

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2009-211480

(P2009-211480A)

(43) 公開日 平成21年9月17日(2009.9.17)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/21 (2006.01)	G06F 17/21 530T	5B075
G06F 17/30 (2006.01)	G06F 17/30 140	5B109
	G06F 17/30 350C	
	G06F 17/21 530E	

審査請求 未請求 請求項の数 39 O L (全 16 頁)

(21) 出願番号 特願2008-54648 (P2008-54648)
 (22) 出願日 平成20年3月5日(2008.3.5)

(71) 出願人 000004237
 日本電気株式会社
 東京都港区芝五丁目7番1号
 (74) 代理人 100079005
 弁理士 宇高 克己
 (72) 発明者 森口 昌和
 東京都港区芝五丁目7番1号 日本電気株式会社内
 (72) 発明者 辰巳 勇臣
 東京都港区芝五丁目7番1号 日本電気株式会社内
 Fターム(参考) 5B075 NR02 PR06 QM08
 5B109 NH02 NH20

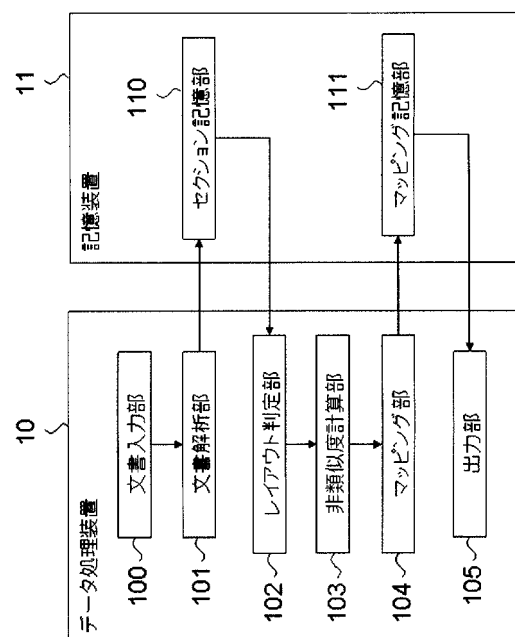
(54) 【発明の名称】 構造化文書処理システム、構造化文書処理方法及び構造化文書処理プログラム

(57) 【要約】

【課題】本発明は、構造化文書をレンダリングせずにセクションのマッピングを行える構造化文書処理技術を提供することにある。また、構造化文書のレイアウト及び構造の両方を考慮したセクションのマッピングを行える構造化文書処理技術を提供することにある。

【解決手段】本発明は、構造化文書を構成している各タグに、レイアウトへの影響度に基づいて重みを割り当て、この割り当てた重み及び各構造化文書内のタグ構造に基づいて、構造化文書間の非類似度を計算する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

複数の構造化文書同士の非類似度を計算する構造化文書処理装置であって、
構造化文書を構成している各タグに、レイアウトへの影響度に基づいて重みを割り当てるレイアウト判定手段と、

前記割り当てた重み及び前記各構造化文書内のタグ構造に基づいて、構造化文書間の非類似度を計算する非類似度計算手段と

を有することを特徴とする構造化文書処理装置。

【請求項 2】

前記計算された非類似度が小さい構造化文書同士を対応付けるマッピング手段を有することを特徴とする請求項 1 記載の構造化文書処理装置。

10

【請求項 3】

マッピング手段は、複数の構造化文書を、非類似度が小さい構造化文書同士から順に対応付けていくことを特徴とする請求項 2 に記載の構造化文書処理装置。

【請求項 4】

マッピング手段は、初めに対応付けた構造化文書を基準にして、各構造化文書を複数のグループに分割し、各グループ内の構造化文書同士を対応付けることを特徴とする請求項 2 又は請求項 3 に記載の構造化文書処理装置。

【請求項 5】

前記非類似度計算手段は、前記割り当てた重み及び構造化文書を細分化したセクション内のタグ構造に基づいて、各構造化文書のセクション間の非類似度を計算することを特徴とする請求項 1 から請求項 4 のいずれかに記載の構造化文書処理装置。

20

【請求項 6】

入力した構造化文書を解析してセクションを抽出する文書解析手段を有することを特徴とする請求項 5 に記載の構造化文書処理装置。

【請求項 7】

レイアウト判定手段は、重みの算出基準として、構造化文書に予め定められているレイアウトの定義を用いることを特徴とする請求項 1 から請求項 6 のいずれかに記載の構造化文書処理装置。

【請求項 8】

レイアウト判定手段は、重みの算出基準として、ユーザが予め定めた定義を用いることを特徴とする請求項 1 から請求項 6 のいずれかに記載の構造化文書処理装置。

30

【請求項 9】

レイアウト判定手段は、重みの算出基準として、予めレンダリングした際に各タグのレイアウトへの影響度を測定し、それに基づいてシステムが定めた定義を用いることを特徴とする請求項 1 から請求項 6 のいずれかに記載の構造化文書処理装置。

【請求項 10】

非類似度計算手段は、非類似度の計算に、タグ構造をツリー型に変換した構造体を用いることを特徴とする請求項 1 から請求項 9 のいずれかに記載の構造化文書処理装置。

【請求項 11】

非類似度計算手段は、非類似度の計算に、タグ構造を単一バイトの文字列に変換した構造体を用いることを特徴とする請求項 1 から請求項 9 のいずれかに記載の構造化文書処理装置。

40

【請求項 12】

非類似度計算手段は、非類似度の計算に、構造体の編集距離を用いる請求項 10 又は請求項 11 に記載の構造化文書処理装置。

【請求項 13】

構造化文書のレイアウトと構造に基づいて構造化文書に対応付ける構造化文書処理システムであって、

前記構造化文書を構成する各タグにレイアウトへの影響度に基づいて重みを割り当てる

50

レイアウト判定手段と、

前記割り当てられた重み及び前記構造化文書のタグ構造に基づいて、構造化文書を細分化したセクション間の非類似度を計算する非類似度計算手段と、

前記計算された非類似度に基づいて構造化文書同士を対応付け、この構造化文書同士の対応付けを示す表示情報を生成するマッピング手段と

前記生成した表示情報を、通信ネットワークを介して情報端末に送信する情報配信手段と

を有することを特徴とする構造化文書処理システム。

【請求項 14】

複数の構造化文書同士の非類似度を計算する構造化文書処理方法であって、

10

構造化文書を構成している各タグに、レイアウトへの影響度に基づいて重みを割り当てるレイアウト判定ステップと、

前記割り当てた重み及び前記各構造化文書内のタグ構造に基づいて、構造化文書間の非類似度を計算する非類似度計算ステップと

を有することを特徴とする構造化文書処理方法。

【請求項 15】

前記計算された非類似度が小さい構造化文書同士を対応付けるマッピングステップを有することを特徴とする請求項 14 に記載の構造化文書処理方法。

【請求項 16】

マッピングステップは、複数の構造化文書を、非類似度が小さい構造化文書同士から順に対応付けていくことを特徴とする請求項 15 に記載の構造化文書処理方法。

20

【請求項 17】

マッピングステップは、初めに対応付けた構造化文書を基準にして、各構造化文書を複数のグループに分割し、各グループ内の構造化文書同士を対応付けることを特徴とする請求項 15 又は請求項 16 に記載の構造化文書処理方法。

【請求項 18】

前記非類似度計算ステップは、前記割り当てた重み及び構造化文書を細分化したセクション内のタグ構造に基づいて、各構造化文書のセクション間の非類似度を計算することを特徴とする請求項 14 から請求項 17 のいずれかに記載の構造化文書処理方法。

【請求項 19】

30

入力した構造化文書を解析してセクションを抽出する文書解析ステップを有することを特徴とする請求項 18 に記載の構造化文書処理方法。

【請求項 20】

レイアウト判定ステップは、重みの算出基準として、構造化文書に予め定められているレイアウトの定義を用いることを特徴とする請求項 14 から請求項 19 のいずれかに記載の構造化文書処理方法。

【請求項 21】

レイアウト判定ステップは、重みの算出基準として、ユーザが予め定めた定義を用いることを特徴とする請求項 14 から請求項 19 のいずれかに記載の構造化文書処理方法。

【請求項 22】

40

レイアウト判定ステップは、重みの算出基準として、予めレンダリングした際に各タグのレイアウトへの影響度を測定し、それに基づいてシステムが定めた定義を用いることを特徴とする請求項 14 から請求項 19 のいずれかに記載の構造化文書処理方法。

【請求項 23】

非類似度計算ステップは、非類似度の計算に、タグ構造をツリー型に変換した構造体を用いることを特徴とする請求項 14 から請求項 22 のいずれかに記載の構造化文書処理方法。

【請求項 24】

非類似度計算ステップは、非類似度の計算に、タグ構造を単一バイトの文字列に変換した構造体を用いることを特徴とする請求項 14 から請求項 22 のいずれかに記載の構造化文

50

書処理方法。

【請求項 25】

非類似度計算ステップは、非類似度の計算に、構造体の編集距離を用いる請求項 23 又は請求項 24 に記載の構造化文書処理方法。

【請求項 26】

マッピングステップにおいて構造化文書同士を対応付けた結果を、通信ネットワークを介して送信する情報送信ステップを有することを特徴とする請求項 15 から請求項 25 のいずれかに記載の構造化文書処理方法。

【請求項 27】

複数の構造化文書同士の非類似度を計算するプログラムであって、前記プログラムは、
情報処理装置に、

構造化文書を構成している各タグに、レイアウトへの影響度に基づいて重みを割り当てるレイアウト判定処理と、

前記割り当てた重み及び前記各構造化文書内のタグ構造に基づいて、構造化文書間の非類似度を計算する非類似度計算処理と
を実行させることを特徴とするプログラム。

【請求項 28】

前記計算された非類似度が小さい構造化文書同士を対応付けるマッピング処理を実行させることを特徴とする請求項 27 に記載のプログラム。

【請求項 29】

マッピング処理は、複数の構造化文書を、非類似度が小さい構造化文書同士から順に対応付けていくことを特徴とする請求項 28 に記載のプログラム。

【請求項 30】

マッピング処理は、初めに対応付けた構造化文書を基準にして、各構造化文書を複数のグループに分割し、各グループ内の構造化文書同士を対応付けることを特徴とする請求項 28 又は請求項 29 に記載のプログラム。

【請求項 31】

前記非類似度計算処理は、前記割り当てた重み及び構造化文書を細分化したセクション内のタグ構造に基づいて、各構造化文書のセクション間の非類似度を計算することを特徴とする請求項 27 から請求項 30 のいずれかに記載のプログラム。

【請求項 32】

入力した構造化文書を解析してセクションを抽出する文書解析処理を実行させることを特徴とする請求項 31 に記載のプログラム。

【請求項 33】

レイアウト判定処理は、重みの算出基準として、構造化文書に予め定められているレイアウトの定義を用いることを特徴とする請求項 27 から請求項 32 のいずれかに記載のプログラム。

【請求項 34】

レイアウト判定処理は、重みの算出基準として、ユーザが予め定めた定義を用いることを特徴とする請求項 27 から請求項 32 のいずれかに記載のプログラム。

【請求項 35】

レイアウト判定処理は、重みの算出基準として、予めレンダリングした際に各タグのレイアウトへの影響度を測定し、それに基づいてシステムが定めた定義を用いることを特徴とする請求項 27 から請求項 32 のいずれかに記載のプログラム。

【請求項 36】

非類似度計算処理は、非類似度の計算に、タグ構造をツリー型に変換した構造体を用いることを特徴とする請求項 27 から請求項 35 のいずれかに記載のプログラム。

【請求項 37】

非類似度計算処理は、非類似度の計算に、タグ構造を単一バイトの文字列に変換した構造体を用いることを特徴とする請求項 27 から請求項 35 のいずれかに記載のプログラム。

10

20

30

40

50

【請求項 38】

非類似度計算処理は、非類似度の計算に、構造体の編集距離を用いる請求項 36 又は請求項 37 に記載のプログラム。

【請求項 39】

マッピング処理における対応付けの結果を、通信ネットワークを介して送信する情報送信処理を実行させることを特徴とする請求項 28 から請求項 38 のいずれかに記載のプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、構造化文書処理する構造化文書処理システム、構造化文書処理方法及び構造化文書処理プログラムに関し、特に HTML (Hyper Text Markup Language) のような特定の文書型定義 (DTD: Document Type Definition) に基づいてレイアウトを形成する構造化文書において、特定の条件で分割された領域を、レイアウトと構造の両方が類似した別の領域とマッピングする構造化文書処理システム、構造化文書処理方法及び構造化文書処理プログラムに関する。

【背景技術】

【0002】

近年、大量の情報を含む Web コンテンツなどの構造化文書から、自動的に文書構造を解析して複数のセクションを抽出して利用するシステムが研究されている。例えば、セクションの中からユーザが必要な情報を、自動あるいは手動で選択して利用する情報提供システムがある。なお、セクションとは、構造化文書を細分化した領域のことで、構造化文書と同様にマークアップ言語 (以下、タグ) で構成される。

【0003】

しかし、その情報提供システムなどのように、常に内容が変化していく Web コンテンツのセクションを利用したアプリケーションでは、時間経過に応じてそのセクションの表示位置が変わったり、消滅したりするため、ユーザに間違ったセクションの情報を提供してしまう場合がある。

【0004】

そこで、特許文献 1 では、構造化文書のレイアウトの特徴に注目し、セクションの抽出順序やセクションの表示座標、および見出しなどのレイアウト情報を用いることで、セクションのレイアウトが変わっても、セクションの位置を推定する方法が提案されている。

【0005】

また、特許文献 2 で提案されている構造化文書同士の類似度を検出する技術をセクションに応用し、変化前と変化後とのセクションの文書構造を比較して、より類似した構造を持つセクション同士をマッピングさせる方法が考えられる。特に、構造化文書の文書構造がツリー型で表現できるため、そのツリーの編集距離を用いて構造化文書の類似性を判定する。

【特許文献 1】特開 2007 293543 号

【特許文献 2】特開 2007 - 052556 号

【発明の開示】

【発明が解決しようとする課題】

【0006】

第 1 の問題点は、特許文献 1 のように座標などのレイアウトの特徴を用いて、より類似したレイアウトのセクション同士をマッピングする場合、構造化文書を少なくとも 1 度はレンダリングして、レイアウト情報を取得しなければならないということである。

【0007】

第 2 の問題点は、レンダリングせずに、特許文献 2 のように文書構造の類似性を判定して、より類似した構造のセクション同士をマッピングする場合、レイアウトを考慮していないために適切なマッピングができないということである。

10

20

30

40

50

【 0 0 0 8 】

本発明が解決しようとする課題は、構造化文書をレンダリングせずにセクションのマッピングを行える構造化文書処理システム、構造化文書処理方法および構造化文書処理プログラムを提供することにある。

【 0 0 0 9 】

また、構造化文書のレイアウト及び構造の両方を考慮したセクションのマッピングを行える構造化文書処理システム、構造化文書処理方法および構造化文書処理プログラムを提供することにある。

【課題を解決するための手段】

【 0 0 1 0 】

10

上記課題を解決するための本発明は、複数の構造化文書同士の非類似度を計算する構造化文書処理装置であって、構造化文書を構成している各タグに、レイアウトへの影響度に基づいて重みを割り当てるレイアウト判定手段と、前記割り当てた重み及び前記各構造化文書内のタグ構造に基づいて、構造化文書間の非類似度を計算する非類似度計算手段とを有することを特徴とする。

【 0 0 1 1 】

20

上記課題を解決するための本発明は、構造化文書のレイアウトと構造に基づいて構造化文書に対応付ける構造化文書処理システムであって、前記構造化文書を構成する各タグにレイアウトへの影響度に基づいて重みを割り当てるレイアウト判定手段と、前記割り当てられた重み及び前記構造化文書のタグ構造に基づいて、構造化文書を細分化したセクション間の非類似度を計算する非類似度計算手段と、前記計算された非類似度に基づいて構造化文書同士に対応付け、この構造化文書同士の対応付けを示す表示情報を生成するマッピング手段と前記生成した表示情報を、通信ネットワークを介して情報端末に送信する情報配信手段とを有することを特徴とする。

【 0 0 1 2 】

30

上記課題を解決するための本発明は、複数の構造化文書同士の非類似度を計算する構造化文書処理方法であって、構造化文書を構成している各タグに、レイアウトへの影響度に基づいて重みを割り当てるレイアウト判定ステップと、前記割り当てた重み及び前記各構造化文書内のタグ構造に基づいて、構造化文書間の非類似度を計算する非類似度計算ステップとを有することを特徴とする。

【 0 0 1 3 】

上記課題を解決するための本発明は、複数の構造化文書同士の非類似度を計算するプログラムであって、前記プログラムは、情報処理装置に、構造化文書を構成している各タグに、レイアウトへの影響度に基づいて重みを割り当てるレイアウト判定処理と、前記割り当てた重み及び前記各構造化文書内のタグ構造に基づいて、構造化文書間の非類似度を計算する非類似度計算処理とを実行させることを特徴とする。

【発明の効果】

【 0 0 1 4 】

40

本発明によると、レンダリングしなくてもセクションのマッピングを行えることにある。その理由は、レイアウトではなくセクションの構成タグを比較対象とするためである。

【 0 0 1 5 】

また、本発明によると、セクションのレイアウト及び構造の両方に基づいたマッピングを行えることにある。その理由は、各タグのレイアウトへの影響度、およびセクションのタグ構造を判定要素とするためである。

【発明を実施するための最良の形態】

【 0 0 1 6 】

本発明の特徴を説明するために、以下において、図面を参照して具体的に述べる。

【 0 0 1 7 】

本発明による構造化文書処理システムの特徴は、レイアウト判定部 1 0 2 と、非類似度計算部 1 0 3 と、マッピング部 1 0 4 とを有する点である。

50

【0018】

レイアウト判定部102は、構造化文書の各タグのレイアウトへの影響度に基づいて、重みをタグに割り当てる。

【0019】

非類似度計算部103は、タグ構造、およびレイアウト判定部102で各タグに割り当てられた重みに基づいて、比較するセクション同士の非類似度を計算する。

【0020】

マッピング部104は、非類似度計算部103で算出した非類似度に基づいて、適切に構造化文書をマッピングする。

【0021】

図1は、本発明による構造化文書処理システムの構成の一例を示すブロック図である。本実施の形態では、構造化文書処理システムは、ハードウェアで構成することも可能であるが、以下ではプログラムに従って動作するパーソナルコンピュータなどの情報処理端末によって実現する場合を用いて説明する。尚、構造化文書処理システムは、構造化文書を複数のセクションに分割して配信するシステム等のビジネスモデルに適用されてもよい。この場合、構造化文書処理システムは、例えば、構造化文書をレンダリングするソフトウェアを搭載した携帯電話やPDA、パーソナルコンピュータ等のユーザ端末と、構造化文書を処理する構造化文書処理サーバとを含んでもよい。

10

【0022】

図1に示すように、本実施の形態では、構造化文書処理システムは、プログラム制御により動作するデータ処理装置10と、情報を記憶する記憶装置11とを有する。

20

【0023】

データ処理装置10は、具体的には、プログラムに従って動作するパーソナルコンピュータやサーバ等によって実現される。

【0024】

データ処理装置10は、文書入力部100と、文書解析部101と、レイアウト判定部102と、非類似度計算部103と、マッピング部104と、出力部105とを有する。また、記憶装置11は、具体的には、メモリやハードディスク装置等によって実現される。記憶装置11は、セクション記憶部110と、マッピング記憶部111とを有する。

30

【0025】

文書入力部100は、外部から構造化文書を取得し、文書解析部101に出力する機能を備える。例えば、文書入力部100は、ユーザの操作に従って、記憶装置11から構造化文書を読み出し、文書解析部101に出力する。また、例えば、文書入力部100は、インターネット等の通信ネットワークを介して構造化文書（例えば、Webコンテンツなど）を受信し、文書解析部101に出力する。

【0026】

文書解析部101は、文書入力部100から取得した構造化文書を解析して、複数のセクションを抽出し、セクション記憶部110に記憶させる機能を備える。なお、取得した構造化文書をそのまま単一セクションとしてセクション記憶部110に記憶させてもよい。

40

【0027】

レイアウト判定部102は、セクション記憶部110からセクションを取得し、各セクションを構成するタグのレイアウトへの影響度に基づいて、各タグに重みを割り当てる機能を備える。また、レイアウト判定部102は、タグに重みを割り当てたセクションを、非類似度計算部に出力する機能を備える。なお、ここで説明するセクションには、セクションの集合であるセクショングループ（構造化文書そのものを含む）を含んでもよい。

【0028】

例えば、レイアウト判定部102は、タグのレイアウト定義をDTD (Document Type Definition) から取得し、各タグのレイアウトへの影響度をブロック要素（見出しや段落など、レイアウトを構成する基本要素）およびインライン要素

50

(強調やリンクなど、表示情報に役割や機能を与える要素)の2種類に分類する。そして、レイアウト判定部102は、ブロック要素に重みを大きく与え、一方インライン要素には重みを小さく与えることによって、レイアウトへの影響度に基づいた重み付けをする。

【0029】

非類似度計算部103は、レイアウト判定部102から比較するセクション(少なくとも2以上のセクション)を取得して、セクションのタグ構造に基づいて、セクション間の構造の類似性を示す非類似度を計算する機能を備える。また、非類似度計算部103は、セクション間の非類似度を、マッピング部104に出力する機能を備える。

【0030】

例えば、非類似度計算部103は、各セクションのタグ構造をツリー型に変換して、ツリー同士の編集距離を計算する。その際、レイアウト判定部102で計算した重みをツリーの各ノードの編集コストとすることによって、レイアウトへの影響度およびセクションの構造の両方に基づいた非類似度を計算する。

【0031】

マッピング部104は、非類似度計算部103から各セクションの非類似度を取得し、最も非類似度が小さい、即ち最も類似しているセクションの組み合わせをマッピングする(対応付ける)機能を備える。また、マッピング部104は、セクションのマッピング結果を、マッピング記憶部111に記憶させる機能を備える。例えば、マッピング部104は、構造化文書Daと構造化文書Dbとの比較において、最も非類似度が小さいセクションの組み合わせから順に、DaとDbとのセクションをすべてマッピングする。また、セクション数の違いからマッピングできずに残ってしまったセクションは、空セクションとマッピングする。なお、マッピング部104は、セクションの相対位置を考慮し、マッピングしたセクションを基準に、構造化文書内のセクションの集合を2つのグループに分けて、その各グループ内でマッピングを行ってもよい。また、マッピング部104は、セクションの絶対位置を考慮し、構造化文書内のセクションの階層構造に基づいて、セクションの集合を複数のグループに分け、その各グループ内でマッピングを行ってもよい。

【0032】

出力部105は、マッピング記憶部111が記憶しているセクションのマッピング結果を表示情報として外部に出力する機能を備える。例えば、出力部105は、マッピング記憶部111から、ユーザが指定したセクションと、そのセクションとマッピングされているセクションとを抽出し、液晶表示部やディスプレイ装置等の表示装置に表示させたり、通信ネットワークを介して情報端末に送信したりする。

【0033】

記憶装置11は、セクション記憶部110と、マッピング記憶部111とを含む。

【0034】

セクション記憶部110は、文書解析部101が解析した構造化文書のセクションを、構造化文書毎に記憶する。

【0035】

マッピング記憶部111は、マッピング部104が計算したセクションのマッピング結果を記憶する。

【0036】

次に、動作について説明する。図2は、構造化文書処理システムがレイアウトと文書構造に基づいてセクションをマッピングする処理の一例を示す流れ図である。

【0037】

まず、文書入力部100は、構造化文書を取得する。例えば、文書入力部100は、記憶装置11に格納されている構造化文書を読み出す。また、例えば、文書入力部100は、通信ネットワークを介して構造化文書(例えば、Webコンテンツなど)を受信する。

【0038】

次に、文書解析部101は、文書入力部100から取得した構造化文書を解析して、複数のセクションを抽出し、構造化文書毎にセクション記憶部110に記憶させる(ステッ

10

20

30

40

50

ブ S 1 1)。

【 0 0 3 9 】

続いて、レイアウト判定部 1 0 2 は、セクション記憶部 1 1 0 からセクションを取得し、タグのレイアウトへの影響度に基づいて、各タグに重みを割り当てる (ステップ S 1 2)。

【 0 0 4 0 】

次に、非類似度計算部 1 0 3 は、レイアウト判定部 1 0 2 から各構造化文書のセクションを取得し、タグに割り当てられた重みとタグの構造とに基づいて、比較元のセクションとこれ以外のセクションとの間の各非類似度を計算する (ステップ S 1 3)。

【 0 0 4 1 】

続いて、マッピング部 1 0 4 は、非類似度計算部 1 0 3 からセクションの非類似度を取得して、その非類似度に基づいて複数の構造化文書のセクション同士をマッピングし、マッピング記憶部 1 1 1 に記憶させる (ステップ S 1 4)。

【 0 0 4 2 】

また、出力部 1 0 5 は、マッピング記憶部 1 1 1 が記憶するセクションのマッピング結果を出力する。

【 0 0 4 3 】

以上のように、本実施の形態によれば、レイアウトへの影響度に基づいてタグに重みを割り当て、その重みとセクションのタグ構造に基づいて非類似度を計算し、非類似度に基づいてセクションをマッピングすることにより、レンダリングせずにレイアウトや文書構造に基づいた適切なセクションのマッピングが可能となる。

【 0 0 4 4 】

例えば、ある 2 つの構造化文書の比較において、微妙にレイアウトが異なっているセクションがある場合、あるセクションを構成するタグの種類や構造が、他のセクションと比べて最も類似していれば、それをマッピングすることができる。また、定性的にしか把握できなかった類似性を、非類似度という形で定量的に把握することができる。

【 0 0 4 5 】

次に、本発明による構造化文書処理システムの具体的な実施例について説明する。

【 0 0 4 6 】

まず、構造化文書処理システムの第 1 の実施例について説明する。なお、本実施例における構造化文書処理システムは、第 1 の実施の形態で示した構造化文書処理システムに相当する。また、本実施例では、データ処理装置がパーソナルコンピュータであり、データ記憶装置が磁気ディスク装置であるものとする。

【 0 0 4 7 】

パーソナルコンピュータ (データ処理装置) は、文書入力手段、文書解析手段、レイアウト判定手段、非類似度計算手段、マッピング手段、及び出力手段として機能する中央演算装置を含む。また、磁気ディスク装置 (記憶装置) は、パーソナルコンピュータによって解析または計算されたセクション情報や非類似度情報を記憶する。なお、データ処理装置は、サーバや携帯電話等でもよく、端末の種類によらない。また、本実施例では、構造化文書の例として、Web コンテンツを対象とする。例えば、パーソナルコンピュータは、インターネットを介して Web コンテンツを受信する。

【 0 0 4 8 】

本実施例では、まず、中央演算装置は、Web コンテンツを受信して、受信した Web コンテンツを解析してセクションを抽出する。そして、中央演算装置は、抽出したセクション情報を磁気ディスク装置に記憶させる。なお、受信した Web コンテンツをそのまま単一セクションとしてもよい。また、構造化文書は、レンダリングするために作成されたものであれば、HTML や XML などの種類に寄らない。本実施例では、構造化文書として HTML を扱う。

【 0 0 4 9 】

図 3 は、1 つの構造化文書の文書構造例を示す図であり、図 4 は、2 つの構造化文書が

10

20

30

40

50

解析されて複数のセクションが抽出された後のセクションのレイアウト構成を示す図である。本実施例では、中央演算装置は、図3のように、構造化文書D1から複数のセクションを抽出し、図4のように、セクションの抽出順序に基づいてそれぞれセクション1~8のように番号を割り振る。もう一方の構造化文書D2でも同様にセクションを抽出し、そのセクション情報を磁気ディスク装置に記憶させる。

【0050】

次に、中央演算装置は、磁気ディスク装置からセクション情報を取得し、各セクションのタグにレイアウトへの影響度に基づいて重みを割り当てる。本実施例では、中央演算装置は、各タグに割り当てる重みを、DTDに基づいて、ブロック要素（見出しや段落など、レイアウトを構成する基本要素）およびインライン要素（強調やリンクなど、表示情報に役割や機能を与える要素）の2つの種類に分けて計算する。なお、磁気ディスク装置から取得するセクションは、構造化文書を解析したセクションに限らず、他のセクションや、構造化文書そのものでもよい。また、取得するセクションあるいは構造化文書は複数でもよい。

10

【0051】

図5は、図3の構造化文書D1のタグをブロック要素およびインライン要素別に分け、重みとしてそれぞれ100と1を割り当てた例を示す説明図である。なお、図5の重みの値は、レイアウトへの影響度を強く評価したいため、ブロック要素を示す“div”等に対して「100」、インライン要素を示す“h”、“a”、“img”等に対して「1」のように、重みの差を大きくしたが、例えば、レイアウトへの影響度を緩く評価するならば、ブロック要素を「3」、インライン要素を「1」のように、重みの差を小さくしてもよい。また、レイアウトではなく、タグの構造の類似性を強く評価したいならば、ブロック要素およびインライン要素を共に1にして、重みの差を無くしてもよい。また、レイアウトへの影響度を定めるタグの種類は、ブロック要素およびインライン要素という定義を使用せずに、DTDの別の定義や、CSS（Cascading Style Sheets）などのDTD以外の構造化文書内の要素の表示を定義したレイアウト定義に従ってもよい。また、ユーザが予め定義したのもよい。また、予めレンダリングした際に各タグのレイアウトへの影響度を計算し、その記憶している結果に基づいてもよい。

20

【0052】

次に、中央演算装置は、セクションのタグ構造に基づいて、比較するセクションとの編集距離を計算する。本実施例では、タグの構造をツリー型に変換し、ツリー構造の編集距離を計算する。さらに、ノード一つの編集コストを各タグの重みとし、それに基づいて編集距離から非類似度を計算する。

30

【0053】

図6は、図5のセクション1をツリー型に変換した例を示す説明図である。なお、図6のRはツリーの根（Root）を示している。また、図6では、重みの値が大きいノードは表示サイズを大きく、逆に重みの値が小さいノードは小さく表現している。

【0054】

図7は、図6のセクション1のツリーを比較元として、別のセクションであるセクション2、セクション3、セクション4のそれぞれとの編集距離および非類似度の計算例を示す説明図である。本実施例では、ツリーの編集距離の計算において、レイアウトへの影響度を中心に計算するため、レイアウト要素（ブロック要素やインライン要素など）が同じならば、タグ名に関係しない。また、本実施例では、置換を使用せずに、置換 = 削除 + 挿入と見なす。さらに、本実施例では、レイアウトにツリーの兄弟順序も影響すると考え、ツリーの兄弟要素を入れ替えての構造一致は許可しない。

40

【0055】

例えば、セクション2のツリー構造を示した(a)のツリーとの編集距離の計算において、比較元のツリーからインライン要素“a”と“img”との2つを削除（del）した構造と同等になるため、編集距離が2となり、重みに基づいた非類似度は2となる。

【0056】

50

同様に、セクション3のツリー構造を示した(b)のツリーは、比較元のツリーから、インライン要素を1つ削除、1つ挿入(insert)、ブロック要素を1つ削除した構造と同等になるため、編集距離は3となり、重みに基づいた非類似度は102となる。

【0057】

さらに、セクション4のツリー構造を示した(c)のツリーは、比較元のツリーから、インライン要素を2つ削除、2つ挿入した構造と同等になるため、編集距離は4となり、重みに基づいた非類似度は4となる。

【0058】

なお、類似判定を厳密にするため、同じ要素でもタグ名が異なっていれば編集コストに一定の重みを与えて置換を使用してもよい。また、類似判定を緩和するため、ツリーの兄弟要素を入れ替えての構造一致を許可してもよい。また、タグの構造を用いるのならば、ツリー型に変換しなくてもよい。例えば、ブロック要素を1、及びインライン要素を0として、二進法的な表現にタグ構造を変換し、変換した文字列の編集距離を求めてもよい。また、中央演算装置は、非類似度を計算した時点で処理を終えてもよい。この場合、磁気ディスク装置に記憶される情報は、非類似度となる。

10

【0059】

次に、中央演算装置は、最も非類似度が小さいセクションの組み合わせから順にマッピングする。本実施例では、マッピングしたセクションを基準にセクションを2つのグループに分割し、それぞれのグループでセクションをマッピングするという処理を繰り返す。

【0060】

図8は、図4の構造化文書D1とD2とのマッピング例を示す説明図である。例えば、まず、最も非類似度が小さいセクション4(D1)とセクション4(D2)との組み合わせをマッピングする。このマッピングによって、構造化文書が2つのグループに分割され、それぞれD1ではG11とG12、D2ではG21とG22となる。

20

【0061】

次に、それぞれのグループで最も非類似度が小さいセクションの組み合わせをマッピングする。例えば、図9のように、G11:G21では、セクション2(D1)とセクション3(D2)、およびG12:G22では、セクション8(D1)とセクション9(D2)との組み合わせをマッピングする。

【0062】

そして、D1ではグループG11がG111とG112に、G12がG121とG122に、一方D2ではG21がG211とG212に、G22がG221とG222に分割される。なお、G122、G212及びG222はセクションが存在しない空グループである。

30

【0063】

続いて、今までと同様にそれぞれのグループで最も非類似度が小さいセクションの組み合わせをマッピングする。例えば、図10のように、G111:G211では、セクション1(D1)とセクション1(D2)、及びG121:G221では、セクション5(D1)とセクション7(D2)の組み合わせをマッピングする。

【0064】

そして、D1ではグループG111がG1111とG1112に、G121がG1211とG1212に、一方D2ではG211がG2111とG2112に、G221がG2211とG2212に分割される。ここで、G112:G212では、セクション3(D1)の組み合わせの相手がD2に存在しないため、存在しないセクションとマッピングされる。

40

【0065】

以下も同様に分割されたグループでのマッピングを繰り返し、最終的に図11のようなマッピング結果になる。

【0066】

なお、図12のように、セクションの文書内における階層構造を利用してグループ分割

50

してもよい。

【0067】

例えば、深さ2まで探索すると、まず、図13(a)のように、深さ1の階層において、グループのマッピングを行う。次に、図13(b)のように、深さ2の階層において、グループのマッピングを行う。その後、それぞれの分割されたグループ内において、セクションのマッピングを行う。

【0068】

上述した本発明は、レイアウトを持った構造化文書を複数のセクションに分割して利用するアプリケーションにおいて、レイアウトと文書構造に基づいて類似したセクションに適切にマッピングすることができるため、構造化文書の変化に強いマッピングが可能となる。また、構造化文書の類似性を定量的に把握できるようになる。

10

【0069】

例えば、ブログなどのレイアウト構成がよく変化するWebコンテンツにおいて、特定のセクションの更新情報を管理するアプリケーションでは、他のセクションが削除あるいは追加され、セクションの構成が変化したとしても、その特定のセクションを一貫して管理することができる。

【0070】

また、Webコンテンツ全体の更新情報を管理するアプリケーションでも、どのセクションが削除されたり追加されたりしたかを識別することができる。さらに、全くURLが異なるWebコンテンツ同士でも、その類似性をセクション単位で把握できるため、特定の類似性をもったセクションをすべてのWebコンテンツで非表示にするなどの処理が可能となる。

20

【図面の簡単な説明】

【0071】

【図1】構造化文書処理システムの他の構成例を示すブロック図である。

【図2】構造化文書処理システムがセクションをマッピングする処理の一例を示す流れ図である。

【図3】構造化文書の文書構造の例を示す説明図である。

【図4】セクションのレイアウト構成を示す説明図である。

【図5】ブロック要素およびインライン要素のタグに対して、重みとしてそれぞれ100と1を割り当てた例を示す説明図である。

30

【図6】セクションをツリー型に変換する例を示す説明図である。

【図7】編集距離および非類似度の計算例を示す説明図である。

【図8】構造化文書D1とD2とのマッピング手順を示す説明図である。

【図9】構造化文書D1とD2とのマッピング手順を示す説明図である。

【図10】構造化文書D1とD2とのマッピング手順を示す説明図である。

【図11】構造化文書D1とD2とのマッピングが完了した例を示す説明図である。

【図12】構造化文書D1の階層構造例を示す説明図である。

【図13】構造化文書D1とD2との階層構造を利用したマッピング例を示す説明図である。

40

【符号の説明】

【0072】

10 データ処理装置

11 記憶装置

100 文書入力部

101 文書解析部

102 レイアウト判定部

103 非類似度計算部

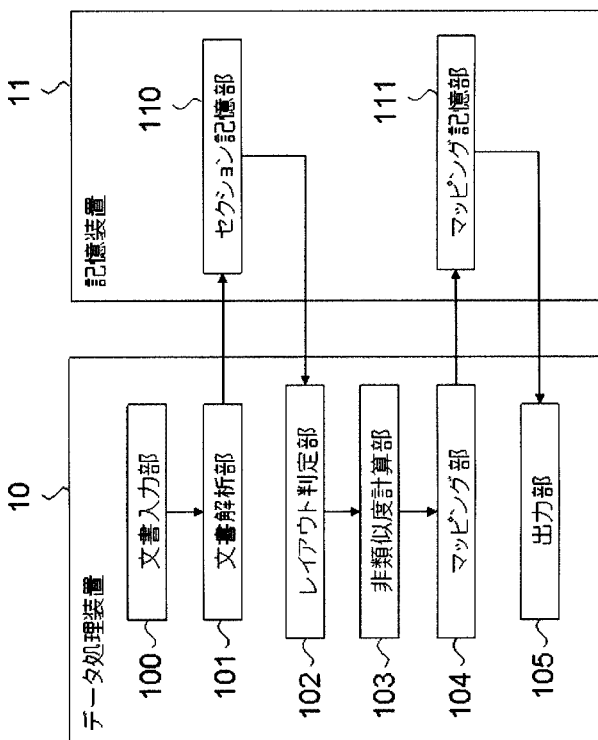
104 マッピング部

105 出力部

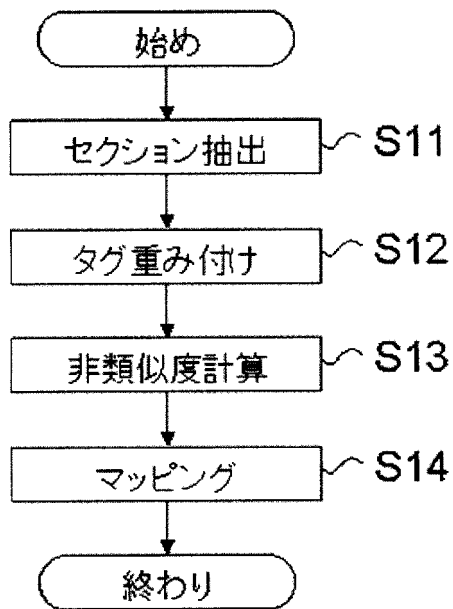
50

- 110 セクション記憶部
- 111 マッピング記憶部

【図1】

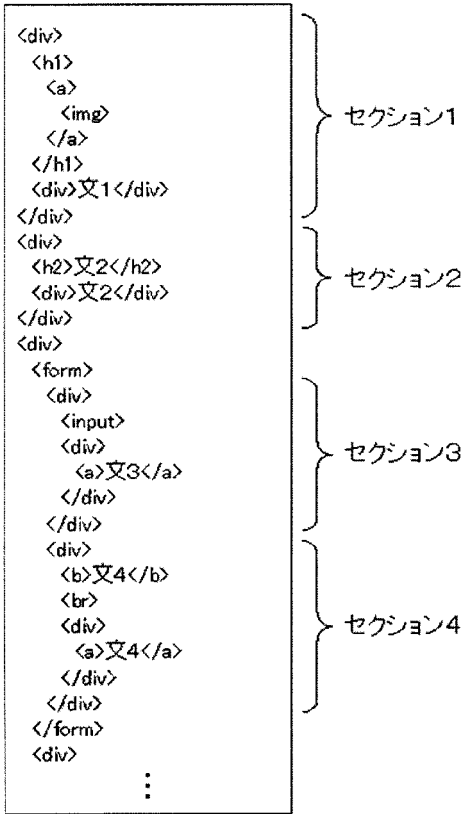


【図2】

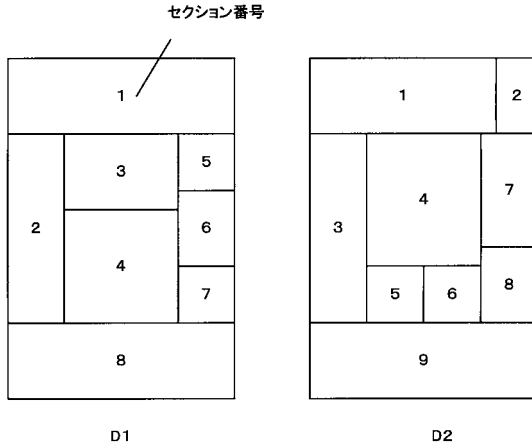


【 図 3 】

構造化文書D1

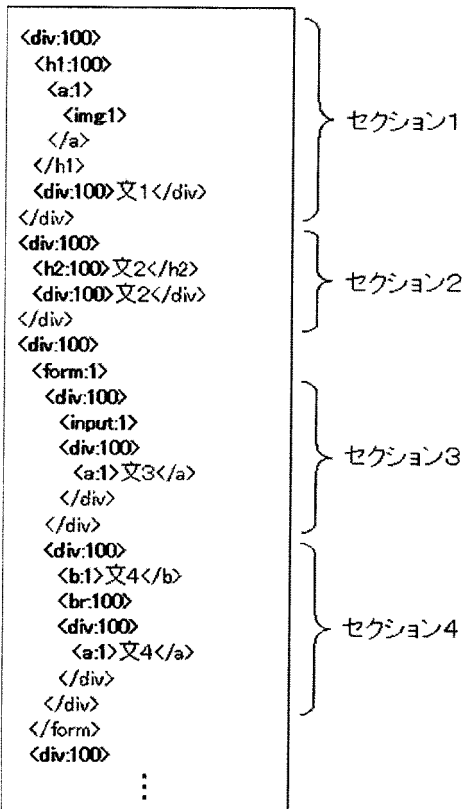


【 図 4 】

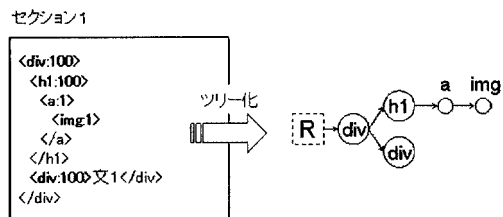


【 図 5 】

構造化文書D1



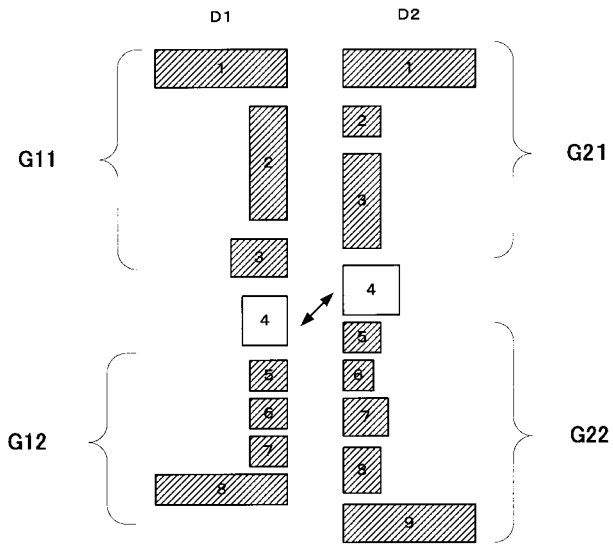
【 図 6 】



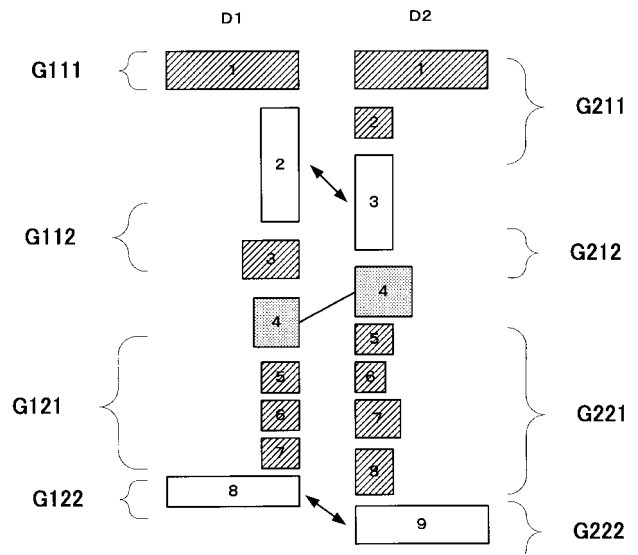
【 図 7 】

比較元セクション	非類似度	比較先セクション
	2 inline.del:2	(a)
	102 inline.del:1 inline.ins:1 block.del:1	(b)
	4 inline.del:2 inline.ins:2	(c)

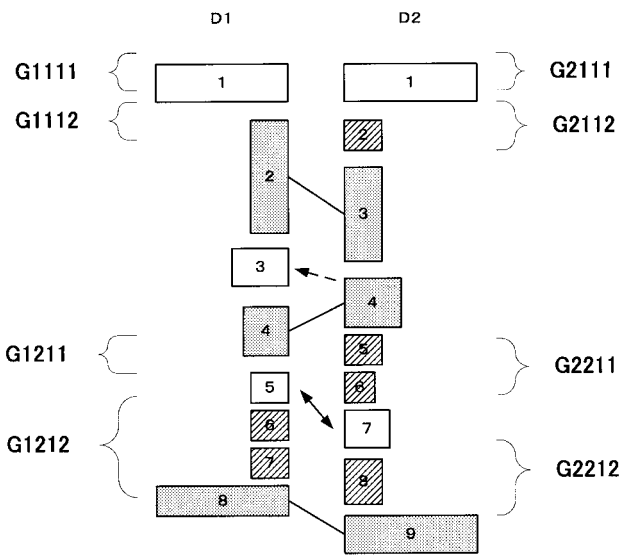
【 図 8 】



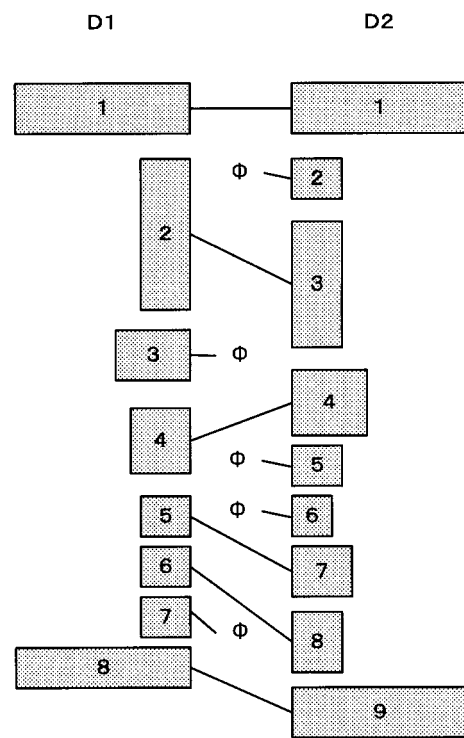
【 図 9 】



【 図 1 0 】

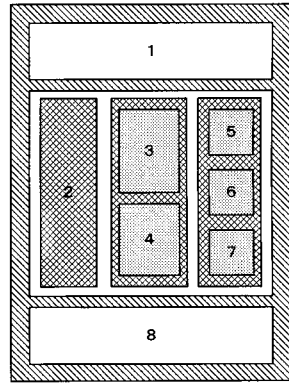
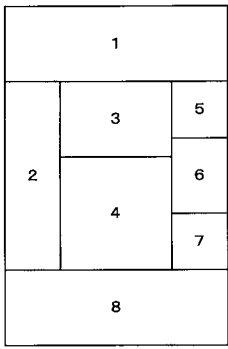


【 図 1 1 】



【図 1 2】

D1



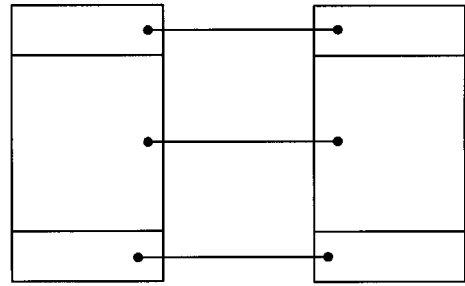
階層化

【図 1 3】

D1

D2

(a) 深さ2



(b) 深さ2

