

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5007405号
(P5007405)

(45) 発行日 平成24年8月22日(2012.8.22)

(24) 登録日 平成24年6月8日(2012.6.8)

(51) Int.Cl. F I
G06F 13/00 (2006.01) G O 6 F 13/00 5 5 0 A
G06F 17/21 (2006.01) G O 6 F 17/21 5 0 1 T

請求項の数 12 (全 50 頁)

<p>(21) 出願番号 特願2008-42722 (P2008-42722) (22) 出願日 平成20年2月25日(2008.2.25) (65) 公開番号 特開2009-199512 (P2009-199512A) (43) 公開日 平成21年9月3日(2009.9.3) 審査請求日 平成22年9月27日(2010.9.27)</p> <p>(出願人による申告)平成19年度独立行政法人情報通信研究機構、研究テーマ「軽度脳障害者のための情報セラピーインタフェースの研究開発」に関する委託研究、産業技術力強化法第19条の適用を受ける特許出願</p>	<p>(73) 特許権者 393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2 (74) 代理人 100115749 弁理士 谷川 英和 (72) 発明者 米澤 朋子 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内 (72) 発明者 光永 法明 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内 (72) 発明者 宮下 敬宏 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p style="text-align: right;">最終頁に続く</p>
---	---

(54) 【発明の名称】 情報処理装置、およびプログラム

(57) 【特許請求の範囲】

【請求項1】

1以上のウェブページを格納し得るウェブページ格納部と、
 分野を示す情報である分野情報と、ページ内のレイアウトに関する情報であるレイアウト情報との組を、2組以上格納し得るレイアウト情報格納部と、
 前記1以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2以上取得するページ情報取得部と、
 前記1以上のウェブページの内容から、当該1以上のウェブページの分野を決定し、分野情報を取得する分野情報決定部と、
 前記ページ情報取得部が取得した2以上のページ情報を出力するページ情報出力部を具備し、
 前記ページ情報出力部は、
 前記分野情報決定部が取得した分野情報と対になるレイアウト情報を取得するレイアウト情報取得手段と、
 前記レイアウト情報に従って、前記ページ情報取得部が取得した2以上のページ情報をレイアウトするレイアウト手段と、
 前記レイアウト手段がレイアウトした2以上のページ情報を出力するページ情報出力手段とを具備する情報処理装置。

【請求項2】

前記ページ情報取得部は、

前記 1 以上のウェブページから情報を取得するための情報であるページ取得ルール情報を、2 以上格納しているページ取得ルール情報格納手段と、

前記 1 以上のウェブページに対応するページ取得ルール情報、または前記分野情報決定部が取得した分野情報と対になるページ取得ルール情報を用いて、前記 1 以上のウェブページから、ページ単位の情報であるページ情報を、2 以上取得するページ情報取得手段とを具備する請求項 1 記載の情報処理装置。

【請求項 3】

前記分野情報決定部は、

分野情報と、当該分野情報が示す分野に対応する 1 以上の用語を含む用語情報を 1 以上格納している用語格納手段と、

10

前記 1 以上のウェブページから、前記 1 以上の用語の出現の度合いに関する情報である出現度情報を取得する出現度情報取得手段と、

前記出現度情報を用いて、前記 1 以上のウェブページの分野を決定し、分野情報を取得する分野情報決定手段を具備する請求項 1 または請求項 2 記載の情報処理装置。

【請求項 4】

前記分野情報決定部は、

前記 1 以上のウェブページから、1 以上のデータタイプのデータ量またはデータサイズに関する情報であるデータ情報を取得するデータ情報取得手段と、

前記データ情報を用いて、前記 1 以上のウェブページの分野を決定し、分野情報を取得する分野情報決定手段を具備する請求項 1 または請求項 2 記載の情報処理装置。

20

【請求項 5】

1 以上のウェブページを格納し得るウェブページ格納部と、

前記 1 以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2 以上取得するページ情報取得部と、

前記ページ情報取得部が取得した 2 以上のページ情報を出力するページ情報出力部を具備し、

前記ページ情報取得部は、

前記 1 以上のウェブページの各々に対して、ウェブページの構造またはタグまたは内容を用いて、主となる情報を取得する主情報取得手段と、

前記主情報取得手段が取得した情報から、ページ情報を取得するページ情報取得手段を具備する情報処理装置。

30

【請求項 6】

1 以上のウェブページを格納し得るウェブページ格納部と、

前記 1 以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2 以上取得するページ情報取得部と、

前記ページ情報取得部が取得した 2 以上のページ情報を出力するページ情報出力部を具備し、

前記ページ情報取得部は、

2 以上のウェブページの中の一のウェブページを決定し、当該一のウェブページと、他の 1 以上の各ウェブページの組のうち、1 以上の組に対して、内容の共通の度合いを示す情報である共通度を、1 以上、取得する共通度取得手段と、

40

前記 2 以上のウェブページから 2 以上のページ情報を取得するページ情報取得手段を具備し、

前記ページ情報取得手段は、

前記共通度に応じて、前記他の 1 以上の各ウェブページからの情報の取得方法が異なる情報処理装置。

【請求項 7】

前記共通度取得手段は、

前記 2 以上の各ウェブページのデータ量を算出し、データ量が最も大きいウェブページに対する、他の各ウェブページの共通度を取得し、

50

前記ページ情報取得手段は、

前記共通度が予め格納されている第一の閾値より小さい場合には、データ量が最も大きいウェブページに対して、前記他のウェブページを結合し、

前記共通度が前記第一の閾値より大きい場合には、各ウェブページを構成する一まとまりの情報であるブロックを取得し、前記データ量が最も大きいウェブページのブロック毎に、前記他のウェブページのブロックのうちで、予め格納されている第二の閾値より大きい共通度を有するブロックを抽出し、当該抽出したブロックを、前記データ量が最も大きいウェブページの前記ブロックと次のブロックの間に挿入し、出力する全ページ情報を構成し、

前記全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得する請求項6記載の情報処理装置。

10

【請求項8】

前記ページ情報取得手段は、

前記共通度が前記第一の閾値より大きい場合、前記他のウェブページのブロック中で、前記データ量が最も大きいウェブページの全ブロックのいずれに対しても、前記第二の閾値より大きな閾値を有しないブロックについて、その前に配置されているブロックで、前記データ量が最も大きいウェブページのいずれかのブロックに対して前記第二の閾値より大きな閾値を有するとして、前記データ量が最も大きいウェブページ内に挿入された箇所の直後に挿入する処理をさらに行い、出力する全ページ情報を構成し、前記全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得する請求項7記載の情報処理装置。

20

【請求項9】

前記ページ情報取得部は、

ユーザから受け付けたキーワードを用いて、1以上のサーバ装置からウェブページを検索し、当該検索したウェブページから、上位のn（nは2以上の自然数）のウェブページを取得する請求項5から請求項8いずれか記載の情報処理装置。

【請求項10】

記憶媒体に、

分野を示す情報である分野情報と、ページ内のレイアウトに関する情報であるレイアウト情報との組を、2組以上格納しており、

コンピュータを、

1以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2以上取得するページ情報取得部と、

前記1以上のウェブページの内容から、当該1以上のウェブページの分野を決定し、分野情報を取得する分野情報決定部と、

前記ページ情報取得部が取得した2以上のページ情報を出力するページ情報出力部として機能させるためのプログラムであって、

前記ページ情報出力部は、

前記分野情報決定部が取得した分野情報と対になるレイアウト情報を取得するレイアウト情報取得手段と、

前記レイアウト情報に従って、前記ページ情報取得部が取得した2以上のページ情報をレイアウトするレイアウト手段と、

前記レイアウト手段がレイアウトした2以上のページ情報を出力するページ情報出力手段とを具備するものとして、コンピュータを機能させるプログラム。

30

40

【請求項11】

コンピュータを、

1以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2以上取得するページ情報取得部と、

前記ページ情報取得部が取得した2以上のページ情報を出力するページ情報出力部として機能させるプログラムであって、

前記ページ情報取得部は、

50

前記 1 以上のウェブページの各々に対して、ウェブページの構造またはタグまたは内容を用いて、主となる情報を取得する主情報取得手段と、
前記主情報取得手段が取得した情報から、ページ情報を取得するページ情報取得手段を具備するものとして、コンピュータを機能させるプログラム。

【請求項 1 2】

コンピュータを、

1 以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2 以上取得するページ情報取得部と、

前記ページ情報取得部が取得した 2 以上のページ情報を出力するページ情報出力部として機能させるプログラムであって、

前記ページ情報取得部は、

2 以上のウェブページの中の一のウェブページを決定し、当該一のウェブページと、他の 1 以上の各ウェブページの組のうち、1 以上の組に対して、内容の共通の度合いを示す情報である共通度を、1 以上、取得する共通度取得手段と、

前記 2 以上のウェブページから 2 以上のページ情報を取得するページ情報取得手段を具備し、

前記ページ情報取得手段は、

前記共通度に応じて、前記他の 1 以上の各ウェブページからの情報の取得方法が異なるものとして、コンピュータを機能させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、1 以上のウェブページから、ページ単位の情報に分割した電子ブックのコンテンツを得る情報処理装置等に関するものである。

【背景技術】

【0002】

従来、より幅広いユーザ層が快適にウェブページ等の閲覧対象情報を閲覧できるように配慮した情報閲覧システムがあった（例えば、特許文献 1 参照）。本システムにおいて、ウェブページを TV モニタに表示させるために、見出しとなる表示要素のサイズを拡大させるなどの処理を行っていた。

【特許文献 1】特開 2006 - 279887 号公報（第 1 頁、第 1 図等）

【発明の開示】

【発明が解決しようとする課題】

【0003】

しかしながら、従来の情報閲覧システムにおいては、ウェブページから、複数ページの集合である電子ブックのコンテンツを自動構成することはできなかった。

【課題を解決するための手段】

【0004】

本第一の発明の情報処理装置は、1 以上のウェブページを格納し得るウェブページ格納部と、前記 1 以上のウェブページから、ページ単位の情報であるページ情報を、2 以上取得するページ情報取得部と、前記ページ情報取得部が取得した 2 以上のページ情報を出力するページ情報出力部を具備する情報処理装置である。

【0005】

かかる構成により、ウェブページから、複数ページの集合である電子ブックのコンテンツを自動構成することができる。

【0006】

また、本第二の発明の情報処理装置は、第一の発明に対して、分野を示す情報である分野情報と、ページ内のレイアウトに関する情報であるレイアウト情報との組を、2 組以上格納しているレイアウト情報格納部と、前記 1 以上のウェブページの内容から、当該 1 以上のウェブページの分野を決定し、分野情報を取得する分野情報決定部とをさらに具備し

10

20

30

40

50

、前記ページ情報出力部は、前記分野情報決定部が取得した分野情報と対になるレイアウト情報を取得するレイアウト情報取得手段と、前記レイアウト情報に従って、前記ページ情報取得部が取得した2以上のページ情報をレイアウトするレイアウト手段と、前記レイアウト手段がレイアウトした2以上のページ情報を出力するページ情報出力手段とを具備する情報処理装置である。

【0007】

かかる構成により、ウェブページの内容に適したレイアウトで、電子ブックのコンテンツを自動構成することができる。

【0008】

また、本第三の発明の情報処理装置は、第二の発明に対して、前記ページ情報取得部は、分野情報と、前記1以上のウェブページから情報を取得するための情報であるページ取得ルール情報との組を、2以上格納しているページ取得ルール情報格納手段と、前記分野情報決定部が取得した分野情報と対になるページ取得ルール情報を用いて、前記1以上のウェブページから、ページ単位の情報であるページ情報を、2以上取得するページ情報取得手段とを具備する情報処理装置である。

10

【0009】

かかる構成により、ウェブページの分野を適切に判断した上で、ウェブページの内容に適したレイアウトで、電子ブックのコンテンツを自動構成することができる。

【0010】

また、本第四の発明の情報処理装置は、第二、第三いずれかの発明に対して、前記分野情報決定部は、分野情報と、当該分野情報が示す分野に対応する1以上の用語を含む用語情報を1以上格納している用語格納手段と、前記1以上のウェブページから、前記1以上の用語の出現の度合いに関する情報である出現度情報を取得する出現度情報取得手段と、前記出現度情報を用いて、前記1以上のウェブページの分野を決定し、分野情報を取得する分野情報決定手段を具備する情報処理装置である。

20

【0011】

かかる構成により、例えば、ニュース記事のウェブページから、新聞風のレイアウトのコンテンツが構成でき、ブログのウェブページから、日記風のレイアウトのコンテンツが構成できる。

【0012】

また、本第五の発明の情報処理装置は、第二、第三いずれかの発明に対して、前記分野情報決定部は、前記1以上のウェブページから、1以上のデータタイプのデータ量またはデータサイズに関する情報であるデータ情報を取得するデータ情報取得手段と、前記データ情報を用いて、前記1以上のウェブページの分野を決定し、分野情報を取得する分野情報決定手段を具備する情報処理装置である。

30

【0013】

かかる構成により、例えば、子供向けのウェブページから、絵本風のレイアウトのコンテンツが構成できる。

【0014】

また、本第六の発明の情報処理装置は、第一の発明に対して、前記ページ情報取得部は、前記1以上のウェブページの各々に対して、ウェブページの構造またはタグまたは内容を用いて、主となる情報を取得する主情報取得手段と、前記主情報取得手段が取得した情報から、ページ情報を取得するページ情報取得手段を具備する情報処理装置である。

40

【0015】

かかる構成により、ウェブページに記載されている情報のうち、重要と思われる情報のみを用いて、電子ブックのコンテンツを自動構成することができる。

【0016】

また、本第七の発明の情報処理装置は、第一の発明に対して、前記ページ情報取得部は、2以上のウェブページの各々から、内容の共通の度合いを示す情報である共通度を取得する共通度取得手段と、前記2以上のウェブページから2以上のページ情報を取得する手

50

段であり、前記共通度に応じて、取得方法が異なるページ情報取得手段を具備する情報処理装置である。

【0017】

かかる構成により、複数のウェブページの内容を用いて、優れた電子ブックのコンテンツを自動構成することができる。

【0018】

また、本第八の発明の情報処理装置は、第七の発明に対して、前記共通度取得手段は、前記2以上の各ウェブページのデータ量を算出し、データ量が最も大きいウェブページに対する、他の各ウェブページの共通度を取得し、前記ページ情報取得手段は、前記共通度が予め格納されている第一の閾値より小さい場合には、データ量が最も大きいウェブページに対して結合し、前記共通度が前記第一の閾値より大きい場合には、各ウェブページを構成する一まとまりの情報であるブロックを取得し、前記データ量が最も大きいウェブページのブロック毎に、前記他の各ウェブページのブロックのうちで、予め格納されている第二の閾値より大きい共通度を有するブロックを抽出し、当該抽出したブロックを、前記データ量が最も大きいウェブページの前記ブロックと次のブロックの間に挿入し、出力する全ページ情報を構成し、前記全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得する情報処理装置である。

10

【0019】

かかる構成により、複数のウェブページの内容を用いて、近似する内容の情報が近いところに配置され、整理された電子ブックのコンテンツを自動構成することができる。

20

【0020】

また、本第九の発明の情報処理装置は、第八の発明に対して、前記ページ情報取得手段は、前記他の各ウェブページのブロック中で、前記データ量が最も大きいウェブページの全ブロックのいずれに対しても、前記第二の閾値より大きな閾値を有しないブロックについて、その前に配置されているブロックで、前記データ量が最も大きいウェブページのいずれかのブロックに対して前記第二の閾値より大きな閾値を有するとして、前記データ量が最も大きいウェブページ内に挿入された箇所の直後に挿入する処理をさらに行い、出力する全ページ情報を構成し、前記全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得する情報処理装置である。

【0021】

30

かかる構成により、複数のウェブページの内容を用いて、近似する内容の情報が近いところに配置され、整理された電子ブックのコンテンツを自動構成することができる。

【0022】

また、本第十の発明の情報処理装置は、第七から第十いずれかの発明に対して、前記ページ情報取得部は、ユーザから受け付けたキーワードを用いて、1以上のサーバ装置からウェブページを検索し、当該検索したウェブページから、上位の n (n は2以上の自然数)のウェブページを取得する情報処理装置である。

【0023】

かかる構成により、ユーザが欲する情報を、複数のウェブページから、自動的に選択し、有用な電子ブックのコンテンツを自動構成することができる。

40

【発明の効果】

【0024】

本発明による情報処理システムによれば、ウェブページから、複数ページの集合である電子ブックのコンテンツを自動構成することができる。

【発明を実施するための最良の形態】

【0025】

以下、情報処理システム等の実施形態について図面を参照して説明する。なお、実施の形態において同じ符号を付した構成要素は同様の動作を行うので、再度の説明を省略する場合がある。

【0026】

50

(実施の形態1)

本実施の形態において、1以上のウェブページから、ページ単位の情報に分割した電子ブックのコンテンツを得る情報処理システム等について説明する。また、本実施の形態において、ウェブページから分野(ニュース、子供向けなど)を決定し、当該分野に対応するテンプレート(後述するレイアウト情報と同意義)を取得し、当該テンプレートに従って、ウェブページを電子ブックのコンテンツに変換する情報処理システム等について説明する。また、ウェブページのURLによって、異なるテンプレートを取得し、当該テンプレートに従って、ウェブページを電子ブックのコンテンツに変換する情報処理システム等について説明する。また、本実施の形態において、ウェブページの分野やウェブページのURLによって、ページ情報の取得方法が異なる情報処理システム等について説明する。さらに、本実施の形態において、ウェブページの主要な情報(後述する主情報)を取得し、主情報のみを用いて、ウェブページから電子ブックのコンテンツを取得する情報処理システム等について説明する。

10

【0027】

図1は、本実施の形態における情報処理システム1の概念図である。情報処理システム1は、情報処理装置11、1以上のサーバ装置12を具備する。情報処理装置11は、ウェブページから電子ブックのコンテンツを得る処理を行う装置である。また、サーバ装置12は、ウェブページを格納しており、情報処理装置11等の他の装置からの要求に応じて、ウェブページを送信する装置である。

【0028】

図2は、本実施の形態における情報処理システム1のブロック図である。情報処理装置11は、指示受付部110、ウェブページ格納部111、ウェブページ受信部112、ウェブページ蓄積部113、レイアウト情報格納部114、分野情報決定部115、ページ情報取得部116、ページ情報出力部117を具備する。

20

【0029】

分野情報決定部115は、用語格納手段1151、出現度情報取得手段1152、分野情報決定手段1153を具備する。

【0030】

ページ情報取得部116は、ページ取得ルール情報格納手段1161、主情報取得手段1162、ページ情報取得手段1163を具備する。

30

【0031】

ページ情報出力部117は、レイアウト情報取得手段1171、レイアウト手段1172、ページ情報出力手段1173を具備する。

【0032】

サーバ装置12は、ウェブページ記憶部121、ウェブページ送信部122を具備する。

【0033】

指示受付部110は、指示を受け付ける。この指示とは、例えば、ページ情報を出力する指示である。また、指示とは、他の指示でも良い。また、ページ情報出力指示は、電子ブックのコンテンツの元になるウェブページを特定する情報(URLやURIなど)を含んでも良いし、元になるウェブページを検索するためのキーワードなどを含んでも良い。また、受付とは、ユーザからの入力の受け付け、他の装置からの受信などを含む概念である。指示の入力手段は、キーボードやマウスやメニュー画面によるもの等、何でも良い。指示受付部110は、キーボード等の入力手段のデバイスドライバーや、メニュー画面の制御ソフトウェア等で実現され得る。

40

【0034】

ウェブページ格納部111は、1以上のウェブページを格納し得る。ウェブページは、通常、他の装置から取得されたウェブページである。ウェブページ格納部111は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。ウェブページ格納部111にウェブページが記憶される過程は問わない。例えば、記録媒体を介してウェブ

50

ブページがウェブページ格納部 1 1 1 で記憶されるようになってよく、通信回線等を介して送信されたウェブページがウェブページ格納部 1 1 1 で記憶されるようになってよく、あるいは、入力デバイスを介して入力されたウェブページがウェブページ格納部 1 1 1 で記憶されるようになってよい。

【 0 0 3 5 】

ウェブページ受信部 1 1 2 は、ウェブページをサーバ装置 1 2 から受信する。ウェブページ受信部 1 1 2 は、通常、無線または有線の通信手段で実現されるが、放送を受信する手段で実現されても良い。

【 0 0 3 6 】

ウェブページ蓄積部 1 1 3 は、ウェブページ受信部 1 1 2 が受信したウェブページを、ウェブページ格納部 1 1 1 に蓄積する。ウェブページ蓄積部 1 1 3 は、通常、MPU やメモリ等から実現され得る。ウェブページ蓄積部 1 1 3 の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

【 0 0 3 7 】

レイアウト情報格納部 1 1 4 は、分野を示す情報である分野情報と、電子ブックのページ内のレイアウトに関する情報であるレイアウト情報との組を、2組以上格納している。また、レイアウト情報格納部 1 1 4 は、ウェブページを特定する情報であるウェブページ特定情報(URL やURI など)とレイアウト情報との組を、2組以上格納している。また、分野情報とレイアウト情報との組、またはウェブページ特定情報とレイアウト情報との組を、以下、適宜、格納レイアウト情報という。分野とは、ニュース、ブログ、子供向け、技術解説などである。分野情報とは、分野に対応する情報である。分野情報は、分野がニュースの場合は「1」、ブログの場合は「2」、子供向けの場合は「3」、技術解説の場合は「4」などである。レイアウト情報とは、ページ情報をレイアウトするための、いわゆる定義データでも良いし、ページ情報をレイアウトするプログラムまたはプログラム名などでも良い。レイアウト情報格納部 1 1 4 は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。レイアウト情報格納部 1 1 4 にレイアウト情報が記憶される過程は問わない。例えば、記録媒体を介してレイアウト情報がレイアウト情報格納部 1 1 4 で記憶されるようになってよく、通信回線等を介して送信されたレイアウト情報がレイアウト情報格納部 1 1 4 で記憶されるようになってよく、あるいは、入力デバイスを介して入力されたレイアウト情報がレイアウト情報格納部 1 1 4 で記憶されるようになってよい。

【 0 0 3 8 】

分野情報決定部 1 1 5 は、1以上のウェブページの内容から、当該1以上のウェブページの分野を決定し、分野情報を取得する。なお、ウェブページの内容とは、例えば、html やxmlなどの記述言語で記述されたウェブページの情報のことである。分野情報決定部 1 1 5 は、例えば、ニュース記事に現れる用語を多数含むウェブページの分野を「ニュース」と決定する。また、分野情報決定部 1 1 5 は、例えば、ウェブページのタイトルに「ニュース」を含むウェブページの分野を「ニュース」と決定する。また、分野情報決定部 1 1 5 は、例えば、口語調の文が多い(口語調の文に出現する単語(例えば、「おれ」「~だし」「だね」など)を多数含む)ウェブページ、日付が先頭から所定の範囲に記載されているウェブページ、主語がない文が多いなどのウェブページの分野を「ブログ」と決定する。また、分野情報決定部 1 1 5 は、例えば、ひらがなの割合が多いウェブページの分野を「子供向け」と決定する。また、分野情報決定部 1 1 5 は、例えば、画像の数や割合が多いウェブページの分野を「子供向け」と決定する。分野情報決定部 1 1 5 は、通常、MPU やメモリ等から実現され得る。分野情報決定部 1 1 5 の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

【 0 0 3 9 】

用語格納手段 1 1 5 1 は、分野情報と、当該分野情報が示す分野に対応する1以上の用

10

20

30

40

50

語を含む用語情報を1以上格納している。用語情報は、例えば、(ニュース、「事件」「総理」「政治」「経済」・・・)、(ブログ、「天気」「晴れ」「曇り」「雨」・・・)などである。なお、ここでの用語情報のデータ構造は(分野情報、1以上の用語)である。用語格納手段1151は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。用語格納手段1151に用語情報が記憶される過程は問わない。例えば、記録媒体を介して用語情報が用語格納手段1151で記憶されるようになってよく、通信回線等を介して送信された用語情報が用語格納手段1151で記憶されるようになってよく、あるいは、入力デバイスを介して入力された用語情報が用語格納手段1151で記憶されるようになってよい。

【0040】

出現度情報取得手段1152は、1以上のウェブページから、1以上の用語の出現の度合いに関する情報である出現度情報を取得する。出現度情報は、用語の出現数でも良いし、出現の割合(出現数/全用語数)でも良いし、タイトルの用語として使われているか否でも良いし、tf・idfにより取得される情報である。その他、出現度情報は、用語の出現の度合いに関する情報であれば良い。出現度情報取得手段1152は、ある用語Aを、1以上のウェブページに対して、パターンマッチングを行い、出現度情報を取得しても良いし、形態素解析を行い、形態素を抽出し、その品詞などを用いて、出現度情報を取得しても良い。出現度情報を取得するアルゴリズムは問わない。出現度情報取得手段1152は、通常、MPUやメモリ等から実現され得る。出現度情報取得手段1152の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

【0041】

分野情報決定手段1153は、出現度情報を用いて、1以上のウェブページの分野を決定し、分野情報を取得する。分野情報決定手段1153は、分野情報が示す分野に対応する1以上の用語の出現度情報が予め決められた条件を満たす(例えば、閾値以上、など)場合は、当該分野情報を取得する。分野情報決定手段1153は、ウェブページのタイトルで使われている用語に対応する分野情報を取得しても良い。分野情報決定手段1153は、通常、MPUやメモリ等から実現され得る。分野情報決定手段1153の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

【0042】

ページ情報取得部116は、1以上のウェブページから、ページ単位の情報であるページ情報を、2以上取得する。ページ情報とは、電子ブック用のページの情報である。例えば、1ページに入るデータ量が決まっており、ページ情報取得部116は、1以上のウェブページから、その1ページ分のデータごとに、2ページ以上、情報を区切って取得する。また、ページ情報の取得とは、1以上のウェブページから、取得した情報を用いて、ページ単位の情報を構成することである、とも言える。ページ情報取得部116の具体的な処理アルゴリズムの例は、後述する。ページ情報取得部116は、通常、MPUやメモリ等から実現され得る。ページ情報取得部116の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

【0043】

ページ取得ルール情報格納手段1161は、分野情報と、1以上のウェブページから情報を取得するための情報であるページ取得ルール情報との組を、2以上格納している。ページ取得ルール情報格納手段1161は、ウェブページ特定情報とページ取得ルール情報との組を、2以上格納していても良い。ページ取得ルール情報とは、取得する情報に対応するタグ、取得する情報のデータタイプ(静止画、テキスト列など)、除く情報(タグ、データ)を指定する情報などにより構成されても良い。ページ取得ルール情報は、プログラムに埋め込まれていても、プログラム自体でも良い。ページ取得ルール情報は、正規表現で記載されていても、特定の構造、文言からなる自然言語で記載されていても良い。ペ

10

20

30

40

50

ページ取得ルール情報の記載方法は問わない。ページ取得ルール情報の具体例については、後述する。ページ取得ルール情報格納手段1161は、揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。ページ取得ルール情報格納手段1161にページ取得ルール情報が記憶される過程は問わない。例えば、記録媒体を介してページ取得ルール情報がページ取得ルール情報格納手段1161で記憶されるようになってよく、通信回線等を介して送信されたページ取得ルール情報がページ取得ルール情報格納手段1161で記憶されるようになってよく、あるいは、入力デバイスを介して入力されたページ取得ルール情報がページ取得ルール情報格納手段1161で記憶されるようになってよくよい。

【0044】

主情報取得手段1162は、1以上のウェブページの各々に対して、ウェブページの構造またはタグまたは内容を用いて、主となる情報（以下、適宜、主情報という。）を取得する。主情報取得手段1162は、例えば、右側にレイアウトされるビットマップをバナー広告であると判断し、主情報から除く処理を行う。また、主情報取得手段1162は、例えば、タグ「<div class="topicsindex">」に対応するテキスト列を、主情報として取得する。その他、主情報の取得方法は、種々あり得る。なお、主情報の取得方法を、ページ取得ルール情報としても良い。かかる場合、主情報取得手段1162は、ページ情報取得手段1163と同様の処理となる。主情報取得手段1162は、通常、MPUやメモリ等から実現され得る。主情報取得手段1162の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

【0045】

ページ情報取得手段1163は、分野情報決定部115が取得した分野情報と対になるページ取得ルール情報を用いて、1以上のウェブページから、ページ単位の情報であるページ情報を、2以上取得する。また、ページ情報取得手段1163は、ウェブページに対応するページ取得ルール情報を用いて、1以上のウェブページから、ページ単位の情報であるページ情報を、2以上取得しても良い。また、ページ情報取得手段1163は、主情報取得手段1162が取得した情報から、2以上のページ情報を取得しても良い。ページ情報取得手段1163は、通常、MPUやメモリ等から実現され得る。ページ情報取得手段1163の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

【0046】

ページ情報出力部117は、ページ情報取得部116が取得した2以上のページ情報を出力する。ここで、出力とは、ディスプレイへの表示、プロジェクターを用いた投影、プリンタへの印字、音出力、外部の装置（例えば、電子ブック装置）への送信、記録媒体（例えば、電子ブック装置に装着して用いられる着脱可能な可搬型の記録媒体）への蓄積、他の処理装置や他のプログラム等への処理結果の引渡し等を含む概念である。ページ情報出力部117は、ディスプレイやスピーカー等の出力デバイスを含むと考えると含まないと考えても良い。ページ情報出力部117は、出力デバイスのドライバーソフトまたは、出力デバイスのドライバーソフトと出力デバイス等で実現され得る。

【0047】

レイアウト情報取得手段1171は、分野情報決定部115が取得した分野情報と対になるレイアウト情報を、レイアウト情報格納部114から取得する。また、レイアウト情報取得手段1171は、ページ情報を構成する元となる処理対象のウェブページに対応するレイアウト情報を、レイアウト情報格納部114から取得しても良い。レイアウト情報取得手段1171は、通常、MPUやメモリ等から実現され得る。レイアウト情報取得手段1171の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

【0048】

レイアウト手段1172は、レイアウト情報取得手段1171が取得したレイアウト情

10

20

30

40

50

報に従って、ページ情報取得部 116 が取得した 2 以上のページ情報をレイアウトする。レイアウト手段 1172 が、ページ情報をレイアウトした結果、電子ブックの 2 以上のページが構成される。レイアウト手段 1172 は、通常、MPU やメモリ等から実現され得る。レイアウト手段 1172 の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアは ROM 等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

【0049】

ページ情報出力手段 1173 は、レイアウト手段 1172 がレイアウトした 2 以上のページ情報を出力する。ページ情報出力手段 1173 は、ディスプレイやスピーカー等の出力デバイスを含むと考えても含まないと考えても良い。ページ情報出力手段 1173 は、出力デバイスのドライバーソフトまたは、出力デバイスのドライバーソフトと出力デバイス等で実現され得る。

10

【0050】

ウェブページ記憶部 121 は、1 以上のウェブページを格納している。ウェブページ記憶部 121 は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。

【0051】

ウェブページ送信部 122 は、ウェブページ記憶部 121 のウェブページを、情報処理装置 11 に送信する。ウェブページ送信部 122 は、通常、情報処理装置 11 からの要求を受け付け、当該要求に対応して、ウェブページを、情報処理装置 11 に送信する。ウェブページ送信部 122 は、通常、無線または有線の通信手段で実現されるが、放送手段で実現されても良い。

20

【0052】

次に、情報処理システム 1 の動作について説明する。まず、情報処理装置 11 の動作例について、図 3 から図 5 のフローチャートを用いて説明する。

【0053】

(ステップ S301) ウェブページ受信部 112 は、ウェブページを受信したか否かを判断する。ウェブページを受信すればステップ S302 に行き、ウェブページを受信しなければステップ S303 に行く。

【0054】

(ステップ S302) ウェブページ蓄積部 113 は、ステップ S301 で受信されたウェブページを、ウェブページ格納部 111 に蓄積する。ステップ S301 に戻る。

30

【0055】

(ステップ S303) 指示受付部 110 は、ページ情報出力指示を受け付けたか否かを判断する。ページ情報出力指示を受け付ければステップ S304 に行き、ページ情報出力指示を受け付けなければステップ S301 に戻る。なお、ページ情報出力指示には、1 以上のウェブページ特定情報を含む。1 以上のウェブページ特定情報は、例えば、トップページの URL であり、例えば、トップページからリンクされているウェブページも含めたウェブページ群から、ページ情報が構成され得る。

【0056】

(ステップ S304) 分野情報決定部 115 は、ウェブページ格納部 111 から、ページ情報出力指示により特定される 1 以上のウェブページを読み出す。

40

【0057】

(ステップ S305) 分野情報決定部 115 は、ステップ S304 で読み出された 1 以上のウェブページから、分野情報を取得する。分野情報を取得する処理（以下、適宜「分野決定処理」という）の詳細例について、図 4 のフローチャートを用いて説明する。

【0058】

(ステップ S306) ページ情報取得部 116 は、ステップ S304 で読み出された 1 以上のウェブページから、ページ情報を取得する処理を行う。ページ情報取得処理の詳細例について、図 5 のフローチャートを用いて説明する。

50

【0059】

(ステップS307)レイアウト情報取得手段1171は、ステップS305で取得された分野情報に対応するレイアウト情報を、レイアウト情報格納部114から読み出す。

【0060】

(ステップS308)レイアウト手段1172は、ステップS306で取得されたページ情報(ここでは、レイアウトされる前の情報)を、ステップS307で読み出されたレイアウト情報に従ってレイアウトし、ページ情報をメモリ上に得る。ここで、レイアウト手段1172は、例えば、ウェブページ内のタグを解釈し、ページ単位に区切る(例えば、改ページコードを挿入する)処理を行う。ウェブページ内のタグを解釈する処理は、いわゆるウェブブラウザの処理である。

10

【0061】

(ステップS309)ページ情報出力手段1173は、ステップS308でレイアウトし、構成された2以上のページ情報を出力する。ステップS301に戻る。

【0062】

なお、図3のフローチャートにおいて、分野情報を自動取得したが、分野情報は、ユーザが入力し、指示受付部110が受け付けても良い。

【0063】

また、図3のフローチャートにおいて、分野情報に応じたレイアウト情報を用いて、ページ情報を構成したが、ウェブページに対応するレイアウト情報を用いて、ページ情報を構成しても良い。

20

【0064】

さらに、図3のフローチャートにおいて、電源オフや処理終了の割り込みにより処理は終了する。

【0065】

次に、ステップS305の分野決定処理の詳細例について、図4のフローチャートを用いて説明する。

【0066】

(ステップS401)出現度情報取得手段1152は、カウンタ*i*に1を代入する。

【0067】

(ステップS402)出現度情報取得手段1152は、*i*番目の分野情報が、用語格納手段1151に格納されているか否かを判断する。*i*番目の分野情報が格納されていればステップS403に行き、格納されていなければステップS411に行く。

30

【0068】

(ステップS403)出現度情報取得手段1152は、*i*番目の分野情報に対応する1以上の用語を、用語格納手段1151から読み出す。

【0069】

(ステップS404)出現度情報取得手段1152は、カウンタ*j*に1を代入する。

【0070】

(ステップS405)出現度情報取得手段1152は、ステップS403で取得した用語の中に、*j*番目の用語が存在するか否かを判断する。*j*番目の用語が存在すればステップS406に行き、*j*番目の用語が存在しなければステップS409に行く。

40

【0071】

(ステップS406)出現度情報取得手段1152は、ステップS405で取得した*j*番目の用語に対する出現度情報を取得する。出現度情報取得手段1152は、例えば、*j*番目の用語の、読み出されている1以上のウェブページにおける出現頻度を取得する。また、出現度情報取得手段1152は、例えば、*j*番目の用語の、読み出されている1以上のウェブページにおける出現割合(出現頻度/全ワード数)を取得する。また、出現度情報取得手段1152は、例えば、*j*番目の用語についての $tf \cdot idf$ を取得する。 $tf \cdot idf$ は、 tf (用語の出現頻度)と idf (逆出現頻度)の二つの指標で計算される。 $tf \cdot idf$ は、公知技術であるので詳細な説明を省略する。

50

【0072】

(ステップS407)出現度情報取得手段1152は、ステップS406で取得された出現度情報を、バッファに一時格納する。

【0073】

(ステップS408)出現度情報取得手段1152は、カウンタjを1、インクリメントする。ステップS405に戻る。

【0074】

(ステップS409)出現度情報取得手段1152は、ステップS407でバッファに一時格納された1以上の出現度情報を用いて、i番目の分野情報についてのスコアを算出し、当該算出結果を、i番目の分野情報と対応付けて、メモリ上に一時格納する。出現度情報取得手段1152は、例えば、ステップS407でバッファに一時格納された1以上の出現度情報の和、または、平均値を算出し、i番目の分野情報についてのスコアとする。

10

【0075】

(ステップS410)出現度情報取得手段1152は、カウンタiを1、インクリメントする。ステップS402に戻る。

【0076】

(ステップS411)分野情報決定手段1153は、ステップS409で算出した、i番目の分野情報についてのスコアのうちで、最も大きなスコアに対応する分野情報を取得し、メモリ上に配置する。上位処理にリターンする。

20

【0077】

なお、図4のフローチャートにおいて、分野情報ごとにスコアを算出するアルゴリズムが異なっても良い。例えば、分野が「ニュース」(分野情報が、例えば「1」)の場合、ウェブページのタイトルで使われている用語の中に「ニュース」が存在すれば、スコアを非常に大きな数にしても良い。かかる場合、ウェブページのタイトルで使われている用語の中に「ニュース」が存在すれば、当該ウェブページの分野は、「ニュース」である、と判断されることを示す。また、例えば、分野が「ブログ」(分野情報が、例えば「2」)の場合、「天気」「晴れ」「曇り」「雨」などの用語の出現頻度を取得し、出現頻度の合計をスコアとしても良い。かかる場合、用語格納手段1151に、例えば、分野情報と、スコアを算出するアルゴリズムの情報(例えば、スコア算出処理を行うプログラム名や関数名など)が対で格納されている。

30

【0078】

次に、ステップS306のページ情報取得処理の詳細例について、図5のフローチャートを用いて説明する。

【0079】

(ステップS501)主情報取得手段1162は、ステップS304で読み込まれた1以上のウェブページの各々に対して、ウェブページの構造またはタグまたは内容を用いて、主情報を取得する。主情報を取得するアルゴリズムは種々あり、公知技術である。また、その例は、上述した。

【0080】

(ステップS502)ページ情報取得手段1163は、ステップS305で取得された分野情報と対になるページ取得ルール情報を、ページ取得ルール情報格納手段1161から読み出す。

40

【0081】

(ステップS503)ページ情報取得手段1163は、ステップS502で読み出したページ取得ルール情報を用いて、1以上のウェブページ(正確には、ここでは、ステップS501で取得された主情報)から、ページ単位の情報であるページ情報を、2以上取得する。

【0082】

なお、図5のフローチャートのステップS503における処理は、ページ取得ルール情

50

報に従った種々の処理が考えられる。その処理例については後述する。

【0083】

また、図5のフローチャートにおいて、ステップS501の主情報取得処理は、必須ではない。

【0084】

また、図5のフローチャートのステップS502において、ウェブページ識別情報と対になるページ取得ルール情報を、ページ取得ルール情報格納手段1161から読み出しても良い。

【0085】

さらに、図5のフローチャートにおいて、上記の情報処理装置11の処理において、ページ情報を構成する元の情報をすべて取得してから、レイアウトした。しかし、一つのまとまりのある情報(サブ情報)を取得した後、レイアウトし、その後さらに、次のサブ情報を取得しても良い。つまり、サブ情報ごとに、取得とレイアウトを繰り返すアルゴリズムでも良い。かかる場合、レイアウトに応じた、また、空きスペースに応じたサブ情報の取得が可能となる。

10

【0086】

次に、サーバ装置12の動作について説明する。サーバ装置12のウェブページ送信部122は、情報処理装置11からの要求に対応して、ウェブページ記憶部121から1以上のウェブページを読み出し、情報処理装置11に送信する。

【0087】

以下、本実施の形態における情報処理システム1の具体的な動作について説明する。情報処理システム1の概念図は図1である。

20

【0088】

今、情報処理装置11のレイアウト情報格納部114には、図6に示すレイアウト情報管理表(レイアウト情報管理表A)と、図7に示すレイアウト情報管理表(レイアウト情報管理表B)が格納されている。

【0089】

レイアウト情報管理表Aは、分野情報がニュースのウェブページに対応するレイアウト情報である。レイアウト情報管理表Aは、「ID」「主情報の種類」「ページ数」「領域」「文字サイズ」の属性値を有するレコードを1以上格納している。「ID」は、レコードを識別する属性であり、表管理のために存在する。「主情報の種類」は、ページ情報取得部116が取得した主情報の種類を示す情報が格納される。「ページ数」は、レイアウトされた場合のページ番号に関する情報が格納される。「領域」は、ページ内の領域を示す情報が格納される。「領域」内の情報のタイプには、ここでは、2種類ある。一つは、座標情報の集合である。「領域」の値が座標情報の集合である場合、領域は、それらを結ぶ矩形であることを示す。他は、始点と幅、高さの情報で定義される。かかる場合、領域は、始点と幅、高さの情報により示される長方形である。また、対応する情報がテキストの場合、「文字サイズ」の属性値により、文字サイズが示される。なお、レイアウト情報格納部114には、別途、分野情報がニュースの場合に、例えば、段組数「3」、文字列は縦書き、など、他の情報も格納されている、とする。また、図6において、「ID=6、7」の領域の情報と、「ID=8,9」の領域の情報は、交互に適用され、2つの記事で、2ページ目以降の各ページを構成することとなる。さらに、図6において、「ID=6」の領域の情報「始点(0,0)」に下線が付与されている(フラグが付与されている、ことと同意義)のは、「始点(0,0)」は、前の記事がはみ出してきた場合に、下方にずれることを意味する。「ID=7」の領域の情報「(0,0)(0,297)」も同様である。かかることにより、後述する図8(b)の81の領域が出来る場合がある。図8の81の領域が出来る場合とは、図8(a)の記事(1ページ目の記事)が2ページ目にはみ出してきた場合である。

30

40

【0090】

レイアウト情報管理表Bは、分野がブログのウェブページに対応するレイアウト情報で

50

ある。レイアウト情報管理表 B は、「ID」「主情報の種類」「配置情報」「文字サイズ」「フォント」「文字色」の属性値を有するレコードが 1 以上格納される。「配置情報」は、ここでは、各主情報の種類ごとの情報の配置順序を示している。また、レイアウト情報格納部 114 には、別途、分野情報がブログの場合に、例えば、段組数「1」、文字列は縦書き、など、他の情報も格納されている、とする。

【0091】

なお、レイアウト情報管理表 A を用いて、ニュースのウェブページをレイアウトした場合、図 8 に示すようなレイアウトとなる。図 8 において、数値は、レイアウト情報管理表 A のレコードの「ID」に対応する。また、図 8 (a) は、1 ページ目のページ情報のレイアウトの概要図、図 8 (b) は、2 ページ目以降のページ情報のレイアウトの概要図である。

10

【0092】

また、レイアウト情報管理表 B を用いて、ブログのウェブページをレイアウトした場合、図 9 に示すようなレイアウトとなる。図 9 において、数値は、レイアウト情報管理表 B のレコードの「ID」に対応する。また、図 9 (a) は、1 ページ目のページ情報のレイアウトの概要図、図 9 (b) は、2 ページ目のページ情報のレイアウトの概要図である。なお、図 9 において、「1」「2」「3」に対応する情報が順に横に、縦書きで配置されることを示しており、そのデータ量に応じて、横(左)にずれていく、ことは言うまでもない。

【0093】

20

図 10、図 11 は、ページ取得ルール情報格納手段 1161 に格納されているページ取得ルール情報管理表である。ページ取得ルール情報管理表は、「分野情報」「ID」「ページ取得ルール」の属性値を有するレコードを 1 以上格納している。「ページ取得ルール」は、「ルール」と「主情報の種類」の属性を有する。「ルール」は、ここでは、自然言語で記載しているが、プログラミング言語で記載しても良いし、プログラムの実行モジュール名などでも良い。「ルール」がプログラムの実行モジュール名である場合、実行モジュールは、ページ取得ルール情報格納手段 1161 に格納されている、とする。

【0094】

図 10、図 11 におけるルールは、例えば、HTML や XML などのウェブページに記載されている情報の中の、特定のタグや、属性や、タグまたは属性に対応する値が条件となっている。また、例えば、HTML や XML などのウェブページに記載されている情報の中の、アンカーを辿っていき、さらに下位のウェブページを取得して、当該下位のウェブページから、特定のタグや、属性や、タグまたは属性に対応する値を条件として情報を取得するルールである。さらに、図 10、図 11 におけるルールは、例えば、情報が配置される領域の大きさと、配置される候補となるデータのサイズ(データ量)とを比較して、合致する(収まりが良い場合で、完全一致するとは限らない)候補のデータを取得するルールである。さらに、図 10、図 11 において、主情報の種類ごとに取得ルールを設けているが、主情報に対して一つのルールでも良い。

30

【0095】

また、用語格納手段 1151 は、図 12 に示す用語情報管理表を格納している。用語情報管理表は、「分野情報」「用語」「スコア」の属性値を有するレコードを 1 以上格納している。「用語」は、単なる単語と、出現領域(例えば、"title")を特定した単語と、用語のタイプ(例えば、「\$日付タイプ」)などがある。「\$日付タイプ」とは、「2008/1/29」「2008年1月29日」「1月29日(火)」「1/29(火)」「平成20年1月29日」などの日付の表記となる文字列を言う。図 12 において、「"ニュース" in title」のレコードは、<title>タグ内に「ニュース」という文字列を含む場合に、ヒットするレコードである。「スコア」とは、ヒットした場合(用語を含むなど)に、カウントされる数値である。

40

【0096】

かかる状況において、ニュースのウェブページから電子ブックのコンテンツを自動生成

50

する具体例 1 と、ブログのウェブページから電子ブックのコンテンツを自動生成する具体例 2 の、2 つの具体例について説明する。

(具体例 1)

【0097】

ニュースのウェブページから電子ブックのコンテンツを自動生成する具体例について説明する。まず、情報処理装置 11 の指示受付部 110 は、あるウェブページの出力指示をユーザから受け付けた、とする。すると、情報処理装置 11 は、サーバ装置 12 に対して、そのウェブページの送信を要求する。そして、サーバ装置 12 のウェブページ送信部 122 は、受け付けた要求に対応するウェブページを、ウェブページ記憶部 121 から読み出し、情報処理装置 11 に送信する。

10

【0098】

情報処理装置 11 のウェブページ受信部 112 は、要求したサイトのウェブページを受信する。そして、ウェブページ蓄積部 113 は、受信されたウェブページを、ウェブページ格納部 111 に蓄積する。なお、ウェブページ格納部 111 は、ここでは、主メモリでも良い。そして、情報処理装置 11 の図示しない処理部は、受信されたウェブページを解釈し、例えば、図 13 に示すようなサイトを画面上に表示する。図 13 において、131 は、電子ブックのコンテンツを自動生成するツールバーである。

【0099】

次に、ユーザは、ツールバー 131 をマウスで押下した、とする。すると、指示受付部 110 は、ページ情報出力指示を受け付ける。つまり、ツールバー 131 の押下は、ページ情報出力指示である。

20

【0100】

次に、分野情報決定部 115 は、現在表示中のウェブページのスクリプトを取得する。かかる、スクリプト(ここでは、HTML で記述された文字列)を、図 14 に示す。

【0101】

分野情報決定部 115 は、図 14 の情報を用いて、以下のように、分野情報を取得する。つまり、出現度情報取得手段 1152 は、図 12 の用語情報管理表の分野情報「ニュース」を取得する。そして、出現度情報取得手段 1152 は、「ニュース」に対応する 1 番目の用語「トピックス」を取得し、図 14 の情報から、用語「トピックス」を検索し、用語「トピックス」の出現回数「1」を取得する。そして、出現度情報取得手段 1152 は、出現回数「1」と、用語情報管理表のスコア「1」を乗算し、「1」を、用語「トピックス」のスコアとして得る。同様に、出現度情報取得手段 1152 は、用語「経済」のスコアとして「1」を得る。また、同様に、用語「社会」のスコアとして「1」を得る。また、同様に、用語「スポーツ」のスコアとして「1」を得る。さらに、出現度情報取得手段 1152 は、用語「"ニュース" in title」から、<title>タグに対応する情報「x サイトニュース」を取得する。そして、出現度情報取得手段 1152 は、情報「x サイトニュース」に用語「ニュース」が含まれていると判断し、「"ニュース" in title」に合致するとする。そして、出現度情報取得手段 1152 は、「"ニュース" in title」と対になるスコア「10」を得る。以上の処理を繰り返し、出現度情報取得手段 1152 は、「用語」のすべての属性値に対するスコアを算出し、それらの和を取得する、とする。ここで、出現度情報取得手段 1152 は、分野情報「ニュース」に対するスコア「18」を得た、とする。

30

40

【0102】

次に、図 14 の情報について、分野情報「ブログ」に対するスコアを算出する。つまり、出現度情報取得手段 1152 は、用語「\$ 日付タイプ」を用いて、図 14 の情報の中に日付タイプの情報が存在するか否かを判断する。ここで、図 14 の情報の中に日付タイプの情報が存在するとして、出現度情報取得手段 1152 は、用語「\$ 日付タイプ」に対するスコア「2」を得る。次に、図 14 の情報の中に用語「俺」が存在するか否かを判断する。ここで、図 14 の情報の中に用語「俺」が存在しない、とする。そして、次に、図 14 の情報の中に用語「おれ」が存在するか否かを判断する。ここで、図 14 の情報の中に

50

用語「おれ」も存在しない、とする。用語が存在しない場合、加算されるスコアは「0」である。以上の処理を繰り返し、出現度情報取得手段1152は、分野情報「ブログ」に対するスコア「2」を得た、とする。

【0103】

そして、次に、分野情報決定手段1153は、上記で算出した、1番目、2番目の分野情報についてのスコアのうちで、最も大きなスコア「18」に対応する分野情報「ニュース」を取得し、メモリ上に配置する。以上の処理により、図14のウェブページの分野情報がニュースである、と決定された。

【0104】

次に、ページ情報取得部116は、図14のウェブページから、以下のように、ページ情報を取得する。つまり、ページ情報取得手段1163は取得された分野情報「ニュース」と対になるすべてのページ取得ルール情報を、図10、図11のページ取得ルール情報管理表から読み出す。そして、ページ情報取得手段1163は、以下のように、ページ取得ルール情報に従って、1以上のページ情報を取得する。

【0105】

まず、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=1」(1)のルール「class=tab?on」に対応するidの値を取得に従って、タグ「<li class="tab0 on">」を検索し、当該タグに対応するidの値"topics"を、タグ「」から取得する。

【0106】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=1」(2)のルールに従い、idの値"topics"を含む<div>内の4以上の項目「XX県知事にABC氏が当選へ」「サッカー 日本対ブラジル」「女優YYさん 結婚」「NY株 反発」・・・を、読み出す。

【0107】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=1」(3)のルールに従い、「」に合致する項目の文字列「XX県知事にABC氏が当選へ」を取得する。そして、ページ情報取得手段1163は、主情報の種類「メイン記事の見出し」とその情報「XX県知事にABC氏が当選へ」を対にして、バッファに一時格納する。

【0108】

次に、レイアウト情報取得手段1171は、取得された分野情報「ニュース」に対応するレイアウト情報(図6)を、レイアウト情報格納部114から読み出す。そして、その中で、主情報の種類「メイン記事の見出し」と対になる、ページ数「1」、領域「(0,0)(0,250)(40,0)(40,250)」、文字サイズ「24」を取得する。

【0109】

そして、レイアウト手段1172は、取得されたページ情報「XX県知事にABC氏が当選へ」を、取得されたレイアウト情報に従ってレイアウトし、ページ情報をメモリ上に得る。ここで、メイン記事の見出しは、図15の151のようにレイアウトされる、こととなる。

【0110】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=2」(1)のルール「メイン記事の見出しに対応するアンカーを取得」に対応するアンカーの値「f/topics/top/1/?1201573096」を取得する。

【0111】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=2」

10

20

30

40

50

(2)のルール「アンカーの先のウェブページを取得」に従って、当該アンカー「f/topics/top/1/?1201573096」が示すウェブページを、サーバ装置12から取得する。

【0112】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=2」(3)のルール「[記事全文]に対応するウェブページを取得」に従って、[記事全文]に対応するアンカーが示すウェブページを、サーバ装置12から取得する。

【0113】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=2」(4)のルール「ウェブページ中の<div>に対応する文字列を取得」に従って、取得したウェブページから、記事の本文となる文字列「任期満了に伴うXX県知事選は27日投票、即日開票の結果、無所属ABC氏が・・・・・・XYZ氏は、あと少しのところまで届かなかった。」を取得し、主情報の種類「メイン記事の内容」と対応付けて、バッファに格納する。

10

【0114】

次に、レイアウト情報取得手段1171は、取得された分野情報「ニュース」であり、主情報の種類「メイン記事の内容」と対になる、ページ数「1」、領域「(40,0)(40,250)(0,250)(0,297)(110,297)(110,0)」、文字サイズ「12」を取得する。

【0115】

そして、レイアウト手段1172は、取得されたページ情報「任期満了に伴うXX県知事選は27日投票、即日開票の結果、無所属ABC氏が・・・・・・XYZ氏は、あと少しのところまで届かなかった。」を、取得されたレイアウト情報に従ってレイアウトし、ページ情報をメモリ上に得る。ここで、メイン記事の見出しは、図15の152のようにレイアウトされる、こととなる。

20

【0116】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=3」(1)のルール「のgifデータを取得」に従って、「<img src=http://photo.gif」を用いて、「photo.gif」を、サーバ装置12から取得する。そして、ページ情報取得手段1163は、画像「photo.gif」を、主情報の種類「メイン画像」と対応付けて、バッファに格納する。

30

【0117】

次に、レイアウト情報取得手段1171は、取得された分野情報「ニュース」であり、主情報の種類「メイン画像」と対になる、ページ数「1」、領域「(110,230)(110,297)(210,230)(210,297)」を取得する。

【0118】

そして、レイアウト手段1172は、取得された画像「photo.gif」を、1ページ目の領域「(110,230)(110,297)(210,230)(210,297)」に配置する。ここで、メイン画像は、図15の153のようにレイアウトされる、こととなる。

【0119】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=4」のルールに従って、ページ情報の1ページ目の残る領域に最も合致する(近い)項目を決定する。つまり、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=4」の(1)から(5)により、既に出力している項目「XX県知事にABC氏が当選へ」を除き、3つ以上の項目「サッカー 日本対ブラジル」「女優YYさん 結婚」「NY株 反発」・・・・、から、それぞれの項目のアンカーの先のウェブページの

40

50

記事の情報 (<div> に対応する文字列) を、それぞれ取得する。

【 0 1 2 0 】

次に、分野情報「ニュース」に対応する「ID = 4」の(6)に従い、ページ情報取得手段1163は、各項目に対応する情報のデータ量(例えば、文字数)を取得する。そして、ページ情報取得手段1163は、1ページ目の残る領域のサイズを取得する。残る領域のサイズは、ページのサイズ(例えば、A4、など)と、既に配置されている情報のデータ量から、取得可能である。次に、ページ情報取得手段1163は、「ID = 4」(6)のルールに従い、1ページ目の残りの領域サイズに最も近い項目の文字列(ここでは、例えば、「女優YYさん 結婚」)を取得する。次に、ページ情報取得手段1163は、文字列「女優YYさん 結婚」と、主情報の種類「サブ記事1の見出し」とを対にして、バッファに格納する。

10

【 0 1 2 1 】

次に、レイアウト情報取得手段1171は、取得された分野情報「ニュース」であり、主情報の種類「サブ記事1の見出し」と対になる、ページ数「1」、領域「(110, 0) (110, 230) (140, 230) (140, 0)」、文字サイズ「20」を取得する。

【 0 1 2 2 】

そして、レイアウト手段1172は、取得された文字列「女優YYさん 結婚」を、1ページ目の領域「(110, 0) (110, 230) (140, 230) (140, 0)」に、文字サイズ「20」として、配置する。ここで、メイン画像は、図15の154のようにレイアウトされる、こととなる。

20

【 0 1 2 3 】

次に、分野情報「ニュース」に対応する「ID = 5」(1)に従い、ページ情報取得手段1163は、サブ記事1のタグ<div>に対応する文字列「女優YYさんは、.....」を、取得する。なお、この文字列「女優YYさんは、.....」は、項目「女優YYさん 結婚」のアンカー「f/topics/top/3/?1201574581」を用いて、サーバ装置12から取得されたウェブページの<div>に対応する文字列である。

30

【 0 1 2 4 】

以上の処理により、電子ブックの1ページ目のページ情報が構成できた。次に、情報処理装置11は、2ページ目以降のページ情報の構成を行う。

【 0 1 2 5 】

つまり、分野情報「ニュース」に対応する「ID = 6」(1)に従い、ページ情報取得手段1163は、他の項目「サッカー 日本対ブラジル」に対応する文字列「サッカー 日本対ブラジル」を取得する。そして、ページ情報取得手段1163は、文字列「サッカー 日本対ブラジル」と、主情報の種類「一般記事の見出し」を対にして、バッファに格納する。

【 0 1 2 6 】

次に、レイアウト情報取得手段1171は、取得された分野情報「ニュース」であり、主情報の種類「一般記事の見出し1」と対になる、ページ数「2~」、領域「始点(0, 0) (w, h) = (20, 150)」、文字サイズ「18」を取得する。

40

【 0 1 2 7 】

そして、レイアウト手段1172は、取得された文字列「サッカー 日本対ブラジル」を、2ページ目の「始点(0, 0) (w, h) = (20, 150)」で特定される領域に、文字サイズ「18」として、配置する。ここで、文字列「サッカー 日本対ブラジル」は、図16の161のようにレイアウトされる、こととなる。

【 0 1 2 8 】

次に、分野情報「ニュース」に対応する「ID = 7」(1)に従い、ページ情報取得手段1163は、一般記事の見出し「サッカー 日本対ブラジル」に対応するアンカー「f

50

/topics/top/2/?1201577113"」を取得する。そして、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=7」(2)に従い、当該アンカー「f/topics/top/2/?1201577113"」の先のウェブページをサーバ装置12から取得する。

【0129】

次に、ページ情報取得手段1163は、分野情報「ニュース」に対応する「ID=7」(3)に従い、当該ウェブページの<div>に対応する文字列「サッカー 日本対ブラジルは、29日19時から・・・」を取得する。そして、ページ情報取得手段1163は、文字列「サッカー 日本対ブラジルは、29日19時から・・・」と、主情報の種類「一般記事の内容」を対にして、バッファに格納する。

【0130】

次に、レイアウト情報取得手段1171は、取得された分野情報「ニュース」であり、主情報の種類「一般記事の内容」、特に、「一般記事の内容1」と対になる、ページ数「2~」、領域「(0,0)(0,297)(210,297)(210,198)(90,198)(90,0)」、文字サイズ「12」を取得する。

【0131】

そして、レイアウト手段1172は、取得された文字列「サッカー 日本対ブラジルは、29日19時から・・・」を、2ページ目の領域「(0,0)(0,297)(210,297)(210,198)(90,198)(90,0)」に、文字サイズ「12」として、配置する。ここで、文字列「サッカー 日本対ブラジルは、29日19時から・・・」は、図16の162のようにレイアウトされる、こととなる。

【0132】

同様に、ページ情報取得手段1163、レイアウト情報取得手段1171、レイアウト手段1172の処理により、タイトルの文字列「NY株 反発」、および文字列「週明けの米株式市場は、・・・」が、図16の163、164のように配置される。なお、文字列「NY株 反発」、および文字列「週明けの米株式市場は、・・・」の配置のために、図6の「ID=8」「ID=9」のレイアウト情報が用いられる。

【0133】

次に、ページ情報出力手段1173は、レイアウトし、構成された2以上のページ情報(図15、図16など)を出力する。ここでの出力は、例えば、図示しない電子ブック装置への送信である。また、出力は、例えば、情報処理装置11に装着されている着脱可能な記憶媒体である。

(具体例2)

【0134】

ブログから電子ブックのコンテンツを自動生成する具体例について説明する。まず、情報処理装置11の指示受付部110は、あるウェブページの出力指示をユーザから受け付けた、とする。すると、情報処理装置11は、サーバ装置12に対して、そのウェブページの送信を要求する。そして、サーバ装置12のウェブページ送信部122は、受け付けた要求に対応するウェブページを、ウェブページ記憶部121から読み出し、情報処理装置11に送信する。

【0135】

情報処理装置11のウェブページ受信部112は、要求したサイトのウェブページを受信する。そして、ウェブページ蓄積部113は、受信されたウェブページを、ウェブページ格納部111に蓄積する。なお、ウェブページ格納部111は、ここでは、主メモリでも良い。そして、情報処理装置11の図示しない処理部は、受信されたウェブページを解釈し、例えば、図17に示すようなサイトを画面上に表示する。

【0136】

次に、ユーザは、ツールバー171をマウスで押下した、とする。すると、指示受付部110は、ページ情報出力指示を受け付ける。つまり、ツールバー171の押下は、ページ情報出力指示である。また、ユーザは、ツールバー171の押下により、出力されるダイアログ(図示しない)に対して、分野情報「ブログ」を入力した、とする。なお、具体

10

20

30

40

50

例 1 においては、分野情報を自動決定したが、本例のように、分野情報が入力されても良い。そして、情報処理装置 11 の指示受付部 110 は、分野情報「ブログ」を受け付け、メモリ上に配置する。

【0137】

次に、ページ情報取得部 116 は、現在表示中のウェブページのスク립トを取得する。かかる、スク립ト（ここでは、HTML で記述された文字列）を、図 18 に示す。

【0138】

次に、ページ情報取得部 116 は、図 18 のウェブページから、以下のように、ページ情報を取得する。つまり、ページ情報取得手段 1163 は、入力された分野情報「ブログ」と対になるすべてのページ取得ルール情報を、図 11 のページ取得ルール情報管理表から読み出す。そして、ページ情報取得手段 1163 は、以下のように、ページ取得ルール情報に従って、1 以上のページ情報を取得する。

10

【0139】

まず、ページ情報取得手段 1163 は、分野情報「ブログ」に対応する「ID = 8」（1）のルール「<h2> タグに対応する文字列を取得」に従って、タグ「<h2>」を検索し、当該タグに対応する値「2008年1月27日(日)」を取得する。そして、ページ情報取得手段 1163 は、文字列「2008年1月27日(日)」と、主情報の種類「日付」とを対応付けて、バッファに格納する。

【0140】

次に、ページ情報取得手段 1163 は、分野情報「ブログ」に対応する「ID = 9」（1）のルール「<h3> タグに対応する文字列を取得」に従って、文字列「スキー場にて」を取得する。そして、ページ情報取得手段 1163 は、文字列「スキー場にて」と、主情報の種類「見出し」とを対応付けて、バッファに格納する。

20

【0141】

また、ページ情報取得手段 1163 は、分野情報「ブログ」に対応する「ID = 10」（1）のルール「<div class="entry-body-text"> タグに対応する文字列を取得」に従って、文字列「XXスキー場の雪は最高。 」を取得する。そして、文字列「XXスキー場の雪は最高。 」と、主情報の種類「本文」とを対応付けて、バッファに格納する。

【0142】

さらに、ページ情報取得手段 1163 は、分野情報「ブログ」に対応する「ID = 8」（1）のルール「<h2> タグに対応する文字列を取得」に従って、タグ「<h2>」を検索し、当該タグに対応する値「2008年1月26日(土)」を取得する。そして、ページ情報取得手段 1163 は、文字列「2008年1月26日(土)」と、主情報の種類「日付」とを対応付けて、バッファに格納する。

30

【0143】

次に、ページ情報取得手段 1163 は、分野情報「ブログ」に対応する「ID = 9」（1）のルール「<h3> タグに対応する文字列を取得」に従って、文字列「バスでの移動」を取得する。そして、文字列「バスでの移動」と、主情報の種類「見出し」とを対応付けて、バッファに格納する。

40

【0144】

次に、ページ情報取得手段 1163 は、分野情報「ブログ」に対応する「ID = 10」（1）のルール「<div class="entry-body-text"> タグに対応する文字列を取得」に従って、文字列「XXスキー場に行くために、バスで深夜から移動した。とても寒く、 」を取得する。そして、文字列「XXスキー場に行くために、バスで深夜から移動した。とても寒く、 」と、主情報の種類「本文」とを対応付けて、バッファに格納する。

【0145】

以上の処理を繰り返し、例えば、バッファ内に図 19 の情報を得る。なお、上記は、説明の簡素化のためにルールを、直ちに、ウェブページの情報に適用したが、通常、ページ

50

情報取得手段 1163 は、ウェブページの中のブロック（一つの区切り）単位の情報を取得し、当該情報に対して、順にルールを適用し、合致するルールに対応する主情報を、前記情報と対にして格納する。

【0146】

次に、レイアウト情報取得手段 1171 は、分野情報「ブログ」に対応するレイアウト情報（図7）を、レイアウト情報格納部 114 から読み出す。

【0147】

そして、レイアウト手段 1172 は、図7のレイアウト情報を用いて、日付、見出し、本文を順に配置し、かつ、それぞれの主情報の種類に応じた文字サイズ、フォント、文字色で、図19の情報を配置する。そして、レイアウト手段 1172 は、図20のページ情報を得て、メモリ上に配置する。そして、2ページ以降のページ情報も、1ページ目と同様に、日付、見出し、本文を順に配置する。

10

【0148】

以上、本実施の形態によれば、ウェブページから、複数ページの集合である電子ブックのコンテンツを自動構成することができる。

【0149】

また、本実施の形態によれば、ウェブページの内容に適したレイアウトで、電子ブックのコンテンツを自動構成することができるため、非常に読みやすいコンテンツを提供できる。例えば、ニュース記事のウェブページから、新聞風のレイアウトのコンテンツが構成でき、ブログのウェブページから、日記風のレイアウトのコンテンツが構成できる。

20

【0150】

また、本実施の形態によれば、ウェブページに記載されている情報のうち、重要と思われる情報のみを用いて、電子ブックのコンテンツを自動構成することができる。例えば、ウェブページに表示される広告などの情報を除いて、電子ブックのコンテンツを自動構成することができる。

【0151】

なお、本実施の形態によれば、レイアウト情報の定義方法、ページ取得ルール情報の定義情報は問わない。これらの情報は、データとして定義されていても良いし、プログラムにより実現されていても良い。かかることは、他の実施の形態においても同様である。

【0152】

また、本実施の形態によれば、主情報取得手段 1162 を有さなくても良い。かかることも、他の実施の形態においても同様である。

30

【0153】

また、本実施の形態において、ウェブページの分野やウェブページ特定情報に応じて、異なるレイアウトで、電子ブックのコンテンツを構成した。しかし、かかる分野情報やウェブページ特定情報を用いずに、常に同一のレイアウトで、電子ブックのコンテンツを構成しても良い。かかることも、他の実施の形態においても同様である。

【0154】

また、本実施の形態において、出力中のウェブページを元として、電子ブックのコンテンツを構成した。しかし、情報処理装置 11 が、電子ブックのコンテンツを構成する元となるウェブサイトの URL を 1 以上、記憶媒体に格納しており、例えば、予め設定された時刻（例えば、毎日 7:00）に、情報処理装置 11 が動作を開始するようにしても良い。かかる場合、情報処理装置 11 は、記憶媒体に格納された 1 以上の URL（ニュース記事が記載されたウェブページの URL）から取得されるウェブサイトから、電子ブックのコンテンツを構成し、例えば、電子ブックプレーヤーに送信する、などの処理を行っても良い。かかる処理を行うことにより、ユーザは、自分の電子ブックプレーヤーに、毎朝、出勤前に、ニュースの情報を新聞に類似したレイアウトのコンテンツを入手でき、通勤中に、見ることができる。かかることは、他の実施の形態においても同様である。

40

【0155】

さらに、本実施の形態における処理は、ソフトウェアで実現しても良い。そして、この

50

ソフトウェアをソフトウェアダウンロード等により配布しても良い。また、このソフトウェアをCD-ROMなどの記録媒体に記録して流布しても良い。なお、このことは、本明細書における他の実施の形態においても該当する。なお、本実施の形態における情報処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータを、1以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2以上取得するページ情報取得部と、前記ページ情報取得部が取得した2以上のページ情報を出力するページ情報出力部として機能させるためのプログラム、である。

【0156】

また、上記プログラムにおいて、前記ページ情報出力部は、ページ内のレイアウトに関する情報であり、前記1以上のウェブページに応じたレイアウト情報を、記憶媒体から取得するレイアウト情報取得手段と、前記レイアウト情報に従って、前記ページ情報取得部が取得した2以上のページ情報をレイアウトするレイアウト手段と、前記レイアウト手段がレイアウトした2以上のページ情報を出力するページ情報出力手段を具備するものとして、コンピュータを機能させることは好適である。

10

【0157】

また、上記プログラムにおいて、コンピュータを、前記1以上のウェブページの内容から、当該1以上のウェブページの分野を決定し、分野情報を取得する分野情報決定部をさらに具備し、前記レイアウト情報格納部は、分野を示す情報である分野情報と、ページ内のレイアウトに関する情報であるレイアウト情報との組を、2組以上格納しており、前記ページ情報出力部のレイアウト情報取得手段は、前記分野情報決定部が取得した分野情報と対になるレイアウト情報を取得するものとして機能させることは好適である。

20

【0158】

また、上記プログラムにおいて、前記ページ情報取得部は、前記1以上のウェブページから情報を取得するための情報であるページ取得ルール情報を、2以上格納しているページ取得ルール情報格納手段と、前記1以上のウェブページに対応するページ取得ルール情報、または前記分野情報決定部が取得した分野情報と対になるページ取得ルール情報を用いて、前記1以上のウェブページから、ページ単位の情報であるページ情報を、2以上取得するページ情報取得手段とを具備するものとして、コンピュータを機能させることは好適である。

30

【0159】

また、上記プログラムにおいて、前記分野情報決定部は、分野情報と、当該分野情報が示す分野に対応する1以上の用語を含む用語情報を1以上格納している用語格納手段と、前記1以上のウェブページから、前記1以上の用語の出現の度合いに関する情報である出現度情報を取得する出現度情報取得手段と、前記出現度情報を用いて、前記1以上のウェブページの分野を決定し、分野情報を取得する分野情報決定手段を具備するものとして、コンピュータを機能させることは好適である。

【0160】

また、上記プログラムにおいて、前記ページ情報取得部は、前記1以上のウェブページの各々に対して、ウェブページの構造またはタグまたは内容を用いて、主となる情報を取得する主情報取得手段と、前記主情報取得手段が取得した情報から、ページ情報を取得するページ情報取得手段を具備するものとして、コンピュータを機能させることは好適である。

40

【0161】

(実施の形態2)

本実施の形態において、1以上のウェブページから、ページ単位の情報に分割した電子ブックのコンテンツを得る情報処理システム等について説明する。また、本実施の形態において、ウェブページから分野(ニュース、子供向けなど)を決定し、当該分野に対応するテンプレートを取得し、当該テンプレートに従って、ウェブページを電子ブックのコンテンツに変換する情報処理システム等について説明する。また、ウェブページのURLに

50

よって、異なるテンプレートを取得し、当該テンプレートに従って、ウェブページを電子ブックのコンテンツに変換する情報処理システム等について説明する。また、本実施の形態において、ウェブページの分野やURLによって、ページ情報の取得方法が異なる情報処理システム等について説明する。さらに、本実施の形態において、特定のデータタイプ（例えば、静止画）のウェブページ内での割合から、ウェブページの分野やレイアウトなどを決定する情報処理システム等について説明する。なお、例えば、静止画の割合が多い場合、分野を「子供向け」と決定し、絵本のページ情報を自動生成する。

【0162】

本実施の形態における情報処理システム2の概念図は図1である。図21は、本実施の形態における情報処理システム2のブロック図である。情報処理システム2は、情報処理装置21、サーバ装置12を具備する。

10

【0163】

情報処理装置21は、指示受付部110、ウェブページ格納部111、ウェブページ受信部112、ウェブページ蓄積部113、レイアウト情報格納部114、分野情報決定部215、ページ情報取得部116、ページ情報出力部117、を具備する。

【0164】

分野情報決定部215は、データ情報取得手段2151、分野情報決定手段2152を具備する。

【0165】

データ情報取得手段2151は、1以上のウェブページから、1以上のデータタイプのデータ量またはデータサイズに関する情報であるデータ情報を取得する。データタイプとは、例えば、静止画、動画、テキスト、グラフィックス、ひらがなのテキスト、漢字のテキストなどである。データ情報とは、データの割合、データ数などである。なお、上記の、1以上のウェブページとは、ページ情報を構成する元になるウェブページであるが、例えば、実施の形態1で述べた主情報取得手段1162が処理した後に取得された主情報でも良い。データ情報取得手段2151は、通常、MPUやメモリ等から実現され得る。データ情報取得手段2151の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

20

【0166】

分野情報決定手段2152は、データ情報を用いて、1以上のウェブページの分野を決定し、分野情報を取得する。分野情報決定手段2152は、例えば、データタイプ、データ情報、分野情報を有する分野情報決定ルールを1以上保持しており、データ情報取得手段2151が取得したデータタイプとデータ情報に対応する分野情報を取得する。分野情報決定手段2152は、通常、MPUやメモリ等から実現され得る。分野情報決定手段2152の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

30

【0167】

次に、情報処理システム2の動作について、情報処理システム1の動作と比較すると、分野決定処理が異なる。そこで、情報処理システム2の分野決定処理について、図22のフローチャートを用いて説明する。なお、図22のフローチャートにおいて、図4のフローチャートと同一の処理について、説明を省略する。

40

【0168】

（ステップS2201）データ情報取得手段2151は、i番目の分野情報に対応するデータタイプの情報を取得する。データタイプの情報とは、例えば、データタイプが静止画なら「1」、動画なら「2」、ひらがなのテキストなら「3」などである。データタイプの情報は、i番目の分野情報に対応付けて、予め記憶媒体に格納されている。

【0169】

（ステップS2202）データ情報取得手段2151は、ステップS2201で取得したデータタイプの情報に対応するデータタイプのデータを、ページ情報の変換元となるウ

50

ウェブページから取得し、当該取得したデータから、データ情報（例えば、割合）を取得する。

【0170】

（ステップS2203）データ情報取得手段2151は、データ情報を用いてスコアを算出し、当該スコアを、バッファに一時格納する。

【0171】

以下、本実施の形態における情報処理システム2の具体的な動作について説明する。情報処理システム2の概念図は図1である。

【0172】

今、情報処理装置21のレイアウト情報格納部114には、図23に示すレイアウト情報管理表（レイアウト情報管理表C）が格納されている。レイアウト情報管理表Cは、分野情報「絵本」に対応するレイアウト情報の管理表である。また、レイアウト情報管理表Cは、「ID」「主情報の種類」「レイアウト情報」の属性値を有するレコードを1以上格納している。「レイアウト情報」は、ここでは「配置情報」「属性」を有する。「配置情報」は、ここでは、ページ情報内における、情報の配置を決定するための情報であり、「始点=(10,10)」は、(10,10)を始点として、対応する情報が配置されることを示す。「配置情報」が「全面」とは、対応する情報が、ページ全体に配置されることを示す。かかる場合、対応する情報は、拡大、または縮小され得る。「配置情報」が「下づめ」とは、対応する情報が、ページ内の、下側に詰めて配置されることを示す。「属性」は、文字の属性であり、ここでは、文字サイズのみを規定している。ただし、文字色やフォントなどの他の文字属性を規定しても良い。なお、レイアウト情報管理表Cには、規定されていないが、文字は、横書きで配置される、とする。また、ページ取得ルール情報格納手段1161は、図24に示すページ取得ルール情報管理表を格納している。

【0173】

さらに、分野情報決定部215は、「ウェブページ内のビットマップデータの割合が30%以上の場合は、分野情報が「絵本」である」と判断する、分野情報決定判断ルールを格納している、とする。

【0174】

かかる状況において、ユーザが、情報処理装置21に、図25のウェブページを表示させた、とする。図25のウェブページを表示させるための処理については、実施の形態1で述べた。

【0175】

次に、ユーザは、電子ブックのコンテンツを自動生成するツールバーを押下した、とする。すると、指示受付部110は、ページ情報出力指示を受け付ける。つまり、ツールバーの押下は、ページ情報出力指示である。

【0176】

次に、分野情報決定部215は、現在表示中のウェブページのスクリプトを取得する。かかる、スクリプト（ここでは、HTMLで記述された文字列）を、図26に示す。

【0177】

分野情報決定部215は、図26の情報を用いて、以下のように、分野情報を取得する。つまり、分野情報決定部215は、図26のスクリプトを解析し、ビットマップのサイズ(S1)を取得する。また、ページ全体のサイズ(S2)を取得する。そして、分野情報決定部215は、「 $S1/S2 \geq 0.3$ 」を満たすか否かを判断する。ここで、分野情報決定部215は、「 $S1/S2 \geq 0.3$ 」を満たす、と判断する。そして、分野情報決定部215は、上記の分野情報決定判断ルールを適用し、分野情報が「絵本」と判断する。

【0178】

次に、ページ情報取得部116は、図26のウェブページから、以下のように、ページ情報を取得する。つまり、ページ情報取得手段1163は取得された分野情報「絵本」と

10

20

30

40

50

対になるすべてのページ取得ルール情報を、図 2 4 のページ取得ルール情報管理表から読み出す。そして、ページ情報取得手段 1 1 6 3 は、以下のように、ページ取得ルール情報に従って、1 以上のページ情報を取得する。

【 0 1 7 9 】

まず、ページ情報取得手段 1 1 6 3 は、分野情報「絵本」に対応する「ID = 1」(1) のルール「<title> タグに対応する文字列を取得」に従って、タグ「<title>」を検索し、当該タグに対応する文字列「元気なポチくん」を取得する。そして、ページ情報取得手段 1 1 6 3 は、図 2 4 の主情報の種類「タイトル」と対応付けて、取得した文字列「元気なポチくん」を、バッファに格納する。

【 0 1 8 0 】

次に、ページ情報取得手段 1 1 6 3 は、分野情報「絵本」に対応する「ID = 1」(2) のルール「<frame src=" " > に従ってデータを取得」に従って、犬のビットマップデータを取得する。つまり、ページ情報取得手段 1 1 6 3 は、「src=" "」で指定されたデータ取得先からデータを取得する。そして、ページ情報取得手段 1 1 6 3 は、図 2 4 の主情報の種類「絵」と対応付けて、取得したビットマップデータを、バッファに格納する。

【 0 1 8 1 】

次に、ページ情報取得手段 1 1 6 3 は、分野情報「絵本」に対応する「ID = 1」(3) のルール「<body> タグ内の文字列を取得」に従って、文字列「うちの愛犬のポチは、ほんとうに元気です。 」を取得する。そして、ページ情報取得手段 1 1 6 3 は、図 2 4 の主情報の種類「文章」と対応付けて、取得した文字列を、バッファに格納する。以上の処理により、図 2 7 に示す情報が、バッファに格納された、こととなる。

【 0 1 8 2 】

次に、レイアウト情報取得手段 1 1 7 1 は、取得された分野情報「絵本」に対応するレイアウト情報(図 2 3) を、レイアウト情報格納部 1 1 4 から読み出す。

【 0 1 8 3 】

次に、レイアウト手段 1 1 7 2 は、バッファ上に取得された文字列「元気なポチくん」を、主情報の種類「タイトル」に従ったレイアウト(「始点 = (1 0 , 1 0)」「属性 = 1 4 p t」) で配置する。その結果、1 4 p t の文字サイズで、文字列「元気なポチくん」がページ内の左上に配置される。なお、ここでは、ページは、左上が原点(0 , 0) であり、右側に行けば x の値が + となり、下側に行けば y の値がプラスとなる、とする。

【 0 1 8 4 】

また、レイアウト手段 1 1 7 2 は、バッファ上のビットマップを、主情報の種類「絵」に従ったレイアウトで配置する。その結果、ページの全面に、ビットマップが拡大され、配置される。

【 0 1 8 5 】

さらに、レイアウト手段 1 1 7 2 は、バッファ上の文字列「うちの愛犬のポチは、ほんとうに元気です。 」を、主情報の種類「文章」に従ったレイアウト(下づめ「属性 = 2 0 p t」) で配置する。

【 0 1 8 6 】

また、レイアウト手段 1 1 7 2 は、ビットマップ(絵) が最背面にくるように、主情報を配置する。さらに、レイアウト手段 1 1 7 2 は、ビットマップの色に対して、文字が映える文字色を決定し、当該文字色に変換して、配置することは好適である。つまり、レイアウト手段 1 1 7 2 は、以下の文字色変換手段(図示しない) を有する。文字色変換手段は、文字列を配置する際に、文字列を配置する領域のビットマップの情報を取得する。そして、文字色変換手段は、取得したビットマップの領域の色の平均値を算出する。次に、文字色変換手段は、当該平均値を用いて、その補色となる色を算出する。次に、文字色変換手段は、補色の色に文字列の色属性を変更する。以上の処理により、見やすい絵本が自動生成できる。

10

20

30

40

50

【0187】

なお、上記の処理において、ページ情報取得手段1163が取得したビットマップは、文字列を有するビットマップである場合もある。かかる場合に対応して、ページ情報取得手段1163は、ビットマップを取得した際には、文字認識処理を行い、内部の文字列を取得することは好適である。そして、ページ情報取得手段1163は、取得できた文字列が一定以上の割合である場合、当該ビットマップを文字列と判断することは好適である。

【0188】

そして、ページ情報出力手段1173は、構成されたページ情報を入力する。以上の処理により、情報処理装置21は、図28に示すページ情報が出力される。図28に示すページ情報は、例えば、絵本の第1ページである。なお、ここでの出力態様は、上述したように、電子ブックプレーヤーへの送信、記憶媒体への蓄積など、種々、考えられる。

10

【0189】

以上、本実施の形態によれば、ウェブページから、複数ページの集合である電子ブックのコンテンツを自動構成することができる。

【0190】

また、本実施の形態によれば、ウェブページの内容に適したレイアウトで、電子ブックのコンテンツを自動構成することができるため、非常に読みやすいコンテンツを提供できる。例えば、子供向けのウェブページから、絵本のレイアウトのコンテンツが構成できる。

【0191】

20

また、本実施の形態によれば、ウェブページに記載されている情報のうち、重要と思われる情報のみを用いて、電子ブックのコンテンツを自動構成することができる。例えば、ウェブページに表示される広告などの情報を除いて、電子ブックのコンテンツを自動構成することができる。

【0192】

さらに、本実施の形態における情報処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータを、1以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2以上取得するページ情報取得部と、前記ページ情報取得部が取得した2以上のページ情報を出力するページ情報出力部として機能させるためのプログラム、である。

30

【0193】

また、上記プログラムにおいて、前記ページ情報出力部は、ページ内のレイアウトに関する情報であり、前記1以上のウェブページに応じたレイアウト情報を、記憶媒体から取得するレイアウト情報取得手段と、前記レイアウト情報に従って、前記ページ情報取得部が取得した2以上のページ情報をレイアウトするレイアウト手段と、前記レイアウト手段がレイアウトした2以上のページ情報を出力するページ情報出力手段を具備するものとして、コンピュータを機能させることは好適である。

【0194】

また、上記プログラムにおいて、前記分野情報決定部は、前記1以上のウェブページから、1以上のデータタイプのデータ量またはデータサイズに関する情報であるデータ情報を取得するデータ情報取得手段と、前記データ情報を用いて、前記1以上のウェブページの分野を決定し、分野情報を取得する分野情報決定手段を具備するものとして、コンピュータを機能させることは好適である。

40

【0195】

(実施の形態3)

本実施の形態において、1以上のウェブページから、ページ単位の情報に分割した電子ブックのコンテンツを得る情報処理システム等について説明する。また、本実施の形態において、複数のウェブページから情報を取得し、1以上のページ情報を構成する処理について説明する。かかる処理は、他のウェブページとの共通度合いに応じて、合成されるアルゴリズムが異なる場合についても説明する。

50

【 0 1 9 6 】

本実施の形態における情報処理システム3の概念図は図1である。図29は、本実施の形態における情報処理システム3のブロック図である。情報処理システム3は、情報処理装置31、サーバ装置12を具備する。

【 0 1 9 7 】

情報処理装置31は、ウェブページ格納部111、ウェブページ受信部112、ウェブページ蓄積部113、レイアウト情報格納部114、ページ情報取得部316、ページ情報出力部117、を具備する。

【 0 1 9 8 】

ページ情報取得部316は、ウェブページ取得手段3161、共通度取得手段3162、ページ情報取得手段3163を具備する。

10

【 0 1 9 9 】

ウェブページ取得手段3161は、ページ情報を構成する元になる、2以上のウェブページを、1以上のサーバ装置12から取得する。具体的には、例えば、ウェブページ取得手段3161は、ユーザから入力されたキーワードを取得する。そして、ウェブページ取得手段3161は、キーワードを用いて、関連キーワードを1以上、上位からm個取得する。なお、キーワードに対する関連キーワードを取得する処理は公知技術であるので、説明を省略する。次に、キーワードおよび1以上(m個)の関連キーワードをキーとして、ウェブページを検索する、かかる技術は、いわゆるウェブの検索エンジンの技術である。そして、次に、ウェブページ取得手段3161は、検索されたウェブページから、上位n個のURLを取得する。次に、ウェブページ取得手段3161は、上位n個のURLを用いて、n個のウェブページを取得する。ウェブページ取得手段3161は、通常、MPUやメモリ等から実現され得る。ウェブページ取得手段3161の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

20

【 0 2 0 0 】

共通度取得手段3162は、2以上のウェブページの各々から、内容の共通の度合いを示す情報である共通度を取得する。ここで、2以上のウェブページは、通常、ウェブページ取得手段3161が取得したウェブページであるが、予め決められていても良い。共通度取得手段3162は、さらに具体的には、2以上の各ウェブページのデータ量を算出し、データ量が最も大きいウェブページに対する、他の各ウェブページの共通度を取得しても良い。また、共通度取得手段3162は、データ量が最も大きいウェブページに1以上のウェブページが挿入された後の、挿入対象のウェブページに対する、他の各ウェブページの共通度を取得しても良い。かかる場合、「データ量が最も大きいウェブページ」には、挿入対象のウェブページも含む、とする。なお、共通度とは、2つのウェブページの記載内容の共通度合いに関する情報である。共通度は、例えば、2つのウェブページに、共通して出現する用語(単語)の数、2つのウェブページに、共通して出現する用語(単語)の、全用語に対する割合などである。その他、共通度の算出方法は何でも良い。共通度取得手段3162は、通常、MPUやメモリ等から実現され得る。共通度取得手段3162の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

30

40

【 0 2 0 1 】

ページ情報取得手段3163は、2以上のウェブページから2以上のページ情報を取得する。なお、ページ情報取得手段3163は、共通度に応じて、取得アルゴリズムが異なることは好適である。例えば、ページ情報取得手段3163は、他の各ウェブページの共通度が予め格納されている第一の閾値より小さい場合には、データ量が最も大きいウェブページに対して、前記他の各ウェブページを結合する。なお、「結合」とは、データ量が最も大きいウェブページの後ろ、または前(通常、後ろ)に、他の1以上のウェブページを追記することを言う。この1以上の他のウェブページは、データ量が最も大きいウェブページに対する、共通度が予め格納されている第一の閾値より小さいウェブページである

50

。そして、ページ情報取得手段3163は、共通度が第一の閾値より大きい場合には、各ウェブページを構成する一まとまりの情報であるブロックを取得し、データ量が最も大きいウェブページ（挿入対象のウェブページも含む。）のブロック毎に、他の各ウェブページのブロックのうちで、予め格納されている第二の閾値より大きい共通度を有するブロックを抽出し、当該抽出したブロックを、データ量が最も大きいウェブページの前記ブロック（共通度が大きいブロック）と次のブロックの間に挿入し、出力する全ページ情報を構成する。そして、ページ情報取得手段3163は、全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得する。なお、ページ情報取得手段3163は、1ページのデータ量（行数など）を保持している、とする。

【0202】

さらに、ページ情報取得手段3163は、以下のような処理を加えても良い。つまり、ページ情報取得手段3163は、他の各ウェブページのブロック中で、データ量が最も大きいウェブページの全ブロックのいずれに対しても、第二の閾値より大きな閾値を有しないブロックについて、その前に配置されているブロックで、データ量が最も大きいウェブページのいずれかのブロックに対して前記第二の閾値より大きな閾値を有するとして、前記データ量が最も大きいウェブページ内に挿入された箇所の直後に挿入する処理を行う。そして、ページ情報取得手段3163は、出力する全ページ情報を構成し、全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得する。ページ情報取得手段3163は、通常、MPUやメモリ等から実現され得る。ページ情報取得手段3163の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア（専用回路）で実現しても良い。

【0203】

次に、情報処理システム3の動作について図30のフローチャートを用いて説明する。図30のフローチャートにおいて、図3のフローチャートと同一の処理について、説明を省略する。

【0204】

（ステップS3001）ページ情報取得部316は、ページ情報を構成する元になる、2以上のウェブページを、1以上のサーバ装置12から取得する。かかるウェブページ取得処理の動作例の詳細について、図31のフローチャートを用いて説明する。

【0205】

（ステップS3002）ページ情報取得部316は、ページ取得処理を行う。ページ取得処理の動作例の詳細について、図32のフローチャートを用いて説明する。

【0206】

なお、図30のフローチャートにおいて、電源オフや処理終了の割り込みにより処理は終了する。

【0207】

次に、ステップS3001のウェブページ取得処理の動作例の詳細について、図31のフローチャートを用いて説明する。

【0208】

（ステップS3101）ウェブページ取得手段3161は、ユーザから入力されたキーワードを取得する。

【0209】

（ステップS3102）ウェブページ取得手段3161は、キーワードを用いて、関連キーワードを1以上、上位からm個取得する。なお、mは、1以上の自然数である。なお、関連キーワードを取得する処理は、概念辞書を用いて、キーワードの同義語などを取得しても良いし、過去にユーザが検索エンジンに入力したキーワードセットから取得しても良い。かかる関連キーワードを取得する処理は、公知技術であるので、詳細な説明を省略する。

【0210】

（ステップS3103）ウェブページ取得手段3161は、キーワードおよび1以上（

10

20

30

40

50

m個)の関連キーワードをキーとして、ウェブページを検索する。なお、ウェブページ取得手段3161は、例えば、公知の検索エンジンの構成を有する。

【0211】

(ステップS3104)ウェブページ取得手段3161は、ステップS3103で検索されたウェブページのURLの中で、上位n個のURLを取得する。上位n個のURLとは、いわゆる検索エンジンが出力した、上位n個のウェブページのURLである。

【0212】

(ステップS3105)ウェブページ取得手段3161は、ステップS3104で取得した上位n個のURLを用いて、n個のウェブページを取得する。上位処理にリターンする。

10

【0213】

次に、ステップS3002のページ取得処理の動作例の詳細について、図32のフローチャートを用いて説明する。

【0214】

(ステップS3201)共通度取得手段3162は、カウンタiに1を代入する。

【0215】

(ステップS3202)共通度取得手段3162は、i番目のウェブページが存在するか否かを判断する。i番目のウェブページが存在すればステップS3203に行き、i番目のウェブページが存在しなければステップS3208に行く。

【0216】

(ステップS3203)共通度取得手段3162は、i番目のウェブページのデータ量を算出し、メモリ上に配置する。

20

【0217】

(ステップS3204)共通度取得手段3162は、i番目のウェブページの全用語を取得する。共通度取得手段3162は、例えば、以下のようにウェブページの全用語を取得する。つまり、共通度取得手段3162は、例えば、ウェブページからタグ(例えば、<HTML>、<HEAD>、</HEAD>など)を消去し、制御文字(例えば、¥、"、{、}など)をスペースに置き換えるなどの処理をする。次に、共通度取得手段3162は、日本語の文字列を形態素解析し、名詞や動詞などの自立語のみを取得し、各形態素を用語とする。また、英文字列や数値列などは、区切り(スペース、リターンなど)ごとに用語として取得する。

30

【0218】

(ステップS3205)共通度取得手段3162は、ステップS3204で取得した1以上の用語の集合である用語群から、(用語、出現回数)の組を1組以上取得する。なお、共通度取得手段3162は、ステップS3204で取得した1以上の用語の集合をソートし、(用語、出現回数)の組を1組以上取得することは好適である。

【0219】

(ステップS3206)共通度取得手段3162は、i番目のウェブページをまとまりのある1以上のブロックに分割し、ブロックごとにメモリ上に配置する。なお、共通度取得手段3162は、例えば、ウェブページの区切りとなるタグ(例えば、<HR>、<DIV>、など)の情報を予め保持しており、当該区切りとなるタグを用いて(パターンマッチングして)、ブロックの区切りを認識し、ウェブページから1以上のブロックを取り出す。

40

【0220】

(ステップS3207)共通度取得手段3162は、カウンタiを1、インクリメントする。ステップS3202に戻る。

【0221】

(ステップS3208)ページ情報取得手段3163は、2以上のウェブページの中で、データ量が最大のウェブページを決定し、取得する。

【0222】

50

(ステップS3209) ページ情報取得手段3163は、カウンタjに1を代入する。
【0223】

(ステップS3210) ページ情報取得手段3163は、j番目の他のウェブページが存在するか否かを判断する。j番目の他のウェブページが存在すればステップS3211に行き、j番目の他のウェブページが存在しなければステップS3216に行く。なお、他のウェブページとは、データ量が最大のウェブページを除くウェブページのことである。

【0224】

(ステップS3211) ページ情報取得手段3163は、データ量が最大のウェブページの(用語、出現回数)の組の集合と、j番目の他のウェブページの(用語、出現回数)の組の集合とを用いて、2つのウェブページの共通度を取得する。ここで、共通度は、共通する用語の数でも良いし、共通する用語に対応する出現回数の和(1つの共通する用語に対して、2つの出現回数が加算される)でも良いし、共通する用語の全用語数に対する割合(用語の種類だけで見ても良いし、用語の出現回数を考慮しても良い)などでも良い。なお、データ量が最大のウェブページは、挿入対象のウェブページである、ことは好適である。挿入対象のウェブページは、元のデータ量が最大のウェブページに対して、1以上のウェブページの内容が挿入されたウェブページである。

10

【0225】

(ステップS3212) ページ情報取得手段3163は、共通度が予め格納されている閾値(第一の閾値)より大きいかなんかを判断する。共通度が第一の閾値より大きい場合ステップS3213に行き、大きくない場合ステップS3215に行く。

20

【0226】

(ステップS3213) ページ情報取得手段3163は、ウェブページの合成処理を行う。合成処理について、図33のフローチャートを用いて説明する。

【0227】

(ステップS3214) ページ情報取得手段3163は、カウンタjを1、インクリメントする。ステップS3210に戻る。

【0228】

(ステップS3215) ページ情報取得手段3163は、j番目の他のウェブページを、挿入対象のウェブページの最後に付加する。ステップS3214に行く。

30

【0229】

(ステップS3216) ページ情報取得手段3163は、構成したページ情報に対して、1ページごとに分割する処理を行う。分割処理とは、単に、1ページごとに区切る(例えば、改ページのコードを挿入など)ことである。分割処理は、いわゆるワードプロセッサのページ単位に区切る処理と同様でも良い。上位処理にリターンする。

【0230】

次に、ステップS3213の合成処理について、図33のフローチャートを用いて説明する。

【0231】

(ステップS3301) ページ情報取得手段3163は、カウンタiに1を代入する。

40

【0232】

(ステップS3302) ページ情報取得手段3163は、ブロックの情報の挿入対象のウェブページについて、i番目のブロックが存在するか否かを判断する。i番目のブロックが存在すればステップS3303に行き、i番目のブロックが存在しなければステップS3313に行く。なお、ここで、ブロックの情報の挿入対象のウェブページとは、元の最大のデータ量のウェブページに対して、1以上の他のウェブページのブロックが挿入されたウェブページ、または元の最大のデータ量のウェブページである。

【0233】

(ステップS3303) ページ情報取得手段3163は、ブロックの情報の挿入対象のウェブページの中の、i番目のブロックの情報を取得する。

50

【 0 2 3 4 】

(ステップ S 3 3 0 4) ページ情報取得手段 3 1 6 3 は、カウンタ j に 1 を代入する。

【 0 2 3 5 】

(ステップ S 3 3 0 5) ページ情報取得手段 3 1 6 3 は、処理対象の他のウェブページについて、 j 番目のブロックが存在するか否かを判断する。 j 番目のブロックが存在すればステップ S 3 3 0 6 に行き、 j 番目のブロックが存在しなければステップ S 3 3 1 2 に行く。

【 0 2 3 6 】

(ステップ S 3 3 0 6) ページ情報取得手段 3 1 6 3 は他のウェブページの中の、 j 番目のブロックの情報を取得する。

10

【 0 2 3 7 】

(ステップ S 3 3 0 7) ページ情報取得手段 3 1 6 3 は、ステップ S 3 3 0 3 で取得したブロックの情報と、ステップ S 3 3 0 6 で取得したブロックの情報とを比較し、共通度を取得する。なお、この共通度は、上述したウェブページの共通度を取得する処理と同様でも良い。また、他のアルゴリズムで共通度を取得しても良い。

【 0 2 3 8 】

(ステップ S 3 3 0 8) ページ情報取得手段 3 1 6 3 は、ステップ S 3 3 0 7 で取得した共通度が閾値(第二の閾値)より大きいか否かを判断する。大きい場合はステップ S 3 3 0 9 に行き、大きくない場合はステップ S 3 3 1 1 に行く。

【 0 2 3 9 】

(ステップ S 3 3 0 9) ページ情報取得手段 3 1 6 3 は、他のウェブページの j 番目のブロックの情報を、挿入対象のウェブページの i 番目のブロックと $(i + 1)$ 番目のブロックの間に挿入する。

20

【 0 2 4 0 】

(ステップ S 3 3 1 0) ページ情報取得手段 3 1 6 3 は、他のウェブページから j 番目のブロックの情報を削除する。

【 0 2 4 1 】

(ステップ S 3 3 1 1) ページ情報取得手段 3 1 6 3 は、カウンタ j を 1、インクリメントする。ステップ S 3 3 0 5 に戻る。

【 0 2 4 2 】

(ステップ S 3 3 1 2) ページ情報取得手段 3 1 6 3 は、カウンタ i を 1、インクリメントする。ステップ S 3 3 0 2 に戻る。

30

【 0 2 4 3 】

(ステップ S 3 3 1 3) ページ情報取得手段 3 1 6 3 は、他のウェブページのブロックのうち、挿入されずに残ったブロックの情報を、挿入対象のウェブページの最後に挿入する。上位処理にリターンする。

【 0 2 4 4 】

以下、本実施の形態における情報処理システム 3 の具体的な動作について説明する。情報処理システム 3 の概念図は図 1 である。

【 0 2 4 5 】

ユーザは、`Latex` の `tabular` 環境について勉強するために、情報処理装置 3 1 に対して、「`Latex tabular`」のキーワードを有するページ情報出力指示を入力した、とする。すると、ページ情報取得部 3 1 6 は、以下のように 3 つのウェブページを取得する。

40

【 0 2 4 6 】

つまり、ウェブページ取得手段 3 1 6 1 は、ユーザから入力されたキーワード「`Latex tabular`」を取得する。次に、ウェブページ取得手段 3 1 6 1 は、キーワードを用いて、関連キーワードを検索し、上位から 2 個「`tex latex` コマンド」を取得する。そして、ウェブページ取得手段 3 1 6 1 は、キーワードおよび 1 以上 (m 個) の関連キーワード「`Latex tabular tex latex` コマンド」をキー

50

として、ウェブページを検索する。

【0247】

次に、ウェブページ取得手段3161は、検索されたウェブページのURLの中で、上位3個のURLを取得する。そして、ウェブページ取得手段3161は、取得した上位3つのURLを用いて、3つのウェブページを取得する。かかるウェブページを、図34の(a)から(c)に示す。かかるウェブページは、出力されている場合のイメージである。

【0248】

次に、ページ情報取得部316は、以下のように3つのウェブページを合成する、ページ取得処理を行う。

【0249】

まず、共通度取得手段3162は、1番目のウェブページ(図34(a))について、データ量「10065byte」を算出する。次に、共通度取得手段3162は、1番目のウェブページの全用語(「表組」「TeX」「表」・・・「表組」・・・「tabular」・・・)を取得する。

【0250】

次に、共通度取得手段3162は、取得した1以上の用語の集合である用語群から、(用語、出現回数)の組を1組以上取得する。ここでは、共通度取得手段3162は、(表組, 2)(TeX, 3)(tabular, 34)(表, 17)・・・、を取得した、とする。

【0251】

次に、共通度取得手段3162は、1番目のウェブページ(図34(a))をまとまりのある1以上のブロックに分割し、ブロックごとにメモリ上に配置する。ここで、共通度取得手段3162は、ウェブページの区切りとなるタグ(例えば、<HR>、<DIV>、など)の情報を予め保持しており、当該区切りとなるタグを用いて(パターンマッチングして)、ブロックの区切りを認識し、ウェブページから1以上のブロックを取り出す。

【0252】

同様に、共通度取得手段3162は、2番目のウェブページ(図34(b))について、データ量「2355byte」を算出する。次に、共通度取得手段3162は、2番目のウェブページの全用語(「LaTeX」「以下」・・・「開発」・・・「tabular」・・・)を取得する。

【0253】

次に、共通度取得手段3162は、取得した1以上の用語の集合である用語群から、(用語、出現回数)の組、(LaTeX, 1)(以下, 1)・・・(tabular, 2)・・・、を取得した、とする。

【0254】

次に、共通度取得手段3162は、2番目のウェブページ(図34(b))をまとまりのある1以上のブロックに分割し、ブロックごとにメモリ上に配置する。

【0255】

同様に、共通度取得手段3162は、3番目のウェブページ(図34(c))について、データ量「7522byte」を算出する。次に、共通度取得手段3162は、2番目のウェブページの全用語(「shortstack」・・・「tabular」・・・)を取得する。

【0256】

次に、共通度取得手段3162は、取得した1以上の用語の集合である用語群から、(用語、出現回数)の組、(shortstack, 1)・・・(tabular, 32)・・・、を取得した、とする。

【0257】

次に、共通度取得手段3162は、3番目のウェブページ(図34(c))をまとまりのある1以上のブロックに分割し、ブロックごとにメモリ上に配置する。

【0258】

10

20

30

40

50

ページ情報取得手段 3 1 6 3 は、2 以上のウェブページの中で、データ量が最大のウェブページを 1 番目のウェブページ (図 3 4 (a)) と決定し、取得する。

【 0 2 5 9 】

次に、ページ情報取得手段 3 1 6 3 は、1 番目の他のウェブページ (図 3 4 (c)) を取得する。なお、ここでは、他のウェブページについて、データ量が多い方から処理する、とする。

【 0 2 6 0 】

次に、ページ情報取得手段 3 1 6 3 は、データ量が最大のウェブページの (用語、出現回数) の組の集合と、1 番目の他のウェブページ (図 3 4 (c)) の (用語、出現回数) の組の集合とを用いて、2 つのウェブページの共通度を、「1 9 8」と取得した、とする。10

【 0 2 6 1 】

次に、ページ情報取得手段 3 1 6 3 は、共通度が予め格納されている閾値 (第一の閾値) 「1 0 0」より大きいかなかを判断する。ここで、2 つのウェブページの共通度「1 9 8」は、第一の閾値「1 0 0」より大きいので、ページ情報取得手段 3 1 6 3 は、図 3 3 のフローチャートを用いて説明したように、1 番目の他のウェブページ (図 3 4 (c)) の各ブロックを、共通度が第二の閾値より大きい、データ量が最大のウェブページのブロックの直後に挿入していく。

【 0 2 6 2 】

そして、ページ情報取得手段 3 1 6 3 は、1 番目の他のウェブページ (図 3 4 (c)) のブロックのうち、データ量が最大のウェブページに挿入されなかったブロックを、データ量が最大のウェブページの最後に挿入していく。20

【 0 2 6 3 】

次に、ページ情報取得手段 3 1 6 3 は、以下のように、2 番目の他のウェブページ (図 3 4 (b)) の合成処理を行う。つまり、ページ情報取得手段 3 1 6 3 は、挿入対象のウェブページの (用語、出現回数) の組の集合と、2 番目の他のウェブページ (図 3 4 (b)) の (用語、出現回数) の組の集合とを用いて、2 つのウェブページの共通度を、「6 8」と取得した、とする。

【 0 2 6 4 】

次に、ページ情報取得手段 3 1 6 3 は、共通度が予め格納されている閾値 (第一の閾値) 「1 0 0」より大きいかなかを判断する。ここで、2 つのウェブページの共通度「6 8」は、第一の閾値「1 0 0」より小さいので、ページ情報取得手段 3 1 6 3 は、図 3 2 のフローチャートを用いて説明したように、1 番目の他のウェブページ (図 3 4 (c)) の各ブロックを、挿入対象のウェブページ (図 3 4 (a) と (c) のブロックが合成された後のウェブページ) の最下位に付加する。30

【 0 2 6 5 】

以上の処理により、3 つのウェブページが合成された。合成後のイメージ図は、図 3 5 である。図 3 5 において、「1」は図 3 4 (a) のウェブページを構成する 1 以上のブロック、「2」は図 3 4 (c) のウェブページを構成する 1 以上のブロック、「3」は図 3 4 (b) のウェブページ (全ブロック) を示す。そして、図 3 6 は、合成後のウェブページである。40

【 0 2 6 6 】

次に、図 3 6 のウェブページに対して、ページ情報取得手段 3 1 6 3 は、1 ページごとに分割する処理を行う。ページ情報取得手段 3 1 6 3 は、例えば、HTML で記述された図 3 6 のウェブページをフラットな text ファイル (ビットマップなどを含み得る) にして、ページ単位に区切る処理を行う。そして、ページ情報取得手段 3 1 6 3 は、2 以上のページ情報を得て、メモリ上に配置する。

【 0 2 6 7 】

次に、レイアウト情報取得手段 1 1 7 1 は、レイアウト情報を、レイアウト情報格納部 1 1 4 から読み出す。そして、レイアウト手段 1 1 7 2 は、取得されたページ情報 (図 3 50

6の情報)を、読み出されたレイアウト情報に従ってレイアウトし、1以上のページ情報をメモリ上に得る。次に、ページ情報出力手段1173は、メモリ上の1以上のページ情報を出力する。ここでの出力は、例えば、電子ブックプレーヤーへの送信や、着脱可能な可搬型の記憶媒体への蓄積などである。

【0268】

以上、本実施の形態によれば、複数のウェブページの内容を用いて、電子ブックのコンテンツを自動構成することができる。また、本実施の形態によれば、複数のウェブページの内容を用いて、近似する内容の情報が近いところに配置され、整理された電子ブックのコンテンツを自動構成することができる。さらに、本実施の形態によれば、かかる構成により、ユーザが欲する情報を、複数のウェブページから、自動的に選択し、有用な電子ブックのコンテンツを自動構成することができる。

10

【0269】

なお、本実施の形態における情報処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータを、1以上のウェブページから、電子ブックのページ単位の情報であるページ情報を、2以上取得するページ情報取得部と、前記ページ情報取得部が取得した2以上のページ情報を出力するページ情報出力部として機能させるためのプログラム、である。

【0270】

また、上記プログラムにおいて、前記ページ情報出力部は、ページ内のレイアウトに関する情報であり、前記1以上のウェブページに応じたレイアウト情報を、記憶媒体から取得するレイアウト情報取得手段と、前記レイアウト情報に従って、前記ページ情報取得部が取得した2以上のページ情報をレイアウトするレイアウト手段と、前記レイアウト手段がレイアウトした2以上のページ情報を出力するページ情報出力手段を具備するものとして、コンピュータを機能させることは好適である。

20

【0271】

また、上記プログラムにおいて、前記ページ情報取得部は、2以上のウェブページの中の一のウェブページを決定し、当該一のウェブページと、他の1以上の各ウェブページの組のうち、1以上の組に対して、内容の共通の度合いを示す情報である共通度を、1以上、取得する共通度取得手段と、前記2以上のウェブページから2以上のページ情報を取得するページ情報取得手段を具備し、前記ページ情報取得手段は、前記共通度に応じて、前記他の1以上の各ウェブページからの情報の取得方法が異なるものとして、コンピュータを機能させるためのプログラム、であることは好適である。

30

【0272】

また、上記プログラムにおいて、前記共通度取得手段は、前記2以上の各ウェブページのデータ量を算出し、データ量が最も大きいウェブページに対する、他の各ウェブページの共通度を取得し、前記ページ情報取得手段は、前記共通度が予め格納されている第一の閾値より小さい場合には、データ量が最も大きいウェブページに対して、前記他のウェブページを結合し、前記共通度が前記第一の閾値より大きい場合には、各ウェブページを構成する一まとまりの情報であるブロックを取得し、前記データ量が最も大きいウェブページのブロック毎に、前記他のウェブページのブロックのうちで、予め格納されている第二の閾値より大きい共通度を有するブロックを抽出し、当該抽出したブロックを、前記データ量が最も大きいウェブページの前記ブロックと次のブロックの間に挿入し、出力する全ページ情報を構成し、前記全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得するものとして、コンピュータを機能させるためのプログラム、であることは好適である。

40

【0273】

また、上記プログラムにおいて、前記ページ情報取得手段は、前記共通度が前記第一の閾値より大きい場合、前記他のウェブページのブロック中で、前記データ量が最も大きいウェブページの全ブロックのいずれに対しても、前記第二の閾値より大きな閾値を有しないブロックについて、その前に配置されているブロックで、前記データ量が最も大きいウ

50

ウェブページのいずれかのブロックに対して前記第二の閾値より大きな閾値を有するとして、前記データ量が最も大きいウェブページ内に挿入された箇所の直後に挿入する処理をさらに行い、出力する全ページ情報を構成し、前記全ページ情報を、1ページ単位に区切り、2以上のページ情報を取得するものとして、コンピュータを機能させるためのプログラム、であることは好適である。

【0274】

また、上記プログラムにおいて、前記ページ情報取得部は、ユーザから受け付けたキーワードを用いて、1以上のサーバ装置12からウェブページを検索し、当該検索したウェブページから、上位のn(nは2以上の自然数)のウェブページを取得するものとして、コンピュータを機能させるためのプログラム、であることは好適である。

10

【0275】

また、図37は、本明細書で述べたプログラムを実行して、上述した実施の形態の情報処理装置等を実現するコンピュータの外観を示す。上述の実施の形態は、コンピュータハードウェア及びその上で実行されるコンピュータプログラムで実現され得る。図37は、このコンピュータシステム340の概観図であり、図38は、コンピュータシステム340の内部構成を示す図である。

【0276】

図37において、コンピュータシステム340は、FD(Floppy Disk(登録商標))ドライブ3411、CD-ROM(Compact Disk Read Only Memory)ドライブ3412を含むコンピュータ341と、キーボード342と、マウス343と、モニタ344とを含む。

20

【0277】

図38において、コンピュータ341は、FDドライブ3411、CD-ROMドライブ3412に加えて、CPU(Central Processing Unit)3413と、CPU3413、CD-ROMドライブ3412及びFDドライブ3411に接続されたバス3414と、ブートアッププログラム等のプログラムを記憶するためのROM3415と、CPU3413に接続され、アプリケーションプログラムの命令を一時的に記憶するとともに一時記憶空間を提供するためのRAM(Random Access Memory)3416と、アプリケーションプログラム、システムプログラム、及びデータを記憶するためのハードディスク3417とを含む。ここでは、図示しないが、コンピュータ341は、さらに、LANへの接続を提供するネットワークカードを含んでも良い。

30

【0278】

コンピュータシステム340に、上述した実施の形態の情報処理装置等の機能を実行させるプログラムは、CD-ROM3501、またはFD3502に記憶されて、CD-ROMドライブ3412またはFDドライブ3411に挿入され、さらにハードディスク3417に転送されても良い。これに代えて、プログラムは、図示しないネットワークを介してコンピュータ341に送信され、ハードディスク3417に記憶されても良い。プログラムは実行の際にRAM3416にロードされる。プログラムは、CD-ROM3501、FD3502またはネットワークから直接、ロードされても良い。

40

【0279】

プログラムは、コンピュータ341に、上述した実施の形態の情報処理装置等の機能を実行させるオペレーティングシステム(OS)、またはサードパーティープログラム等は、必ずしも含まなくても良い。プログラムは、制御された態様で適切な機能(モジュール)を呼び出し、所望の結果が得られるようにする命令の部分のみを含んでいれば良い。コンピュータシステム340がどのように動作するかは周知であり、詳細な説明は省略する。

【0280】

なお、上記プログラムにおいて、情報を送信する送信ステップや、情報を受信する受信ステップなどでは、ハードウェアによって行われる処理、例えば、送信ステップにおける

50

モデムやインターフェースカードなどで行われる処理（ハードウェアでしか行われない処理）は含まれない。

【0281】

また、上記プログラムを実行するコンピュータは、単数であってもよく、複数であってもよい。すなわち、集中処理を行ってもよく、あるいは分散処理を行ってもよい。

【0282】

また、上記各実施の形態において、一の装置に存在する2以上の通信手段（端末情報送信部、端末情報受信部など）は、物理的に一の媒体で実現されても良いことは言うまでもない。

【0283】

また、上記各実施の形態において、各処理（各機能）は、単一の装置（システム）によって集中処理されることによって実現されてもよく、あるいは、複数の装置によって分散処理されることによって実現されてもよい。

【0284】

本発明は、以上の実施の形態に限定されることなく、種々の変更が可能であり、それらも本発明の範囲内に包含されるものであることは言うまでもない。

【産業上の利用可能性】

【0285】

以上のように、本発明にかかる情報処理システムは、ウェブページから、複数ページの集合である電子ブックのコンテンツを自動構成することができる、という効果を有し、情報処理装置、電子ブックプレーヤー等として有用である。

【図面の簡単な説明】

【0286】

【図1】実施の形態1における情報処理システム概念図

【図2】同情報処理システムのブロック図

【図3】同情報処理装置の動作について説明するフローチャート

【図4】同分野決定処理の動作について説明するフローチャート

【図5】同ページ情報取得処理の動作について説明するフローチャート

【図6】同レイアウト情報管理表を示す図

【図7】同レイアウト情報管理表を示す図

【図8】同レイアウトイメージを示す図

【図9】同レイアウトイメージを示す図

【図10】同ページ取得ルール情報管理表を示す図

【図11】同ページ取得ルール情報管理表を示す図

【図12】同用語情報管理表を示す図

【図13】同ウェブページの例を示す図

【図14】同ウェブページのスクリプトを示す図

【図15】同レイアウト結果を示す図

【図16】同レイアウト結果を示す図

【図17】同ウェブページの例を示す図

【図18】同ウェブページのスクリプトを示す図

【図19】同バッファ上のデータ例を示す図

【図20】同レイアウト結果を示す図

【図21】実施の形態2における情報処理システムのブロック図

【図22】同分野情報決定処理について説明するフローチャート

【図23】同レイアウト情報管理表を示す図

【図24】同ページ取得ルール情報管理表を示す図

【図25】同ウェブページの例を示す図

【図26】同ウェブページのスクリプトを示す図

【図27】同バッファ上のデータ例を示す図

10

20

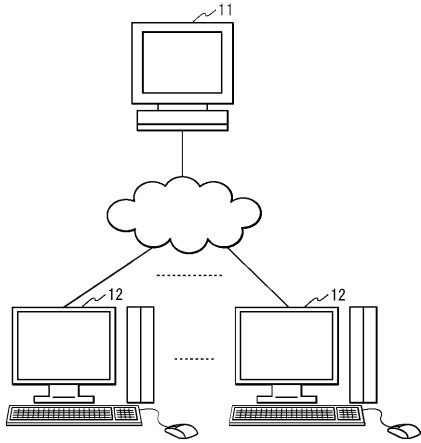
30

40

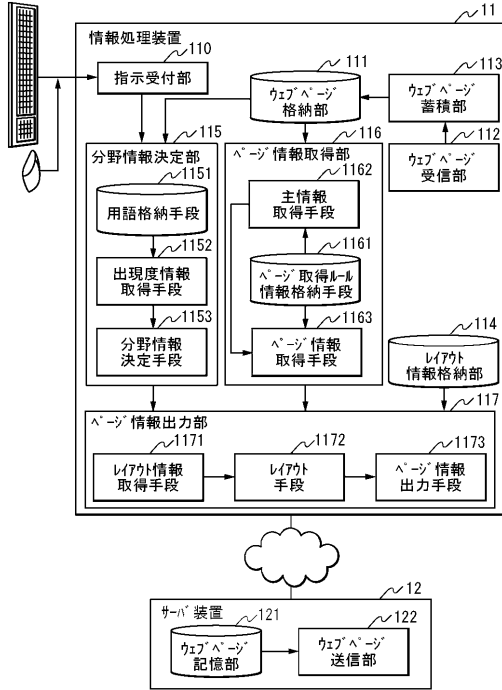
50

【図 2 8】	同レイアウト結果を示す図	
【図 2 9】	実施の形態 3 における情報処理システムのブロック図	
【図 3 0】	同情報処理システムの動作について説明するフローチャート	
【図 3 1】	同ウェブページ取得処理の動作について説明するフローチャート	
【図 3 2】	同ページ取得処理の動作について説明するフローチャート	
【図 3 3】	同合成処理の動作について説明するフローチャート	
【図 3 4】	同ウェブページの例を示す図	
【図 3 5】	同合成後のイメージ図	
【図 3 6】	同合成後のウェブページのイメージ図	
【図 3 7】	同コンピュータシステムの外観図	10
【図 3 8】	同コンピュータシステムのブロック図	
【符号の説明】		
【 0 2 8 7 】		
1、2、3	情報処理システム	
1 1、2 1、3 1	情報処理装置	
1 2	サーバ装置	
1 1 0	指示受付部	
1 1 1	ウェブページ格納部	
1 1 2	ウェブページ受信部	
1 1 3	ウェブページ蓄積部	20
1 1 4	レイアウト情報格納部	
1 1 5、2 1 5	分野情報決定部	
1 1 6、3 1 6	ページ情報取得部	
1 1 7	ページ情報出力部	
1 2 1	ウェブページ記憶部	
1 2 2	ウェブページ送信部	
1 1 5 1	用語格納手段	
1 1 5 2	出現度情報取得手段	
1 1 5 3、2 1 5 2	分野情報決定手段	
1 1 6 1	ページ取得ルール情報格納手段	30
1 1 6 2	主情報取得手段	
1 1 6 3、3 1 6 3	ページ情報取得手段	
1 1 7 1	レイアウト情報取得手段	
1 1 7 2	レイアウト手段	
1 1 7 3	ページ情報出力手段	
2 1 5 1	データ情報取得手段	
3 1 6 1	ウェブページ取得手段	
3 1 6 2	共通度取得手段	

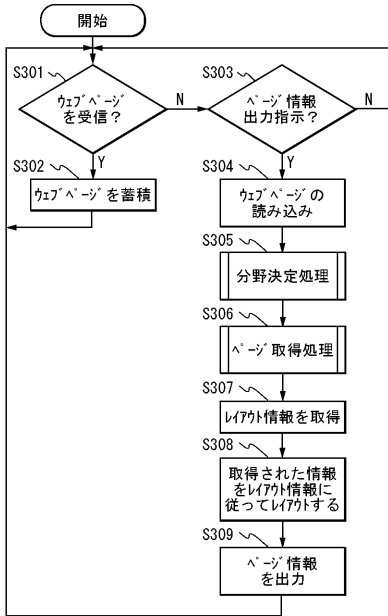
【図1】



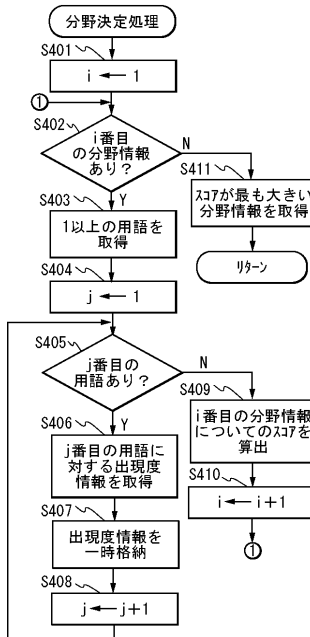
【図2】



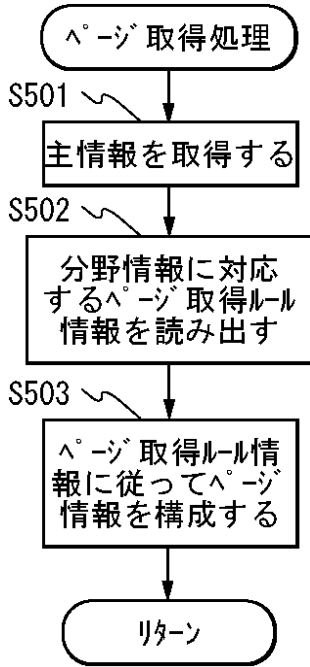
【図3】



【図4】



【図5】



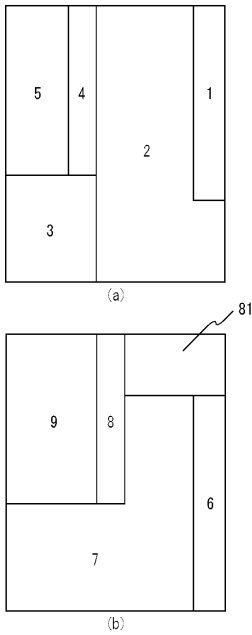
【図6】

ID	主情報の種類	ページ数	領域	文字サイズ
1	メイン記事の見出し	1	(0, 0) (0, 250) (40, 0) (40, 250)	24
2	メイン記事の内容	1	(40, 0) (40, 250) (0, 250) (0, 297) (110, 297) (110, 0)	12
3	メイン画像	1	(110, 230) (110, 297) (210, 230) (210, 297)	—
4	サブ記事1の見出し	1	(110, 0) (110, 230) (140, 230) (140, 0)	20
5	サブ記事1の内容	1	(140, 0) (140, 230) (210, 230) (210, 0)	12
6	一般記事の見出し1	2~	始点 (0, 0) (w, h) = (20, 150)	18
7	一般記事の内容1	2~	(0, 0) (0, 297) (210, 297) (210, 198) (90, 198) (90, 0)	12
8	一般記事の見出し2	2~	始点 (90, 0) (w, h) = (20, 198)	18
9	一般記事の内容2	2~	(110, 0) (110, 198) (210, 198) (210, 0)	12

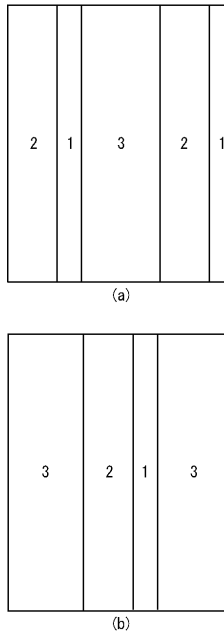
【図7】

ID	主情報の種類	配置情報	文字サイズ	フォント	文字色
1	日付	1	14	ゴシック	赤
2	見出し	2	28	明朝・太	青
3	本文	3	16	明朝	黒

【図8】



【図9】



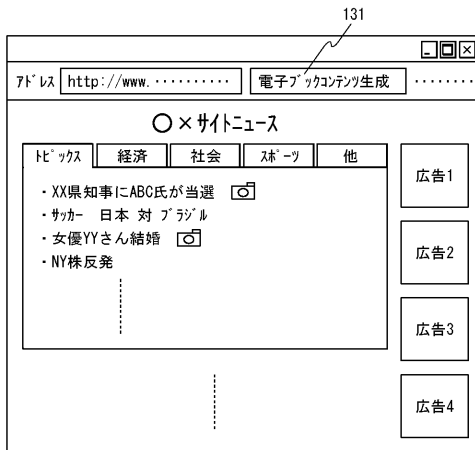
【図10】

分野情報	ID	ページ取得ルール	
		ルール	主情報の種類
ニュース	1	(1) class="tab? on"に対応するidの値を取得 (2) idの値の文字列を含む<div>内の項目に対応)を取得 (3) が付与されている文字列を取得	メイン記事の見出し
	2	(1) メイン記事の見出しに対応するアカーを取得 (2) アカーの先のウェブページを取得 (3) [記事全文]に対応するウェブページを取得 (4) ウェブページ中の<div>に対応する文字列を取得	メイン記事の内容
	3	(1) のgifデータを取得	メイン画像
	4	(1) class="tab? on"に対応するidの値を取得 (2) idの文字列を含む<div>内の項目であり、メイン記事の見出しを除く項目を取得 (3) 各項目に対応するアカーを取得 (4) アカーの先のウェブページを取得 (5) ウェブページ中の<div>に対応する文字列を取得 (6) 文字列中で、そのサイズが1ページ目の残り領域サイズに最も合致する項目(文字列)を取得	サブ記事1の見出し
	5	(1) サブ記事1の<div>に対応する文字列を取得	サブ記事1の内容

【図11】

分野情報	ID	ページ取得ルール	
		ルール	主情報の種類
ニュース	6	(1) 他の項目に対応する文字列を取得	一般記事の見出し
	7	(1) 一般記事の見出しに対応するアカーを取得 (2) アカーの先のウェブページを取得 (3) ウェブページ中の<div>に対応する文字列を取得	一般記事の内容
ブログ	8	(1) <h2>タグに対応する文字列を取得	日付
	9	(1) <h3>タグに対応する文字列を取得	見出し
	10	(1) <div class="entry-body-text" >タグに対応する文字列を取得	本文

【図13】



【図12】

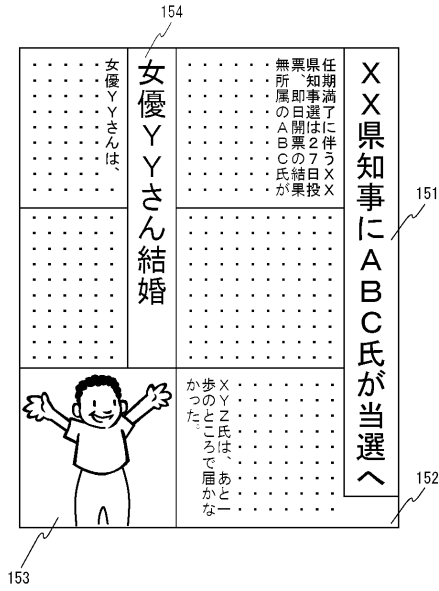
分野情報	用語	スコア
ニュース	トピックス	1
	経済	1
	社会	1
	スポーツ	1
	“ニュース” in title	10
	⋮	⋮
ブログ	\$日付タイプ	2
	俺	1
	おれ	1
	⋮	⋮

【図14】

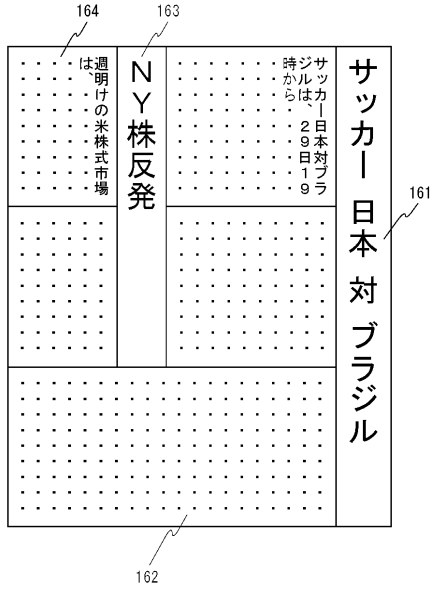
```

<html lang="ja">
<head>
:
:
<title>〇×サイトニュース</title>
:
<hr class="separate">
<div id="division">
<div id="main">
<div id="topicsbox" class="bx">
<div class="hd">
<ul class="tab on0">
<li class="tab0 on"><span><a id="topics" href="r/tp" hidefocus="true">
トピックス</a></span></li>
<li class="tab1"><span><a id="economy" href="r/tcco"
hidefocus="true">経済</a>< span>< li>
<li class="tab2"><span><a id="social" href="r/tent" hidefocus="true">
社会</a></span></li>
<li class="tab3"><span><a id="sports" href="r/tspo" hidefocus="true">
スポーツ</a></span></li>
<li class="tab4 last"><span><a id="others" href="r/toth"
hidefocus="true">他</a></span>< li>
</ul>
</div>
<div id="topicsboxbd">
<div id="topicsfb" class="current">
<div class="topicsindex">
<ul class="emphasis">
<li><a href="f/topics/top/1/?1201573096">XX県知事にABC氏が当選
</a>< img src="http://photo.gif" alt="[photo]" > </li>
<li><a href="f/topics/top/2/?1201577113">サッカー 日本対ブラジル
</a></li>
<li><a href="f/topics/top/3/?1201574581">女優YYさん 結婚</a>< li>
<li><a href="f/topics/top/7/?1201575515">NY株 反発</a></li>
:
:
</ul>
:
:
</html>
    
```

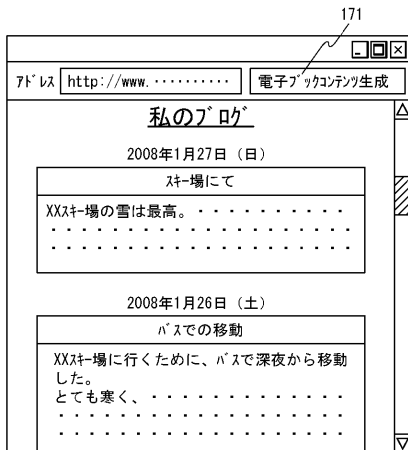
【図15】



【図16】



【図17】



【図18】

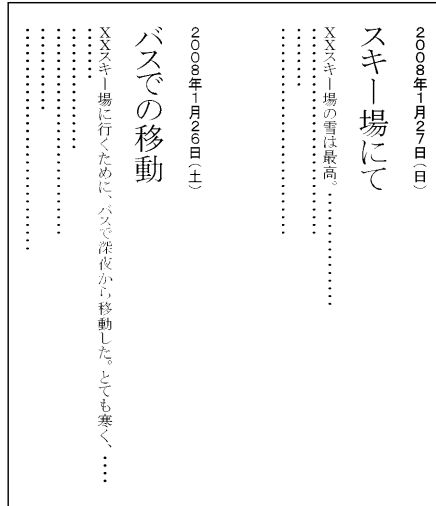
```

<html xmlns="*.....">
<head>
  <title>私のブログ日記</title>
  <h2>2008年1月27日(日)</h2>
  <h3>スキー場にて</h3>
  <div class="entry-body">
    <div class="entry-body-text">
      XXスキー場の雪は最高。.....
    </div>
  </div>
  <h2>2008年1月26日(土)</h2>
  <h3>バスでの移動</h3>
  <div class="entry-body">
    <div class="entry-body-text">
      XXスキー場に行くために、バスで深夜から移動した。とても寒く、...
    </div>
  </div>
</body>
</html>
    
```

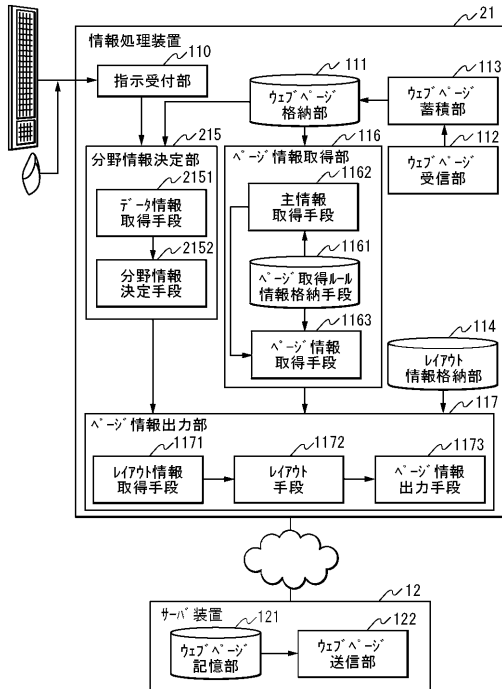
【図19】

ID	主情報の種類	データ
1	日付	2008年1月27日(日)
2	見出し	スキー場にて
3	本文	XXスキー場は最高.....
4	日付	2008年1月26日(土)
5	見出し	バスでの移動
6	本文	XXスキー場に行くために、 バスで深夜から移動した。 とても寒く、.....
.....

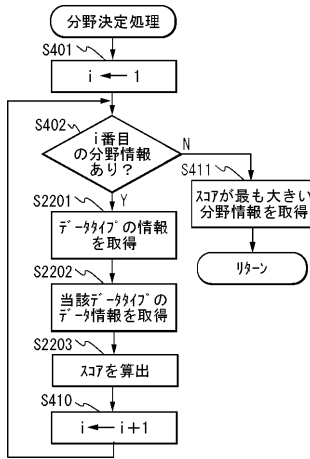
【図20】



【図21】



【図22】



【図23】

ID	主情報の種類	レイアウト情報	
		配置情報	属性
1	タイトル	始点 = (10, 10)	14pt
2	絵	全面	—
3	文章	下づめ	20pt
.....

【図24】

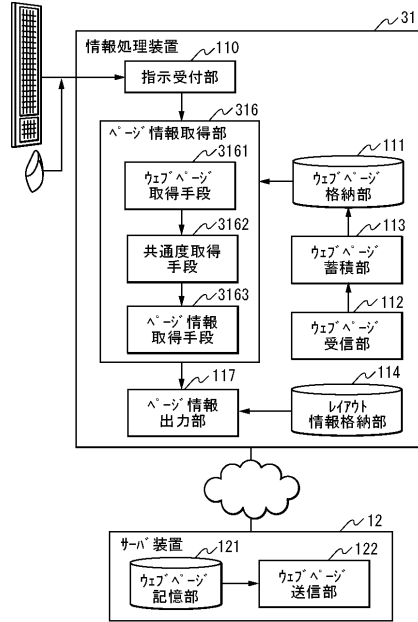
ID	分野情報	ページ取得ルール		
		ルール	主情報の種類	
1	絵本	(1)	<title>が</title>に対応する文字列を取得	タイトル
		(2)	<frame src=" " " >内のデータを取得	絵
		(3)	<body>が</body>内の文字列を取得	文章
...

【図26】

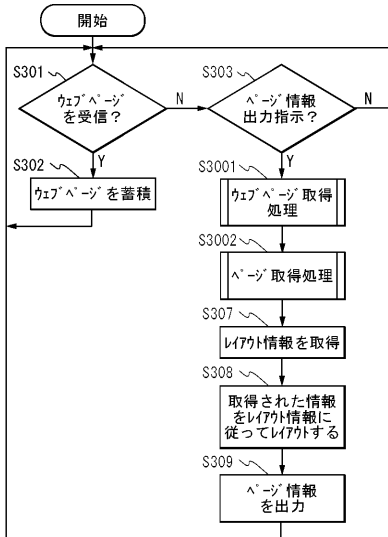
```

<html>
<head>
<title>元気なポチくん</title>
</head>
<meta http-equiv="Content-Type" content="text/html; charset=EUC-JP">
<frameset rows="48,*" frameborder="1" border="1" framespacing="1">
<frame
src="kghead.php?SY=2&MD=0&FM=0&BL=&TP=http://www.e-hon.jp/kushar/kusj1.htm" scrolling="no" marginwidth="0"
marginheight="0" name="kidsgoo_header">
<frame
src="kgbody.php?SY=2&MD=0&FM=0&BL=&TP=http://www.e-hon.jp/kushar/kusj1.htm" scrolling="auto" name="_kgbody">
</frameset>
<noframes>
<body bgcolor="#FFFFFF">
うちの愛犬のポチは、ほんとうに元気です。……………
</body>
</noframes>
</html>
    
```

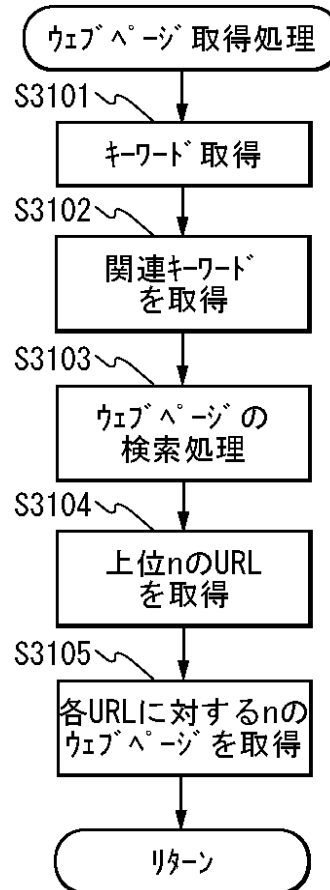
【図29】



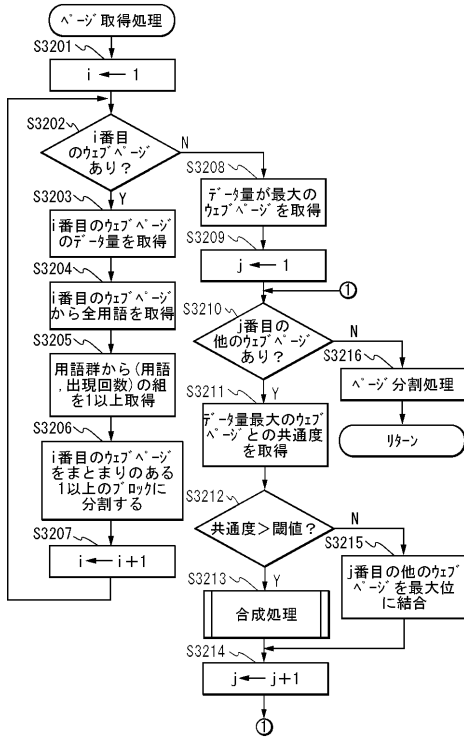
【図30】



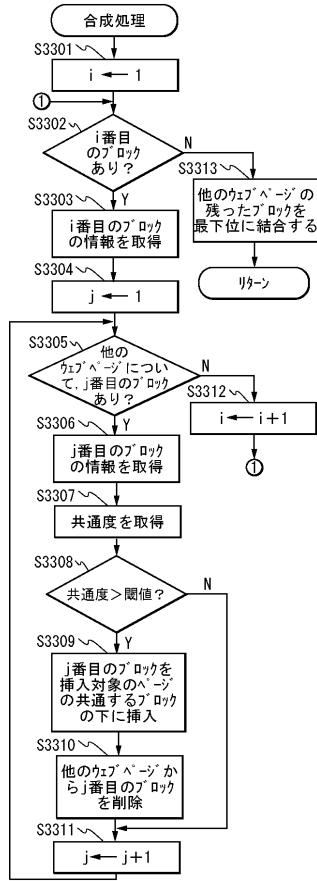
【図31】



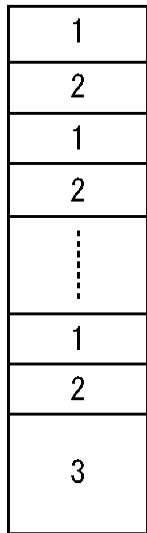
【図 3 2】



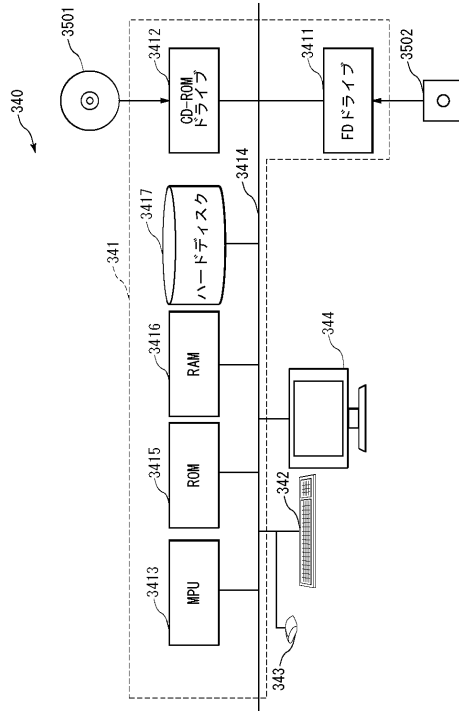
【図 3 3】



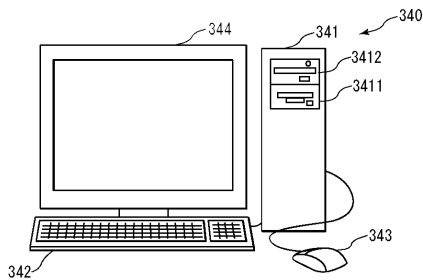
【図 3 5】



【図 3 8】




【図 3 7】



【図 25】



【図 27】

主情報の種類	情報
タイトル	元気なポチくん
絵	
文章	うちの愛犬のポチは、ほんとうに元気です。.....

【図 28】

元気なポチくん



うちの愛犬のポチは、ほんとうに元気
です。.....
.....。

【図 3 4】

:

表組

- `\begin{tabular}{l}`:
項目を左揃え。
- `\begin{tabular}{c}`:
項目を中央揃え。
- `\begin{tabular}{r}`:
項目を右寄せ。

:

:

(a)

:

LaTeXで以下のような表を作ろうと思ったとき、容易に作成できるソフトウェアを開発しました。

理論値	実験値	
	室内	室外
0	0	0.5
10	10.5	9.5
100	90	110

:

「VisualTabular」といいます。

(b)

:

`\shortstack` を使いまして

:

```

\begin{tabular}{rrr}
... & \shortstack{line1¥¥ line2}
& ...
\end{tabular>

```

:

:

:

(c)

【図 3 6】

：

表組

- ・`\begin{tabular}{l}`：
項目を左揃え。
- ・`\begin{tabular}{c}`：
項目を中央揃え。
- ・`\begin{tabular}{r}`：
項目を右寄せ。

`\shortstack` を使いまして

：

```
\begin{tabular}{rrr}
... & \shortstack{line1\\line2}
& ...
\end{tabular}
```

：

：

：

：

：

：

：

LaTeXで以下のような表を作ろうと思ったとき、容易に作成できるソフトウェアを開発しました。

理論値	実験値	
	室内	室外
0	0	0.5
10	10.5	9.5
100	90	110

：

「VisualTabular」といいます。

フロントページの続き

(72)発明者 安部 伸治

京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

審査官 木村 雅也

(56)参考文献 特開2006-279887(JP,A)

山根 康宏, まるごとNokia E61 初版, 日本, 株式会社技術評論社 片岡 巖, 2007年 9月10日, 第1版, 第186頁

(58)調査した分野(Int.Cl., DB名)

G06F 13/00

G06F 17/21