\backslash \cdot [0-9]\{2\})\{0,1\}$$

Figure 1

【図 2】

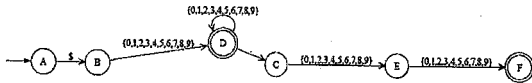


Figure 2

【図 3】

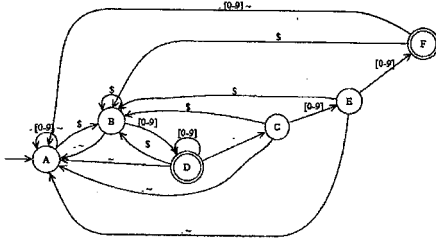


Figure 3

【図 4】

状態	シンボル	0...9	.	\$	...
A	A...A	A	B	A...	
B	D...D	A	B	A...	
C	E...E	A	B	A...	
D	D...D	C	B	A...	
E	F...F	A	B	A...	
F	A...A	A	B	A...	

Figure 4

【図 7】

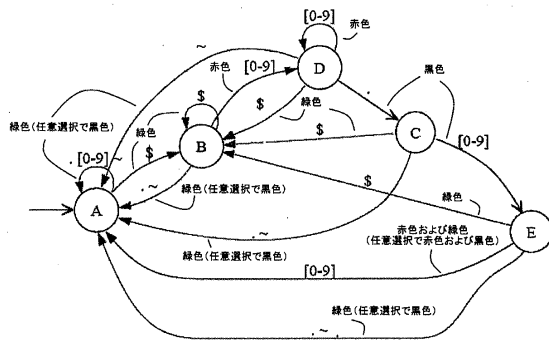


Figure 7

【図 8】

```

function PARTITION( $Q, \Sigma, \delta: Q \times \Sigma \mapsto Q$ ): ( $result, K, \kappa: \Sigma \mapsto K$ )
  WorkList  $\leftarrow \Sigma$ 
  result  $\leftarrow \emptyset$ 
  i  $\leftarrow -1$ 
  while WorkList  $\neq \emptyset$  do
    w  $\leftarrow$  WorkListに属する何らかの要素
    equiv  $\leftarrow \{a \in \Sigma \mid \forall q \in Q \delta(q, a) = \delta(q, w)\}$ 
    i  $\leftarrow i + 1$ 
    K  $\leftarrow K \cup \{i\}$ 
    foreach sym  $\in$  equiv do  $\kappa(sym) \leftarrow i$ 
    result  $\leftarrow result \cup equiv$ 
    WorkList  $\leftarrow WorkList - equiv$ 
  end

```

Figure 8

【図 5】

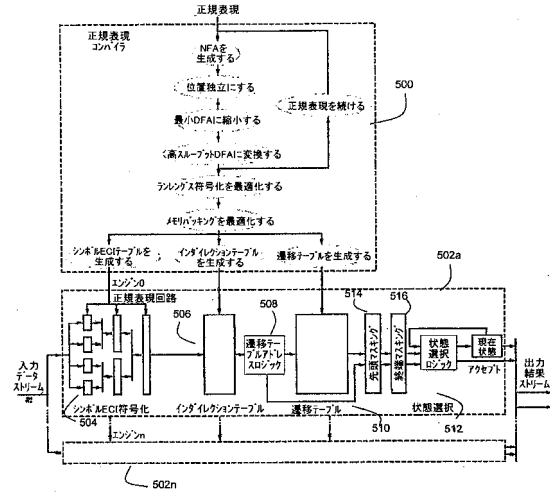


Figure 5

【図 6】

```

function DOUBLESTRIDE( $\delta^k$ ):  $\delta^{2k}$ 
  foreach  $q \in Q$  do
    foreach  $a \in \Sigma^k$  do
      foreach  $b \in \Sigma^k$  do
         $\delta^{2k}(q, ab) \leftarrow \delta^k(\delta^k(q, a), b)$ 
      end
    end
  end

```

Figure 6

【図 9】

シンボル	ECI
\$	0
.	1
0-9	2
その他すべて	3

(a)

状態	0	1	2	3
A	(B,1,0)	(A,1,0)	(A,1,0)	(A,1,0)
B	(B,1,0)	(A,1,0)	(D,0,1)	(A,1,0)
C	(B,1,0)	(A,1,0)	(E,0,0)	(A,1,0)
D	(B,1,0)	(C,0,0)	(D,0,1)	(A,1,0)
E	(B,1,0)	(A,1,0)	(A,0,1)	(A,1,0)

(b)

Figure 9

【図 10】

シンボル	ECI
0	0
0	1
0	2
0	3
1	0
1	1
1	2
1	3
2	0
2	1
2	2
2	3
3	0
3	1
3	2
3	3

(a)

状態	0	1	2	3	4	5	6	7
A	(B,1,0)	(A,1,0)	(D,1,1)	(A,1,0)	(B,1,0)	(A,1,0)	(A,1,0)	(A,1,0)
B	(B,1,0)	(A,1,0)	(D,1,1)	(A,1,0)	(B,1,0)	(C,0,1)	(D,0,1)	(A,1,1)
C	(B,1,0)	(A,1,0)	(D,1,1)	(A,1,0)	(B,1,0)	(A,1,0)	(A,0,1)	(A,1,0)
D	(B,1,0)	(A,1,0)	(D,1,1)	(E,0,0)	(B,1,0)	(C,0,1)	(D,0,1)	(A,1,1)
E	(B,1,0)	(A,1,0)	(D,1,1)	(A,1,0)	(B,1,0)	(A,1,1)	(A,1,1)	(A,1,1)

(b)

Figure 10

## 【図 11】

```

function BASEDFA( $Q, \Sigma, q_0, \delta, A, \rho$ ):  $DFA^1 = (Q', \Sigma', q'_0, K, \kappa, \delta')$ 
 $K \leftarrow \{0\}$ 
foreach  $a \in \Sigma$  do  $\kappa(a) \leftarrow 0$ 
 $WorkList \leftarrow \{q_0\}$ 
while  $WorkList \neq \emptyset$  do
 $w \leftarrow WorkList$ に属する何らかの要素
 $WorkList \leftarrow WorkList - \{w\}$ 
/*  $w$  既に  $\lambda$  閉方された状態の集合
    $q \in w$  構造による */
 $Q' \leftarrow Q' \cup \{w\}$ 
foreach  $a \in \Sigma$  do
 $targs \leftarrow \{t \mid (\exists s \in w) \delta(s, a) = t\}$ 
 $targs \leftarrow LAMBDA CLOSURE(targs)$ 
 $red \leftarrow targs \cap A \neq \emptyset$ 
 $green \leftarrow (\forall s \in w - \{q_0\}) (\forall t \in Q - \{q_0\}) \delta(s, a) \neq t$ 
 $\delta' \leftarrow \delta' \cup \{w \times a \mapsto (targs, green, red)\}$ 
if  $targs \notin Q'$ 
then  $WorkList \leftarrow WorkList \cup \{targs\}$ 
end

```

Figure 11

## 【図 12】

状態	0	1	2	3
A	5(B,1,0)	3(A,1,0)	1(A,1,0)	5(A,1,0)
B			1(D,0,1)	
C			1(E,0,0)	
D		1(C,0,0)	1(D,0,1)	
E		1(A,1,0)	1(A,0,1)	

Figure 12

## 【図 13】

51110 30101 12000 10110 10101 13011 14010 13011 10101 30110 10101

Figure 13

## 【図 16】

ECI	ポインタ	インデックス	カウント	アドレス
0	0	0	1	0
1	0	1	3	1
2	1	1	5	2
3	3	0	1	3

インダイレクションテーブル

アドレス	0	1	2	3
0	5(B,1,0)	3(A,1,0)	4(C,0,0)	
1	5(A,1,0)	1(A,1,0)	2(D,0,1)	
2	3(E,0,0)	4(D,0,1)	5(A,0,1)	
3	5(A,1,0)			

遷移テーブルメモリ

Figure 16

## 【図 17】

シンボル	ECI	先頭 ポインタ	末尾 インデックス	ワード インデックス	カウント	アドレス
S	0	0	0	0	0	0
.	1	0	1	0	1	1
[0-9]	2	1	1	2	1	2
-	3	3	0	0	0	3

インダイレクションテーブル

アドレス	0	1	2	3
0	5(B,1,0)	3(A,1,0)	4(C,0,0)	
1	5(A,1,0)	1(A,1,0)	2(D,0,1)	
2	3(E,0,0)	4(D,0,1)	5(A,0,1)	
3	5(A,1,0)			

遷移テーブルメモリ

Figure 17

## 【図 14】

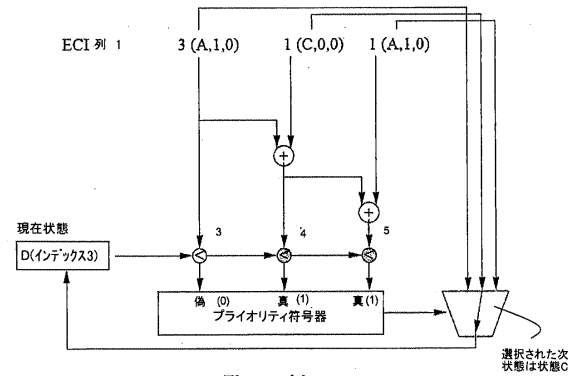


Figure 14

## 【図 15】

ECI	ポインタ	インデックス	カウント	アドレス
0	0	0	1	0
1	0	1	3	1
2	1	1	5	2
3	3	0	1	3

インダイレクションテーブル

アドレス	0	1	2	3
0	5(B,1,0)	3(A,1,0)	1(C,0,0)	
1	1(A,1,0)	1(A,1,0)	1(D,0,1)	
2	1(E,0,0)	1(D,0,1)	1(A,0,1)	
3	5(A,1,0)			

遷移テーブルメモリ

(a)

(b)

Figure 15

## 【図 18】

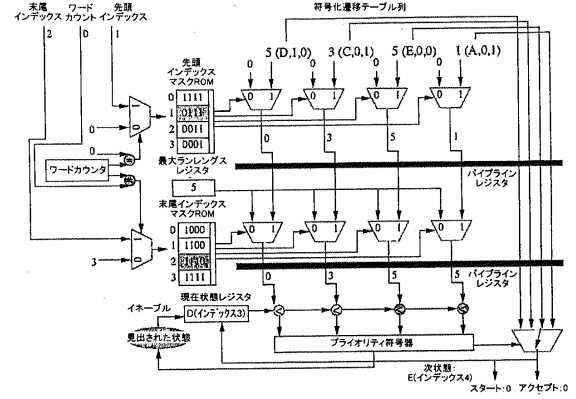


Figure 18

## 【図 19】

状態	0	1	2	3
A	(B,1,0)	(A,1,0)	(A,1,0)	(A,1,0)
C	(B,1,0)	(A,1,0)	(E,0,0)	(A,1,0)
E	(B,1,0)	(A,1,0)	(A,0,1)	(A,1,0)
B	(B,1,0)	(A,1,0)	(D,0,1)	(A,1,0)
D	(B,1,0)	(C,0,0)	(D,0,1)	(A,1,0)

(a)

状態	0	1	2	3
A	(5,B,1,0)	(4,A,1,0)	(1,A,1,0)	(5,A,1,0)
C			(1,E,0,0)	
B			(1,A,0,1)	
D		(1,C,0,0)	(2,D,0,1)	

(b)

Figure 19

## 【図 20】

```

function DIFFMATRIX( $Q, \Sigma, \delta$ ):  $D$ 
foreach  $q \in Q$  do
  foreach  $p \in Q$  do
    foreach  $s \in \Sigma$  do  $D[q][p] \leftarrow D[q][p] + \delta(q, s) \neq \delta(p, s)$ 
end

```

Figure 20

【図 2 1】

```

function RLOPTIM( $Q, q_0, \Sigma, \delta, D$ ):  $DFA' = (Q', \Sigma, q_0', \delta')$ 
   $Used \leftarrow \emptyset$ 
   $cur \leftarrow \text{MAX}(\sum_{i \in Q} D[i])$ 
  while  $Used \neq Q$  do
     $next \leftarrow \text{MAX}(D[cur] - Used)$ 
     $Used \leftarrow Used \cup \{next\}$ 
     $cur \leftarrow next$ 
  end

```

Figure 21

【図 2 2】

シンボル	ECI	先頭 ポインタ	末尾 インデックス	ワード インデックス	カウント	アドレス
\$	0	0	0	0	0	5 (B,1,0)
.	1	0	1	2	0	1 (A,1,0)
[0-9]	2	1	0	0	1	2 (D,0,1)
~	3	2	1	1	0	5 (A,1,0)

インダイレクションテーブル

アドレス	データ
0	5 (B,1,0)
1	1 (A,1,0)
2	2 (D,0,1)
3	5 (A,1,0)

遷移テーブルメモリ

Figure 22

【図 2 3】

```

procedure COLPACK( $R, x$ )
  while  $R \neq \emptyset$  do
     $longest \leftarrow \text{MAXLENGTH}(R)$ 
     $R \leftarrow R - \{longest\}$ 
    PACK( $longest$ )
     $sum \leftarrow x - (longest \bmod x)$ 
     $L_0 \leftarrow \{0\}$ 
    foreach  $col \in R$  do  $L_i \leftarrow \text{MERGE}(L_{i-1}, L_{i-1} + \text{LENGTH}(col))$ 
    sumより大きいすべての要素を $L_i$ から取り除く
     $S \leftarrow \text{Max}(L_n)$ 
    foreach  $col \in S$  do
       $R \leftarrow R - \{col\}$ 
      call PACK( $col$ )
  end

```

Figure 23

【図 2 6】

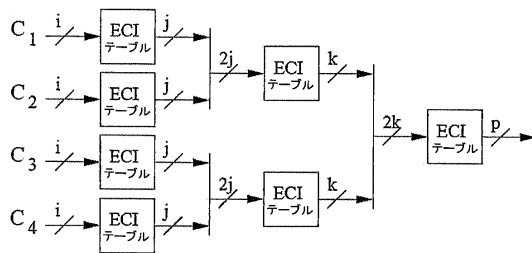


Figure 26

【図 2 4】

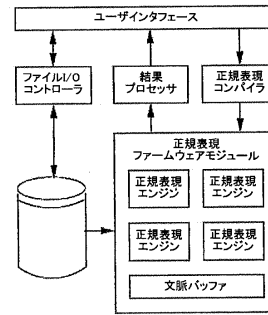


Figure 24

【図 2 5】

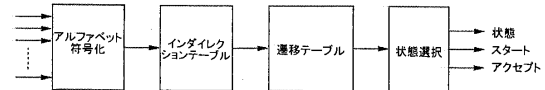


Figure 25

---

フロントページの続き

- (72)発明者 サイトロン, ロン・ケイ  
アメリカ合衆国、ミズーリ・６３１２４、セント・ルイス、グラナダ・ウェイ・４８
- (72)発明者 テイラー, デイビッド・エドワード  
アメリカ合衆国、ミズーリ・６３１１８、セント・ルイス、ミズーリ・アベニュー・３４４８
- (72)発明者 プロデイ, ベンジャミン・カリー  
アメリカ合衆国、ミズーリ・６３１３０、ユニバーシティ・シティ、リーランド・アベニュー・１  
・ダブリュ・７２６

審査官 吉田 誠

- (56)参考文献 特開２００５－２４２９９７（ＪＰ，Ａ）  
特開平１０－２０７９１２（ＪＰ，Ａ）  
特開２００５－２４２６７２（ＪＰ，Ａ）

- (58)調査した分野(Int.Cl.，ＤＢ名)  
Ｇ０６Ｆ １７／３０