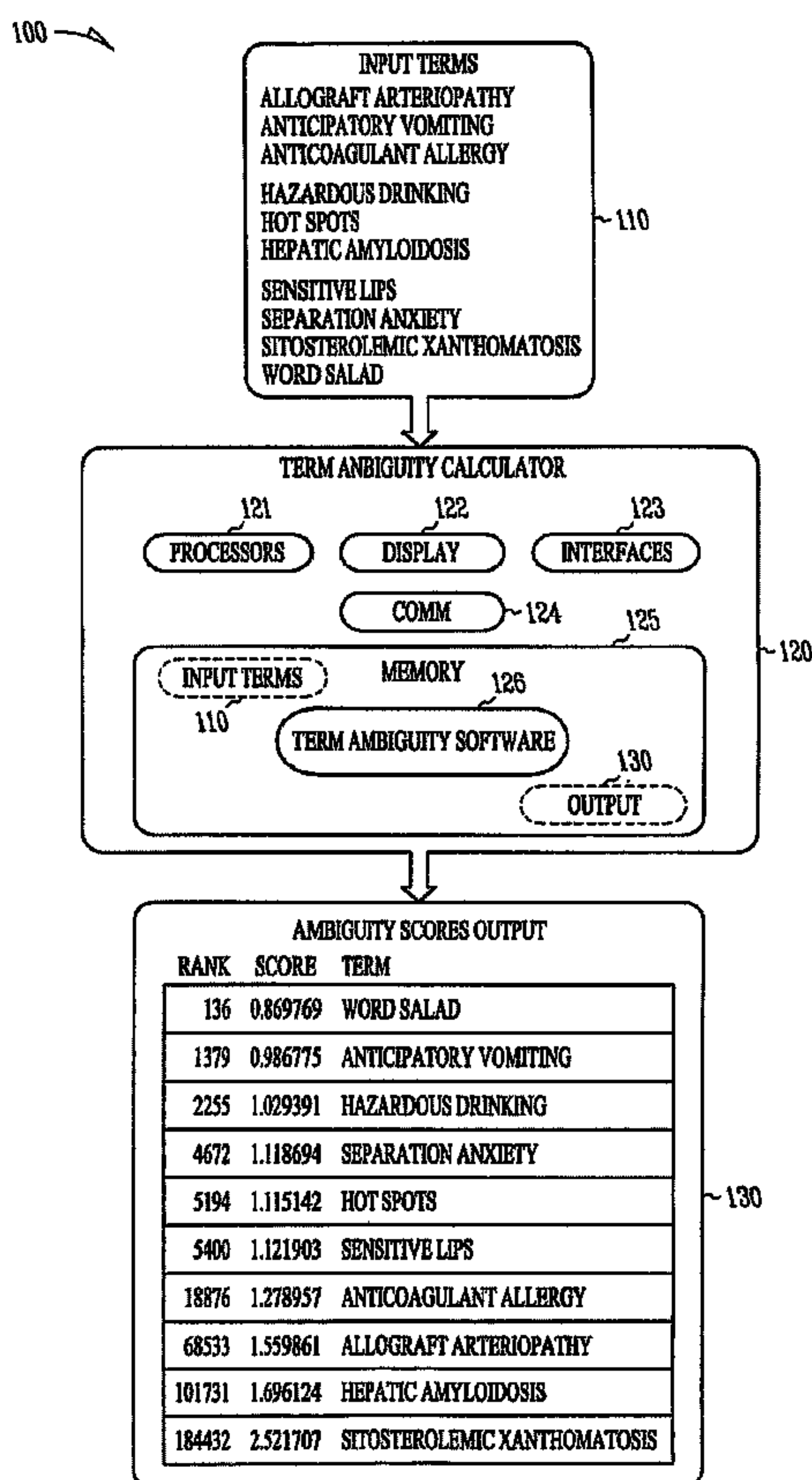




(86) Date de dépôt PCT/PCT Filing Date: 2006/10/04
 (87) Date publication PCT/PCT Publication Date: 2007/04/19
 (45) Date de délivrance/Issue Date: 2016/01/26
 (85) Entrée phase nationale/National Entry: 2008/04/03
 (86) N° demande PCT/PCT Application No.: US 2006/038671
 (87) N° publication PCT/PCT Publication No.: 2007/044350
 (30) Priorité/Priority: 2005/10/04 (US60/723,483)

(51) Cl.Int./Int.Cl. *G06F 17/30* (2006.01)
 (72) Inventeurs/Inventors:
 DOZIER, CHRISTOPHER C., US;
 CHAUDHARY, MARK, US;
 KONDADADI, RAVI, US
 (73) Propriétaire/Owner:
 THOMSON REUTERS GLOBAL RESOURCES, CH
 (74) Agent: MARKS & CLERK

(54) Titre : SYSTEMES, PROCEDES ET LOGICIELS DE LEVE D'AMBIGUITE SUR DES TERMES MEDICAUX
 (54) Title: SYSTEMS, METHODS, AND SOFTWARE FOR ASSESSING AMBIGUITY OF MEDICAL TERMS



(57) Abrégé/Abstract:

Some known medical terms may function as non-medical terms depending on their particular context. Accordingly, the present inventors devised systems, methods, and software that facilitate determining whether a term that is found in a medical corpus is likely to be a medical term when found in another corpus. An exemplary embodiment receives a term and computes an ambiguity score based on language models for a medical and a non-medical corpus.

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau(43) International Publication Date
19 April 2007 (19.04.2007)

PCT

(10) International Publication Number
WO 2007/044350 A3(51) International Patent Classification:
G06F 17/30 (2006.01)

(21) International Application Number:

PCT/US2006/038671

(22) International Filing Date: 4 October 2006 (04.10.2006)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:

60/723,483 4 October 2005 (04.10.2005) US

(71) Applicant (for all designated States except US): THOMSON GLOBAL RESOURCES [IE/CH]; Landis + Gyr-Str. 3, CH-6300 Zug (CH).

(71) Applicants and

(72) Inventors: DOZIER, Christopher, C. [US/US]; 4825 Knox Avenue So., Minneapolis, MN 55419 (US). CHAUDHARY, Mark [IN/US]; P.o. Box 21192, Eagan, MN 55121 (US). KONDADADI, Ravi [IN/US]; 4110 Lexington, Avenue So., #105, Eagan, MN 55123 (US).

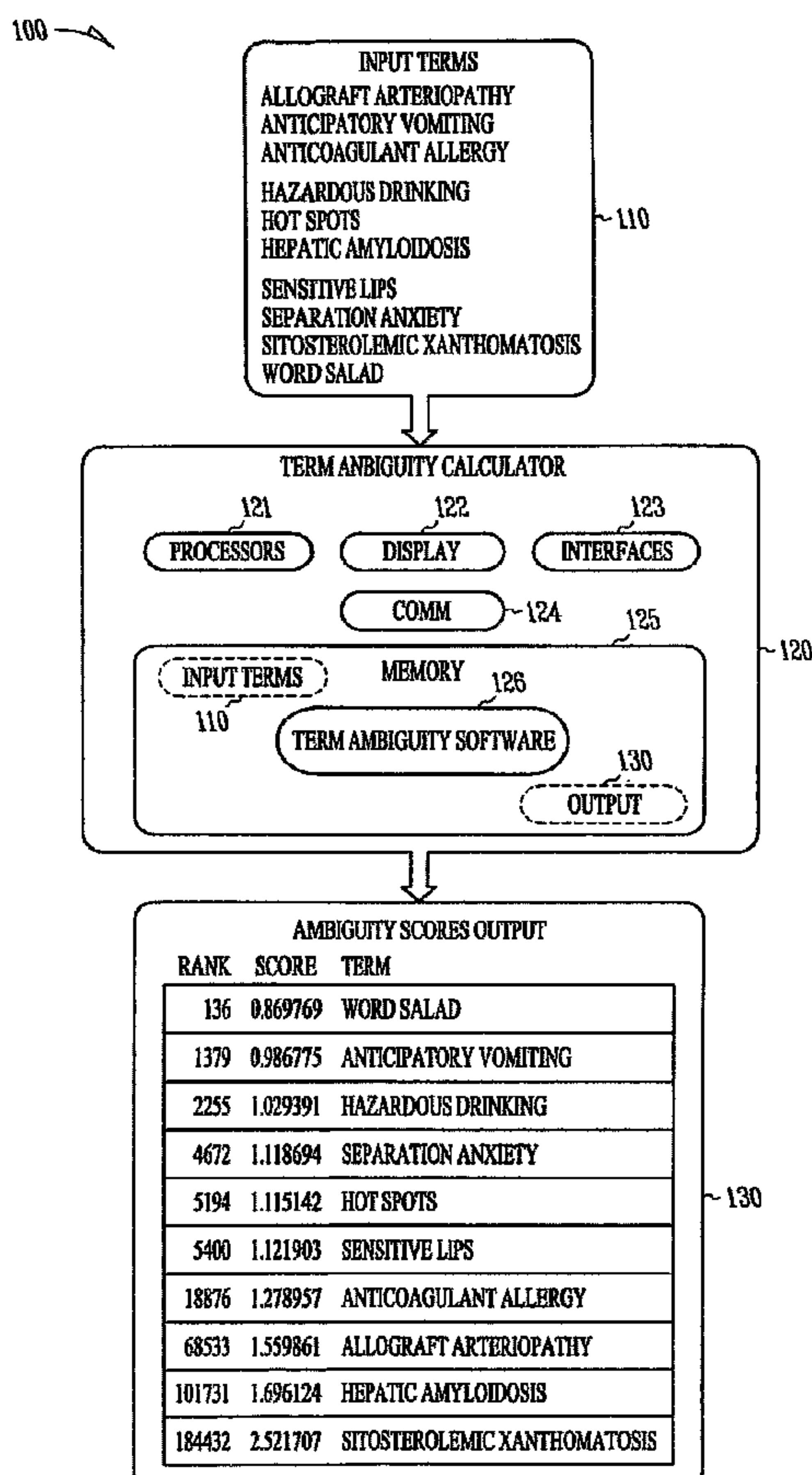
(74) Agents: STEFFEY, Charles, E. et al.; Schwegman Lundberg Woessner & Kluth, PA, P.O. Box 2938, Minneapolis, MN 55402 (US).

(81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LV, LY, MA, MD, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, SV, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

(84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IS, IT, LT, LU, LV, MC, NL, PL, PT, RO, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

[Continued on next page]

(54) Title: SYSTEMS, METHODS, AND SOFTWARE FOR ASSESSING AMBIGUITY OF MEDICAL TERMS



(57) Abstract: Some known medical terms may function as non-medical terms depending on their particular context. Accordingly, the present inventors devised systems, methods, and software that facilitate determining whether a term that is found in a medical corpus is likely to be a medical term when found in another corpus. An exemplary embodiment receives a term and computes an ambiguity score based on language models for a medical and a non-medical corpus.

WO 2007/044350 A3

WO 2007/044350 A3



Published:

— *with international search report*

(88) Date of publication of the international search report:

21 June 2007

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

Systems, Methods, and Software For Assessing Ambiguity of Medical Terms

5

Copyright Notice and Permission

A portion of this patent document contains material subject to copyright protection. The copyright owner has no objection to the facsimile reproduction by anyone of the patent document or the patent disclosure, as it appears in the
10 Patent and Trademark Office patent files or records, but otherwise reserves all copyrights whatsoever. The following notice applies to this document:
Copyright © 2005-2006, Thomson Global Resources.

Technical Field

15 Various embodiments of the present invention concern systems, methods, and software for identifying medical content in documents and linking those documents to other documents based on the medical content.

Background

20 The fantastic growth of the Internet and other computer networks has fueled an equally fantastic growth in the data accessible via these networks. One of the seminal modes for interacting with this data is through the use of hyperlinks within electronic documents.

25 Hyperlinks are user-selectable elements, such as highlighted text or icons, that link one portion of an electronic document to another portion of the same document or to other documents in a database or computer network. With

proper computer equipment and network access, a user can select or invoke a hyperlink and almost instantaneously view the other document, which can be located on virtually any computer system in the world.

Although many hyperlinks are created and inserted into documents manually, recent years have seen development of automated techniques for identifying specific types of document text and linking the identified text using hyperlinks to other related documents. For example, to facilitate legal research, the Westlaw legal research system automatically identifies legal citations and attorney names in text and links the citations to corresponding legal documents in a database and the attorney names to biographical entries in an online directory. For further details, see U.S. Patent 7,003,719 and U.S. Published Patent Application US2003/0135826A1.

Although the automated linking technology in the Westlaw system is highly effective for legal citations and names, the present inventors have identified that this technology is not well suited for other types of content, such as medical terms. For example, the inventors recognize that identifying legal citations and entity names within a text is generally simpler than identifying medical terms because terms may function as medical terms in one context and as non-medical terms in another. Legal citations and person names, on the other hand, generally function as legal citations and person names regardless of context.

Accordingly, the present inventors have identified a need for automated methods identifying whether terms are medical terms or non-medical terms.

25

Summary

To address this and/or other needs, the inventors devised, among other things, systems, methods, and software that facilitate determining whether a term is a medical term or a non-medical term.

2a

According to an aspect of the present invention, there is provided a computer-implemented method for inserting in a document a hyperlink to a medical document associated with a medical term in the document, the method comprising:

receiving by a computer a term included within a document;

determining by the computer an ambiguity score for the term, wherein the ambiguity score is based on first and second language models for respective medical and non-medical corpuses of documents;

determining by the computer, based on the ambiguity score, whether to insert a hyperlink in the document to a medical document associated with the term; and

inserting the hyperlink in the document if the determination to insert the hyperlink is affirmative.

According to another aspect of the present invention, there is provided a computerized system comprising:

an input for receiving a set of terms;

processor for executing code adapted to determine an ambiguity score for a term from the set of terms, wherein the ambiguity score is based on first and second language models for respective medical and non-medical corpuses of documents;

means for determining, based on the ambiguity score, whether to insert a hyperlink in a first document including the term to a medical document associated with the term; and

means for inserting the hyperlink in the first document if the determination to insert the hyperlink is affirmative.

According to a further aspect of the present invention, there is provided a computer-readable medium on which is stored a set of instructions for inserting in a document a hyperlink to a medical document associated with a medical term in the document which, when executed by a computer, perform steps comprising:

receiving by the computer a term included within a document;

determining by the computer an ambiguity score for the term, wherein the ambiguity score is based on first and second language models for respective medical and non-medical corpuses of documents;

determining by the computer, based on the ambiguity score, whether to insert a hyperlink in the document to a medical document associated with the term; and

inserting the hyperlink in the document if the determination to insert the hyperlink is affirmative.

Brief Description of Drawings

Figure 1 is a block diagram of an exemplary system 100 which corresponds to one or more embodiment of the present invention.

Figure 2 is a flow chart of an exemplary method of operating system 100
5 which corresponds to one or more embodiments of the invention.

Detailed Description of Exemplary Embodiments

The following detailed description, which references and incorporates Figures 1 and 2, describes and illustrates one or more exemplary embodiments of
10 the invention. These embodiments, offered not to limit but only to exemplify and teach the invention, are shown and described in sufficient detail to enable those skilled in the art to make and use the invention. Thus, where appropriate to avoid obscuring the invention, the description may omit certain information known to those of skill in the art.

15

Exemplary Computer System Embodying the Invention

Figure 1 shows a diagram of an exemplary computer system 100 incorporating a system, method, and software for assessing the ambiguity of terms, such as medical terms. Though the exemplary system is presented as an
20 interconnected ensemble of separate components, some other embodiments implement their functionality using a greater or lesser number of components. Moreover, some embodiments intercouple one or more the components through wired or wireless local- or wide-area networks. Some embodiments implement one or more portions of system 100 using one or more mainframe computers or
25 servers.) Thus, the present invention is not limited to any particular functional partition.

Generally, system 100 includes input terms 110, term-ambiguity calculator 120, and ambiguity scores output 130.

Input terms 110 includes one or more terms, such as a set of terms from a
30 medical database. In the exemplary embodiment, input terms 110 includes terms from the Unified Medical Language System (UMLS). The table below shows

that UMLS includes a great number of terms in disease, injury, medical procedure, body part, and drug categories.

Category	Terms	Concepts
Disease	189,712	69,948
Injury	42,141	28,997
Medical procedure	134,179	72,918
Body part	38,041	22,260
Drugs	244,752	129,959

- 5 In some embodiment, input terms 110 are terms extracted from one or more input documents, such as an electronic judicial opinion, or other type legal document.

Coupled to database 110 is term-ambiguity calculator 120. Calculator 120 includes one or more conventional processors 121, display device 122,
 10 interface devices 123, network-communications devices 124, and memory 125. Memory 125, which can take a variety of forms, such as coded instructions or data on an electrical, magnetic, and/or optical carrier medium, includes term-ambiguity software 126. Term-ambiguity software 126 includes various software and data components, for determining or calculating for each input term
 15 t and ambiguity score, $Score(term)$ defined as

$$Score(term) = \lambda_1 \frac{\log(P(t | News_lang))}{\log(P(t | UMLS_lang))} + \lambda_2 \frac{\log(P(t | Legal_lang))}{\log(P(t | UMLS_lang)}$$

where

20

$$\log(P(t | lang)) = \sum_{i=1}^n \log(P(ngram | lang))$$

and lamda1 and lamda2 are constants, which in some embodiments are used to normalize or smooth the scoring function. In some embodiments, lambda1 and

lamda2 are set to 0.5. The exemplary embodiment uses ngram backoff with Witten Bell smoothing to smooth the language models.

The exemplary scoring function is based on the intuition that medical ngrams, such as "hepatic," occur relatively more often in UMLS than in news or legal and that ngrams such as "drinki" will occur relatively more often in news or legal than in UMLS. Terms having ngrams that are more highly predicted by UMLS than news or legal tend to yield a larger score and thus indicate that the given term is more likely a medical term than not a medical term when found in a news or legal document.

Term-ambiguity calculator 120 outputs a set 130 of one or more ambiguity scores based on the input terms. (Figure 1 shows that the input terms 110 and output scores 130 are also retained in memory 130.) In the exemplary embodiment, the scores are output as a ranked list, with each score associated with corresponding terms. (Note that term may include one or more words.)

The ambiguity scores can be used for a variety of purposes, including for example determining whether it is appropriate to insert a link in a document including a given term back to a UMLS document associated with the term. For example, in the output terms shown the terms having an ambiguity score greater than 1.5 may be considered as clearly being medical terms and thus linked with high confidence back to related UMLS documents. On the other hand, terms such as "word salad" or "anticipatory vomiting" that have lower scores should not generally be linked back to a related UMLS document without contextual corroboration.

Exemplary Operation of System 100

Figure 2 shows a flowchart illustrating an exemplary method of operating system 100. Flow chart includes process blocks 210-230. Though these blocks (and those of other flow charts in this document) are arranged serially in the exemplary embodiment, other embodiments may reorder the blocks, omit one or more blocks, and/or execute two or more blocks in parallel using multiple processors or a single processor organized as two or more virtual machines or subprocessors. Moreover, still other embodiments implement the blocks as one or more specific interconnected hardware or integrated-circuit

modules with related control and data signals communicated between and through the modules. Thus, this and other exemplary process flows in this document are applicable to software, firmware, hardware, and other types of implementations.

5 Block 210 entails receiving a set of terms. In the exemplary embodiment, this entails receiving a set of terms from UMLS or an input news or legal document into memory 126 of term-ambiguity calculator 120. Execution continues at block 220.

10 Block 220 entails determining one or more ambiguity scores for one or more of the input terms. In the exemplary embodiment this entails computing ambiguity scores according to the definition set forth above for $Score(term)$ in equation above, which provides a sum of two conditional probability ratios. Each conditional probability is based on language model of set or corpus of documents. In some embodiments, one of the conditional probability ratios is
15 omitted from the scoring function. Also, in some embodiments, the conditional probability ratios are inverted.

20 Block 230 entails outputting one or more of the determined ambiguity scores. In the exemplary embodiment, this entails outputting in printed or other human readable form; however, in other embodiments, the output may also be used by another machine, component, or software module, or simply retained in memory.

Conclusion

25 The embodiments described above are intended only to illustrate and teach one or more ways of practicing or implementing the present invention, not to restrict its breadth or scope. The actual scope of the invention, which embraces all ways of practicing or implementing the teachings of the invention, is defined only by the following claims and their equivalents.

The embodiments of the invention in which an exclusive property or privilege is claimed are defined as follows:

1. A computer-implemented method for inserting in a document a hyperlink to a medical document associated with a medical term in the document, the method comprising:
 - receiving by a computer a term included within a document;
 - determining by the computer an ambiguity score for the term, wherein the ambiguity score is based on first and second language models for respective medical and non-medical corpuses of documents;
 - determining by the computer, based on the ambiguity score, whether to insert a hyperlink in the document to a medical document associated with the term; and
 - inserting the hyperlink in the document if the determination to insert the hyperlink is affirmative.
2. The computer-implemented method of claim 1, wherein the second language model is based on a legal or general news corpus of documents.
3. The computer-implemented method of claim 1, wherein the ambiguity score is based on a ratio of a probability of the term given a non-medical corpus to a probability of the term given a medical corpus.
4. The computer-implemented method of any one of claims 1 to 3, wherein first language model is based on the Unified Medical Language System (UMLS).
5. The computer-implemented method of any one of claims 1 to 4, further comprising:
 - determining that the ambiguity score meets a specified level.
6. The computer-implemented method of any one of claims 1 to 5, further comprising:
 - linking the term to a related document.

7. The computer-implemented method of any one of claims 1 to 3, wherein the ambiguity score is based on

$$\lambda_1 \frac{\log(P(t | \text{News_lang}))}{\log(P(t | \text{UMLS_lang}))} + \lambda_2 \frac{\log(P(t | \text{Legal_lang}))}{\log(P(t | \text{UMLS_lang}))}$$

where t denotes the term; News_lang denotes a news corpus; UMLS_lang denotes a medical corpus; Legal_lang denotes a legal corpus; λ_1 and λ_2 are constants; and

$$\log(P(t | \text{lang})) = \sum_{i=1}^n \log(P(\text{ngram} | \text{lang}))$$

where lang is a placeholder for the corpus of interest; and n denotes the number of words in the term t .

8. The computer-implemented method of any one of claims 1 to 5, wherein the received term is a term within an electric document.

9. A computerized system comprising:

an input for receiving a set of terms;

processor for executing code adapted to determine an ambiguity score for a term from the set of terms, wherein the ambiguity score is based on first and second language models for respective medical and non-medical corpuses of documents;

means for determining, based on the ambiguity score, whether to insert a hyperlink in a first document including the term to a medical document associated with the term; and

means for inserting the hyperlink in the first document if the determination to insert the hyperlink is affirmative.

10. The computerized system of claim 9, wherein the second language model is based on a legal or general news corpus of documents.

11. The computerized system of claim 9, wherein each ambiguity score is based on a ratio of a probability of the term given a non-medical corpus to a probability of the term given a medical corpus.
12. The computerized system of any one of claims 9 to 11, wherein first language model is based on the Unified Medical Language System (UMLS).
13. The computerized system of any one of claims 9 to 12, further comprising: determining that the ambiguity score meets a specified level.
14. The computerized system of any one of claims 9 to 13, further comprising: linking the term to a related document.
15. The computerized system of any one of claims 9 to 11, wherein the means for determining the ambiguity score is adapted to determine the ambiguity score is based on

$$\lambda_1 \frac{\log(P(t | News_lang))}{\log(P(t | UMLS_lang))} + \lambda_2 \frac{\log(P(t | Legal_lang))}{\log(P(t | UMLS_lang))}$$

where t denotes the term; $News_lang$ denotes a news corpus; $UMLS_lang$ denotes a medical corpus; $Legal_lang$ denotes a legal corpus; λ_1 and λ_2 are constants; and

$$\log(P(t | lang)) = \sum_{i=1}^n \log(P(ngram | lang))$$

where $lang$ is a placeholder for the corpus of interest; and n denotes the number of words in the term t .

16. The computerized system of any one of claims 9 to 13, wherein the received set of terms is a set of terms within an electric document.
17. A computer-readable medium on which is stored a set of instructions for inserting in a document a hyperlink to a medical document associated with a medical term in the document which, when executed by a computer, perform steps comprising:

receiving by the computer a term included within a document;
determining by the computer an ambiguity score for the term, wherein the ambiguity score is based on first and second language models for respective medical and non-medical corpuses of documents;
determining by the computer, based on the ambiguity score, whether to insert a hyperlink in the document to a medical document associated with the term; and
inserting the hyperlink in the document if the determination to insert the hyperlink is affirmative.

18. The computer-readable medium of claim 17, wherein the second language model is based on a legal or general news corpus of documents.

19. The computer-readable medium of claim 17, wherein each ambiguity score is based on a ratio of a probability of the term given a non-medical corpus to a probability of the term given a medical corpus.

20. The computer-readable medium of any one of claims 17 to 19, wherein first language model is based on the Unified Medical Language System (UMLS).

21. The computer-readable medium of any one of claims 17 to 19, further comprising:
determining that the ambiguity score meets a specified level.

22. The computer-readable medium of any one of claims 17 to 19, further comprising:
linking the term to a related document.

100

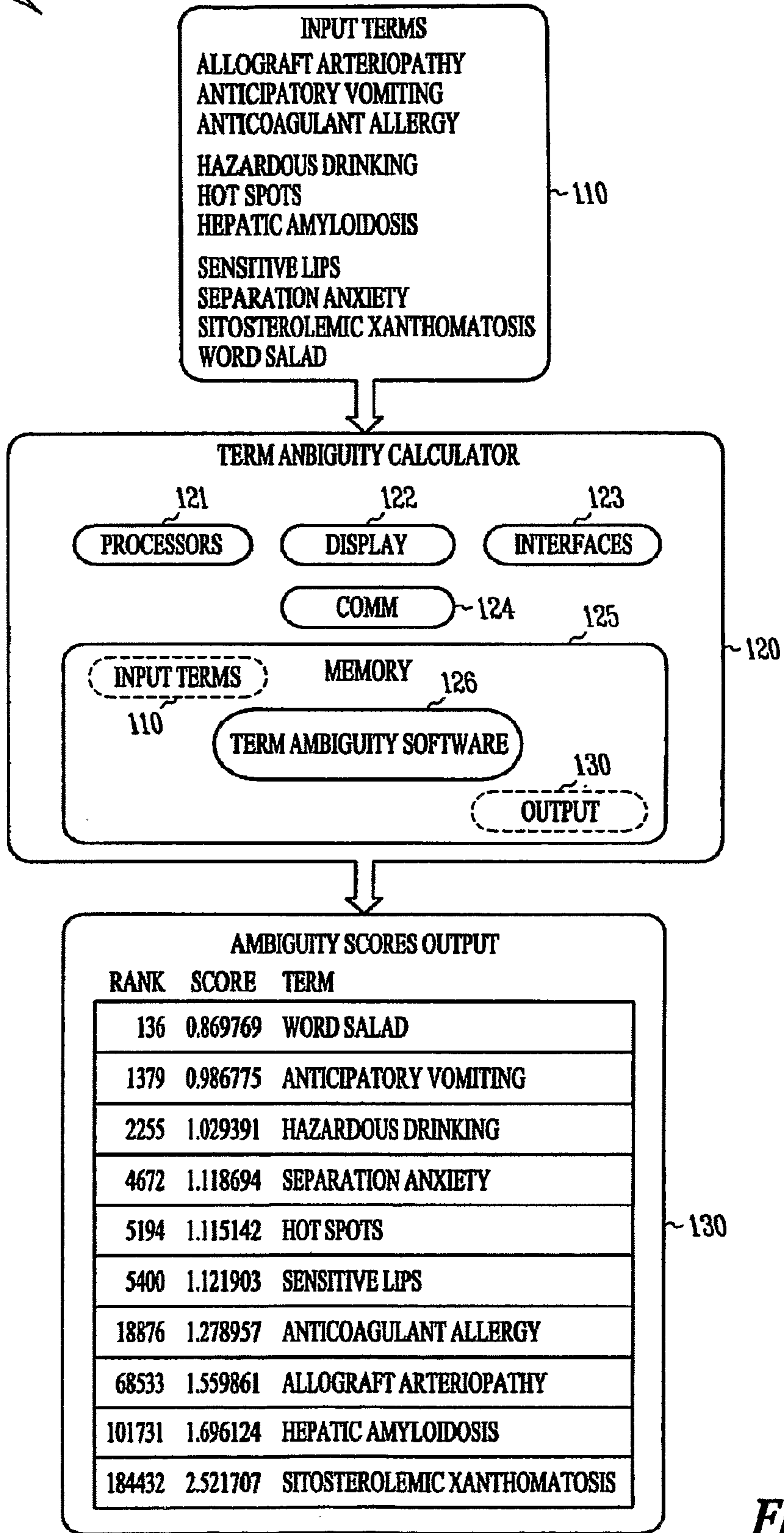


FIG. 1

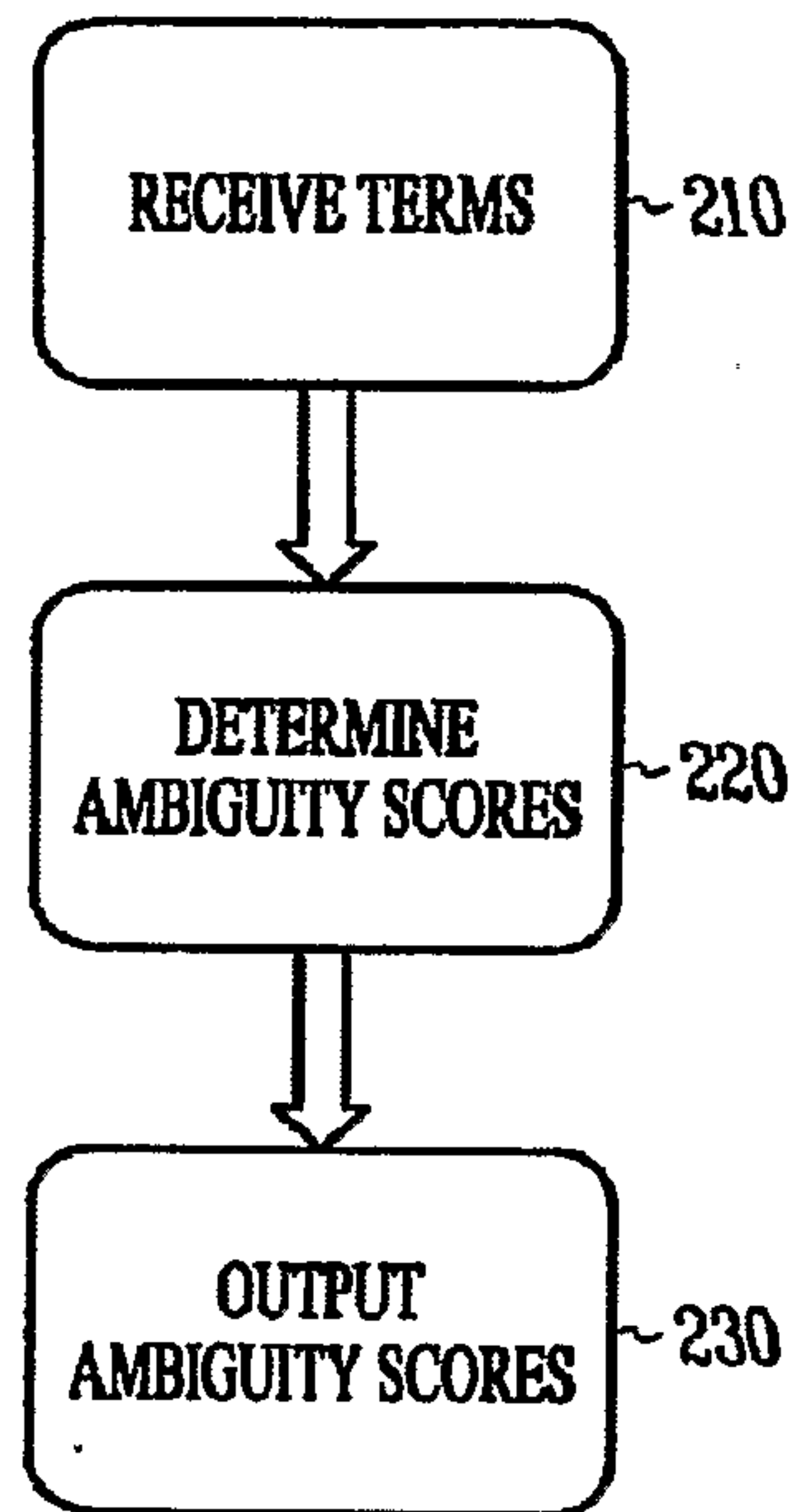


FIG. 2

