

(19)



(11)

EP 2 289 007 B1

(12)

EUROPEAN PATENT SPECIFICATION

(45) Date of publication and mention of the grant of the patent:
22.04.2015 Bulletin 2015/17

(51) Int Cl.:
G06F 17/30^(2006.01)

(21) Application number: **09730808.4**

(86) International application number:
PCT/US2009/036597

(22) Date of filing: **10.03.2009**

(87) International publication number:
WO 2009/126394 (15.10.2009 Gazette 2009/42)

(54) SEARCH RESULTS RANKING USING EDITING DISTANCE AND DOCUMENT INFORMATION

SUCHERGEBNISEINSTUFUNG UNTER VERWENDUNG VON EDITIERDISTANZ- UND DOKUMENTINFORMATIONEN

CLASSEMENT DE RÉSULTATS DE RECHERCHE GRÂCE AU CALCUL D'UNE DISTANCE D'ÉDITION ET À L'EXTRACTION D'INFORMATIONS DOCUMENTAIRES

(84) Designated Contracting States:
AT BE BG CH CY CZ DE DK EE ES FI FR GB GR HR HU IE IS IT LI LT LU LV MC MK MT NL NO PL PT RO SE SI SK TR

- **MEYERZON, Dmitriy**
Redmond, WA 98052-6399 (US)
- **XU, Jun**
Redmond, WA 98052-6399 (US)

(30) Priority: **11.04.2008 US 101951**

(74) Representative: **Grünecker Patent- und Rechtsanwälte**
PartG mbB
Leopoldstraße 4
80802 München (DE)

(43) Date of publication of application:
02.03.2011 Bulletin 2011/09

(73) Proprietor: **Microsoft Technology Licensing, LLC**
Redmond, WA 98052 (US)

(56) References cited:
KR-A- 20030 080 826 US-A1- 2004 141 354
US-A1- 2004 141 354 US-A1- 2006 004 732
US-A1- 2006 069 982 US-A1- 2006 149 723
US-A1- 2006 294 100 US-A1- 2008 005 068
US-B2- 6 738 764

- (72) Inventors:
- **TANKOVICH, Vladimir**
Redmond, WA 98052-6399 (US)
 - **LI, Hang**
Redmond, WA 98052-6399 (US)

EP 2 289 007 B1

Note: Within nine months of the publication of the mention of the grant of the European patent in the European Patent Bulletin, any person may give notice to the European Patent Office of opposition to that patent, in accordance with the Implementing Regulations. Notice of opposition shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

Description

BACKGROUND

5 **[0001]** In a typical search engine service a user can enter a query by selecting the topmost relevant documents out of an indexed collection of URLs (universal resource locators) that match the query. To serve the queries quickly the search engine utilizes one or more methods (e.g., an inverted index data structure) that map keywords to documents. For example, a first step performed by the engine can be to identify the set of candidate documents that contain the keywords specified by the user query. These keywords can be located in the document body or the metadata, or additional

10 **[0002]** In a large index collection the cardinality of the candidate document set can be big, depending on the commonality of the query terms (e.g., potentially millions). Instead of returning the entire set of candidate documents the search engine performs a second step of ranking of the candidate documents with respect to relevance. Typically, the search engine utilizes a ranking function to predict the degree of relevance of a document to a particular query. The ranking function takes multiple features from the document as inputs and computes a number that allows the search engine to sort the documents by predicted relevance.

15 **[0003]** The quality of the ranking function with respect as to how accurately the function predicts relevance of a document is ultimately determined by the user satisfaction with the search results or how many times on average the user finds the answer to the question posed. The overall user satisfaction with the system can be approximated by a single number (or metric), because the number can be optimized by varying the ranking function. Usually, the metrics are computed over a representative set of queries that are selected up front by random sampling of the query logs, and involve assigning relevance labels to each result returned by the engine for each of the evaluation queries. However, these processes for document ranking and relevance are still inefficient in providing the desired results.

20 **[0004]** US 2006/0004732 A1 relates to search engine methods and systems for generating relevant search results and advertisements. This document discloses search engine methods and systems for generating relevant search results and displaying relevant advertisements with those search results. This document discloses methods and systems for modelling and storing data in neutral forms, and then applying topification techniques to the data to generate search results that are relevant to a particular user's search request. This document also applies topification and relevancy methods to associate ads that are relevant to a user with the search results for display.

SUMMARY

25 **[0005]** It is the object of the present invention to integrate uniform resource locator information into the determination of ranking of search results.

30 **[0006]** This object is solved by the subject matter of the independent claims.

[0007] Preferred embodiments are defined by the dependent claims.

35 **[0008]** The following presents a simplified summary in order to provide a basic understanding of some novel embodiments described herein. This summary is not an extensive overview, and it is not intended to identify key/critical elements or to delineate the scope thereof. Its sole purpose is to present some concepts in a simplified form as a prelude to the more detailed description that is presented later.

40 **[0009]** The architecture provides a mechanism for extracting document information from documents received as search results based on a query string and computing an edit distance between a data string and the query string. The data string can be a short and accurate description of the document obtained from document information such as TAUC (title, anchor text, URL (uniform resource locator), and clicks), for example. The edit distance is employed in determining relevance of the document as part of result ranking. The mechanism improves the relevance of search results ranking by employing a set of proximity-related features to detect near-matches of a whole query or part of the query.

45 **[0010]** The edit distance is processed to evaluate how close the query string is to a given data stream that includes the document information. The architecture includes the index-time splitting of compound terms in the URL to allow the more effective discovery of query terms. Additionally, index-time filtering of anchor text is utilized to find the top N anchors of one or more of the document results. Using the TAUC information can be input to a neural network (e.g., 2-layer) to improve relevance metrics for ranking the search results.

50 **[0011]** To the accomplishment of the foregoing and related ends, certain illustrative aspects are described herein in connection with the following description and the annexed drawings. These aspects are indicative, however, of but a few of the various ways in which the principles disclosed herein can be employed and is intended to include all such aspects and equivalents. Other advantages and novel features will become apparent from the following detailed description when considered in conjunction with the drawings.

BRIEF DESCRIPTION OF THE DRAWINGS

[0012]

- 5 FIG. 1 illustrates a computer-implemented relevance system.
 FIG. 2 illustrates a flow chart of an exemplary the matching algorithm for computing edit distance.
 FIG. 3 illustrates processing and generating edit distance values based on a query string and data string using the modified edit distance and matching algorithm.
 FIG. 4 illustrates another example of processing and generating edit distance values based on a query string and data string using the modified edit distance and matching algorithm.
 10 FIG. 5 illustrates a computer-implemented relevance system that employs a neural network to assist in generating a relevance score for the document.
 FIG. 6 illustrates the types of data that can be employed in the document information for determining the edit distance between the query string and the data string.
 15 FIG. 7 illustrates an index-time processing data flow.
 FIG. 8 illustrates a block diagram showing inputs to the neural network from the index process of FIG. 7 for result ranking.
 FIG. 9 illustrates an exemplary system implementation of a neural network, edit distance inputs and raw feature inputs for computing generating search results.
 20 FIG. 10 illustrates a method of determining document relevance of a document result set.
 FIG. 11 illustrates a method of computing relevance of a document.
 FIG. 12 illustrates a block diagram of a computing system operable to execute edit distance processing for search result ranking using TAUC features in accordance with the disclosed architecture.

25 DETAILED DESCRIPTION

[0013] The disclosed architecture improves the relevance of search results ranking by implementing a set of proximity-related features to detect near-matches of a whole query or matches with accurate metadata about the document, such as titles, anchors, URLs, or clicks. For example, consider a query "company store", a document title "company store online" of a first document and a document title "new NEC LCD monitors in company store" of a second document. Assuming other properties is same for both the first and second documents, the architecture assigns a score for a document based on how much editing effort is devoted to make a chosen stream match the query. In this example, the document title is selected for evaluation. The title of first document requires only one delete operation (delete the term "online") to make a full match, while the title of second document requires five deletes (delete the terms "new", "NEC", "LCD", "monitors" and "in"). Thus, the first document is computed to be more relevant.

[0014] The title is one element of TAUC (title, anchor, URL, and clicks) document information for which processing can be applied to some streams of data (e.g., a URL) so that query terms can be found from compound terms. For example, consider, again, the query "company store", and the URL is "www.companystore.com". The result is that the URL is split into four parts (or terms): "www", "company", "store", and "com".

40 **[0015]** Reference is now made to the drawings, wherein like reference numerals are used to refer to like elements throughout. In the following description, for purposes of explanation, numerous specific details are set forth in order to provide a thorough understanding thereof. It may be evident, however, that the novel embodiments can be practiced without these specific details. In other instances, well-known structures and devices are shown in block diagram form in order to facilitate a description thereof.

45 **[0016]** FIG. 1 illustrates a computer-implemented relevance system 100. The system 100 includes a processing component 102 for extracting document information 104 from a document 106 received as search results 108 based on a query string 110. The system 100 can also include a proximity component 112 for computing the edit distance 114 between a data string 116 derived from the document information 104 and the query string 110. The edit distance 114 is employed in determining relevance of the document 106 as part of the search results 108.

50 **[0017]** The document information 104 employed to generate the data string 116 can include title information (or characters), link information (e.g., URL characters), click stream information, and/or anchor text (or characters), for example. The processing component 102 splits compound terms of the document information 104 at index time to compute the edit distance 114. The processing component 102 also filters document information such as anchor text at index time to compute a top-ranked set of anchor text.

55 **[0018]** The computing of the edit distance 114 is based on insertion and deletion of terms to increase proximity (bring closer) between the data string 116 and the query string 110. The computing of the edit distance 114 can also be based on costs associated with insertion and deletion of terms to increase the proximity (bring closer) between the data string 116 and the query string 110.

[0019] Consider a scenario of generating a data string 116 (e.g., TAUC) based on the insertion and/or deletion of terms from the query string 110. This term processing can be performed according to four operations: insert a non-query word into the query string 110; insert a query term into the query string 110; delete a TAUC term from the query string 110; and/or, delete a non-TAUC term from the query string 110.

[0020] The edit distance 114 is based on the insertion and deletion operations, but not substitution. There can be two types of cost defined for insertion. Consider a scenario of generating the data string 116 from the query string 110. In the generation, a word can be inserted into the query string 110, which exists in the original query string 110, then the cost is defined as one; otherwise, the cost is defined as $w_1 (\geq 1)$. Here, w_1 is a weighting parameter that is tuned. For example, if the query string 110 is AB, then the cost of generating the data string of ABC is higher than that of the data string ABA. The intuition is that by inserting "irrelevant words" into the data string 116 makes the entire data string 116 (e.g., TAUC) more irrelevant.

[0021] There can be two types of cost for deletion. Again, consider a scenario of generating the data string 116 from the query string 110. When deleting a term in the query string 110, which term exists in the original data string 116, then the cost is defined as one; otherwise, the cost is defined as $w_2 (\geq 1)$.

[0022] Another type of cost is a position cost. If a deletion or insertion occurs at the first position of the data string 116, then there is an additional cost ($+w_3$). The intuition is that a matching at the beginning of the two strings (query string 110 and data string 116) is given greater importance than matches later in the strings. Consider the following example where the query string 110 is "cnn", and the data string 116 is title = "cnn.com - blur blur". If insertion and deletion occur at the first position, it can significantly reduce the effectiveness of the solution.

[0023] FIG. 2 illustrates a flow chart of an exemplary the modified matching algorithm 200 for computing edit distance. While, for purposes of simplicity of explanation, the one or more methodologies shown herein, for example, in the form of a flow chart or flow diagram, are shown and described as a series of acts, it is to be understood and appreciated that the methodologies are not limited by the order of acts, as some acts may, in accordance therewith, occur in a different order and/or concurrently with other acts from that shown and described herein. For example, those skilled in the art will understand and appreciate that a methodology could alternatively be represented as a series of interrelated states or events, such as in a state diagram. Moreover, not all acts illustrated in a methodology may be required for a novel implementation.

[0024] At 200, elements of the query string and the data (or target) string are enumerated. This is accomplished by setting n to be the length of the query string (where each term in query string is $s[i]$), and setting m to be the length of the target (or data) string (where each term in target string is denoted $t[j]$). At 202, a matrix is constructed that contains $0 \dots m$ rows and $0 \dots n$ columns (where each term in the matrix is denoted as $d[j,i]$). At 204, the first row is initialized with a value that depends on the different cost of deletion and the first column is initialized with a value that depends on the different cost of insertion. At 206, if $n = 0$, return $d[m, 0]$ and exit, and if $m = 0$, return $d[0, n]$ and exit, as indicated at 208. At 210, each character of the query string is examined (i from 1 to n). At 212, each character of the target data string is examined (j from 1 to m). At 214, if the character string in the query string equals the character in the data string, flow is to 216 where the cost is zero and the next matrix cell is processed. In other words, if $s[i]$ equals $t[j]$, the cost is 0 and $d[j,i] = d[j-1,i-1]$.

[0025] If the character in the query string cell is not equal to the character in the data string cell, flow is from 214 to 218 where the current cell is set to the immediate cell above or immediate cell to the left, plus the insertion or deletion cost. In other words, if $s[i]$ is not equal to $t[j]$, set cell $d[j,i]$ of the matrix equal to the minimum of the cell immediately above plus corresponding insertion cost (represented $d[j-1,i] + \text{cost_insertion}$) or the cell immediately to the left plus corresponding deletion cost (represented $d[j,i-1] + \text{cost_deletion}$). At 220, the steps 210, 212, 214, 216 and 218 are iterated to completion. At 222, final cost found in cell $d[m, n]$ is output. Note that both the cost_insertion and the cost_deletion in the example have two kinds of values; for example $w_1=1$, $w_3=4$ for insertion cost, and $w_2=1$, $w_4=26$ for deletion cost.

[0026] In other words, $d[j,i]$ contains the edit distance between strings $s[0..i]$ and $t[0..j]$. $d[0,0] = 0$ by definition (no edits needed to make an empty string equal to empty string). $d[0, y] = d[0,y-1] + (w_2 \text{ or } w_4)$. If it is known how many edits are used to make the string $d[0,y-1]$, then $d[0,y]$ can be calculated as $d[0, y-1] + \text{cost of deleting current character from the target string}$, which cost can be w_2 or w_4 . The cost w_2 is used if the current character is present in both $s[0..n]$, $t[0..m]$; and w_4 , otherwise. $d[x, 0] = d[x-1,0] + (w_1 \text{ or } w_3)$. If it is known how many edits are used to make the string $d[x-1,0]$ then $d[x,0]$ can be calculated as $d[x-1,0] + \text{cost of insertion of the current character from } s \text{ to } t$, which cost can be w_1 or w_3 . The cost w_1 is used if the current character is present in both $s[0..n]$, $t[0..m]$; and w_3 , otherwise.

[0027] For each (j,i) , $d[j,i]$ can be equal to $d[j-1,i-1]$ if $s[i]=t[j]$. The edit distance can be computed between strings $t[j-1]$, $s[i-1]$, and if $s[i]=t[j]$, a common character can be appended to both strings to make the strings equal, without causing edits. Thus, there are three moves employed, where the move that provides the minimal edit distance for current $d[j,i]$ is selected. Put another way,

```

d[j,i] = min(
    d[j-1,i-1] if s[i]=t[j];
    d[j-1,i] + (w1, if s[j] is present in both strings; else, w3);
    d[j,i-1] + (w2, if t[i] is present in both strings; else, w4)
)

```

5
10
15
20
25

[0028] FIG. 3 illustrates processing and generating edit distance values based on a query string and data string using the modified edit distance and matching algorithm. The process involves one or more of left-right, top-down, and diagonal computations. A query string of terms "A B C" is processed against with a target data string of terms "C B A X" (where X denotes a term not in the query string). The process for computing an edit distance can be performed in different ways; however, the specific details for performing a modified version of an edit distance is different as computed according to the disclosed architecture. A 4x5 matrix 300 is constructed based on nxm where n = 3 for the query string and m = 4 for the data string. The query string 302 is placed along the horizontal axis and the target data string 304 is along the vertical axis of the matrix 300.

30
35

[0029] The description will use the matrix 300 denoted with four columns (0-3) and five rows (0-4). Applying the edit distance matching algorithm described in FIG. 2 from left to right beginning in row 0, column 0, the intersecting cell d[0,0] receives "0" since the compare of the empty cell of the query string ABC to the empty cell of the target data string CBAX does not cause insertion or deletion of a term to make the query string the same as the target data string. The "terms" are the same so the edit distance is zero.

40
45

[0030] Moving right to compare the A term of query string 302 to the empty cell of row 0 uses one deletion to make the strings the same; thus, the cell d[0,1] receives a value of "1". Moving right again to the column 2, the compare is now made between terms AB of the query string 302 to the empty cell of the target data string column. Thus, two deletions in the query string 302 are used to make the strings the same resulting in an edit distance of "2" being placed into cell d[0,2]. The same process applies to column 3 where the terms ABC of the query string 302 are compared to the empty cell in target string column, using three deletions to make the strings alike, resulting in an edit distance of "3" in the cell d[0,3].

50
55

[0031] Dropping down to row 1 and continuing left to right, the empty cell of the query string row is compared to the first term C of the target data string 304. One deletion is used to make the strings the same, with an edit distance of "1" in d[1,0]. Moving right to column 1, the compare is made between the A term of the query string 302 and the C term of the target data string 304. A deletion and insertion is used to make the strings alike, thus, a value of "2" is inserted into cell d[1,1]. Skipping to the last cell d[1,3], the matching process for matching ABC to C results in using two deletions for an edit distance of "2" in the cell d[1,3]. Moving to row 4 and column 3 for brevity and to find the overall edit distance, matching terms ABC to terms CBAX results in an edit distance of "8" in cell d[4,3] using insertion/deletion in the first term C of the target string for a value of "2", a value of "0" for the match between the B terms, an insertion/deletion for the match of the third terms C and A for a value of "2", an insertion of the term X for a value of "1" and a value of "3" for position cost, resulting in a final edit distance value of "8" in cell d[4,3].

60
65

[0032] FIG. 4 illustrates another example of processing and generating edit distance values based on a query string and target data string using the modified edit distance and matching algorithm. Here, a matrix 400 is generated for comparing a query string 402 of ABC to a target data string 404 of AB based on weightings for cost_insertion of w1=1, w3=4 for insertion cost, and w2=1 and w4=26 for deletion cost. In other words, working row 0 from left to right, matching term A of the query string 402 to the empty cell before the target string 404 results in one insertion in the target string 404 of the term A for a value of "1" cell d[0,1]. Matching terms AB of the query string 402 to the empty cell before the target string 404 results in two insertions in the target string 404 of the terms AB for a value of "2" cell d[0,2], and matching terms ABC of the query string 402 to the empty cell before the target string 404 results in the two insertions in the target string 404 of the terms AB value plus value w4=26 for the term C for a value of "28" in cell d[0,3], since the term C is not in both strings.

70
75

[0033] Working row 1 from left to right (understanding that d[1,0] = 1), matching term A of the query string 402 to the term A of the target string 404 results in equality in the target string 404 and the query string 402 for a value of "0" in cell d[1,1], by taking the value from d[j-1,i-1] = d[0,0] = "0". Matching terms AB of the query string 402 to the term A of the target string 404 results in one insertion in the target string 404 for the term B for a minimum value of "1" cell d[1,2]. Matching terms ABC of the query string 402 to term A of the target string 404 for the cell d[1,3] results in a minimum value associated with the value of d[j-1,i] = d[0,3] plus w3 for a value of "28" in cell d[1,3] compared to a value of d[j,i-1] = d[1,2] for 1 plus 26 for 27, since the term C is not in both strings, resulting in the minimum value of "27" in d[1,3].

80
85

[0034] Working row 2 from left to right, matching term A of the query string 402 to the terms AB of the target string 404 results in a deletion in the target string 404 for a value of "1" in cell d[2,1]. Matching terms AB of the query string 402 to the terms AB of the target string 404 for the distance in cell d[2,2] results in an equality, thereby pulling the value from d[j-1,i-1] = d[1,1] as the value "0" for cell d[2,2]. Matching terms ABC of the query string 402 to terms AB of the target string 404 for the cell d[2,3] results in a minimum value associated with the value of d[j-1,i] = d[1,3] = 27 plus w3=1

for a value of "28" compared to, since C is not in the target string (also based on a value of $d[i,j-1] = d[2,2] = 0$ plus 26 for 26, since the term C is not in both strings, for the minimum value of "26" in $d[2,3]$).

[0035] FIG. 5 illustrates a computer-implemented relevance system 500 that employs a neural network 502 to assist in generating a relevance score 504 for the document 106. The system 500 includes the processing component 102 for extracting document information 104 from the document 106 received as the search results 108 based on the query string 110, and the proximity component 112 for computing the edit distance 114 between the data string 116 derived from the document information 104 and the query string 110. The edit distance 114 is employed in determining relevance of the document 106 as part of the search results 108.

[0036] The neural network 502 can be employed to receive the document information 104 as an input for computing a relevance score for the document 106. Based solely or in part on the relevance scores for some or all of the search results 108, the documents in the search results 108 can be ranked. The system 500 employs the neural network 502 and codebase to generate the relevance score for ranking of the associated document in the search results 108.

[0037] Following is a description of the edit distance algorithm for calculating the edit distance between the query string and each of the data strings to obtain a TAUC score for each pair.

[0038] Because there is only one title in a document, the TAUC score can be calculated with respect to title as follows:

$$TAUC(Title) = ED(Title)$$

where $TAUC(Title)$ is later used as an input to the neural network after application of a transform function, and $ED(Title)$ is the edit distance of the title.

[0039] There can be multiple instances of anchor text for a document, as well as URLs and clicks (where a click is a previously executed query for which this document was clicked on). The idea is that this document is more relevant for similar queries. At index time, the N anchor texts having the highest frequencies are selected. Then the ED score is calculated for each selected anchor. Finally, the TAUC score is determined for an anchor as follows:

$$TAUC(Anchor) = \text{Min}\{ED(Anchor_i)\} \quad i: \text{top N anchors};$$

[0040] The intuition is that if a good match exists with one of the anchors, then it is sufficient. $TAUC(Anchor)$ is used as a neural network input after applying a transform function.

[0041] Special processing is utilized before calculating the ED for URL strings. At index time URL strings are split into parts using a set of characters as separators. Then terms are found in each part from a dictionary of title and anchor terms. Each occurrence of a term from dictionary is stored in an index with the position measured in characters from the beginning of the URL string.

[0042] At query time all occurrences of the query terms are read from the index stored at index time and the breaks are filled in with "non-query" terms. After this processing the ED is calculated. The result of ED processing is a neural network input, after application of a transform function.

[0043] Another property that can be processed is the number of "clicks" the user enters for a given document content. Each time a user clicks on the document, a stream is entered into a database and associated with the document. This process can also be applied to stream data in the document information text such as short streams of data.

[0044] The index-time URL processing algorithm splits the entire URL into parts using a set of characters as separators. The split function also sets `urlpart.startpos` to a position of part in the source string. The split function performs filtering of insignificant parts of the URL.

[0045] For example, "http://www.companymeeting.com/index.html" is filtered into "companymeeting/index" and split into "companymeeting" and "index".

```
Startpos: 0
```

```
Urlparts = split(url, dictionary)
// find terms in different url parts.
For each (term in dictionary)
{
  Int pos = 0;
```

```

For each(urlpart in urlparts)
{
    pos = urlpart.Find(term, pos);
    while (pos >= 0)
5      {
        // parts_separator is used to distinguish different parts at
query time
        storeOccurrence (term, pos +
urlpart.startpos*parts_separator);
10      pos = url.Find(term, pos + term.length);
    }
}
setIndexStreamLength(parts_separator * urlparts.Count);
}

```

15 **[0046]** Assuming the dictionary contains "company meeting comp", the following keys can be generated: Company: 0; Meeting: 7; and Comp: 0. The total length of the string is parts_separator*2.

20 **[0047]** With respect to query-time processing before ED, at query time the occurrences of the query terms are read, a string of query terms constructed in the order of appearance in the source URL string, and space between the terms filled in with "non-query" word marks. For example, consider a query string of "company policy" and a resulting string of "company" "non-query term" "non-query term".

25 **[0048]** A parts_separator, query term positions, and stream length are determined to know how many parts are in the original URL string and what part contains a given query. Each part without terms is deemed to contain a "non-query term". If a part does not start with a query term, a "non-query term" is inserted before the term. All spaces between query terms are filled with "non-query terms".

30 **[0049]** FIG. 6 illustrates the types of data that can be employed in the document information 104 for determining the edit distance between the query string and the data string. The document information 104 can include TAUC data 602, such as title text 604, anchor text 606, URL 608 text or characters, and click information 610, for example, for processing by the processing component 102 and generation of the data (or target) string 116. The document information 104 can also include click information 610 related to the number of times a user clicks on document content, the type of content the user selects (via the click), the number of clicks on the content, the document in general, etc.

35 **[0050]** FIG. 7 illustrates an index-time processing data flow 700. At the top, document information in the form of the title 604, document anchors 606, click information 610, etc., are received based on document analysis and extraction. The title 604 is processed through a term-splitting algorithm 704 and then to a dictionary 706. The dictionary 706 is a temporary storage of different terms found in the title 604, anchors 606, click information 610, etc. The dictionary 706 is used to split the URL 608 via a URL splitting algorithm 708. The output of the URL splitting algorithm 708 is sent to an indexing process 710 for relevance and ranking processing. The document anchors 606 can also be processed through a filter 712 for the top N anchors. The click information 610 can be processed directly via the indexing process 710. Other document information can be processed accordingly (e.g., term splitting, filtering, etc.).

40 **[0051]** FIG. 8 illustrates a block diagram 800 showing inputs to the neural network from the index process 710 of FIG. 7 for result ranking. The indexing process 710 can be used for computing a URL edit distance (ED) 802 relative to the query string 110, a top-N-anchors ED 804 relative to the query string 110, a title ED 806 relative to the query string 110, a click ED 808 relative to the query string 110, as well as other features 810 not related to edit distance, some or all of which (URL ED 802, top-N-anchors ED 804, title ED 806, click ED 808, and other features 810) can be employed as inputs to the neural network 502, ultimately to find the relevance score for the associated document, and then ranking of the document among other document search results. The neural network 502 can be a 2-layer model that receives at least the TAUC features as raw input features that contribute to identifying relevance of the document. The neural network determines how these features are combined into a single number that can be used for sorting by the search engine.

45 **[0052]** It is to be appreciated that the neural network 502 is just one example of mathematical or computational models that can be employed for the relevance and ranking processing. Other forms of statistical regression can be employed such as naive Bayes, Bayesian networks, decision trees, fuzzy logic models, and other statistical classification models representing different patterns of independence can be employed, where classification is inclusive of methods used to assign rank and/or priority.

50 **[0053]** FIG. 9 illustrates an exemplary system 900 implementation of the neural network 502, edit distance inputs and raw feature inputs for computing generating search results. The set of raw ranking features 810 on the input(s) of the neural network 502 can include a BM25 function 902 (e.g., BM25F), click distance 904, URL depth 906, file types 908, and language match 910. The BM25 components can include body, title, author, anchor text, URL display name, and extracted title, for example.

[0054] FIG. 10 illustrates a method of determining relevance. At 1000, a query string is received as part of a search process. At 1002, document information is extracted from a document returned during the search process. At 1004, a data string is generated from the document information. At 1006, the edit distance is computed between the data string and the query string. At 1008, a relevance score is calculated based on the edit distance.

[0055] Other aspects of the method can include employing term insertion as part of computing the edit distance and assessing an insertion cost for insertion of a term in the query string to generate the data string, the cost represented as a weighting parameter. The method can further comprise employing term deletion as part of computing the edit distance and assessing a deletion cost for deletion of a term in the query string to generate the data string, the cost represented as a weighting parameter. A position cost can be computed as part of computing the edit distance, the position cost associated with term insertion and/or term deletion of a term position in the data string. Additionally, a matching process is performed between characters of the data string and characters of the query string to compute an overall cost of computing the edit distance.

[0056] The splitting compound terms of a URL of the data string can occur at index time. The method can further comprise the filtering of anchor text of the data string to find a top-ranked set of anchor text based on frequency of occurrence in the document and computing an edit distance score for anchor text in the set. The edit distance score, derived from computing the edit distance, can be input into a two-layer neural network after application of a transform function, the score generated based on calculating the edit distance associated with at least one of title information, anchor information, click information, or URL information.

[0057] FIG. 11 illustrates a method of computing relevance of a document. At 1100, a query string is processed as part of a search process to return a result set of documents. At 1102, a data string is generated based on the document information extracted from a document of the result set, the document information includes one or more of title information, anchor text information, click information, and URL information from the document. At 1104, the edit distance is computed between the data string and the query string based on term insertion, term deletion, and term position. At 1106, a relevance score is calculated based on the edit distance, the relevance score used to rank the document in the result set.

[0058] The method can further comprise computing a cost associated with each of the term insertion, term deletion and term position, and factoring the cost into computation of the relevance score, and splitting compound terms of the URL information at index time and filtering the anchor text information at index time to find a top-ranked set of anchor text based on frequency of occurrence of the anchor text in the document. The reading of occurrences of the terms of the query string can be performed to construct a string of query terms in order of appearance in a source URL string and filling space between the terms with word marks.

[0059] As used in this application, the terms "component" and "system" are intended to refer to a computer-related entity, either hardware, a combination of hardware and software, software, or software in execution. For example, a component can be, but is not limited to being, a process running on a processor, a processor, a hard disk drive, multiple storage drives (of optical and/or magnetic storage medium), an object, an executable, a thread of execution, a program, and/or a computer. By way of illustration, both an application running on a server and the server can be a component. One or more components can reside within a process and/or thread of execution, and a component can be localized on one computer and/or distributed between two or more computers.

[0060] Referring now to FIG. 12, there is illustrated a block diagram of a computing system 1200 operable to execute edit distance processing for search result ranking using TAUC features in accordance with the disclosed architecture. In order to provide additional context for various aspects thereof, FIG. 12 and the following discussion are intended to provide a brief, general description of a suitable computing system 1200 in which the various aspects can be implemented. While the description above is in the general context of computer-executable instructions that may run on one or more computers, those skilled in the art will recognize that a novel embodiment also can be implemented in combination with other program modules and/or as a combination of hardware and software.

[0061] Generally, program modules include routines, programs, components, data structures, etc., that perform particular tasks or implement particular abstract data types. Moreover, those skilled in the art will appreciate that the inventive methods can be practiced with other computer system configurations, including single-processor or multiprocessor computer systems, minicomputers, mainframe computers, as well as personal computers, hand-held computing devices, microprocessor-based or programmable consumer electronics, and the like, each of which can be operatively coupled to one or more associated devices.

[0062] The illustrated aspects can also be practiced in distributed computing environments where certain tasks are performed by remote processing devices that are linked through a communications network. In a distributed computing environment, program modules can be located in both local and remote memory storage devices.

[0063] A computer typically includes a variety of computer-readable media. Computer-readable media can be any available media that can be accessed by the computer and includes volatile and non-volatile media, removable and non-removable media. By way of example, and not limitation, computer-readable media can comprise computer storage media and communication media. Computer storage media includes volatile and non-volatile, removable and non-removable media implemented in any method or technology for storage of information such as computer-readable

instructions, data structures, program modules or other data. Computer storage media includes, but is not limited to, RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM, digital video disk (DVD) or other optical disk storage, magnetic cassettes, magnetic tape, magnetic disk storage or other magnetic storage devices, or any other medium which can be used to store the desired information and which can be accessed by the computer.

5 **[0064]** With reference again to FIG. 12, the exemplary computing system 1200 for implementing various aspects includes a computer 1202 having a processing unit 1204, a system memory 1206 and a system bus 1208. The system bus 1208 provides an interface for system components including, but not limited to, the system memory 1206 to the processing unit 1204. The processing unit 1204 can be any of various commercially available processors. Dual micro-processors and other multi-processor architectures may also be employed as the processing unit 1204.

10 **[0065]** The system bus 1208 can be any of several types of bus structure that may further interconnect to a memory bus (with or without a memory controller), a peripheral bus, and a local bus using any of a variety of commercially available bus architectures. The system memory 1206 can include non-volatile memory (NON-VOL) 1210 and/or volatile memory 1212 (e.g., random access memory (RAM)). A basic input/output system (BIOS) can be stored in the non-volatile memory 1210 (e.g., ROM, EPROM, EEPROM, etc.), which BIOS are the basic routines that help to transfer information between elements within the computer 1202, such as during start-up. The volatile memory 1212 can also include a high-speed RAM such as static RAM for caching data.

15 **[0066]** The computer 1202 further includes an internal hard disk drive (HDD) 1214 (e.g., EIDE, SATA), which internal HDD 1214 may also be configured for external use in a suitable chassis, a magnetic floppy disk drive (FDD) 1216, (e.g., to read from or write to a removable diskette 1218) and an optical disk drive 1220, (e.g., reading a CD-ROM disk 1222 or, to read from or write to other high capacity optical media such as a DVD). The HDD 1214, FDD 1216 and optical disk drive 1220 can be connected to the system bus 1208 by a HDD interface 1224, an FDD interface 1226 and an optical drive interface 1228, respectively. The HDD interface 1224 for external drive implementations can include at least one or both of Universal Serial Bus (USB) and IEEE 1394 interface technologies.

20 **[0067]** The drives and associated computer-readable media provide nonvolatile storage of data, data structures, computer-executable instructions, and so forth. For the computer 1202, the drives and media accommodate the storage of any data in a suitable digital format. Although the description of computer-readable media above refers to a HDD, a removable magnetic diskette (e.g., FDD), and a removable optical media such as a CD or DVD, it should be appreciated by those skilled in the art that other types of media which are readable by a computer, such as zip drives, magnetic cassettes, flash memory cards, cartridges, and the like, may also be used in the exemplary operating environment, and further, that any such media may contain computer-executable instructions for performing novel methods of the disclosed architecture.

25 **[0068]** A number of program modules can be stored in the drives and volatile memory 1212, including an operating system 1230, one or more application programs 1232, other program modules 1234, and program data 1236. The one or more application programs 1232, other program modules 1234, and program data 1236 can include the system 100 and associated blocks, the system 500 and associated blocks, the document information 104, TAUC data 602, click information 610, the data flow 700 (and algorithms), and block diagram 800 (and associated blocks).

30 **[0069]** All or portions of the operating system, applications, modules, and/or data can also be cached in the volatile memory 1212. It is to be appreciated that the disclosed architecture can be implemented with various commercially available operating systems or combinations of operating systems.

35 **[0070]** A user can enter commands and information into the computer 1202 through one or more wire/wireless input devices, for example, a keyboard 1238 and a pointing device, such as a mouse 1240. Other input devices (not shown) may include a microphone, an IR remote control, a joystick, a game pad, a stylus pen, touch screen, or the like. These and other input devices are often connected to the processing unit 1204 through an input device interface 1242 that is coupled to the system bus 1208, but can be connected by other interfaces such as a parallel port, IEEE 1394 serial port, a game port, a USB port, an IR interface, etc.

40 **[0071]** A monitor 1244 or other type of display device is also connected to the system bus 1208 via an interface, such as a video adaptor 1246. In addition to the monitor 1244, a computer typically includes other peripheral output devices (not shown), such as speakers, printers, etc.

45 **[0072]** The computer 1202 may operate in a networked environment using logical connections via wire and/or wireless communications to one or more remote computers, such as a remote computer(s) 1248. The remote computer(s) 1248 can be a workstation, a server computer, a router, a personal computer, portable computer, microprocessor-based entertainment appliance, a peer device or other common network node, and typically includes many or all of the elements described relative to the computer 1202, although, for purposes of brevity, only a memory/storage device 1250 is illustrated. The logical connections depicted include wire/wireless connectivity to a local area network (LAN) 1252 and/or larger networks, for example, a wide area network (WAN) 1254. Such LAN and WAN networking environments are commonplace in offices and companies, and facilitate enterprise-wide computer networks, such as intranets, all of which may connect to a global communications network, for example, the Internet.

50 **[0073]** When used in a LAN networking environment, the computer 1202 is connected to the LAN 1252 through a wire

and/or wireless communication network interface or adaptor 1256. The adaptor 1256 can facilitate wire and/or wireless communications to the LAN 1252, which may also include a wireless access point disposed thereon for communicating with the wireless functionality of the adaptor 1256.

[0074] When used in a WAN networking environment, the computer 1202 can include a modem 1258, or is connected to a communications server on the WAN 1254, or has other means for establishing communications over the WAN 1254, such as by way of the Internet. The modem 1258, which can be internal or external and a wire and/or wireless device, is connected to the system bus 1208 via the input device interface 1242. In a networked environment, program modules depicted relative to the computer 1202, or portions thereof, can be stored in the remote memory/storage device 1250. It will be appreciated that the network connections shown are exemplary and other means of establishing a communications link between the computers can be used.

[0075] The computer 1202 is operable to communicate with wire and wireless devices or entities using the IEEE 802 family of standards, such as wireless devices operatively disposed in wireless communication (e.g., IEEE 802.11 over-the-air modulation techniques) with, for example, a printer, scanner, desktop and/or portable computer, personal digital assistant (PDA), communications satellite, any piece of equipment or location associated with a wirelessly detectable tag (e.g., a kiosk, news stand, restroom), and telephone. This includes at least Wi-Fi (or Wireless Fidelity), WiMax, and Bluetooth™ wireless technologies. Thus, the communication can be a predefined structure as with a conventional network or simply an ad hoc communication between at least two devices. Wi-Fi networks use radio technologies called IEEE 802.11x (a, b, g, etc.) to provide secure, reliable, fast wireless connectivity. A Wi-Fi network can be used to connect computers to each other, to the Internet, and to wire networks (which use IEEE 802.3-related media and functions).

[0076] What has been described above includes examples of the disclosed architecture. It is, of course, not possible to describe every conceivable combination of components and/or methodologies, but one of ordinary skill in the art may recognize that many further combinations and permutations are possible. Accordingly, the novel architecture is intended to embrace all such alterations, modifications and variations that fall within the scope of the appended claims. Furthermore, to the extent that the term "includes" is used in either the detailed description or the claims, such term is intended to be inclusive in a manner similar to the term "comprising" as "comprising" is interpreted when employed as a transitional word in a claim.

Claims

1. A computer-implemented method, comprising:

extracting (1002) document information (104) from a document received (1000) as a search result, the search being based on a query string (110), the document information including a uniform resource locator (608), the uniform resource locator including a compound term;
 splitting (704) the compound term into multiple, separate terms;
 finding at least one of the multiple, separate terms in a dictionary (706) of terms;
 generating (1004) a target data string (116) based on the extracted document information, the target data string including one of the at least one multiple, separate terms found in the dictionary; and
 computing (1006) edit distance (114) between the target data string and the query string, the edit distance being employed in determining (1008) relevance of the document as part of result ranking.

2. The method of claim 1, wherein the document information (104) includes at least one of a title (604) information, uniform resource locator (608) information, click (610) information, or anchor (606) text.

3. The method of claims 1 or 2, wherein the compound terms of the document information (104) are split (708) at index time to compute (1008) the edit distance (114) relative to the uniform resource locator (608).

4. The method of one of claims 1 to 3, wherein anchor (606) text of the document information (104) is filtered at index time to compute a top-ranked set of anchor (606) text.

5. The method of one of claims 1 to 4, wherein the document information (104) includes at least one of title (604) characters, anchor (606) characters, click (610) characters, or uniform resource locator (608) characters, which document information is input to a neural network (502) along with raw input features (810) of a BM25F function (902), click distance (904), file type (908), language (910) and uniform resource locator depth (906) to determine the relevance of the document.

6. The method of one of claims 1 to 5, wherein the computing (1006) of the edit distance (114) is based on insertion

and deletion of terms to increase proximity between the target data string (116) and the query string (110).

- 5
7. The method of one of claims 1 to 6, wherein the computing (1006) of the edit distance (114) is based on costs associated with insertion and deletion of terms to increase proximity between the target data string (116) and the query string (110).
- 10
8. The method of claim 1, wherein determining (1008) relevance of the document includes calculating (1008) a relevance score (504) based on the edit distance (114).
- 15
9. The method of claim 8, further comprising employing term insertion as part of computing (1006) the edit distance (114) and assessing an insertion cost for insertion of a term in the query string (110) to generate the target data string (116), the cost represented as a weighting parameter.
- 20
10. The method of claims 8 or 9, further comprising employing term deletion as part of computing (1006) the edit distance (114) and assessing a deletion cost for deletion of a term in the query string (110) to generate the target data string (116), the cost represented as a weighting parameter, or further comprising computing a position cost as part of computing (1006) the edit distance (114), the position cost associated with term insertion and/or term deletion of a term position in the target data string (116).
- 25
11. The method of one of claims 8 to 10, further comprising performing a matching process between characters of the target data string (116) and characters of the query string (110) to compute an overall cost of computing (1006) the edit distance (116), or further comprising splitting (708) the compound terms of the uniform resource locator (608) at index time.
- 30
12. The method of one of claims 8 to 11, further comprising filtering anchor (606) text of the target data string (116) to find a top-ranked set of anchor text based on frequency of occurrence in the document, and/or further comprising computing an edit distance (114) score for anchor (606) text in the set.
- 35
13. The method of one of claims 8 to 12, further comprising inputting a score, derived from computing (1006) the edit distance (114), into a two-layer neural network (502) after application of a transform function, the score generated based on calculating the edit distance associated with at least one of title (604) information, anchor (606) information, click (610) information, or uniform resource locator (608) information, and other raw input features (810).
- 40
14. A system, comprising:
 one or more processors (1204); and
 a memory (1206) coupled to the one or more processors, the memory storing computer-executable instructions which, when executed by the one or more processors, cause the one or more processors to perform the method of one of claims 1 to 13.
- 45
15. One or more computer-readable media (1206, 1210, 1212, 1214, 1218, 1222, 1250) having stored computer-executable instructions for performing the method of one of claims 1 to 13.

Patentansprüche

- 50
1. Computerimplementiertes Verfahren, das umfasst:
 Extrahieren (1002) von Dokumentinformationen (104) aus einem Dokument, das als ein Suchergebnis empfangen wird (1000), wobei die Suche auf einer Abfragezeichenkette (110) basiert, wobei die Dokumentinformationen einen einheitlichen Ressourcenanzeiger bzw. URL (608) enthalten, wobei der URL einen zusammengesetzten Term enthält,
 Aufteilen (704) des zusammengesetzten Terms in mehrere, separate Terme,
 55 Finden wenigstens eines aus den mehreren, separaten Termen in einem Wörterbuch (706) von Termen, Erzeugen (1004) einer Zieldaten-Zeichenkette (116) basierend auf den extrahierten Dokumentinformationen, wobei die Zieldaten-Zeichenkette einen der wenigstens einen mehreren, separaten Terme, die in dem Wörterbuch gefunden wurden, enthält, und

EP 2 289 007 B1

Berechnen (1006) der Editierdistanz (114) zwischen der Zieldaten-Zeichenkette und der Abfragezeichenkette, wobei die Editierdistanz verwendet wird, um die Relevanz des Dokuments als Teil eines Ergebnisrankings zu bestimmen.

- 5 2. Verfahren nach Anspruch 1, wobei die Dokumentinformationen (104) Titel (604)-Informationen, URL (608)-Informationen, Klick (610)-Informationen und/oder Anker (606)-Text enthalten.
3. Verfahren nach Anspruch 1 oder 2, wobei die zusammengesetzten Terme der Dokumentinformationen (104) zur Indexzeit aufgeteilt (708) werden, um die Editierdistanz (114) relativ zu dem URL (608) zu berechnen (1008).
- 10 4. Verfahren nach einem der Ansprüche 1 bis 3, wobei der Anker (606)-Text der Dokumentinformationen (104) zur Indexzeit gefiltert wird, um einen im Ranking ganz oben eingestuften Satz von Anker (606)-Text zu berechnen.
- 15 5. Verfahren nach einem der Ansprüche 1 bis 4, wobei die Dokumentinformationen (104) Titel (604)-Zeichen, Anker (606)-Zeichen, Klick (610)-Zeichen und/oder URL (608)-Zeichen enthalten, wobei die Dokumentinformationen in ein neuronales Netz (502) zusammen mit Roheingabemerkmale (810) einer BM25F-Funktion (902), der Klickdistanz (904), dem Dateityp (908), der Sprache (910) und der URL-Tiefe (906) eingegeben werden, um die Relevanz des Dokuments zu bestimmen.
- 20 6. Verfahren nach einem der Ansprüche 1 bis 5, wobei das Berechnen (1006) der Editierdistanz (114) auf dem Einfügen und Entfernen von Termen zum Vergrößern der Nähe zwischen der Zieldaten-Zeichenkette (116) und der Abfragezeichenkette (110) basiert.
- 25 7. Verfahren nach einem der Ansprüche 1 bis 6, wobei das Berechnen (1006) der Editierdistanz (114) auf den mit dem Einfügen und Entfernen von Termen zum Vergrößern der Nähe zwischen der Zieldaten-Zeichenkette (116) und der Abfragezeichenkette (110) assoziierte Kosten basiert.
8. Verfahren nach Anspruch 1, wobei:
30 das Bestimmen (1008) der Relevanz des Dokuments das Berechnen (1008) einer Relevanzbewertung (504) basierend auf der Editierdistanz (114) umfasst.
9. Verfahren nach Anspruch 8, das weiterhin das Verwenden einer Termeinfügung als Teil des Berechnens (1006) der Editierdistanz (114) und das Einschätzen der Einfügungskosten für das Einfügen eines Terms in die Abfragezeichenkette (110) für das Erzeugen der Zieldaten-Zeichenkette (116) umfasst, wobei die Kosten als ein Gewichtungparameter wiedergegeben werden.
- 35 10. Verfahren nach Anspruch 8 oder 9, das weiterhin das Verwenden einer Termentfernung als Teil des Berechnens (1006) der Editierdistanz (114) und das Einschätzen der Entfernungskosten für das Entfernen eines Terms aus der Abfragezeichenkette (110) für das Erzeugen der Zieldaten-Zeichenkette (116) umfasst, wobei die Kosten als ein Gewichtungparameter wiedergegeben werden, oder
40 weiterhin das Berechnen von Positionskosten als Teil des Berechnens (1006) der Editierdistanz (114) umfasst, wobei die Positionskosten mit der Termeinfügung und/oder der Termentfernung einer Termposition in der Zieldaten-Zeichenkette (116) assoziiert sind.
- 45 11. Verfahren nach einem der Ansprüche 8 bis 10, das weiterhin das Durchführen eines Vergleichsprozesses zwischen Zeichen der Zieldaten-Zeichenkette (116) und Zeichen der Abfragezeichenkette (110) umfasst, um die Gesamtkosten für das Berechnen (1006) der Editierdistanz (116) zu berechnen, oder
50 weiterhin das Aufteilen (708) der zusammengesetzten Terme des URL (608) zur Indexzeit umfasst.
12. Verfahren nach einem der Ansprüche 8 bis 11, das weiterhin das Filtern von Anker (606)-Text der Zieldaten-Zeichenkette (116) umfasst, um einen ganz oben im Ranking eingestuften Ankertextsatz basierend auf der Häufigkeit des Auftretens in dem Dokument zu finden, und/oder
55 weiterhin das Berechnen einer Editierdistanz (114)-Bewertung für den Anker (606)-Text in dem Satz umfasst.
13. Verfahren nach einem der Ansprüche 8 bis 12, das weiterhin das Eingeben einer Bewertung, die aus dem Berechnen (1006) der Editierdistanz (114) abgeleitet wird, in ein zweischichtiges neuronales Netz (502) nach der Anwendung einer Transformationsfunktion umfasst, wobei die Bewertung auf dem Berechnen der Editierdistanz, die mit den

Titel (604)-Informationen, Anker (606)-Informationen, Klick (610)-Informationen und/oder URL (608)-Informationen assoziiert ist, und auf anderen Roheingabemerkmale (810) basiert.

14. System, das umfasst:

einen oder mehrere Prozessoren (1204), und einen Speicher (1206), der mit dem einen oder den mehreren Prozessoren gekoppelt ist, wobei der Speicher computerausführbare Befehle speichert, die bei einer Ausführung durch den einen oder die mehreren Prozessoren veranlassen, dass der eine oder die mehreren Prozessoren das Verfahren nach einem der Ansprüche 1 bis 13 durchführen.

15. Ein oder mehrere computerlesbare Medien (1206, 1210, 1212, 1214, 1218, 1222, 1250), auf denen computerausführbare Befehle zum Durchführen des Verfahrens nach einem der Ansprüche 1 bis 13 gespeichert sind.

Revendications

1. Procédé mis en oeuvre par un ordinateur, comprenant :

l'extraction (1002) d'information de document (104) d'un document reçu (1000) sous forme d'un résultat de recherche, la recherche étant basée sur une chaîne de requête (110), l'information de document comprenant un localisateur uniforme de ressource (608), le localisateur uniforme de ressource comprenant un terme composé ;
la division (704) du terme composé en une multitude de termes séparés ;
la localisation d'au moins un terme de la multitude de termes séparés dans un dictionnaire (706) de termes ;
la génération (1004) d'une chaîne de données cible (116) basée sur l'information de document extraite, la chaîne de données cible comprenant un desdits au moins un terme de la multitude de termes séparés trouvés dans le dictionnaire ; et
le calcul (1006) d'une distance d'édition (114) entre la chaîne de données cible et la chaîne de requête, la distance d'édition étant employée pour la détermination (1008) de la pertinence du document dans le cadre d'un classement de résultats.

2. Procédé selon la revendication 1, dans lequel l'information de document (104) comprend au moins un élément parmi une information de titre (604), une information de localisateur uniforme de ressource (608), une information de clic (610) et un texte d'ancrage (606).

3. Procédé selon la revendication 1 ou 2, dans lequel les termes composés de l'information de document (104) sont divisés (708) au temps d'indexation pour calculer (1008) la distance d'édition (114) relative au localisateur uniforme de ressource (608).

4. Procédé selon l'une des revendications 1 à 3, dans lequel le texte d'ancrage (606) de l'information de document (104) est filtré au temps d'indexation pour calculer un ensemble de texte d'ancrage (606) en tête de classement.

5. Procédé selon l'une des revendications 1 à 4, dans lequel l'information de document (104) comprend au moins un élément parmi des caractères de titre (604), des caractères d'ancrage (606), des caractères de clic (610) et des caractères de localisateur uniforme de ressource (608), ladite information de document étant entrée dans un réseau neuronal (502) conjointement à des caractéristiques d'entrée brutes (810) d'une fonction BM25F (902), une distance de clic (904), un type de fichier (908), une langue (910) et une profondeur de localisateur uniforme de ressource (906) pour déterminer la pertinence du document.

6. Procédé selon l'une des revendications 1 à 5, dans lequel le calcul (1006) de la distance d'édition (114) est basé sur l'insertion et l'effacement de termes pour augmenter la proximité entre la chaîne de données cible (116) et la chaîne de requête (110).

7. Procédé selon l'une des revendications 1 à 6, dans lequel le calcul (1006) de la distance d'édition (114) est basé sur des coûts associés à l'insertion et à l'effacement de termes pour augmenter la proximité entre la chaîne de données cible (116) et la chaîne de requête (110).

8. Procédé selon la revendication 1, dans lequel la détermination (1008) de la pertinence du document comprend le calcul (1008) d'un score de pertinence (504) sur base de la distance d'édition (114).
- 5 9. Procédé selon la revendication 8, comprenant en outre l'emploi d'une insertion de terme dans le cadre du calcul (1006) de la distance d'édition (114) et de l'évaluation d'un coût d'insertion pour insérer un terme dans la chaîne de requête (110) afin de générer la chaîne de données cible (116), le coût étant représenté par un paramètre de pondération.
- 10 10. Procédé selon la revendication 8 ou 9, comprenant en outre l'emploi d'un effacement de terme dans le cadre du calcul (1006) de la distance d'édition (114) et de l'évaluation d'un coût d'effacement pour effacer un terme dans la chaîne de requête (110) afin de générer la chaîne de données cible (116), le coût étant représenté par un paramètre de pondération, ou
- 15 comprenant en outre le calcul d'un coût de position dans le cadre du calcul (1006) de la distance d'édition (114), le coût de position étant associé à l'insertion d'un terme et/ou à l'effacement d'un terme d'une position de terme dans la chaîne de données cible (116).
- 20 11. Procédé selon l'une des revendications 8 à 10, comprenant en outre la mise en oeuvre d'un processus de mise en concordance entre des caractères de la chaîne de données cible (116) et des caractères de la chaîne de requête (110) pour calculer un coût d'ensemble du calcul (1006) de la distance d'édition (116), ou
- 25 comprenant en outre la division (708) des termes composés du localisateur uniforme de ressource (608) au temps d'indexation.
- 30 12. Procédé selon l'une des revendications 8 à 11, comprenant en outre la filtration de texte d'ancrage (606) de la chaîne de données cible (116) pour trouver un ensemble de texte d'ancrage en tête de classement sur base de la fréquence d'occurrence dans le document, et/ou
- 35 comprenant en outre le calcul d'un score de distance d'édition (114) pour un texte d'ancrage (606) dans l'ensemble.
- 40 13. Procédé selon l'une des revendications 8 à 12, comprenant en outre l'entrée d'un score dérivé du calcul (1006) de la distance d'édition (114) dans un réseau neuronal bicouche (502) après application d'une fonction de transformée, le score étant généré sur base du calcul de la distance d'édition associée à au moins un élément parmi de l'information de titre (604), de l'information d'ancrage (606), de l'information de clic (610) et de l'information de localisateur uniforme de ressource (608), ainsi que d'autres caractéristiques d'entrée brutes (810).
- 45 14. Système comprenant :
- un ou plusieurs processeurs (1204) ; et
- une mémoire (1206) couplée auxdits un ou plusieurs processeurs, la mémoire stockant des instructions exécutables par un ordinateur qui, lorsqu'elles sont exécutées par lesdits un ou plusieurs processeurs, commandent
- 50 auxdits un ou plusieurs processeurs de mettre en oeuvre le procédé selon l'une des revendications 1 à 13.
- 55 15. Un ou plusieurs supports lisibles par un ordinateur (1206, 1210, 1212, 1214, 1218, 1222, 1250) sur lesquels sont stockées des instructions exécutables par un ordinateur pour mettre en oeuvre le procédé selon l'une des revendications 1 à 13.

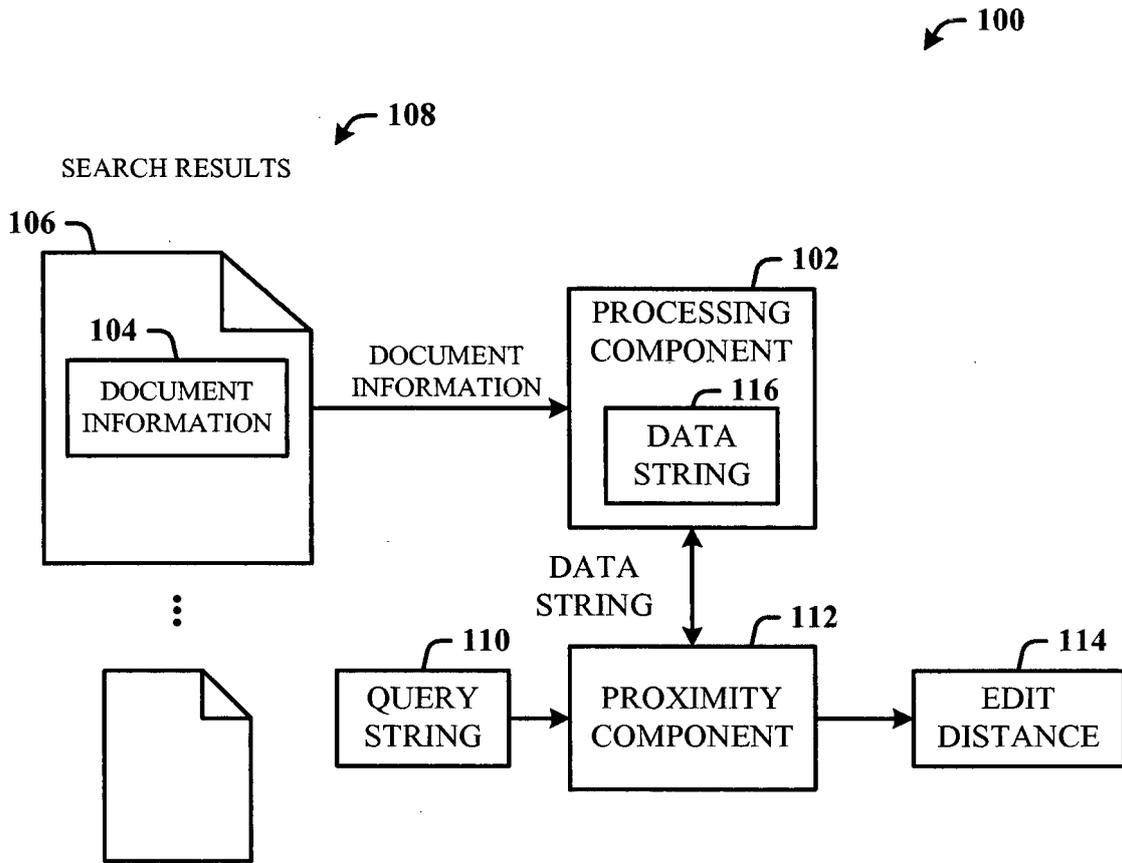


FIG. 1

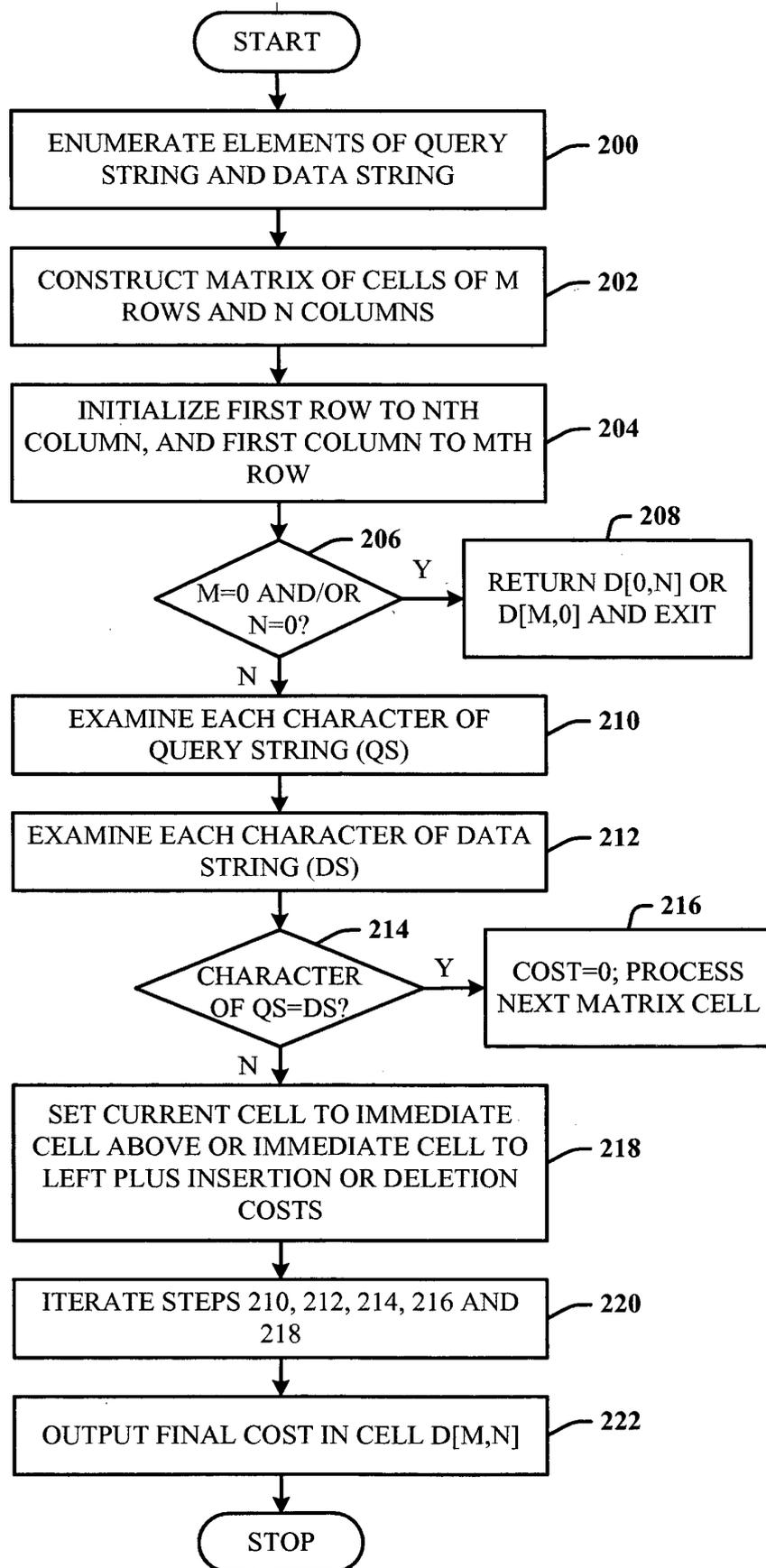


FIG. 2

← 300

304
↓
COLUMNS

	0	1	2	3	
		A	B	C	
0		0	1	2	3
1	C	1	2	3	2
2	B	2	3	2	3
3	A	3	2	3	4
4	X	7	6	7	8

ROWS ← 302

FIG. 3

↖ 400

404 COLUMNS
↓ 0 1 2 3

		A	B	C	
0		0	1	2	28
1	A	1	0	1	27
2	B	2	1	0	26

ROWS

← 402

FIG. 4

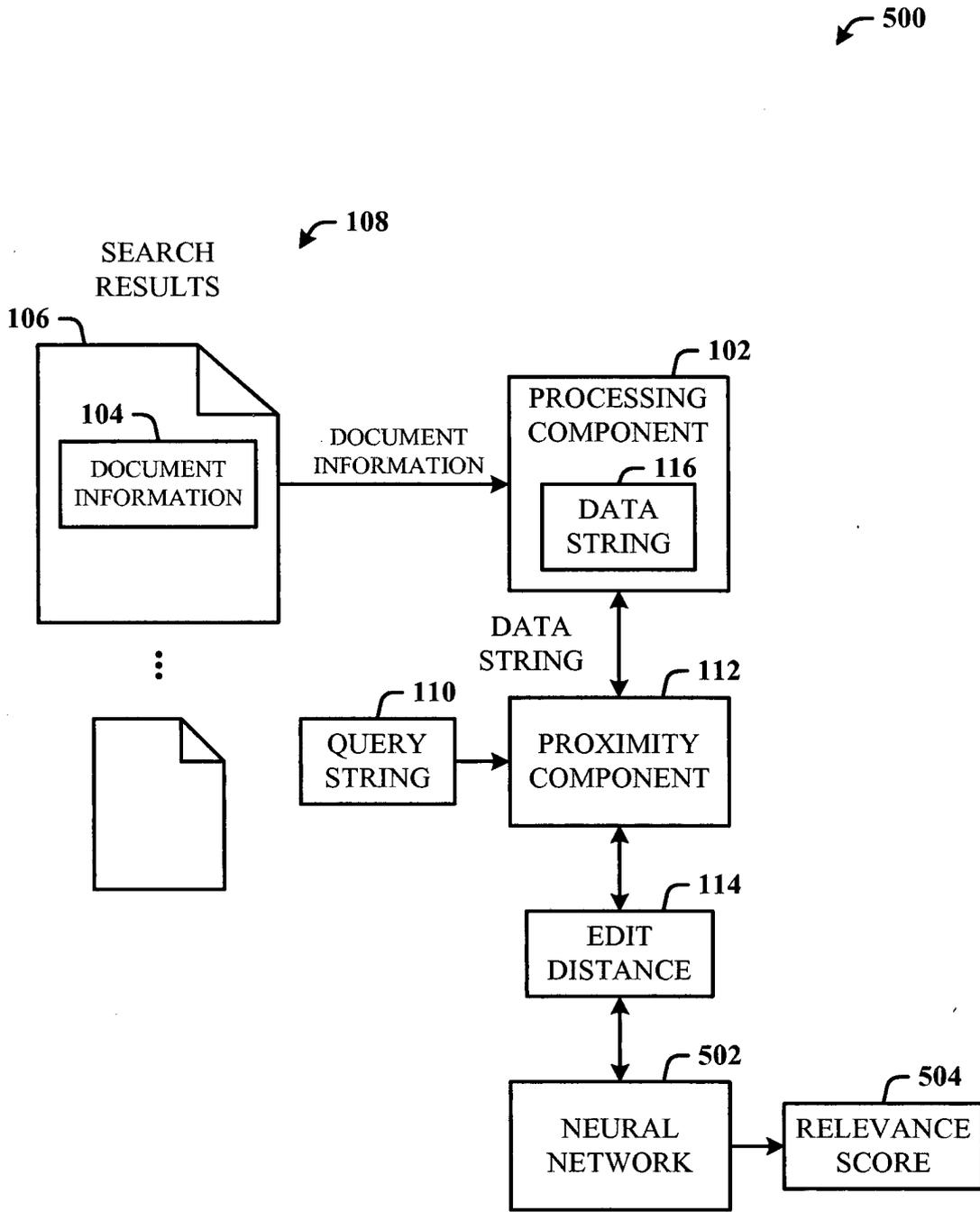


FIG. 5

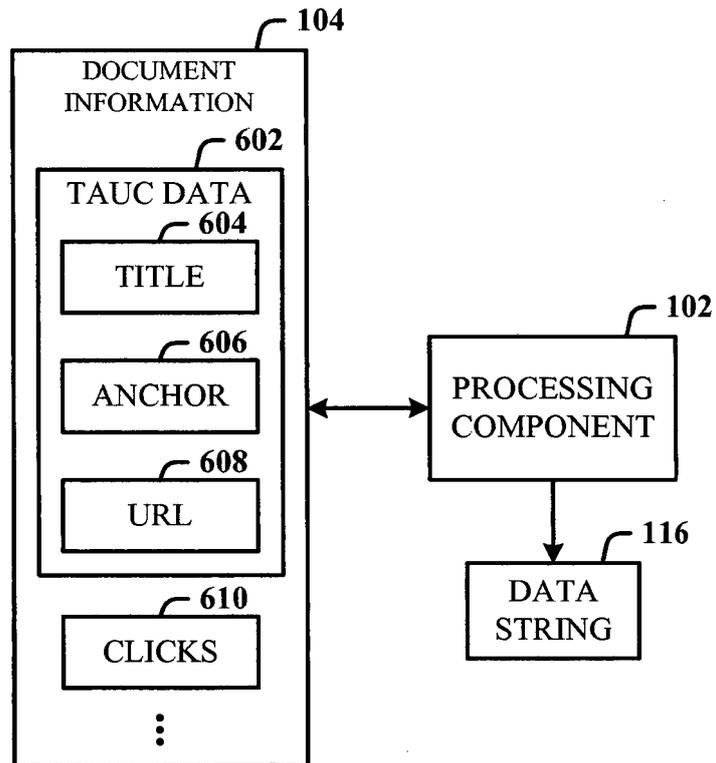


FIG. 6

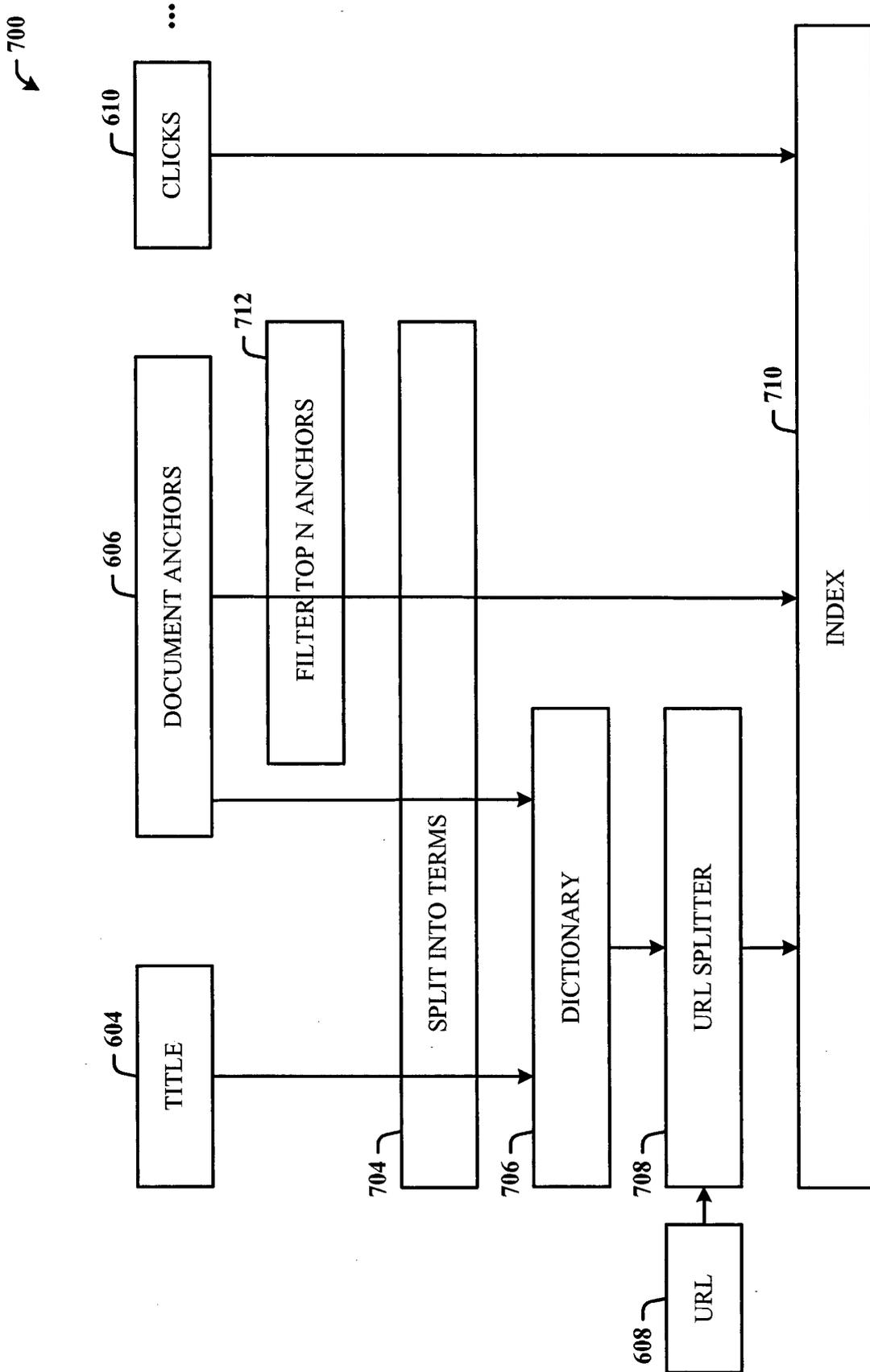


FIG. 7

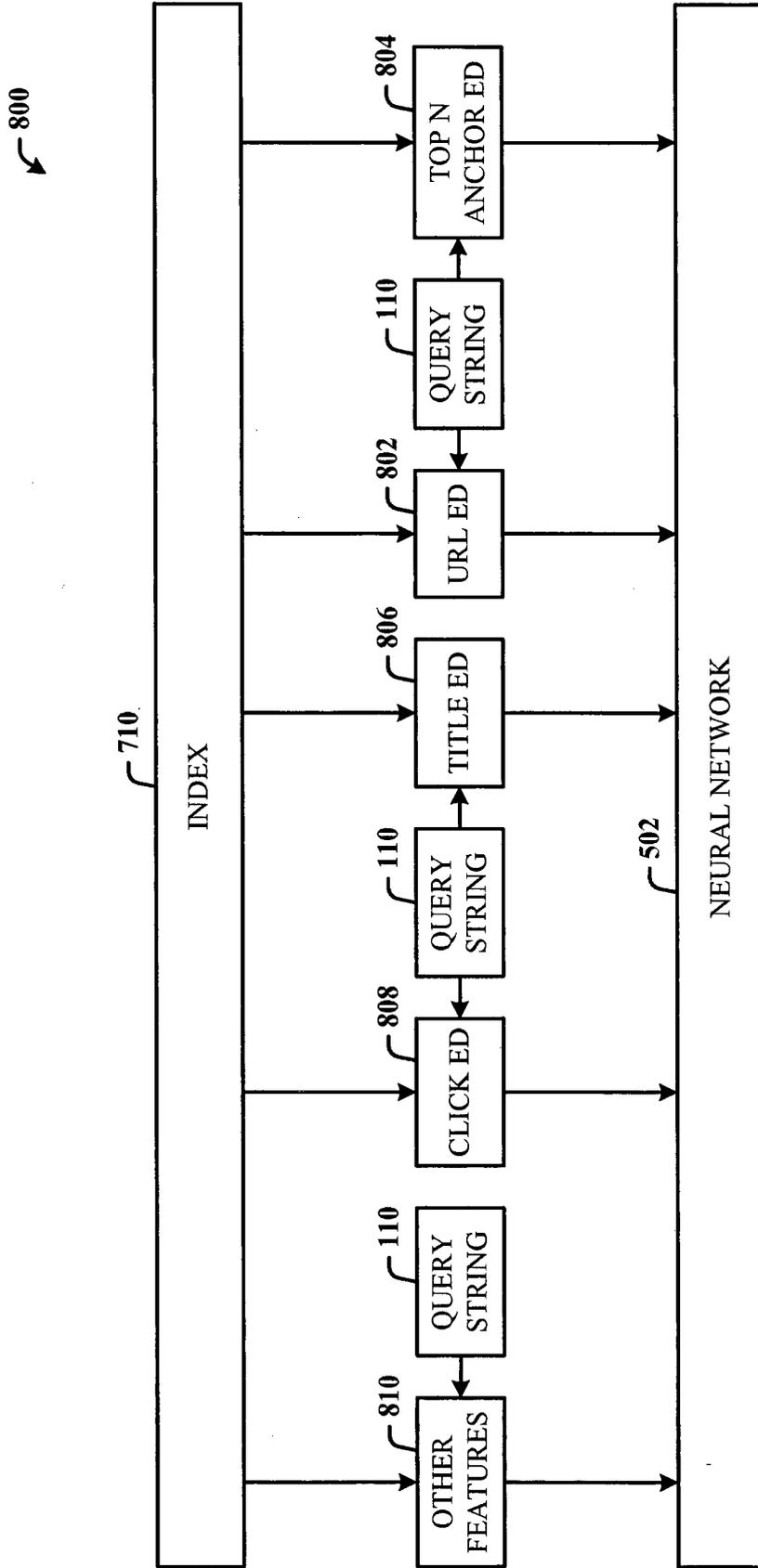


FIG. 8

900

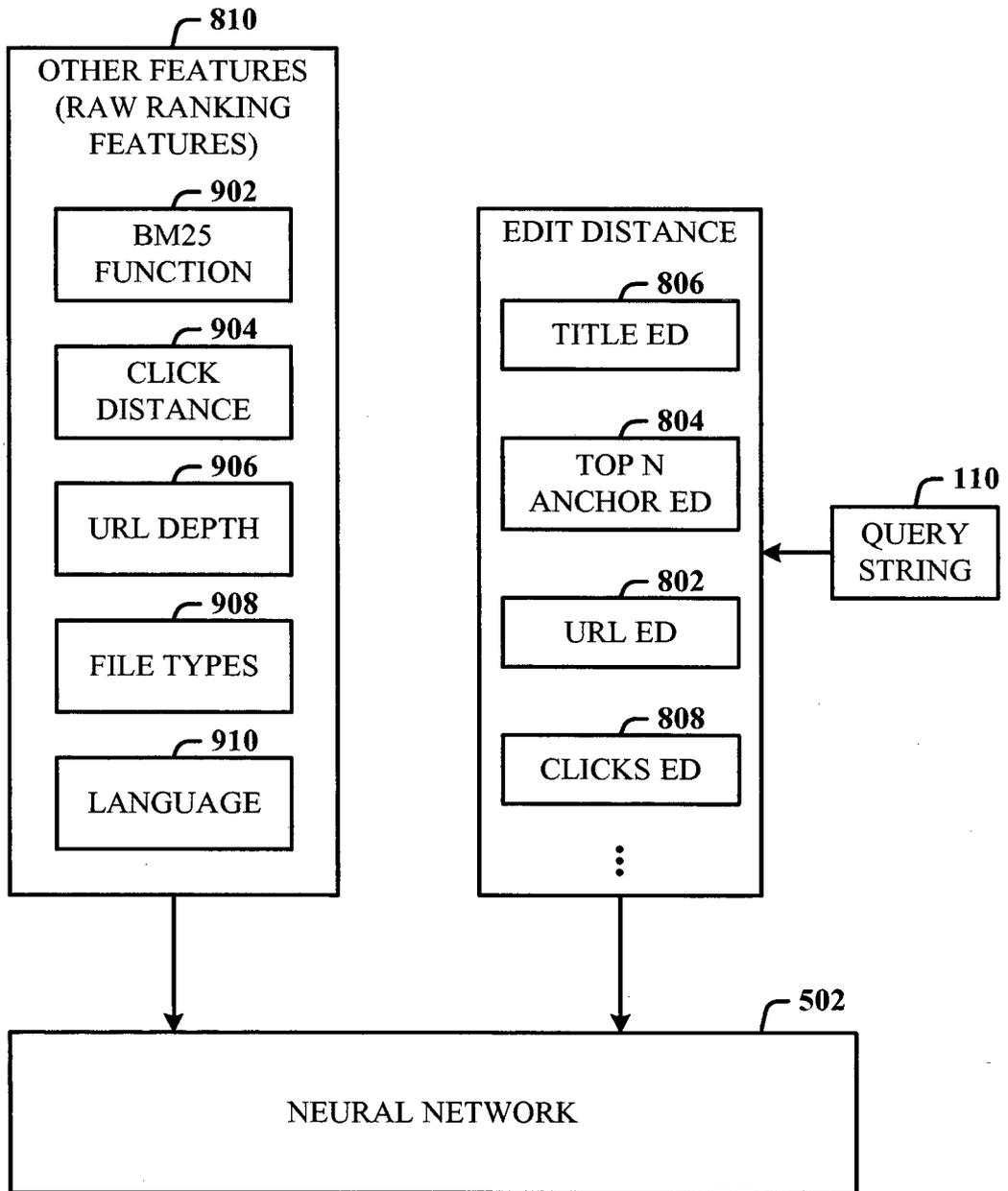


FIG. 9

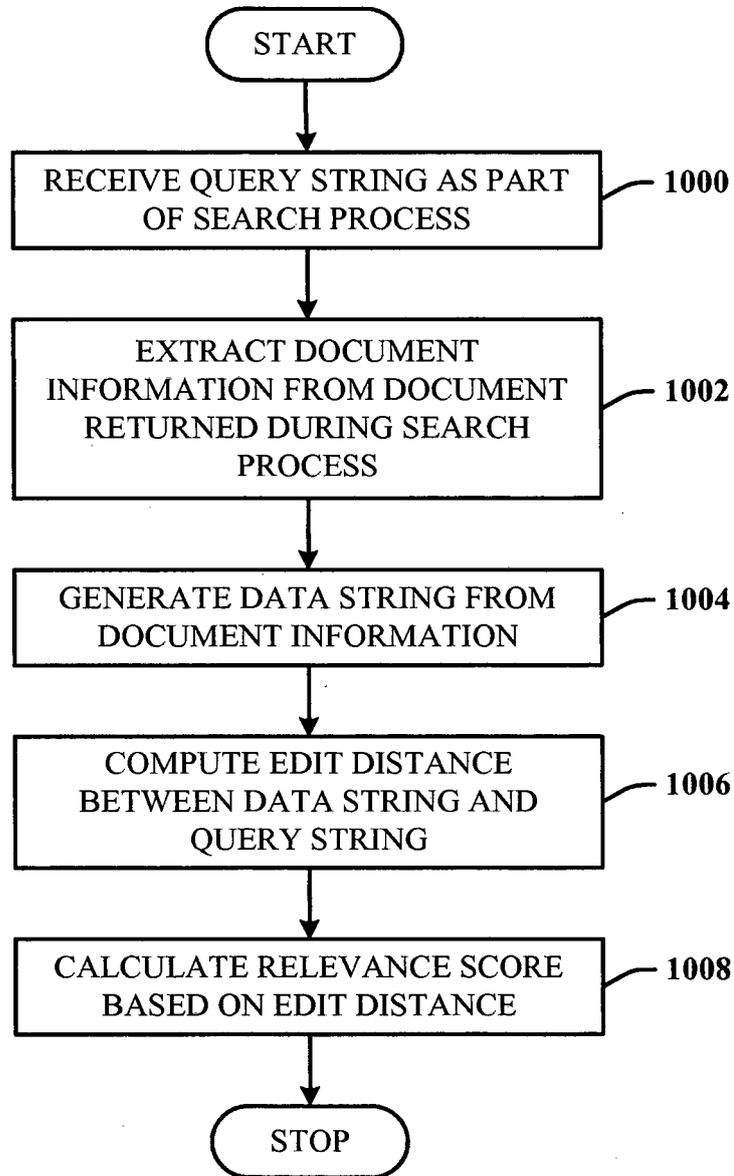


FIG. 10

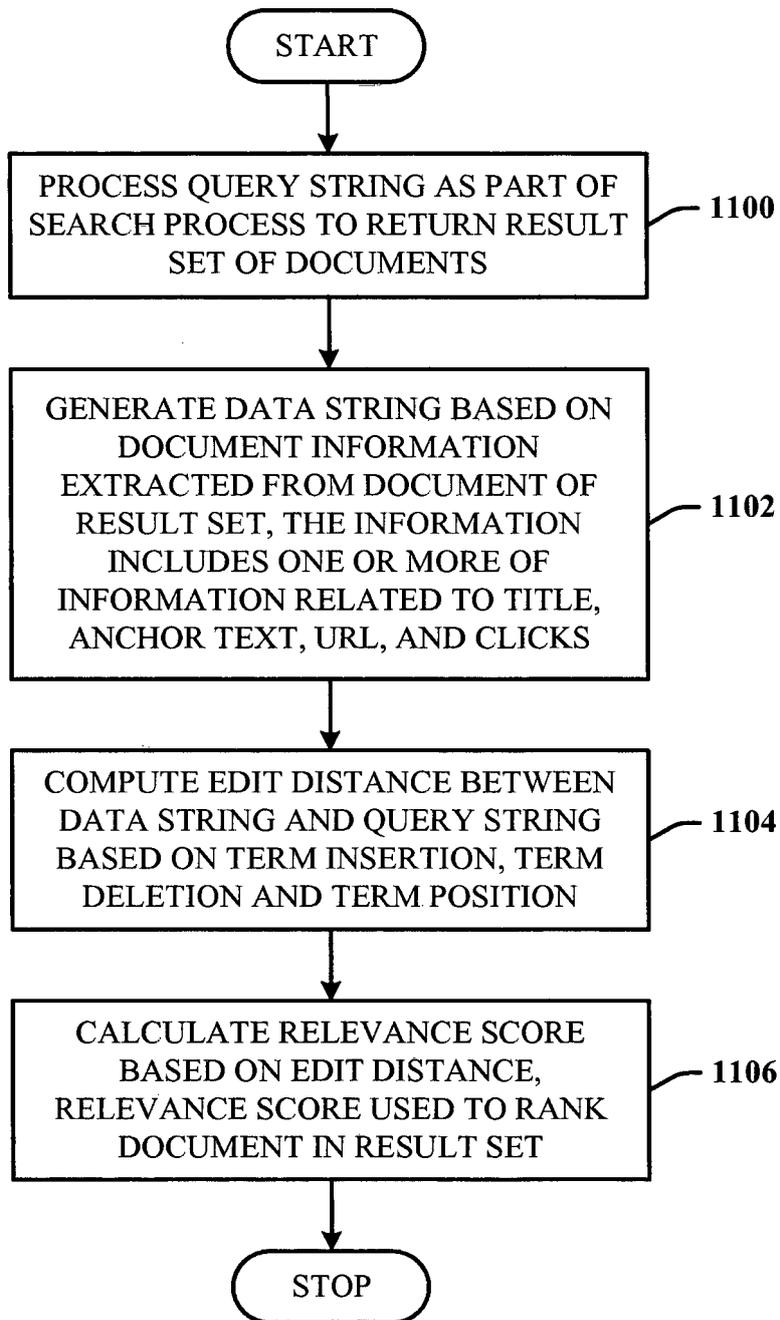


FIG. 11

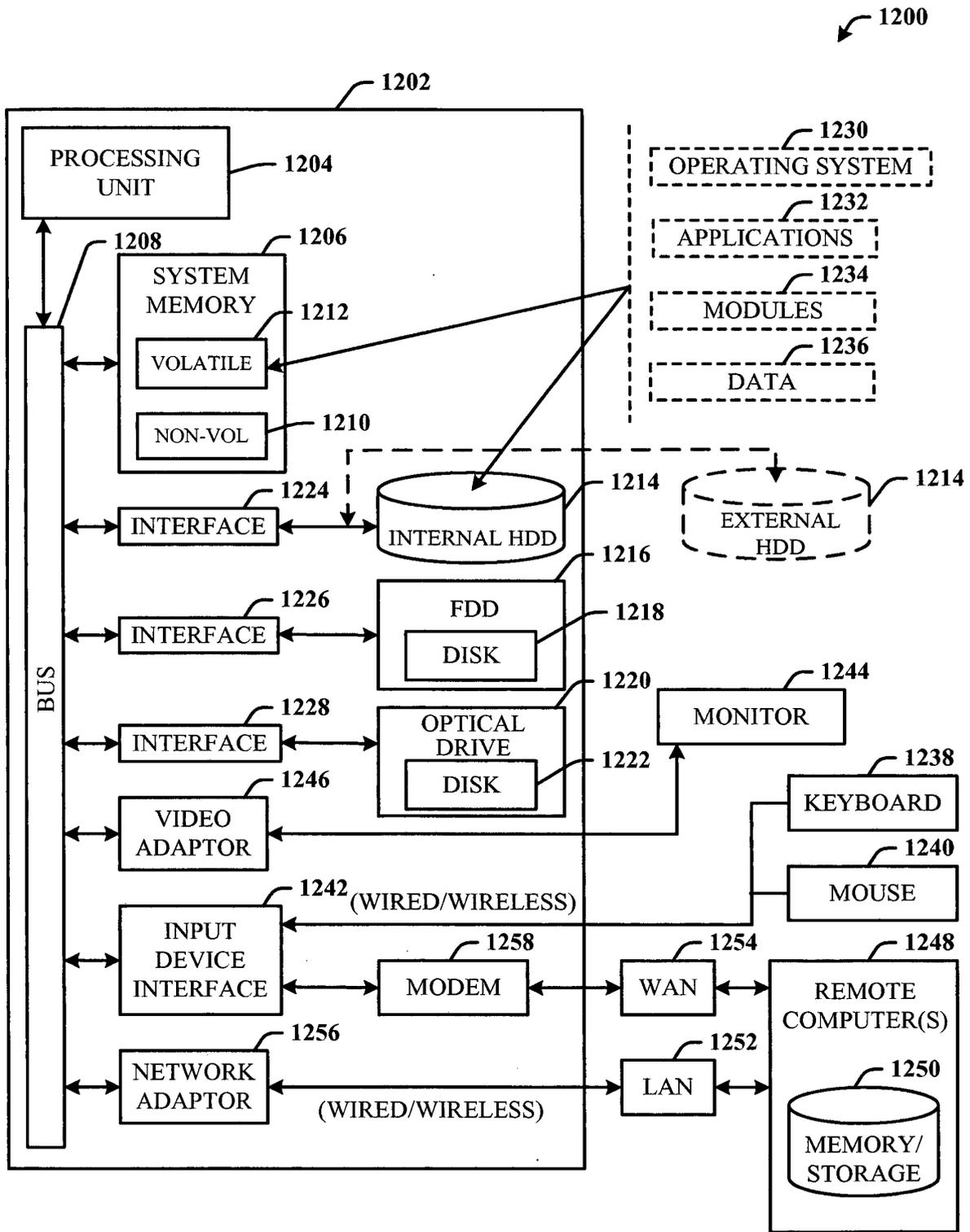


FIG. 12

REFERENCES CITED IN THE DESCRIPTION

This list of references cited by the applicant is for the reader's convenience only. It does not form part of the European patent document. Even though great care has been taken in compiling the references, errors or omissions cannot be excluded and the EPO disclaims all liability in this regard.

Patent documents cited in the description

- US 20060004732 A1 [0004]