



(19) **United States**

(12) **Patent Application Publication**  
ARAI et al.

(10) **Pub. No.: US 2014/0343926 A1**

(43) **Pub. Date: Nov. 20, 2014**

(54) **LANGUAGE MODEL GENERATING DEVICE,  
METHOD THEREOF, AND RECORDING  
MEDIUM STORING PROGRAM THEREOF**

**Publication Classification**

(51) **Int. Cl.**  
*G06F 17/21* (2006.01)  
*G06F 17/27* (2006.01)  
(52) **U.S. Cl.**  
CPC ..... *G06F 17/21* (2013.01); *G06F 17/2755*  
(2013.01)  
USPC ..... **704/9**

(71) Applicant: **NEC Informatec Systems, Ltd.**, Tokyo (JP)

(72) Inventors: **Kazuhiro ARAI**, Tokyo (JP); **Tadashi EMORI**, Tokyo (JP)

(21) Appl. No.: **14/446,540**

(57) **ABSTRACT**

(22) Filed: **Jul. 30, 2014**

A text in a corpus including a set of world wide web (web) pages is analyzed. At least one word appropriate for a document type set according to a voice recognition target is extracted based on an analysis result. A word set is generated from the extracted at least one word. A retrieval engine is caused to perform a retrieval process using the generated word set as a retrieval query of the retrieval engine on the Internet, and a link to a web page from the retrieval result is acquired. A language model for voice recognition is generated from the acquired web page.

**Related U.S. Application Data**

(62) Division of application No. 13/271,424, filed on Oct. 12, 2011, now Pat. No. 8,831,945.

**Foreign Application Priority Data**

(30) Oct. 12, 2010 (JP) ..... 2010-229526

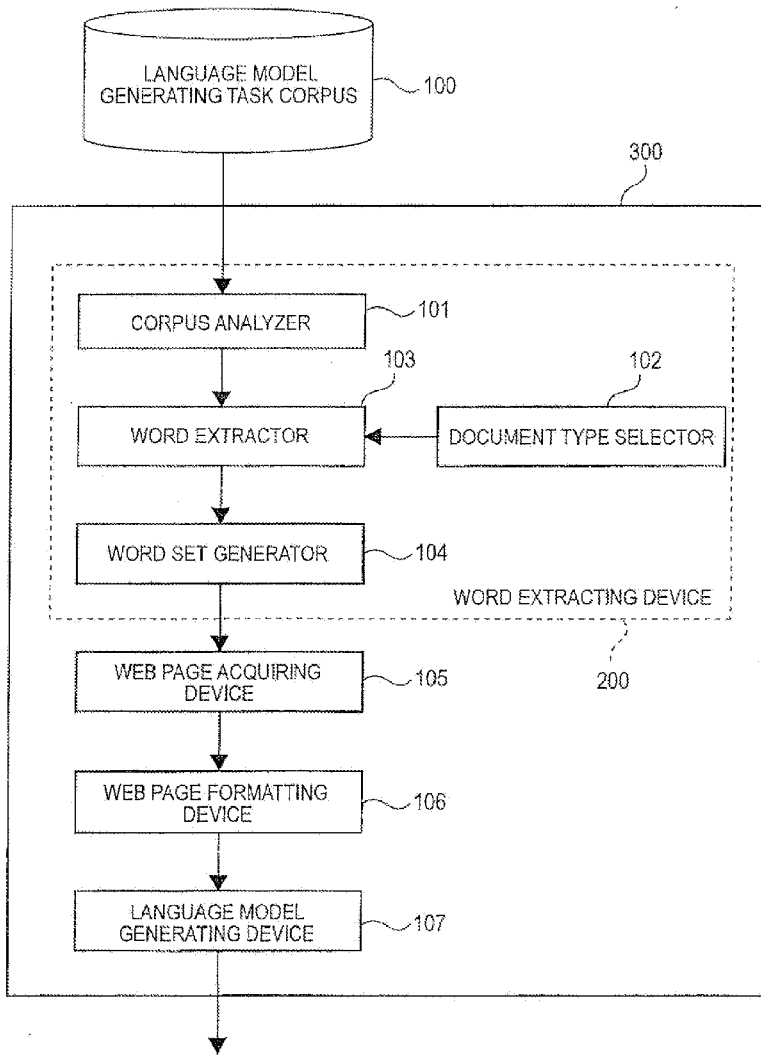


FIG. 1

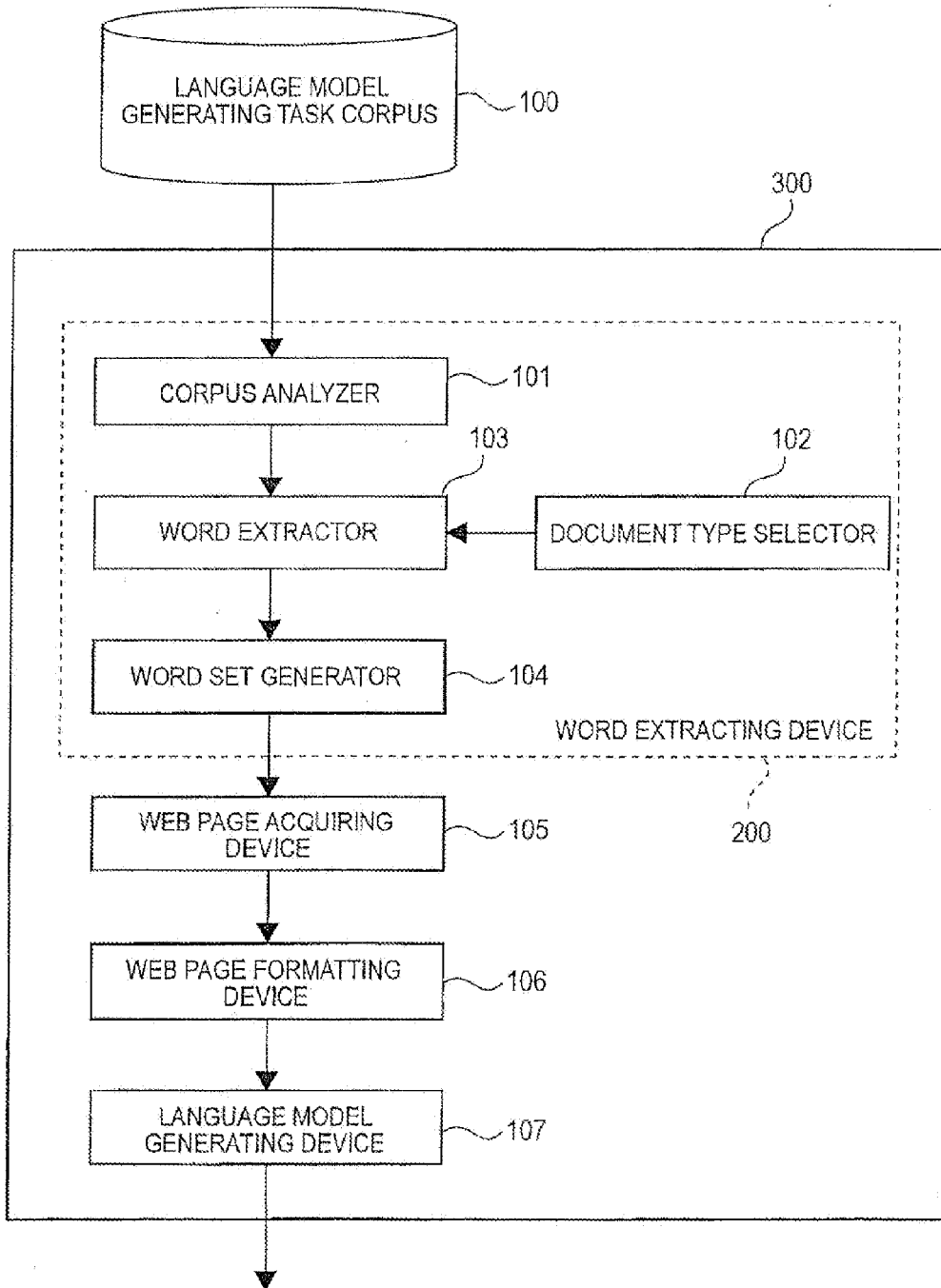
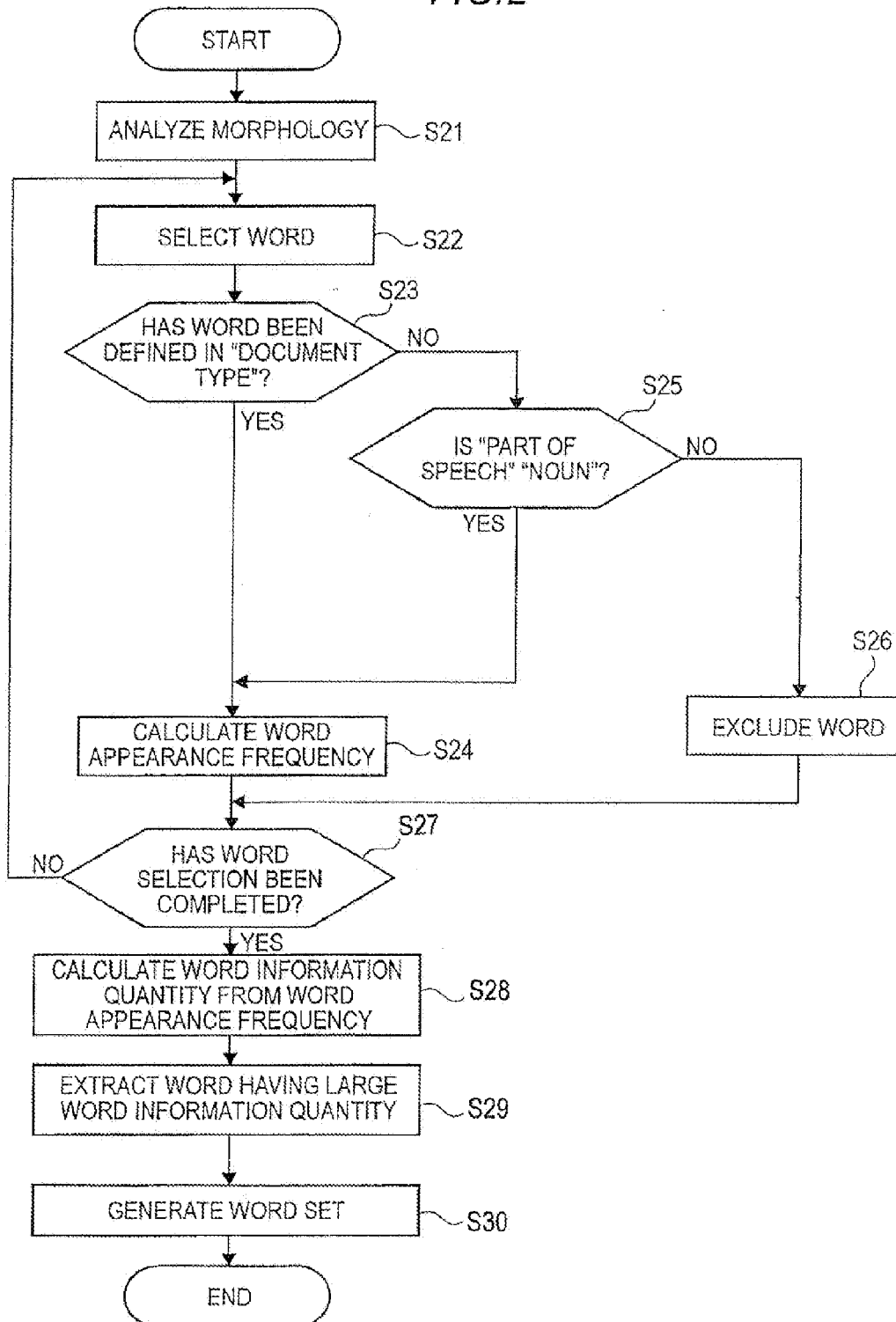


FIG. 2



**FIG.3**

IT IS FINE TODAY.	
NOTATION	PART OF SPEECH
IT	PRONOUN
IS	VERB
FINE	ADJECTIVE
TODAY	NOUN
.	PUNCTUATION MARK
IT WILL BE SUNNY TOMORROW AGAIN.	
NOTATION	PART OF SPEECH
IT	PRONOUN
WILL	MODAL VERB
BE	VERB
SUNNY	ADJECTIVE
TOMORROW	NOUN
AGAIN	ADVERB
.	PUNCTUATION MARK

*FIG.4*

WORDING (NOTATION)	
BRITISH ENGLISH	AMERICAN ENGLISH
COLOUR	COLOR
AEROPLANE	AIRPLANE
LIFT	ELEVATOR
FLYOVER	OVERPASS
FITTED	FIT
SELECTED MEANING	
NOTATION	PART OF SPEECH
SEE	VERB
SAW	VERB
LIVE	VERB, ADJECTIVE
SPOKEN LANGUAGE (NOTATION)	
RIGHT?	
OH, YES	
COULD BE	

*FIG.5*

WORD	APPEARANCE FREQUENCY
SAME PERIOD	32
PRECEDING YEAR	31
DECREASE	19
BUSINESS	16
SALES	16
INCREASE	15
FIELD	14
FORECAST	12
SECURITIES	10
DESCRIPTION	9
FUTURE	8
PROFIT	7
QUATER	7
SERVICE	7

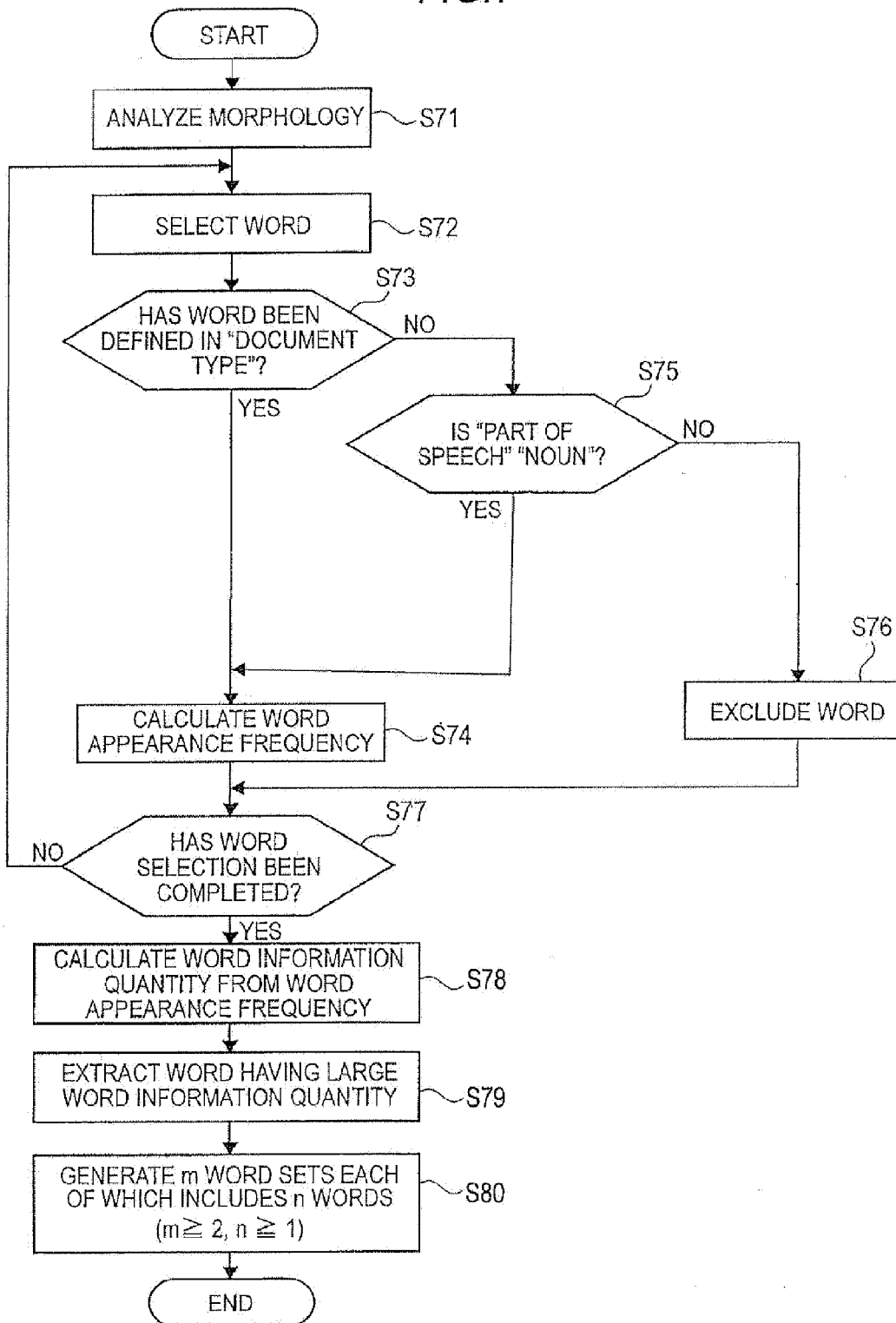
FIG.6

WORD SET	BUSINESS PERFORMANCE, PRECEDING YEAR, SAME PERIOD, FIELD, SALES
----------	--

FIG.8

FIRST WORD SET	BUSINESS PERFORMANCE, PRECEDING YEAR, SERVICE
SECOND WORD SET	FIELD, BUSINESS PERFORMANCE, SALES
THIRD WORD SET	QUARTER, SALES, BUSINESS PERFORMANCE
FOUR WORD SET	SERVICE, SALES, BUSINESS PERFORMANCE
FIVE WORD SET	BUSINESS PERFORMANCE, SERVICE, INVESTMENT

FIG. 7





## LANGUAGE MODEL GENERATING DEVICE, METHOD THEREOF, AND RECORDING MEDIUM STORING PROGRAM THEREOF

**[0001]** This application is based upon and claims the benefit of priority from Japanese patent application No. 2010-229526, filed on Oct. 12, 2010, the disclosure of which is incorporated herein in its entirety by reference.

### TECHNICAL FIELD

**[0002]** The present invention relates to a language model generating device, a method thereof, and a recording medium storing a program thereof, and more particularly to a language model generating device, a method thereof, and a recording medium storing a program thereof that efficiently collect a corpus similar to the content of a voice recognition target from web pages present on the Internet and generate a language model.

### BACKGROUND ART

**[0003]** As a method of complementing word information of a voice recognition target when a language model for voice recognition is generated, there is a method of collecting information similar to the content of the target from web pages on the Internet and generating a language model.

**[0004]** Generally, when the retrieval is performed by a retrieval system that retrieves a world wide web (WWW) based on a designated retrieval word, links of web pages which are arranged according to rank, decided by a predetermined evaluation criterion, are output on a page obtained as a retrieval result. Examples of an evaluation index include; appearance frequency of the retrieval word, metadata of a hypertext markup language (HTML), the number of page links, the presence and absence of a link from a page having many user references, and the like. In order to generate the language model, a web page of a link destination linked from the retrieval result page is acquired. However, as the web page of the link destination, web pages including the content similar to a voice recognition target may be provided, but in most cases, web pages including a plurality of topics or mentioning specialized fields are provided. Thus, if the language model is generated without selecting the web page to acquire, the degree of recognition accuracy of voice recognition is lowered.

**[0005]** For this reason, various techniques for selecting the web page and extracting words involved in selection of the web page have been suggested.

**[0006]** For example, Non-Patent Literature 1 discloses a technique for extracting a word whose part-of-speech information represents a noun from a word string as a result of voice recognition, retrieving a news site on the internet using the extracted word as a retrieval word, and collecting similar web pages. In a technique disclosed in Non-Patent Literature 2, in order to collect a medical related corpus, only the word "medical" is used as the retrieval word, and information is collected up to below two layers of the link destination of the retrieval result. In a technique disclosed in Non-Patent Literature 3, words whose appearance frequencies in the recognition result are ranked within top five are extracted as the retrieval words.

**[0007]** Further, Patent Literature 1 discloses a technique of preventing a concatenation of words including a word with a high appearance frequency from having an unreasonably high language probability when the language model is gen-

erated. Patent Literature 2 discloses a technique of changing retrieval priority according to the background color of a character string inside an image in a system of retrieving information on a network using the character string. In a user interface design tool capable of designing voice recognition, a voice rule synthesis, or the like disclosed in Patent Literature 3, a designer can set a character recognition part and set a recognition mode to "hiragana" or the like.

### CITATION LIST

#### Patent Literature

- [0008]** {PTL 1} JP-A-2002-229588 (paragraph {0103})  
**[0009]** {PTL 2} JP-A-2010-039647 (paragraph {0044})  
**[0010]** {PTL 3} JP-A-09-114623 (paragraph {0107} and FIG. 13)

#### Non-Patent Literature

- [0011]** {NPTL 1} Ito Yusuke, Nishizaki Hiromitsu, Sekiguchi Yoshihiro "Improving recognition performance of spoken documents using similar documents on the Internet", IEIC Technical Report, NLC2005-65, SP2005-98 (2005-12), p.49-53  
**[0012]** {NPTL 2} Nishimura Ryuichi, Nagatomo Kentaro, Komatsu Kumiko, Kuroda Yuka, Lee Akinobu, Saruwatari Hiroshi, Shikano Kiyohiro, "Automatic Language Model Building Tool using Web Texts", Speech Language Information Processing 35-8 (2001.2.3), p.43-47  
**[0013]** {NPTL 3} Takahashi Shin-ya, Morimoto Tsuyoshi, and Irie Yuki, "Analysis of Training Corpus Collected by Web Crawling for Speech Recognition of News Documents", Fukuoka University review of technological sciences 77, Sep., 2006, p.37-44

### SUMMARY OF INVENTION

#### Technical Problem

**[0014]** However, in the technique disclosed in Non-Patent Literature 1, a wide variety of web pages including a word having a plurality of meanings may possibly be retrieved even though the part of speech is narrowed down to a noun. For this reason, when web pages are used as a learning corpus of the language model of voice recognition, the degree of recognition accuracy may be lowered. Further, even though information of top two layers of the retrieval result is collected using the technique of Non-Patent Literature 2, since information of the web pages of the two layer destinations does not necessarily include the desired web page, the degree of recognition accuracy may be lowered. In Non-Patent Literature 3, frequently-appearing words within top five ranks are used as the retrieval words. However, high frequency words may be words that frequently appear in all documents. Thus, when these words are used as the retrieval words, it is difficult to obtain the desired retrieval result. Further, since a synonym may be included due to paraphrasing of a word or fluctuation of notation, the retrieval result may be biased, and it is difficult to collect a sufficient amount of corpuses for desired information from the link destinations of the retrieval result.

**[0015]** Further, Patent Literature 1 does not disclose the technique of reliably selecting a feature word representing a feature of a sentence at the time of word extraction. Patent Literature 2 does not disclose the technique of giving a priority applied at the time of generating the language model into

a character string. In Patent Literature 3, the type of character such as Hiragana or a Chinese character is recognized, but it is difficult to recognize text according to the content of voice recognition target.

[0016] Thus, an exemplary object of the present invention is directed to provide a language model generating device, a method thereof, and a program thereof that efficiently collect a corpus similar to the content of the voice recognition target from the web pages present on the Internet and generate a language model.

#### Solution to Problem

[0017] According to an exemplary first aspect of the present invention, there is provided a language model generating device, which includes:

[0018] a corpus analyzer which analyzes text in a corpus including a set of world wide web (web) pages;

[0019] an extractor which extracts at least one word appropriate for a document type set according to a voice recognition target based on an analysis result by the corpus analyzer;

[0020] a word set generator which generates a word set from the at least one word extracted by the extractor;

[0021] a web page acquiring device which causes a retrieval engine to perform a retrieval process using the word set generated by the word set generator as a retrieval query of the retrieval engine on the Internet and acquires a web page related to the word set from the retrieval result; and

[0022] a language model generating device which generates a language model for voice recognition from the web page.

[0023] According to an exemplary second aspect of the present invention, there is provided a language model generating method, which includes the steps of:

[0024] analyzing a text in a corpus including a set of world wide web (web) pages;

[0025] extracting at least one word appropriate for a document type set according to a voice recognition target based on an analysis result;

[0026] generating a word set from at least one extracted word;

[0027] causing a retrieval engine to perform a retrieval process using the generated word set as a retrieval query of the retrieval engine on the Internet and acquiring a web page related to the word set from the retrieval result; and

[0028] generating a language model for voice recognition from the acquired web page.

[0029] According to an exemplary third aspect of the present invention, there is provided a non-transitory computer-readable recording medium storing a language model generating program used in a computer of a language model generating device and causing the computer to execute:

[0030] a step of analyzing a text in a corpus including a set of world wide web (web) pages;

[0031] a step of extracting at least one word appropriate for a document type set according to a voice recognition target based on an analysis result;

[0032] a step of generating a word set from at least one extracted word;

[0033] a step of causing a retrieval engine to perform a retrieval process using the generated word set as a retrieval query of the retrieval engine on the Internet and acquiring a web page related to the word set from the retrieval result; and

[0034] a step of generating a language model for voice recognition from the acquired web page.

#### Advantageous Effects of Invention

[0035] According to the present invention, the corpus including a web page is analyzed, a word appropriate for a predetermined document type is extracted, a word set generated from the extracted word is used as a retrieval query to cause an Internet information retrieval engine to perform the retrieval, and a language model is generated based on web pages acquired from link destinations of a retrieval result. Thus, the language model can be generated by efficiently collecting the corpus similar to the content of the voice recognition target.

#### BRIEF DESCRIPTION OF DRAWINGS

[0036] FIG. 1A block diagram shows an example of a configuration of a language model generating system according to a first exemplary embodiment of the present invention;

[0037] FIG. 2A flowchart shows an example of a word set generating process according to the first exemplary embodiment of the present invention;

[0038] FIG. 3A diagram shows an example of a morphological analysis result acquired by a corpus analyzing unit according to an exemplary embodiment of the present invention;

[0039] FIG. 4A table shows an example of a “document type” defined or selected by a document type selecting unit according to an exemplary embodiment of the present invention;

[0040] FIG. 5A diagram shows an example of a calculation result of an appearance frequency according to an exemplary embodiment of the present invention;

[0041] FIG. 6A diagram shows an example of a word set according to the first exemplary embodiment of the present invention;

[0042] FIG. 7A flowchart shows an example of a word set generating process according to a second exemplary embodiment of the present invention; and

[0043] FIG. 8A diagram shows an example of a word set according to the second exemplary embodiment of the present invention.

#### DESCRIPTION OF EMBODIMENTS

[0044] Hereinafter, exemplary embodiments of the present invention will be described in detail with reference to the accompanying drawings.

##### First Exemplary Embodiment

[0045] FIG. 1 is a block diagram illustrating an example of a configuration of a language model generating system according to a first exemplary embodiment of the present invention.

[0046] In the language model generating system illustrated in FIG. 1, a language model generating device 300 generates a language model according to a predetermined voice recognition task based on a language model generating task corpus 100.

[0047] The language model generating task corpus 100 is a language resource including various documents (web (world wide web) pages) described in a document description language such as an HTML or an extensible markup language (XML) and used as a learning corpus of the language model generating device 300. In the present embodiment, it is

assumed that the language model generating task corpus 100 includes a linguistic expression similar to a voice recognition target task.

[0048] The language model generating device 300 includes; a word extracting device 200, a web page acquiring unit 105, a web page formatting unit 106, and a language model generating unit 107. The word extracting device 200 includes; a corpus analyzing unit 101, a document type selecting unit 102, a word extracting unit 103, and a word set generating unit 104.

[0049] The corpus analyzing unit 101, of the word extracting device 200, performs an analysis process including morphological analysis of a document on the description content of each document inside the language model generating task corpus 100. For example, for the morphological analysis, an existing morphological analysis technique such as “Tree Tagger” may be used (developed at the Institute for Computational Linguistics of University of Stuttgart; <http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/DecisionTreeTagger.html>). In Japanese, a Mecab system may be used (provided by a joint research between department of informatics, Kyoto University-NTT communication science laboratory, see <http://mecab.sourceforge.net/>) or a Chasen system (provided by a natural Language processing course of Graduate School of Information Science, Nara Institute of Science and Technology, see <http://chasen-legacy.sourceforge.jp/>).

[0050] For example, the morphological analysis result by the corpus analyzing unit 101 is represented as in FIG. 3 which will be described in detail later. The morphological analysis result of a text includes information of each word included in the text, such as “notation”, and “part of speech” as illustrated in FIG. 3.

[0051] The document type selecting unit 102 of the word extracting device 200 sets a rule for word extraction performed by the word extracting unit 103. A “document type” according to the content or the style desired to voice-recognize is defined or selected by an operation of the document type selecting unit 102. An example of the “document type” defined by the document type selecting unit 102 is illustrated in FIG. 4 which will be described in detail later. In the example illustrated in FIG. 4, wording such as a

[0052] British English word and a American English word, “part of speech” selected from a multisense word, and a colloquial document which includes spoken language used in a daily conversation are presented. These documents are examples of what is defined as the “document type”. Further, in order to change the content or the style desired to voice-recognize, the definition of the document type is changed by the document type selecting unit 102.

[0053] The word extracting unit 103 of the word extracting device 200 extracts a feature word, from which a retrieval result associated with the aim of voice recognition is obtained, which is obtained from the analysis results by the corpus analyzing unit 101 according to the rule (document type) set by the document type selecting unit 102.

[0054] For example, for formal documents such as meeting minutes, the word extracting unit 103 extracts a noun and a verb such as “decrease” or “increase” as the feature word rather than a pronoun such as “it” or “we”, and a modal verb such as “can” or “may”. Further, when the meeting minutes includes a British English word and an American

[0055] English word, both of a character string including a British English word and a character string including an

American English word are extracted as the feature words. In the case of daily conversation, a character string having a feature of a colloquial expression such as “right?”, “Oh, yes”, or “Could be” is extracted. Further, a method of selecting a word having a high appearance probability with reference to frequency information of words obtained from a huge amount of text information on a web released by Google may be used (in case of Japanese, see <http://googlejapan.blogspot.com/2007/11/n-gram.HTML>).

[0056] The word set generating unit 104 of the word extracting device 200 generates a combination of a plurality of words extracted by the word extracting unit 103 as a word set. For example, the word set generating unit 104 generates a combination of relevant words such as “same period”, “sales”, or “service” as the word set. The word set generated by the word extracting device 200 in the above described manner is used for web page acquisition by the web page acquiring unit 105.

[0057] The web page acquiring unit 105 has the function of accessing (a set of) web pages through a predetermined retrieval engine. The web page acquiring unit 105 generates a retrieval query for performing a narrowing-down retrieval on the word set generated by the word set generating unit 104, causes the retrieval engine to perform a retrieval process by the query, and acquires a retrieval result. The retrieval result includes links to web pages according to the retrieval query. The web page acquiring unit 105 acquires data of the web pages through the links.

[0058] The web page formatting unit 106 removes a symbol (for example, an angle bracket, a triangular mark, or an asterisk) or a HTML tag (for example, <br> or <head>) which is unnecessary for generation of the language model for voice recognition from the web page data acquired by the web page acquiring unit 105 and so extracts the sentence portion.

[0059] The language model generating unit 107 generates the language model using the sentence portion extracted by the web page formatting unit 106 as the corpus for generating the language model. For example, for the sake of generation of the language model for voice recognition in Japanese, there may be used an existing tool such as a Julius system (Kawahara Lab., Kyoto University, Information-technology Promotion Agency (IPA), Shikano Lab., Nara Institute of Science and Technology, and Julius project team, Nagoya Institute of Technology, <http://julius.sourceforge.jp/>) and an HTK system (University of Cambridge, <http://htk.eng.cam.ac.uk/>). The HTK system is also available in English speech recognition. The generated language model is used for various voice recognition processes.

[0060] In the first exemplary embodiment, the language model generating device 300 is constructed by hardware such as a IC or a plurality of ICs. However, all or some of functions of the language model generating device 300 may be implemented by software such that the program stored in program memory is read and executed by a processor in an information processing device (for example, a personal computer (PC)) including the processor, the program memory, and a work memory. Alternatively, all or some of functions of the language model generating device 300 may be implemented by a combination of hardware and software. Further, an input device such as a keyboard or a mouse may be connected to the document type selecting unit 102 in a wired or wireless manner, and the document type may be defined or selected according to an operation of the input device by a user.

[0061] Subsequently, a description will be made in connection with an operation of the word extracting device 200.

[0062] FIG. 2 is a flowchart illustrating an example of a word set generating process according to the first exemplary embodiment of the present invention.

[0063] In the word extracting process according to the present embodiment, in step S21, first the corpus analyzing unit 101 performs morphological analysis on the description content of the language model generating task corpus 100 and acquires a morphological analysis result. For example, the morphological analysis result includes information representing “notation (character string)”, and “part of speech” of words configuring a text as illustrated in FIG. 3.

[0064] FIG. 3 is a diagram illustrating an example of a morphological analysis result acquired by the corpus analyzing unit 101 according to an exemplary embodiment of the present invention. In the example illustrated in FIG. 3, five words of “it”, “is”, “fine”, “today”, and “.” are detected as the morphological analysis result of a text “It is fine today.”.

[0065] Of these words, each information of notation “it”, and part of speech “pronoun” is detected on the word “it”. Similarly, a morphological analysis result including information of notation “is”, and part of speech “verb” is detected on the word “is”. Further, a morphological analysis result including information of notation “fine”, and part of speech “adjective” is detected on the word “fine”. A morphological analysis result including information of notation “today”, and part of speech “noun” is detected on the word “today”. Information of part of speech “punctuation mark” is acquired on a period “.”.

[0066] In the example illustrated in FIG. 3, seven words of “it”, “will”, “be”, “sunny”, “tomorrow”, “again” and “.” are detected as the morphological analysis result of a text “It will be sunny tomorrow again.”.

[0067] Each information of notation “it”, and part of speech “pronoun” is detected on the word “it”. Similarly, a morphological analysis result including information of notation “will” and part of speech “modal verb” is detected on the word “will”. A morphological analysis result including information of notation “be”, and part of speech “verb” is detected on the word “be”. A morphological analysis result including information of notation “sunny” and part of speech “adjective” is detected on the word “sunny”. A morphological analysis result including information of notation “tomorrow” and part of speech “noun” is detected on the word “tomorrow”. A morphological analysis result including information of notation “again”, and part of speech “adverb” is detected on the word “again”. Information of part of speech “punctuation mark” is acquired on a period “.”.

[0068] In step S22, the word extracting unit 103 selects one word from the morphological analysis result illustrated in FIG. 3 in a predetermined order. The word extracting unit 103 determines whether or not the selected word is to be used as a word set generation target word by a subsequent process. For example, the selection order may be based on an appearance order inside the text.

[0069] In step S23, the word extracting unit 103 determines whether or not “notation”, or “notation” and “part of speech” of the selected word correspond to (match with) information of “notation”, or “notation” and “part of speech” of any character string in the document type defined by the document type selecting unit 102.

[0070] FIG. 4 is a table illustrating an example of the “document type” defined or selected by the document type

selecting unit 102 according to an exemplary embodiment of the present invention. In the example illustrated in FIG. 4, as notations (character strings) for defining the document type of the “wording”, the five pairs of notation character strings are defined. Each of the five pairs of notation character strings consists of a British English word and an American English word. That is, “colour” and “color”, “aeroplane” and “airplane”, “lift” and “elevator”, “flyover” and “overpass”, and “fitted” and “fit” are defined.

[0071] In the example illustrated in FIG. 4, as notations (character strings) and part of speech for defining the “selected meaning” category of the document types, three notation character strings are defined. That is, “saw” as a verb, “see” as a verb, “live” as a verb, and “live” as an adjective are defined, but “saw” as a noun, “see” as a noun aren’t defined. Therefore, the word extracting unit 103 doesn’t extract “saw” as a noun, “see” as a noun.

[0072] Further, in the example illustrated in FIG. 4, as notations (character strings) for defining the “colloquial expression” category of document types, three notation character strings are defined. That is, “Right?”, “Oh, yes”, and “Could be are define.

[0073] In step S23, the word extracting unit 103 determines whether or not “notation”, or “notation” and “part of speech” of the word selected in step S22 correspond to (match with) “notation”, or “notation” and “part of speech” of the character strings defined by the document type selecting unit 102.

[0074] When “notation”, or “notation” and “part of speech” of the selected word correspond to “notation”, or “notation” and “part of speech” of the character string set in the “document type” (Yes in step S23), in step S24, the selected word is set as the word set generating target, and the appearance frequency (the number of appearance times) of the selected word is calculated.

[0075] FIG. 5 is a diagram illustrating an example of a calculation result of the appearance frequency according to an exemplary embodiment of the present invention.

[0076] In the example illustrated in FIG. 5, the “appearance frequency” represents the number of times that the same word is selected in step S22. In this example, a numerical value “32” is given as the appearance frequency of a word having “notation” of “same period”.

[0077] Similarly, “31” is obtained as the appearance frequency of “precedingyear”, and “19” is obtained as the appearance frequency of “decrease”. “16” is obtained as the appearance frequency of “business”, and “16” is obtained as the appearance frequency of “sale” (the rest is similar to above).

[0078] Meanwhile, when “notation”, or “notation” and “part of speech” of the selected word do not correspond to “notation”, or “notation” and “part of speech” of the character strings set in the “document type” (No in step S23), in step S25, the word extracting unit 103 determines whether or not “part-of-speech” information of the selected word represents a noun.

[0079] When the “part-of-speech” information of the selected word does not represent a noun (No in step S25), in step S26, the corresponding word is excluded from the selection target of the word set.

[0080] However, when the “part-of-speech” information of the selected word represents a noun (Yes in step S25), in step S24, the appearance frequency of the corresponding word is calculated as illustrated in FIG. 5.

**[0081]** Thereafter, in step S27, the word extracting unit 103 determines whether or not the determination as to whether or not the selected word is to be used as the word set generation target has been made on all words obtained as the morphological analysis result. When a non-selected word still remains (No in step S27), the process returns to step S22. Another word is selected, and the subsequent process is repeated.

**[0082]** Meanwhile, when the word extracting unit 103 has completed the selection on all words obtained as the morphological analysis result (Yes in step S27), in step S28, the word set generating unit 104 calculates a word information quantity  $I_x$  of each word based on the appearance frequency calculated in step S24 according to Formula 1.

$$I_x = \frac{T_x}{\sum_{x=1} T_x} \times 100 \quad \{\text{Formula 1}\}$$

Here, the calculation of the word information quantity  $I_x$  in step S28 is performed on each word that has not been excluded in step S26.

**[0083]** In Formula 1,  $T_x$  represents the power of the appearance frequency of each word. A word having a large word information quantity  $I_x$  value represents the feature of the language model generating task corpus 100, that is, it is recognized as a word having large similarity to the language model generating task corpus 100.

**[0084]** In step S29, the word set generating unit 104 extracts a predetermined number of words having the word information quantity  $I_x$  value larger than a predetermined threshold value. The words may be extracted in the descending order of the word information quantity  $I_x$ , or a predetermined number of words may be randomly selected.

**[0085]** In step S30, the word set generating unit 104 generates a word set based on the extracted words as illustrated in FIG. 6.

**[0086]** FIG. 6 is a diagram illustrating an example of a word set according to the first exemplary embodiment of the present invention.

**[0087]** In the example illustrated in FIG. 6, a combination of five words of “business performance”, “preceding year”, “same period”, “field”, and “sales” is generated as the word set. The five words are high in the appearance frequency as illustrated in FIG. 5 and recognized as expressly representing the feature of the language model generating task corpus.

**[0088]** The word extracting device 200 performs the above described word set generating process and then transmits the generated word set to the web page acquiring unit 105.

**[0089]** As described above, the web page acquiring unit 105 inputs the word set to the retrieval engine as the retrieval query and obtains the retrieval result relevant to the word set. The web page formatting unit 106 removes any symbols or HTML tags from the web pages acquired by the web page acquiring unit 105 and extracts the sentence portion. The language model generating unit 107 generates the language model using the sentence portion extracted by the web page formatting unit 106 as the language model generating corpus for voice recognition.

**[0090]** As described above, according to the language model generating system according to the present embodiment, the narrowing-down retrieval is performed using words obtained by the word extracting device 200, and thus the

retrieval result list of web pages including linguistic expressions having high similarity or high relevance to the voice recognition target task can be displayed in the upper ranks of the retrieval result. The web page data is acquired using links from the retrieval result list page, and thus the high-accuracy language model generation corpus can be obtained.

**[0091]** Next, a description will be made in connection with a language model generating system according to another exemplary embodiment of the present invention. In the following exemplary embodiment, parts corresponding to the first exemplary embodiment are denoted by corresponding reference numerals, and a detailed description will not be repeated.

#### Second Exemplary Embodiment

**[0092]** A language model generating system according to a second exemplary embodiment of the present invention has the same configuration as in the first exemplary embodiment illustrated in FIG. 1, and thus redundant descriptions will not be repeated.

**[0093]** In the word set generating process according to the first exemplary embodiment, the word set combination of five words is generated as one set as an example. In the present embodiment, the number of words in the word set and the number of sets to generate are adjusted according to a sentence quantity or the content of the language model generating task corpus 100, so that the web page is more efficiently selected.

**[0094]** FIG. 7 is a flowchart illustrating an example of a word set generating process according to the second exemplary embodiment of the present invention.

**[0095]** In a word extracting process according to the present embodiment, in step S71, first the corpus analyzing unit 101 performs morphological analysis on the description content of the language model generating task corpus 100 and acquires a morphological analysis result. For example, the morphological analysis result includes information representing “notation (character string)”, and “part of speech” of words configuring a text as illustrated in FIG. 3, similarly to the first exemplary embodiment.

**[0096]** In step S72, the word extracting unit 103 selects a word from the morphological analysis result illustrated in FIG. 3 in a predetermined order. The word extracting unit 103 determines whether or not the selected word is to be used as a word set generation target word by a subsequent process. For example, the selection order may be based on an appearance order inside the text.

**[0097]** In step S73, the word extracting unit 103 determines whether or not “notation”, or “notation” and “part of speech” of the selected word correspond to (match with) information of “notation”, or “notation” and “part of speech” of any character string in the document type defined by the document type selecting unit 102. The definition of the document type by the document type selecting unit 102 may be made, for example, as illustrated in FIG. 4, similarly to the first exemplary embodiment.

**[0098]** In step S73, the word extracting unit 103 determines whether or not “notation”, or “notation” and “part of speech” of the word selected in step S72 correspond to “notation”, or “notation” and “part of speech” of the character strings defined by the document type selecting unit 102.

**[0099]** When “notation”, or “notation” and “part of speech” of the selected word correspond to “notation”, or “notation” and “part of speech” of the character string set in the “docu-

ment type" (Yes in step S73), in step S74, the selected word is set as the word set generating target, and the appearance frequency (the number of appearance times) of the selected word is calculated. An example of the appearance frequency of each word is illustrated in FIG. 5, similarly to the first exemplary embodiment.

[0100] Meanwhile, when "notation", or "notation" and "part of speech" of the selected word do not correspond to "notation", or "notation" and "part of speech" of the character string set in the "document type" (No in step S73), in step S75, the word extracting unit 103 determines whether or not the "part-of-speech" information of the selected word represents a noun.

[0101] When the "part-of-speech" information of the selected word does not represent a noun (No in step S75), in step S76, the corresponding word is excluded from the selection target of the word set.

[0102] However, when the "part-of-speech" information of the selected word represents a noun (Yes in step S75), in step S74, the appearance frequency of the corresponding word is calculated as illustrated in FIG. 5.

[0103] Thereafter, in step S77, the word extracting unit 103 determines whether or not the determination as to whether or not the selected word is to be used as the word set generation target has been made on all words obtained as the morphological analysis result. When a non-selected word still remains (No in step S77), the process returns to step S72. Another word is selected, and the subsequent process is repeated.

[0104] Meanwhile, when the word extracting unit 103 has completed the selection on all words obtained as the morphological analysis result (Yes in step S77), in step S78, the word set generating unit 104 calculates a word information quantity  $I_x$  of each word based on the appearance frequency calculated in step S24 according to Formula 1, similarly to the first exemplary embodiment.

[0105] In step S79, the word set generating unit 104 extracts a predetermined number of words having a word information quantity  $I_x$  value larger than a predetermined threshold value. The words may be extracted in descending order of the word information quantity  $I_x$ , or a predetermined number of words may be randomly selected.

[0106] In step S80, the word set generating unit 104 generates  $m$  word sets each of which has  $n$  words from the extracted words.

[0107] Values of the variables  $m$  and  $n$  may be set by the user through the input device connected to the word extracting device 200, or the values of the variables  $m$  and  $n$  may be set based on the analysis result of the sentence quantity of the language model generating task corpus 100 by the corpus analyzing unit 101.

[0108] For example, when the sentence quantity of the language model generating corpus 100 is small, since the number of words extracted in step S80 decreases, a narrowing-down retrieval range gets narrow. Thus, it is difficult to acquire web pages similar to the voice recognition target, and a degree of accuracy of voice recognition is not improved. In this case, for example, five word sets each of which has three words are generated as illustrated in FIG. 8 so that a variation can be given to the retrieval. As a result, the narrowing-down retrieval range can get broaden.

[0109] FIG. 8 is a diagram illustrating an example of a word set according to the second exemplary embodiment of the present invention.

[0110] In the example illustrated in FIG. 8, the combination of three words "business performance", "preceding year", and "service" is generated as word set 1. Further, word set 2 includes three words "field", "business performance", and "sales", and word set 3 includes three words "quarter", "sales", and "business performance". Further, word set 4 includes three words "service", "sales", and "business performance", and word set 5 includes three words "quarter", "service", and "investment". For example, the words included in the word sets 1 to 5 are randomly selected from among words having the large word information quantity extracted in step S80. Alternatively, the user may designate a selection method in advance.

[0111] As a result, the narrowing-down retrieval range can be broadened, and web pages causing an improved degree of accuracy of voice recognition can be acquired.

[0112] The word extracting device 200 performs the above described word set generating process and then transmits the generated word set to the web page acquiring unit 105.

[0113] As described above, the web page acquiring unit 105 inputs the word set to the retrieval engine as the retrieval query and obtains the retrieval result relevant to the word set. The web page formatting unit 106 removes any symbols or HTML tags from the web pages acquired by the web page acquiring unit 105 and extracts the sentence portion. The language model generating unit 107 generates the language model using the sentence portion extracted by the web page formatting unit 106 as the language model generating corpus for voice recognition.

[0114] As described above, according to the language model generating system according to the present embodiment, the narrowing-down retrieval is performed using a plurality of word sets obtained by the word extracting device 200. Thus, even when the sentence quantity included in the language model generating task corpus 100 is small, the high-accuracy language model generation corpus can be obtained.

[0115] That is, according to the language model generating system according to the present embodiment, an appropriate word can be extracted even from a small amount of language model generating task corpus 100, and a language corpus similar to the language model generating task corpus 100 and including the content of the voice recognition target can be efficiently collected from web pages on the Internet.

[0116] According to the exemplary embodiments of the present invention, information necessary for generating the language model in the voice recognition technique can be efficiently collected sufficiently from the Internet, and the language model having a high degree of accuracy of voice recognition can be generated.

[0117] According to the exemplary embodiments of the present invention, the language model similar to the target task or the language model capable of supporting even the utterance style such as the regional dialect can be efficiently generated only by preparing a small amount of corpus data including the content desired to voice-recognize. Further, the cost for preparing a large amount of data necessary for generating the high-accuracy language model can be reduced.

[0118] In the exemplary embodiments of the present invention, the word information quantity  $I_x$  representing the feature of the sentence is calculated, and the feature word in each text is extracted based on the word information quantity  $I_x$ . Further, the word representing the feature of the sentence desired to voice-recognize is preferentially extracted, and

then the language model is generated. Thus, the high-accuracy language model can be generated.

**[0119]** Further, in the word set generating process, as for the word that is not defined in the document type, a noun having the same meaning of the word is extracted. Thus, disambiguation can be made at the time of web retrieval, and the desired information can be efficiently acquired.

**[0120]** The word extracting device **200** and the language model generating device **300** according to the exemplary embodiments of the present invention are implemented by hardware. However, the functions of the word extracting device **200** and the language model generating device **300** according to the exemplary embodiments of the present invention may be implemented such that a program causing a computer to function as the word extracting device **200** and the language model generating device **300** is stored in a computer-readable storage medium such as a hard disc device, semiconductor memory or the like, and the program is read and executed by the computer.

**[0121]** By implementing the word extracting device **200** according to the exemplary embodiments of the present invention by hardware and combining the word extracting device **200** with an existing language model generating means, a function equivalent to the language model generating device **300** may be provided.

**[0122]** In the above embodiments, the language model generating task corpus **100** is defined as a set of web documents, but the content of the corpus is not limited thereto.

**[0123]** The above embodiments are exemplary embodiments of the present invention. However, the scope of the present invention is not limited to the above embodiments, and various modifications can be made within the scope not departing from the gist of the present invention. For example, the present invention can apply to Spanish, Chinese or the like.

**[0124]** All or some of the above exemplary embodiments can be described as in the following notes, but the present invention is not limited to the following supplementary notes.

**[0125]** (Supplementary note 1) A language model generating device, comprising:

**[0126]** a corpus analyzer which analyzes text in a corpus including a set of world wide web (web) pages;

**[0127]** an extractor which extracts at least one word appropriate for a document type set according to a voice recognition target based on an analysis result by the corpus analyzer;

**[0128]** a word set generator which generates a word set from the at least one word extracted by the extractor;

**[0129]** a web page acquiring device which causes a retrieval engine to perform a retrieval process using the word set generated by the word set generator as a retrieval query of the retrieval engine on the Internet and acquires a web page related to the word set from the retrieval result; and

**[0130]** a language model generating device which generates a language model for voice recognition from the web page.

**[0131]** (Supplementary note 2) The language model generating device according to Supplementary note 1,

**[0132]** further comprising a selector which sets a character string defining the document type,

**[0133]** wherein the extractor compares each word obtained from the analysis result with the character string and extracts the word when the word corresponds to the character string.

**[0134]** (Supplementary note 3) The language model generating device according to Supplementary note 2,

**[0135]** wherein the character string includes information of notation and a part of speech of the character string,

**[0136]** the corpus analyzer performs morphological analysis of the text in the corpus and detects information of notation and a part of speech for each word in the text, and

**[0137]** the extracting unit compares the notation, or the notation and part of speech of each word in the text with the notation, or the notation and part of speech of the character string and extracts a word corresponding to the notation, or the notation and part of speech of the character string.

**[0138]** (Supplementary note 4) The language model generating device according to Supplementary note 3,

**[0139]** wherein the extractor determines whether or not a part of speech of a word not corresponding to the character string is a noun and excludes the word from an extraction target when the part of speech of the word is not the noun.

**[0140]** (Supplementary note 5) The language model generating device according to any one of Supplementary notes 1 to 4,

**[0141]** wherein the word set generator calculates a word information quantity, representing similarity to the corpus on each word extracted by the extractor and generates the word set from a word whose word information quantity value is at the predetermined value or more.

**[0142]** (Supplementary note 6) The language model generating device according to any one of Supplementary notes 1 to 5,

**[0143]** wherein the word set generator generates a predetermined number of word sets each of which includes a predetermined number of words from words extracted by the extractor.

**[0144]** (Supplementary note 7) The language model generating device according to any one of Supplementary notes 1 to 6,

**[0145]** wherein the word set generator calculates a word information quantity, representing similarity to the corpus on each word extracted by the extractor and generates the predetermined number of word sets each of which includes the predetermined number of words randomly from words whose word information quantity value is at a predetermined value or more.

**[0146]** (Supplementary note 8) A language model generating method, comprising the steps of:

**[0147]** analyzing a text in a corpus including a set of world wide web (web) pages;

**[0148]** extracting at least one word appropriate for a document type set according to a voice recognition target based on an analysis result;

**[0149]** generating a word set from at least one extracted word;

**[0150]** causing a retrieval engine to perform a retrieval process using the generated word set as a retrieval query of the retrieval engine on the Internet and acquiring a web page related to the word set from the retrieval result; and

**[0151]** generating a language model for voice recognition from the acquired web page.

**[0152]** (Supplementary note 9) The language model generating method according to Supplementary note 8,

**[0153]** further comprising a step of setting a character string defining the document type,

[0154] wherein the at least one word is extracted such that each word obtained from the

[0155] analysis result is compared with the character string, and the word is extracted when the word corresponds to the character string.

[0156] (Supplementary note 10) The language model generating method according to Supplementary note 9,

[0157] wherein the character string includes information of notation and a part of speech of the character string,

[0158] the corpus is analyzed such that morphological analysis of the text in the corpus is performed, and information of notation and a part of speech is detected for each word in the text, and

[0159] the at least one word is extracted such that the notation, or the notation and part of speech of each word in the text are compared with the notation, or the notation and part of speech of the character string, and a word corresponding to the notation, or the notation and part of speech of the character string is extracted.

[0160] (Supplementary note 11) The language model generating method according to Supplementary note 10,

[0161] wherein the at least one word is extracted such that it is determined whether or not a part of speech of a word not corresponding to the character string is a noun, and the word is excluded from an extraction target when the part of speech of the word is not the noun.

[0162] (Supplementary note 12) The language model generating method according to any one of Supplementary notes 8 to 11,

[0163] wherein the word set is generated such that a word information quantity representing similarity to the corpus is calculated on each extracted word, and the word set is generated from a word whose word information quantity value is a predetermined value or more.

[0164] (Supplementary note 13) The language model generating method according to any one of Supplementary notes 8 to 12,

[0165] wherein the word set is generated such that a predetermined number of word sets each of which includes a predetermined number of words are generated from extracted words.

[0166] (Supplementary note 14) The language model generating method according to any one of Supplementary notes 8 to 13,

[0167] wherein the word set is generated such that a word information quantity representing similarity to the corpus is calculated on each extracted word, and the predetermined number of word sets each of which includes the predetermined number of words are generated randomly from words whose word information quantity value is at a predetermined value or more.

[0168] (Supplementary note 15) A non-transitory computer-readable recording medium storing a language model generating program used in a computer of a language model generating device and causing the computer to execute:

[0169] a step of analyzing a text in a corpus including a set of world wide web (web) pages;

[0170] a step of extracting at least one word appropriate for a document type set according to a voice recognition target based on an analysis result;

[0171] a step of generating a word set from at least one extracted word;

[0172] a step of causing a retrieval engine to perform a retrieval process using the generated word set as a retrieval

query of the retrieval engine on the Internet and acquiring a web page related to the word set from the retrieval result; and [0173] a step of generating a language model for voice recognition from the acquired web page.

[0174] (Supplementary note 16) The non-transitory computer-readable recording medium according to Supplementary note 15,

[0175] the program further causing the computer to execute a step of setting a character string defining the document type,

[0176] wherein in the step of extracting the at least one word, each word obtained from the analysis result is compared with the character string, and the word is extracted when the word corresponds to the character string.

[0177] (Supplementary note 17) The non-transitory computer-readable recording medium according to Supplementary note 16,

[0178] wherein the character string includes the information of notation and a part of speech of the character string,

[0179] in the step of analyzing the text in the corpus, morphological analysis of the text in the corpus is performed, and information of notation and a part of speech is detected for each word in the text, and

[0180] in the step of extracting the at least one word, the notation, or the notation and part of speech of each word in the text are compared with the notation, or the notation and part of speech of the character string, and a word corresponding to the notation, or the notation and part of speech of the character string is extracted.

[0181] (Supplementary note 18) The non-transitory computer-readable recording medium according to Supplementary note 17,

[0182] wherein in the step of extracting the at least one word, it is determined whether or not a part of speech of a word not corresponding to the character string is a noun, and the word is excluded from an extraction target when the part of speech of the word is not a noun.

[0183] (Supplementary note 19) The non-transitory computer-readable recording medium according to any one of Supplementary notes 15 to 18,

[0184] wherein in the step of generating the word set, a word information quantity representing similarity to the corpus is calculated on each extracted word, and the word set is generated from a word whose word information quantity value is at a predetermined value or more.

[0185] (Supplementary note 20) The non-transitory computer-readable recording medium according to any one of Supplementary notes 15 to 19,

[0186] wherein in the step of generating the word set, a predetermined number of word sets each of which includes a predetermined number of words are generated from extracted words .

[0187] (Supplementary note 21) The non-transitory computer-readable recording medium according to any one of Supplementary notes 15 to 20,

[0188] wherein in the step of generating the word set, a word information quantity representing similarity to the corpus is calculated on each extracted word, and the predetermined number of word sets each of which includes the predetermined number of words are generated randomly from words whose word information quantity value is at a

1. A word retrieval device, comprising:

an analyzer configured to carry out morphological analysis of text;



an extractor configured to extract at least one word representing a feature of the text from words of an morphological analysis result by the analyzer, using word information quantity on the words;

a retrieval device configured to retrieve text related to the at least one word from a web page using the at least one word as a retrieval query.

2. The word retrieval device according to claim 1, wherein the word information quantity is represented by  $I_x$ , where  $T_x$  represents a power of an appearance frequency of each word and  $I_x$  is defined as follows:

$$I_x = \frac{T_x}{\sum_{x=t} T_x} \times 100$$

3. The word retrieval device according to claim 1, the morphological analysis includes detection of information of notation or the notation and a part of speech for each word in the text.

4. The word retrieval device according to claim 1, further comprising a selector which sets a character string defining a document type,

wherein the extractor compares each word obtained from the morphological analysis result with the character string and extracts the word when the word corresponds to the character string.

5. The word retrieval device according to claim 4, wherein the character string includes information of notation and a part of speech of the character string, and the extractor compares the notation or the notation and part of speech of each word in the text, with the notation or the notation and part of speech of the character string and extracts a word corresponding to the notation or the notation and part of speech of the character string.

6. The word retrieval device according to claim 5, wherein the extractor determines whether or not a part of speech of a word not corresponding to the character string is a noun and excludes the word from an extraction target when the part of speech of the word is not the noun.

7. A word retrieval method of a word retrieval device, the method comprising:

carrying out morphological analysis of text;  
extracting at least one word representing a feature of the text from words of an morphological analysis result by the analyzer, using word information quantity on the words;

retrieving text related to the at least one word from a web page using the at least one word as a retrieval query.

8. The word retrieval method according to claim 7, wherein the word information quantity is represented by  $I_x$ , where  $T_x$  represents a power of an appearance frequency of each word and  $I_x$  is defined as follows:

$$I_x = \frac{T_x}{\sum_{x=t} T_x} \times 100$$

9. The word retrieval method according to claim 7, the morphological analysis includes detection of information of notation or the notation and a part of speech for each word in the text.

10. The word retrieval method according to claim 7, further comprising:

setting a character string defining a document type;  
comparing each word obtained from the morphological analysis result with the character string and extracting the word when the word corresponds to the character string.

11. The word retrieval method according to claim 10, wherein the character string includes information of notation and a part of speech of the character string, the comparing being comparing of the notation or the notation and part of speech of each word in the text, with the notation or the notation and part of speech of the character string,

the extracting being extracting of a word corresponding to the notation or the notation and part of speech of the character string.

12. The word retrieval method according to claim 11, further comprising:

determining whether or not a part of speech of a word not corresponding to the character string is a noun, and excluding the word from an extraction target when the part of speech of the word is not the noun.

13. A non-transitory computer-readable recording medium storing a word retrieving program used in a computer of a word retrieval device and causing the computer to execute a method comprising:

carrying out morphological analysis of text;  
extracting at least one word representing a feature of the text from words of an morphological analysis result by the analyzer, using word information quantity on the words;  
retrieving text related to the at least one word from a web page using the at least one word as a retrieval query.

14. The non-transitory computer-readable recording medium according to claim 13,

wherein the word information quantity is represented by  $I_x$ , where  $T_x$  represents a power of an appearance frequency of each word and  $I_x$  is defined as follows:

$$I_x = \frac{T_x}{\sum_{x=t} T_x} \times 100$$

15. The non-transitory computer-readable recording medium according to claim 13, the morphological analysis includes detection of information of notation or the notation and a part of speech for each word in the text.

16. The non-transitory computer-readable recording medium according to claim 13, further comprising:

setting a character string defining a document type;  
comparing each word obtained from the morphological analysis result with the character string and extracting the word when the word corresponds to the character string.

17. The non-transitory computer-readable recording medium according to claim 16,

wherein the character string includes information of notation and a part of speech of the character string,

the comparing being comparing of the notation or the notation and part of speech of each word in the text, with the notation or the notation and part of speech of the character string,

the extracting being extracting of a word corresponding to the notation or the notation and part of speech of the character string.

**18.** The non-transitory computer-readable recording medium according to claim 17, further comprising:

determining whether or not a part of speech of a word not corresponding to the character string is a noun, and excluding the word from an extraction target when the part of speech of the word is not the noun.

\* \* \* \* \*