

(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2011-123740

(P2011-123740A)

(43) 公開日 平成23年6月23日(2011.6.23)

(51) Int.Cl.
G06F 13/00 (2006.01)F I
G06F 13/00 540Aテーマコード (参考)
5B084

審査請求 未請求 請求項の数 8 O L (全 14 頁)

(21) 出願番号 特願2009-281880 (P2009-281880)
(22) 出願日 平成21年12月11日 (2009.12.11)(71) 出願人 306037311
富士フイルム株式会社
東京都港区西麻布2丁目26番30号
(74) 代理人 100083116
弁理士 松浦 憲三
(72) 発明者 福島 敏貢
東京都港区赤坂9丁目7番3号 富士フイルム株式会社内
Fターム(参考) 5B084 AA12 AB04 AB06 CB12 CB22
DB02 DC02 DC03

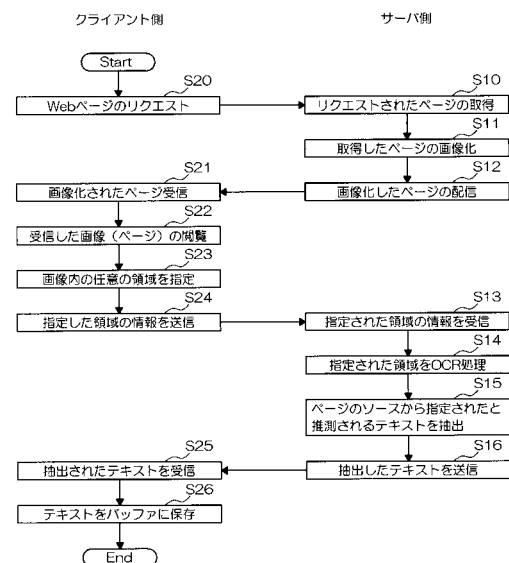
(54) 【発明の名称】 閲覧システム、サーバ、テキスト抽出方法及びプログラム

(57) 【要約】

【課題】画像化したウェブページを端末に送信し、端末装置でウェブページを閲覧する場合において、端末装置に表示された画像内の文字を正確に抽出することができる。

【解決手段】サーバ10は、インターネットからウェブページを取得し(ステップS10)、取得したウェブページから画像を生成し(ステップS11)、画像をクライアント端末20へ送信する(ステップS12)。クライアント端末20は、画像を受信し(ステップS21)、表示部23へ表示し(ステップS22)、矩形領域を指定し(ステップS23)、その情報をサーバ10へ送信する(ステップS24)。サーバ10は、画像から矩形領域の画像を切り出し、OCR処理によりテキストを認識し(ステップS14)、HTMLファイルのソースから認識されたテキストと最も一致度の高いテキストを抽出し(ステップS15)、クライアント端末20へ送信する(ステップS16)。

【選択図】 図4



【特許請求の範囲】

【請求項 1】

表示手段が設けられた端末装置と、前記端末装置と接続されたサーバとで構成された閲覧システムであって、

前記端末装置は、

前記サーバから送信された画像データを受信する端末側受信手段と、

前記受信された画像データに基づいて前記表示手段に画像を表示させる表示制御手段と

、前記表示手段に表示された画像の中の所定の領域を選択する選択手段と、

前記選択された所定の領域の情報を前記サーバへ送信する端末側送信手段と、を備え、

前記サーバは、

ウェブページのソースを取得する取得手段と、

前記取得されたウェブページのソースに基づいて当該ウェブページの画像データを生成する画像生成手段と、

前記生成された画像データを前記端末装置に送信するサーバ側送信手段と、

前記端末装置から送信された所定の領域の情報を受信するサーバ側受信手段と、

前記受信された所定の領域の情報と前記生成された画像データとに基づいて、前記所定の領域の画像からOCR処理により文字を認識する文字認識手段と、

前記OCR処理により認識された文字と推定される文字列を前記取得されたウェブページのソースから抽出する文字列抽出手段と、を備え、

前記サーバ側送信手段は、前記抽出された文字列を前記端末装置に送信し、

前記端末側受信手段は、前記送信された文字列を受信することを特徴とする閲覧システム。

【請求項 2】

前記サーバは、前記所定の領域が閾値以上であるか否かを判断する判断手段を備え、

前記所定の領域が閾値以上であると判断されなかった場合には、前記サーバ側送信手段は、前記OCR処理により認識された文字列を送信することを特徴とする請求項 1 に記載の閲覧システム。

【請求項 3】

前記端末側送信手段は、前記所定の領域の情報として当該所定の領域の座標の情報を前記サーバへ送信し、

前記文字認識手段は、前記生成された画像データと、前記所定の領域の座標の情報とから前記所定の領域の画像を切り出し、当該切り出された所定の領域の画像から文字を認識することを特徴とする請求項 1 又は 2 に記載の閲覧システム。

【請求項 4】

前記文字列抽出手段は、前記OCR処理により認識された文字をキーと前記取得されたソースに含まれるテキストとを比較し、前記OCR処理により認識された文字と最も一致度の高い文字列を抽出することを特徴とする請求項 1、2 又は 3 に記載の閲覧システム。

【請求項 5】

前記端末装置は、前記受信した文字列を記憶する記憶手段を備えたことを特徴とする請求項 1 から 4 のいずれかに記載の閲覧システム。

【請求項 6】

請求項 1 から 5 のいずれかに記載の閲覧システムを構成するサーバ。

【請求項 7】

携帯端末からウェブページの閲覧要求を受け付けるステップと、

前記受け付けられた閲覧要求に基づいてウェブページのソースを取得するステップと、

前記取得されたウェブページのソースに基づいて当該ウェブページの画像データを生成するステップと、

前記端末装置から所定の領域の情報を受信するステップと、

前記受信した所定の領域の情報と前記生成された画像データとに基づいて、前記所定の

10

20

30

40

50

領域の画像からOCR処理により文字を認識するステップと、

前記取得されたソースから前記OCR処理により認識された文字と推定される文字列を抽出するステップと、

前記抽出された文字列を前記端末装置に送信するステップと、
を含むことを特徴とするテキスト抽出方法。

【請求項 8】

請求項 7 に記載のテキスト抽出方法を演算装置に実行させることを特徴とするプログラム。

【発明の詳細な説明】

【技術分野】

10

【0001】

本発明は閲覧システム、サーバ、テキスト抽出方法及びプログラムに係り、特に携帯端末でウェブページが閲覧可能な閲覧システム、サーバ、テキスト抽出方法及びプログラムに関する。

【背景技術】

【0002】

近年、携帯電話にフルブラウザが搭載されることが多くなり、携帯電話からPC用のウェブページを閲覧することが可能となっている。しかしながら、携帯電話でPC用のウェブページを閲覧する場合には、画面が小さいため、ページのレイアウトが崩れて閲覧しづらい等といった問題が起こる場合がある。また、企業のイントラページなどは、安全性を確保するため、アクセスが制限され、携帯電話からは閲覧することができない。

20

【0003】

このような問題を解決するための方法として、サーバでウェブページやイントラページを画像化して携帯電話へ配信するというシステムが考えられる。

【0004】

引用文献 1 には、ウェブページをサーバ側でレンダリングし、画像に変換したページをクライアントに配信するシステムが記載されている。

【0005】

引用文献 2 には、クライアント装置のウェブブラウザからOCR処理の対象とする領域を指定し、サーバでOCR処理を行うシステムが記載されている。

30

【0006】

引用文献 3 には、画像データを文字認識（OCR（Optical Character Reader）処理）にかけ、テキストを抽出し、さらに抽出したテキストデータを構文意味解析処理にかけることにより文章のエラーを検出し、修正を行うことで文字（文章）の認識精度を高めるシステムが記載されている。

【先行技術文献】

【特許文献】

【0007】

【特許文献 1】特開 2004 - 220260 号公報

【特許文献 2】特開 2005 - 327258 号公報

【特許文献 3】特開 2006 - 350663 号公報

40

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかしながら、特許文献 1 に記載の発明では、クライアントに配信するウェブページは画像化されているため、テキスト領域を選択してコピーするといった操作ができなかった。

【0009】

特許文献 2 に記載の発明では、OCR処理により画像データからテキストデータを得ることはできるが、テキストデータの精度を向上させる方法については記載されていない。

50

【 0 0 1 0 】

特許文献 3 に記載の発明では、OCR 処理の精度が低い場合には、構文意味解析ができず、正しいテキストデータが得られないという問題がある。また、構文意味解析ができた場合であっても、得られたテキストデータが画像データに実際に含まれるテキストデータとならないという問題がある。

【 0 0 1 1 】

本発明はこのような事情に鑑みてなされたもので、画像化したウェブページを端末に送信し、端末装置でウェブページを閲覧する場合において、端末装置に表示された画像内の所定の領域に含まれる文字を正確に抽出することができる閲覧システム、サーバ、テキスト抽出方法及びプログラムを提供することを目的とする。

10

【課題を解決するための手段】

【 0 0 1 2 】

請求項 1 に記載の閲覧システムは、表示手段が設けられた端末装置と、前記端末装置と接続されたサーバとで構成された閲覧システムであって、前記端末装置は、前記サーバから送信された画像データを受信する端末側受信手段と、前記受信された画像データに基いて前記表示手段に画像を表示させる表示制御手段と、前記表示手段に表示された画像の中の所定の領域を選択する選択手段と、前記選択された所定の領域の情報を前記サーバへ送信する端末側送信手段と、を備え、前記サーバは、ウェブページのソースを取得する取得手段と、前記取得されたウェブページのソースに基づいて当該ウェブページの画像データを生成する画像生成手段と、前記生成された画像データを前記端末装置に送信するサーバ側送信手段と、前記端末装置から送信された所定の領域の情報を受信するサーバ側受信手段と、前記受信された所定の領域の情報と前記生成された画像データとに基づいて、前記所定の領域の画像から OCR 処理により文字を認識する文字認識手段と、前記 OCR 処理により認識された文字と推定される文字列を前記取得されたウェブページのソースから抽出する文字列抽出手段と、を備え、前記サーバ側送信手段は、前記抽出された文字列を前記端末装置に送信し、前記端末側受信手段は、前記送信された文字列を受信することを特徴とする。

20

【 0 0 1 3 】

請求項 1 に記載の閲覧システムによれば、サーバでは、ウェブページのソースが取得され、取得されたウェブページのソースに基づいて当該ウェブページの画像データが生成され、生成された画像データが端末装置に送信される。端末装置では、送信された画像データが受信され、受信された画像データに基づいて表示手段に画像が表示され、表示手段に表示された画像の中の所定の領域が選択され、選択された所定の領域の情報がサーバへ送信される。サーバでは、端末装置から送信された所定の領域の情報が受信され、受信された所定の領域の情報と生成された画像データとに基づいて所定の領域の画像から OCR 処理により文字が認識され、OCR 処理により認識された文字と推定される文字列が取得されたソースから抽出され、抽出された文字列が端末装置に送信される。携帯端末では、サーバから送信された文字列が受信される。これにより、OCR 処理のミスにより間違ったテキストが認識された場合においても、そのミスを補完し、選択した領域に含まれる正確なテキストデータを得ることができる。例えば、下線付き文字や表の一部等 OCR 処理の精度が低い場合においても、正確なテキストデータを得ることができる。

30

40

【 0 0 1 4 】

請求項 2 に記載の閲覧システムは、請求項 1 に記載の閲覧システムにおいて、前記サーバは、前記所定の領域が閾値以上であるか否かを判断する判断手段を備え、前記所定の領域が閾値以上であると判断されなかった場合には、前記サーバ側送信手段は、前記 OCR 処理により認識された文字列を送信することを特徴とする。

【 0 0 1 5 】

請求項 2 に記載の閲覧システムによれば、サーバでは、所定の領域が閾値以上であるか否かが判断され、所定の領域が閾値以上であると判断されなかった場合には、OCR 処理により認識された文字列が端末装置へ送信される。これにより、効率よく、かつ精度よく

50

選択した領域に含まれるテキストデータを得ることができる。

【0016】

請求項3に記載の閲覧システムは、請求項1又は2に記載の閲覧システムにおいて、前記端末側送信手段は、前記所定の領域の情報として当該所定の領域の座標の情報を前記サーバへ送信し、前記文字認識手段は、前記生成された画像データと、前記所定の領域の座標の情報とから前記所定の領域の画像を切り出し、当該切り出された所定の領域の画像から文字を認識することを特徴とする。

【0017】

請求項3に記載の閲覧システムによれば、所定の領域の情報として所定の領域の座標の情報が端末装置からサーバへ送信されると、サーバでは、生成された画像データと、所定の領域の座標の情報とから所定の領域の画像が切り出され、切り出された所定の領域の画像から文字が認識される。これにより、処理能力の高いサーバで重い処理、すなわち座標に従い指定された領域の画像を抽出する処理を行い、処理能力の低い端末装置で行う処理は、処理コストの小さい矩形領域の座標の送信のみとすることができる。

10

【0018】

請求項4に記載の閲覧システムは、請求項1、2又は3に記載の閲覧システムにおいて、前記文字列抽出手段は、前記OCR処理により認識された文字をキーと前記取得されたソースに含まれるテキストとを比較し、前記OCR処理により認識された文字と最も一致度の高い文字列を抽出することを特徴とする。

【0019】

請求項4に記載の閲覧システムによれば、文字列抽出手段では、OCR処理により認識された文字をキーと取得されたソースに含まれるテキストとが比較され、OCR処理により認識された文字と最も一致度の高い文字列が抽出される。これにより、ソースから選択した領域に含まれるテキストデータを抽出することができる。

20

【0020】

請求項5に記載の閲覧システムは、請求項1から4のいずれかに記載の閲覧システムにおいて、前記端末装置は、前記受信した文字列を記憶する記憶手段を備えたことを特徴とする。

【0021】

請求項5に記載の閲覧システムによれば、端末装置では、サーバから送信された文字列が記憶手段に記憶される。これにより、サーバから送信されたテキストを、任意のテキストフィールドへの貼り付けなどに利用することができる。すなわち、クライアント端末で選択された領域の画像に含まれるテキストのコピーと同等の効果を得ることができる。

30

【0022】

請求項6に記載のサーバは、請求項1から5のいずれかに記載の閲覧システムを構成する。

【0023】

請求項7に記載のテキスト抽出方法は、携帯端末からウェブページの閲覧要求を受け付けるステップと、前記受け付けられた閲覧要求に基づいてウェブページのソースを取得するステップと、前記取得されたウェブページのソースに基づいて当該ウェブページの画像データを生成するステップと、前記端末装置から所定の領域の情報を受信するステップと、前記受信した所定の領域の情報と前記生成された画像データとに基づいて、前記所定の領域の画像からOCR処理により文字を認識するステップと、前記取得されたソースから前記OCR処理により認識された文字と推定される文字列を抽出するステップと、前記抽出された文字列を前記端末装置に送信するステップと、を含むことを特徴とする。

40

【0024】

請求項8に記載のプログラムは、請求項7に記載のテキスト抽出方法を演算装置に実行させることを特徴とする。

【発明の効果】

【0025】

50

本発明によれば、画像化したウェブページを端末に送信し、端末装置でウェブページを閲覧する場合において、端末装置に表示された画像内の所定の領域に含まれる文字を正確に抽出することができる。

【図面の簡単な説明】

【0026】

【図1】本発明が適用された閲覧システム1の概略図である。

【図2】閲覧システム1を構成するサーバの概略図である。

【図3】閲覧システム1を構成するクライアント端末の概略図である。

【図4】閲覧システム1のクライアント端末がテキストデータをコピーする取得する処理の流れを示すフローチャートである。

10

【図5】クライアント端末に表示される閲覧用画像の一例である。

【図6】OCR処理を説明するための図である。

【図7】テキスト抽出処理を説明するための図である。

【図8】一致度が最も高いテキストを抽出する方法を説明するための図である。

【図9】テキスト送信処理を説明するための図である。

【図10】本発明が適用された閲覧システム2のクライアント端末がテキストデータをコピーする取得する処理の流れを示すフローチャートである。

【図11】閲覧システム2のテキスト抽出処理について説明するための図である。

【発明を実施するための形態】

【0027】

20

< 第1の実施の形態 >

閲覧システム1は、主として、サーバ10と、クライアント端末20とで構成される。サーバ10と接続されるクライアント端末20は1台でも良いし、複数でもよい。

【0028】

サーバ10は、図2に示すように、主として、CPU11と、データ取得部12と、画像生成部13と、OCR処理部14と、テキスト抽出部15と、通信部16とで構成される。

【0029】

CPU11は、サーバ10の全体の動作を統括制御する制御手段として機能するとともに、各種の演算処理を行う演算手段として機能する。CPU11は、制御プログラムであるファームウェア、ウェブページを表示するためのプログラムであるブラウザ、制御に必要な各種データ等を記憶するメモリ領域を有する。また、CPU11は、CPU11の作業用領域として利用されるとともに、表示用の画像データなどの一時記憶領域として利用されるメモリ領域を有する。

30

【0030】

データ取得部12は、インターネット31と接続されており、クライアント端末20から要求されたウェブページのコンテンツ等をインターネット31を介して取得する。また、データ取得部12は、文書データベース(DB)32と接続されており、クライアント端末20から要求された文書ファイルなどの各種データを文書DB32から取得する。

【0031】

40

画像生成部13は、データ取得部12が取得したコンテンツ、文書データから画像(以下、閲覧用画像という)を生成する。画像生成部13は、生成した閲覧用画像をCPU11のメモリ領域に記憶する。

【0032】

OCR処理部14は、入力された画像に含まれる文字を識別して文書に変換する。OCR処理自体は一般的な技術であるため、詳細な説明は省略する。

【0033】

テキスト抽出部15は、CPU11により取得されたウェブページのソースから、OCR処理部14が取得したテキストと最も一致度が高いテキストを抽出する。また、テキスト抽出部15は、CPU11により取得された文書データから、OCR処理部14が取得

50

したテキストと最も一致度が高いテキストを抽出する。テキスト抽出部 15 の処理の詳細については、後に詳述する。

【0034】

通信部 16 は、閲覧用画像等をクライアント端末 20 へ送信する。また、通信部 16 は、クライアント端末 20 から送信されたウェブページ閲覧要求等を受信する。

【0035】

クライアント端末 20 は、例えば小型ノートパソコンや携帯電話等であり、図 1 に示すように、ネットワークを介してサーバ 10 と接続される。クライアント端末 20 は、図 3 に示すように、主として、CPU 21 と、入力部 22 と、表示部 23 と、表示制御部 24 と、通信部 25 とで構成される。なお、クライアント端末 20 は、小型ノートパソコンや携帯電話に限定されるものではなく、ウェブブラウザを動作させ得る情報端末であればどのような端末でもよい。

10

【0036】

CPU 21 は、クライアント端末 20 の全体の動作を統括制御するとともに、各種の演算処理を行う演算手段として機能する。CPU 21 は、クライアント端末 20 のクライアント端末情報や、各種制御に必要なプログラム等が記憶されるメモリ領域を有する。また、CPU 21 は、サーバ 10 から送信された各種データを一時的に記憶するバッファを有する。

【0037】

入力部 22 は、ユーザが各種指示を入力するためのものであり、テンキー、十字キー等で構成される。

20

【0038】

表示部 23 は、例えば、カラー表示が可能な液晶ディスプレイである。なお、表示部 23 は、カラー表示に限定されず、白黒表示でもよい。また、表示部 23 は、液晶ディスプレイに限定されず、有機 EL 等を用いてもよい。

【0039】

表示制御部 24 は、サーバ 10 から送信された閲覧用画像を表示部 23 に表示させる。

【0040】

通信部 25 は、サーバ 10 から送信された閲覧用画像、テキストデータ等を受信する。また、通信部 25 は、ウェブページ閲覧要求、領域の情報等をサーバ 10 へ送信する。

30

【0041】

上記のように構成された閲覧システム 1 の作用について説明する。閲覧システム 1 では、クライアント端末 20 にウェブページ（又は文書データ）の画像が表示され、クライアント端末 20 により所定の領域が選択されると、その領域内のテキストをコピーすることができる。図 4 は、クライアント端末 20 が表示部 23 に表示されたウェブページ内のテキストをコピーする処理の流れを示すフローチャートである。

【0042】

クライアント端末 20 の CPU 21 は、メモリ領域に記憶されたウェブブラウザを起動する。入力部 22 により閲覧したいウェブページの情報（URL 等）が入力されると、CPU 21 は、これを受け付けてサーバ 10 へリクエストを送信する（ステップ S20）。

40

【0043】

サーバ 10 の CPU 11 は、リクエストを受信するとデータ取得部 12 に指示を出し、データ取得部 12 はインターネットからリクエストされたウェブページを取得する（ステップ S10）。この場合には、サーバ 10 はプロキシとして動作し、外部のサーバからコンテンツ（例えば、ウェブページの HTML ファイル）を取得する。CPU 11 は、取得したコンテンツをバッファに記憶する。なお、サーバ 10 は、ウェブサーバとしても機能しても良く、この場合にはサーバ 10 の図示しないメモリに記憶されているコンテンツを取得する。

【0044】

データ取得部 12 は取得したコンテンツを画像生成部 13 に出力し、画像生成部 13 は

50

コンテンツから閲覧用画像を生成する（ステップS 1 1）。ウェブページのH t m l ファイルを取得した場合には、画像生成部 1 3 は、H t m l ファイルを解析し、解析結果に基づいて文字や画像を適切に配置した結果を画像化（レンダリング）し、gif、jpeg等の画像ファイルとして保存する。

【0045】

画像生成部 1 3 は生成した閲覧用画像をC P U 1 1へ出力し、C P U 1 1は閲覧用画像をクライアント端末 2 0へ送信する（ステップS 1 2）。

【0046】

クライアント端末 2 0のC P U 2 1は、サーバ 1 0から送信された閲覧用画像を受信し（ステップS 2 1）、表示制御部 2 4へ出力する。表示制御部 2 4は、受信した画像を表示部 2 3へ表示させる（ステップS 2 2）。これにより、図 5 に示すように、クライアント端末 2 0にリクエストしたウェブページの画像が表示され、ユーザがウェブページを閲覧可能となる。

【0047】

表示部 2 3に閲覧用画像が表示された状態で、入力部 2 2によりテキストを抽出（コピー）したい領域の指定が行われる（ステップS 2 3）。領域の指定は、例えば、ユーザが入力部 2 2の十字キー等でカーソルを移動させ、領域の始点及び終点の位置を選択入力することにより行われる。入力部 2 2による入力結果がC P U 2 1で検出されると、C P U 2 1は、図 5 に示すように、始点と終点とにより形成される矩形領域が指定されたと認識する。なお、領域の指定は、この形態に限らず、始点と終点の座標の値を直接入力する等の様々な方法により行うことができる。

【0048】

C P U 2 1は、認識した矩形領域の情報をサーバ 1 0へ送信する（ステップS 2 4）。矩形領域の情報としては、矩形領域の始点及び終点の座標が考えられる。図 5 に示す場合には、閲覧用画像の左上を原点（X 座標、Y 座標共に 0）とし、右方向を+ X 方向、下方向を+ Y 方向として座標が指定される。ただし、座標の指定方法はこれに限定されるものではない。C P U 2 1は、矩形領域の情報として、閲覧用画像から矩形領域を切り出し、切り出された画像を矩形領域の情報として送信するようにしてもよい。

【0049】

サーバ 1 0のC P U 1 1は、クライアント端末 2 0から送信された矩形領域の情報を受信する（ステップS 1 3）。C P U 1 1は、矩形領域の情報をO C R 処理部 1 4へ出力する。

【0050】

O C R 処理部 1 4は、矩形領域の情報に基づいて矩形領域に含まれる文字を認識する（ステップS 1 4）。矩形領域の情報として矩形領域の始点及び終点の座標が入力された場合には、O C R 処理部 1 4は、画像生成部 1 3から閲覧用画像を取得し、閲覧用画像と座標とから矩形領域の画像を切り出す。本実施の形態では、O C R 処理部 1 4は、図 5 の点線で囲まれた領域の画像を矩形領域の画像として切り出す。

【0051】

そして、O C R 処理部 1 4は、切り出した画像をO C R 処理することにより、矩形領域に含まれる文字を認識する。図 6 に示すように、O C R 処理部 1 4は、矩形領域に含まれる「ベルリンで開催された世界陸上をはじめ、週末のスポーツイベント結果ほか、今注目すべき選手についてご紹介」という文字をO C R 処理し、「ベルリンで開催された世界陸上をばじ助、週末のスポーツイベント結果ほか、いま注目すべき選手1についてご紹介。」という認識結果を得る。

【0052】

矩形領域の情報として閲覧用画像から切り出された画像が入力された場合には、O C R 処理部 1 4は、座標情報から画像を抽出する処理は不要であり、入力された画像を直接O C R 処理し、文字を認識する。閲覧システムの実施形態としては、一般的にクライアント端末とサーバではサーバの処理能力のほうが高いため、クライアント端末では処理コスト

10

20

30

40

50

の小さい矩形領域の座標の送信のみを行い、サーバで座標に従い指定された領域の画像を抽出する処理を行う方が好ましい。

【 0 0 5 3 】

OCR処理部14は、得られた認識結果をテキストデータとしてテキスト抽出部15に出力する。テキスト抽出部15は、バッファに記憶されたHTMLファイルを取得し、図7に示すように、HTMLファイルのソースに含まれるテキストの中から入力されたテキストデータと推定されるテキストを抽出する(ステップS15)。ステップS15の処理は、例えば、入力されたテキストデータをキーとして、ソース内から最も一致度の高いテキストを抽出することにより行われる。本実施の形態では、ページのソースとしてHTMLファイルを用いたが、HTMLファイルに限られるものではなく、クライアント端末20に送信した閲覧用画像の基となるウェブページをレンダリングするために必要な情報であればどのようなものでもよい。

10

【 0 0 5 4 】

最も一致度の高いテキストを抽出する方法について、図8を用いて説明する。OCR処理部14により「ABC」というテキストが認識された場合には、テキスト抽出部15は、「ABC」というテキストとソースとを順番に比較し、一致度を算出する。例えば、「ABC」というテキストとソース内のテキスト「AVA」との一致度は33%であり、「ABC」というテキストとソース内のテキスト「VAB」との一致度は0%であり、「ABC」というテキストとソース内のテキスト「ABA」との一致度は66%であり、「ABC」というテキストとソース内のテキスト「EAC」との一致度は33%である。一致度が最も高いのは、「ABC」というテキストとソース内のテキスト「ABA」とを比較した場合であるため、テキスト抽出部15は、ソース内のテキスト「ABA」を抽出する。

20

【 0 0 5 5 】

図7に示す場合には、テキスト抽出部15は、ステップS14で認識されたテキスト「ベルリンで開催された世界陸上をばじ助、週末のスポーツイベント結果ほか、いま注目すべ舌選手1ごついてご紹介。」をキーとして、ソース内から最も一致度の高いテキストの抽出を行う。その結果、テキスト抽出部15は、「ベルリンで開催された世界陸上をはじめ、週末のスポーツイベント結果ほか、いま注目すべき選手についてご紹介。」というテキストを抽出する。

30

【 0 0 5 6 】

そして、テキスト抽出部15は、抽出されたテキストをクライアント端末20で指定された矩形領域に含まれるテキストと判定する。クライアント端末20で指定された矩形領域に含まれるテキストは、必ずソース内に含まれるテキストである。したがって、ソース内に含まれるテキストからOCR処理の結果得られたテキストを推測して抽出することにより、OCR処理のミスにより間違ったテキストが認識された場合においても、そのミスを補完し、正しいテキストを抽出することができる。

【 0 0 5 7 】

なお、本実施の形態では、ステップS15において、ステップS10で取得され、バッファに記憶されたHTMLファイルを用いたが、ステップS15の処理の前に改めてHTMLファイルを取得してもよい。また、ステップS15においては、ソースに含まれるテキスト全てを抽出対象としても良いし、ソースがHTMLファイルでメタ情報(タグ)が含まれている場合等であれば、タグを除いたレンダリングの対象となるテキストのみを抽出対象としても良い。

40

【 0 0 5 8 】

テキスト抽出部15は、抽出したテキストをCPU11に出力し、図9に示すように、CPU11はテキストをクライアント端末20へ送信する(ステップS16)。クライアント端末20のCPU21は、サーバ10から送信されたテキストを受信し(ステップS25)、受信したテキストをCPU21内のバッファに記憶する(ステップS26)。バッファに保存したテキストは、例えば任意のテキストフィールドへの貼り付けなどに利用

50

することなどが考えられる。

【0059】

本実施の形態によれば、ウェブページや文書データを画像化してクライアント端末に表示させる場合に、クライアント端末に表示された画像の一部を選択することにより、選択した領域に含まれる正確なテキストデータを得ることができる。そして、得られたテキストデータを記憶することにより、クライアント端末で選択された領域の画像に含まれるテキストをコピーすることと同等の効果を得ることができる。

【0060】

従来のシンクライアント型ブラウザでは、クライアント端末で閲覧されるウェブページは画像化されているため、ウェブページに含まれるテキストをコピーすることはできなかった。しかしながら、OCR処理とソートからのテキスト抽出とを組み合わせることにより、シンクライアント型ブラウザを用いる場合においても所望のテキストのコピーアンドペーストが可能となる。

10

【0061】

また、本実施の形態によれば、下線付き文字や表の一部等OCR処理の精度が低い場合においても、正確なテキストデータをコピーすることができる。例えば、ステップS23で図5の一点鎖線で囲んだ領域が矩形領域として選択された場合には、ステップS14のOCR処理において、行間の線が原因で上段のテキストは正確な認識結果は得られない。しかしながら、図7に示すようにソースと比較することにより、「各党の政権公約比較」「安全保障」及び「候補者情報」「マニフェスト」「選挙ニュース」というテキストを抽出することができる。

20

【0062】

なお、本実施の形態では、図4に示すようにウェブページを閲覧する場合を例に作用を説明したが、ウェブページの閲覧のみでなく、文書データを閲覧する場合においても同様の方法により、選択した矩形領域内のテキストを抽出することができる。

【0063】

<第2の実施の形態>

第1の実施の形態は、OCR処理のミスにより間違ったテキストが認識された場合においても、そのミスを補完し、正しいテキストを抽出するため、ソースに含まれるテキストの中からテキストを抽出する処理を行なったが、必ずしもソースからのテキスト抽出処理が必要とは限らない。例えば、単語等テキストの長さが短い場合には、OCR処理の制度が高いため、処理結果が正しい場合も多い。

30

【0064】

第2の実施の形態は、クライアント端末で選択された矩形領域の大きさ、即ちテキストの長さに応じてテキスト抽出処理をするかしないかを異ならせる形態である。以下、第2の実施の形態に係る閲覧システム2について説明する。なお、閲覧システム2の構成は閲覧システム1と同様であるため、説明を省略する。また、第1の実施の形態と同一の部分については、同一の符号を付し、詳細な説明を省略する。

【0065】

図10は、閲覧システム2において、クライアント端末20により選択された領域内のテキストをコピーする処理の流れを示すフローチャートである。

40

【0066】

クライアント端末20のCPU21は、メモリ領域に記憶されたウェブブラウザを起動する。入力部22により閲覧したいウェブページの情報(URL等)が入力されると、CPU21は、これを受け付けてサーバ10へリクエストを送信する(ステップS20)。

【0067】

サーバ10のCPU11は、リクエストを受信するとデータ取得部12に指示を出し、データ取得部12はインターネットからリクエストされたウェブページを取得する(ステップS10)。データ取得部12は取得したコンテンツを画像生成部13に出力し、画像生成部13はコンテンツから閲覧用画像を生成する(ステップS11)。画像生成部13

50

は生成した閲覧用画像をCPU11へ出力し、CPU11は閲覧用画像をクライアント端末20へ送信する(ステップS12)。

【0068】

クライアント端末20のCPU21は、サーバ10から送信された閲覧用画像を受信し(ステップS21)、表示制御部24へ出力する。表示制御部24は、受信した画像を表示部23へ表示させる(ステップS22)。これにより、クライアント端末20にリクエストしたウェブページの画像が表示され、ユーザがウェブページを閲覧可能となる。

【0069】

表示部23に閲覧用画像が表示された状態で、テキストを抽出(コピー)したい矩形領域の指定が行われる(ステップS23)。指定された矩形領域の情報はCPU21で検出され、CPU21は、認識した矩形領域の情報をサーバ10へ送信する(ステップS24)。

【0070】

サーバ10のCPU11は、クライアント端末20から送信された矩形領域の情報を受信する。CPU11は、受信された矩形領域の情報に基づいて、矩形領域の大きさ(面積)を算出する(ステップS17)。

【0071】

CPU11は、矩形領域の情報をOCR処理部14へ出力する。OCR処理部14は、矩形領域の情報に基づいて矩形領域に含まれる文字を認識する(ステップS14)。

【0072】

CPU11はステップS13で受信された矩形領域の大きさが閾値以上であるか否かを判断する(ステップS18)。なお、閾値は、予め設定された任意の値であり、CPU11のメモリ領域に記憶されている。閾値は、必要に応じてクライアント端末20等から変更することもできる。閾値としては、OCR処理により正しい結果が得られる最大の長さ(単語レベルの長さ)のテキストが含まれるような面積とすることが望ましい。

【0073】

矩形領域の大きさが閾値以上である場合(ステップS18でYES)は、クライアント端末20により指定された領域に含まれるテキストは文章等の長いテキストであると推定される。テキストが長い場合には、OCR処理の精度は低く、正確に文字が認識できない場合が多い。したがって、OCR処理部14は得られた認識結果をテキストデータとしてテキスト抽出部15に出力し、テキスト抽出部15はバッファに記憶されたHTMLファイルのソースに含まれるテキストの中から入力されたテキストデータと推定されるテキストを抽出する(ステップS15)。テキスト抽出部15は抽出されたテキストをCPU11に出力し、CPU11はテキストをクライアント端末20へ送信する(ステップS19)。これにより、OCR処理のミスにより間違っ

【0074】

矩形領域の大きさが閾値以上でない場合(ステップS17でNO)は、クライアント端末20により指定された領域に含まれるテキストは単語レベルであると推定される。単語であれば、OCR処理の精度がある程度期待できる。また、短いテキストをソースから抽出することで、間違っ

【0075】

ステップS18~S19の処理について、図11を用いて具体的に説明する。閾値が「50」である場合に、ステップS17で算出された面積が「200」である場合には、算出された面積「200」は閾値「50」より大きいため、HTMLファイルのソースに含まれるテキストの中から正しいと推定されるテキストを抽出し、その結果をクライアント端末20で指定された矩形領域に含まれるテキストと判定する。それに対し、ステップS17で算出された面積が「10」である場合には、算出された面積「10」は閾値「50

10

20

30

40

50

」より小さいため、テキスト抽出は行わず、OCR処理により得られた結果をクライアント端末20で指定された矩形領域に含まれるテキストと判定する。

【0076】

クライアント端末20のCPU21は、サーバ10から送信されたテキストを受信し(ステップS25)、受信したテキストをCPU21内のバッファに記憶する(ステップS26)。バッファに保存したテキストは、例えば任意のテキストフィールドへの貼り付けなどに利用することなどが考えられる。

【0077】

本実施の形態によれば、矩形領域の大きさに応じて送信するテキストの抽出方法を変えることにより、効率、精度の良い処理を行うことができる。

10

【0078】

なお、上記第1、第2の実施の形態では、サーバとクライアント端末とを有するシステムを例に説明したが、本発明は、システムに限らず、外部の装置へ画像を配信するサーバとして提供することもできる。また、サーバ、クライアント端末に適用するプログラムとして提供することもできる。

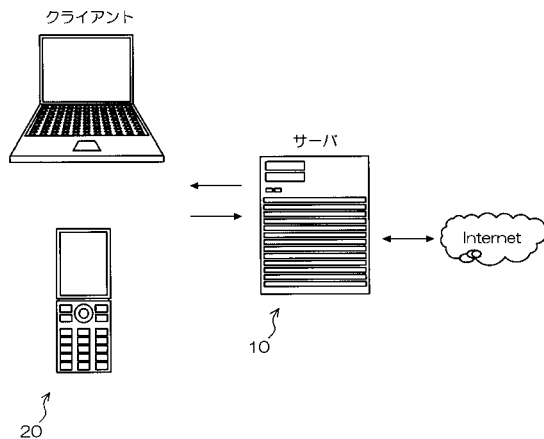
【符号の説明】

【0079】

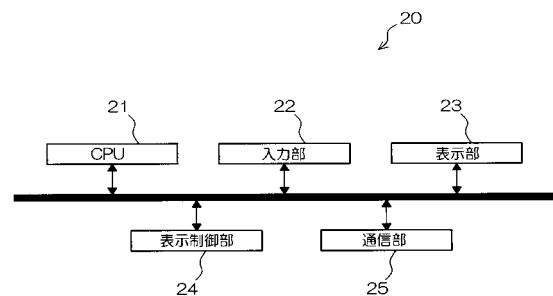
1、2：閲覧システム、10：サーバ、11：CPU、12：データ取得部、13：画像生成部、14：OCR処理部、15：テキスト抽出部、16：通信部、20：クライアント端末、21：CPU、22：入力部、23：表示部、24：表示制御部、25：通信部

20

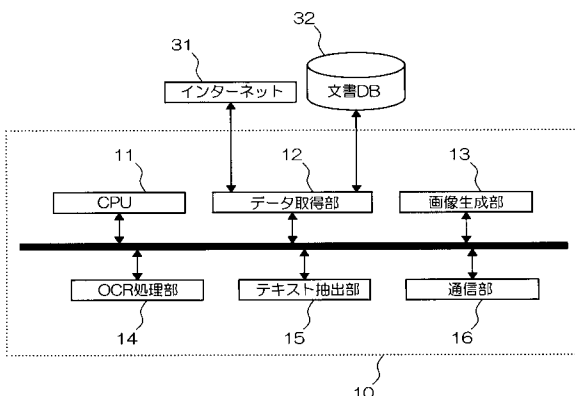
【図1】



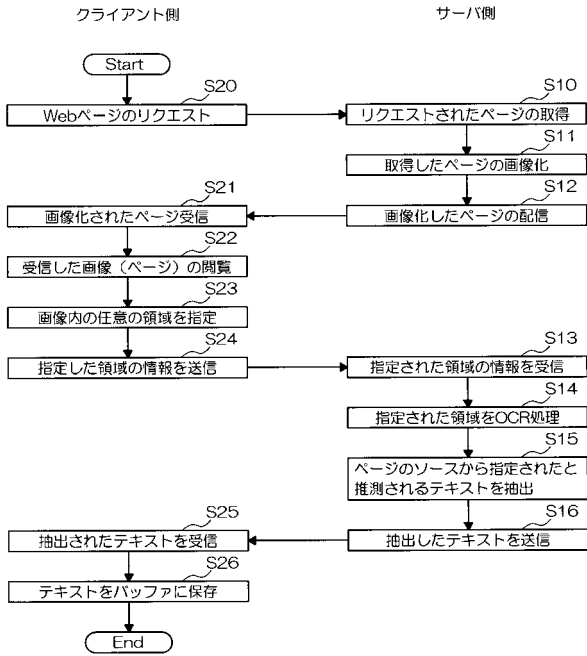
【図3】



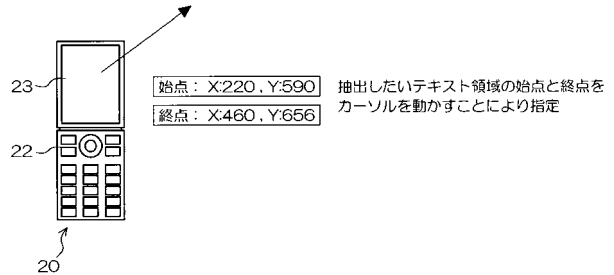
【図2】



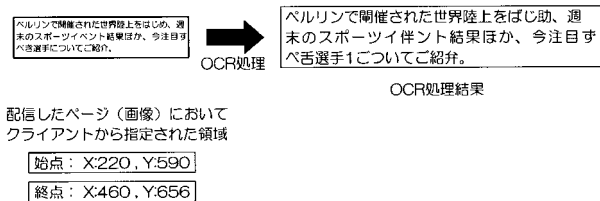
【図 4】



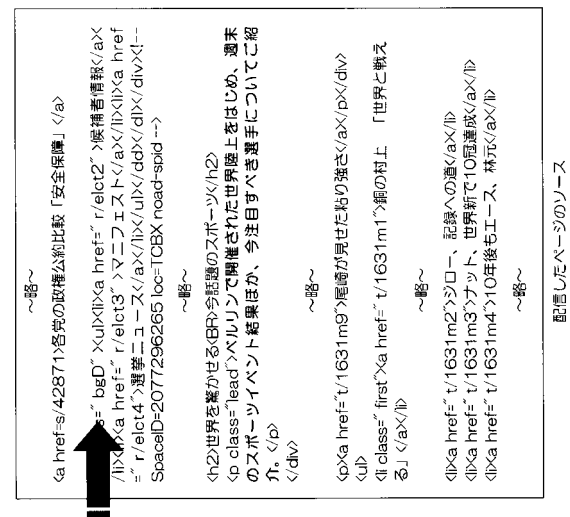
【図 5】



【図 6】



【図 7】

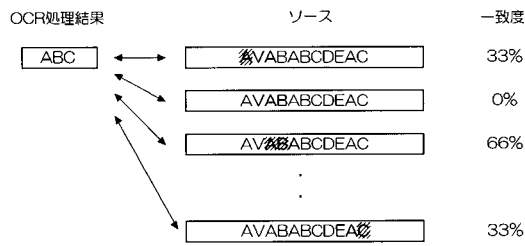


ベルリンで開催された世界陸上をはじめ、週末のスポーツイベント結果ほか、今注目すべき選手1についてご紹介。

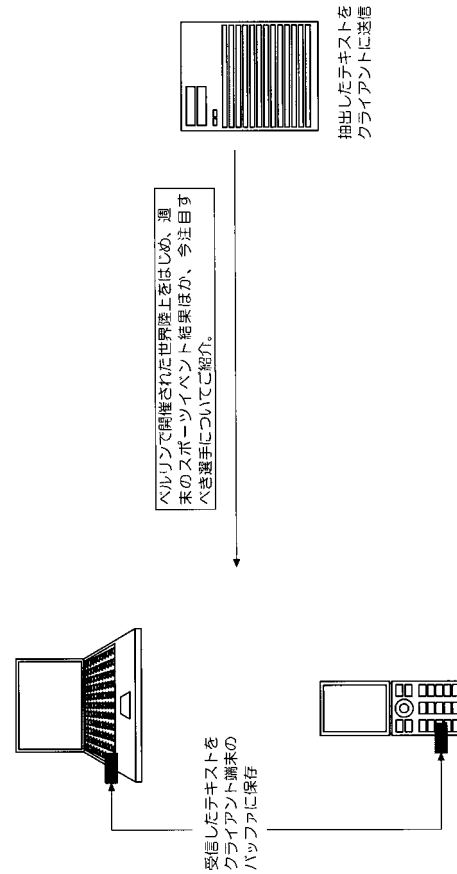
OCR処理結果

配信したページのソース内テキストからOCR処理結果のテキスト(文章)と一致度が高いテキストを、クライアントから指定された領域(領域に含まれる)テキストとして抽出する。

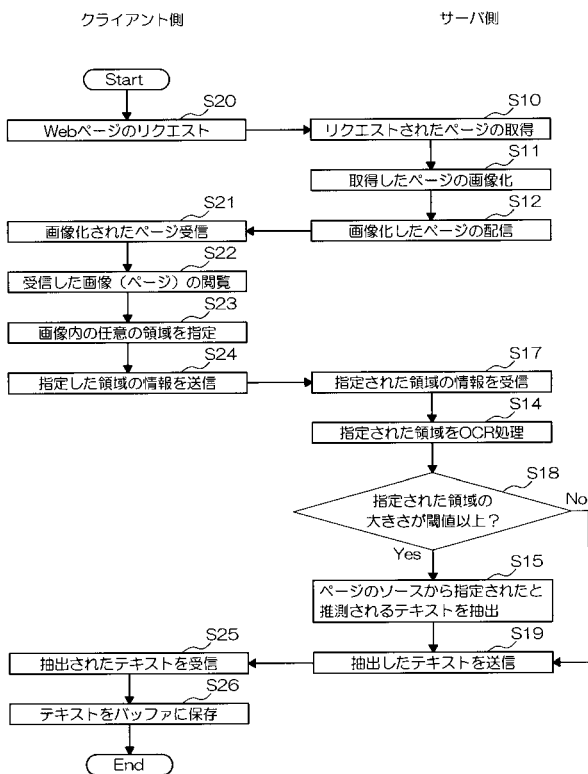
【図 8】



【図 9】



【図 10】



【図 11】

