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

(51) International Patent Classification:
G06F 17/30 (2006.01)

(21) International Application Number:
PCT/IB2010/054156

(22) International Filing Date:
15 September 2010 (15.09.2010)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
TO2009A000704 16 September 2009 (16.09.2009)   IT

(71) Applicants (for all designated States except US): IN-
TELLISEMANTIC SRL [IT/IT]; Via Giaglione, 1-7,
I-10126 Torino (IT). POLITECNICO DI TORINO [IT/
IT]; Corso Duca Degli Abruzzi, 24, I-10129 Torino (IT).

(72) Inventors; and
(75) Inventors/Applicants (for US only): CORNO, Fulvio
[IT/IT]; c/o Politecnico Di Torino, Corso Duca Degli
Abruzzi, 24, I-10129 Torino (IT). PELLEGRINO, Paolo
[IT/IT]; c/o Politecnico Di Torino, Corso Duca Degli
Abruzzi, 24, I-10129 Torino (IT). CIARAMELLA, Al-
berto [IT/IT]; c/o Intellisemantic Srl, Via Giaglione, 1-7,
I-10126 Torino (IT).

(74) Agent: ROBBA, Pierpaolo; Interpatent S.r.l., Via
Caboto, 35, I-10129 Torino (IT).

(81) Designated States (unless otherwise indicated, for every
kind of national protection available): AE, AG, AL, AM,
AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ,
CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO,
DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT,
HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP,
KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD,
ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI,
NO, NZ, OM, PE, PG, PH, PL, PT, RO, RS, RU, SC, SD,
SE, SG, SK, SL, SM, ST, SV, SY, TH, 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, LR, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG,
ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ,
TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK,
EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU,
LV, MC, MK, MT, NL, NO, PL, PT, RO, SE, SI, SK,
SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ,
GW, ML, MR, NE, SN, TD, TG).

Published:
—  with international search report (Art. 21(3))

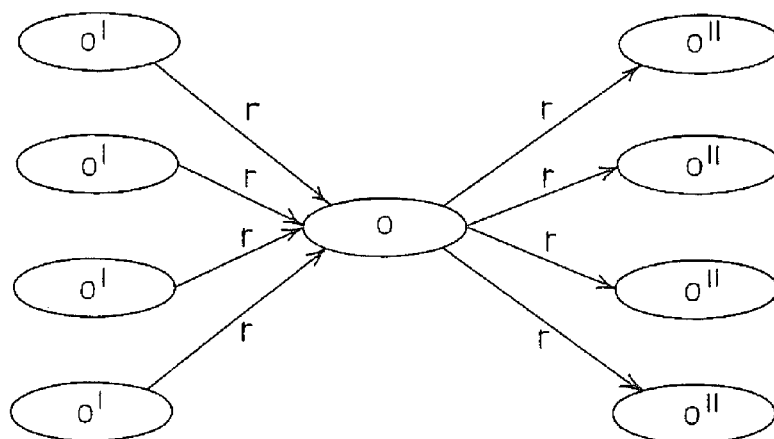(54) Title: SYSTEM AND METHOD FOR CONTENT CLASSIFICATION



FIG. 6

(57) Abstract: The present invention relates to a system for automatic classification of resources, in particular multimedia re-
sources, comprising at least one computer system (12, 20a) having stored therein resources and metadata associated to said re-
sources, said metadata including a plurality of elements, each of said elements comprising a different metadata field. The system
further comprises at least one set of semantic processes arranged for managing one metadata field and for generating in output, by
using as a reference a constant information including at least a semantic lexical network, a plurality of category weights that are
representative of the classification of the metadata field. Present invention also relates to a method for automatic classification of
resources.

## System and method for content classification

**Technical Field**

The present invention relates, in general, to the management of contents, in particular multimedia contents.

5   More in particular, the present invention relates to the automated classification of such contents.

**Background Art**

As known, methods exist for automated classification of textual documents or information (web pages, technical reports, papers, news, etc).

10   The known methods for textual classification, however, are ineffective for classifying non-textual information (songs, images, etc) and known art has used in such contexts a content-based approach (i.e., analysis of the actual images, recordings, videos, etc) or a resource description (metadata) based approach.

Applicants, in general, have noted that such known methods are applicable in very

15   selected contexts, and that no methods have been proved applicable in the general cases, i.e., for automated classification of textual and non-textual information or resources.

As known, the largest multimedia public archives over the Internet (e.g., YouTube, Flickr) use a socially based approach to content description and classification.

In these cases, where most of the content is uploaded by users (UGC, User Generated

20   Content), a description of each new resource (video, image, etc.) is provided by the original uploader, i.e., the user, in a very simplified form (if compared with "professional" multimedia archives, such as clipart libraries, on-line museums, etc).

The description of the resource is composed of some textual fields (such as title, description, and most importantly, tags) specified by the uploader, and may be enriched by

25   other information or resources provided by "viewers" of the resource (e.g., comments, additional tags, new linked contents, bookmarks to the resource, etc).

From publication "Classification of Multi-Media Content (Video's on YouTube) Using Tags and Focal Points" by Ankur Satyendrakumar Sharma and Mohamed Elidrisi published, on the date of filing of present application, at the address "www-

30   users.cs.umn.edu/~ankur/projects.html", a method for classifying Multi-Media contents that uses Tags as input information, is known.

Because tags, as already outlined, may comprise words inserted by users that have uploaded the resource or by other users and may contain English words or words in other

languages than English, slang, acronyms, spelling errors, proper nouns, or anything else, it is apparent that such a prior art suffers at least two problems:

- using only one input information (i.e. tags);

- not managing the inaccuracy of the input information.

Applicants have, in summary, noted that the automated classification of contents is differently managed as a function of the content type and of the context and that no method exists applicable in general cases.

Present invention tries to solve said problems and, in order to better define the content of the invention, in the following description it is assumed that a set of terms are used according to the following definitions, in any case commonly used in the field:

- **Resource** - a resource is any atomic unit of content, that can be identified, classified and searched on a given, preferably on-line, archive.

Preferably the present description refers to multimedia resources, i.e., non-textual resources, but resource may be intended more generally as referring to any kind of resource.

- **Resource Description** - description of a resource is the set of metadata associated to the resource. Such metadata may be, for instance, in the form of text, or tags, or categories. Such metadata may be provided by the original author, or by some website visitor.

For instance, resource description can comprise one or more of the following elements or element types or metadata types or metadata fields, mostly of textual nature:

- Title - a short (1-2 lines) fragment of text, assigned by the author, synthetically identifying the resource;

- Description - a longer text (some paragraphs) describing in more detail the content of the resource;

- Site Category - one or more categories selected among a list of categories; it means that the author must select, from a predefined (by archive or website creators) list of categories, the one or the ones most relevant to the resource. This information is not always reliable, since it comes from the user, and the predefined list of categories is very often ambiguous and vague (not following information architecture state-of-the-art principles);

- Tags - a set of words (uncontrolled keywords) describing the resource. Tags may contain English words or words in other languages than English, slang, acronyms, spelling errors, proper nouns, or anything else. Such tags are selected by the uploader or user, often by picking among the "most popular" tags in a given domain. In some systems tags may

contain internal spaces (such as "Deep Purple"), in other systems they can't (and then one would have "Deep" and "Purple" as separate tags). In some systems other users can add their own (personal) tags to a resource that somebody else has previously uploaded;

- User - an information about which user has uploaded the resource is also a resource description or metadata;

- Comments - comments are, usually, a paragraph or more of text, and may contain textual information useful to identify the resource; in general they are a totally uncontrolled source of information because they are added by other users to the resource;

- Bookmarks - bookmarks are, usually, personal or favorite resources; for instance, other users can add a resource to their bookmarks. In some systems, adding a resources to one's favorites requires the user to select some (personal) tags to classify it.

- **Target categories** - target categories, indicated in the following description as C(i), are application dependent conceptual classes, used to identify and group sets of resources with similar contents.

Preferably, target categories comprise a set of predetermined categories; such a set, according to the preferred embodiment is an input information for the method as disclosed in present invention.

On the basis of the above definition, Applicants have noted that a common problem of the resource description or metadata is that there is no control, no standard, no guarantee about the quality of metadata, nor about their actual relevance to the resource being described.

Therefore, a problem exists on identifying the right target category or categories C(i) associated to a certain resource by analysing metadata.

**Disclosure of the Invention**

Present invention intends to solve the above problem.

In other words, the goal of the present invention is a system and method arranged for automatically identifying the correct category or a set of correct categories C(i) corresponding to a given resource, in particular a multimedia resource, by analyzing only the textual metadata of the resource.

As a matter of facts, such a goal is very useful in all cases where a search is made, for example over the Internet, for rapidly identifying and collect the right information.

According to the present invention, such a problem is solved by means of a system for content management having the features set forth in the claims that follow.

The present invention also relates to a method for content management as well as to a computer program product loadable in the memory of at least one computer unit and including software code portions for performing the steps of the method of the invention when the product is run on at least one computer unit. As used here, the reference to such a computer program product is meant as equivalent to the reference to computer readable medium containing instructions for controlling a system or a device so as to co-ordinate execution of the method according to the invention. Reference to "at least one computer unit" is meant to highlight the possibility for the method of the invention to be carried out in a decentralized manner over a plurality of computer units. Claims are an integral part of the teaching of the present invention.

According to a feature of a preferred embodiment present invention discloses a system for automatic classification of resource descriptions or metadata wherein a computer system is arranged, by using a set of semantic processes, to recognise senses associated to the metadata and to associate category weights to the metadata.

According to a further feature of present invention the set of semantic processes comprises at least one pre-processor block arranged to find the widest possible senses representing all the possible meanings of the metadata, at least one expander block arranged to identify senses that are recurring in the widest possible senses as expanded senses, and to isolate and delete senses that are marginal, and at least one matching block arranged to compare the expanded senses to constant sets of senses corresponding to target categories.

According to another feature of present invention the system comprises an inhibition process arranged for inhibiting the output of the matching block, on the basis of statistical information, so that no classification is made of the metadata field.

**Brief Description of Drawings**

These and further features and advantages of the present invention will appear more clearly from the following detailed description of a preferred embodiment, provided by way of non-limiting examples with reference to the attached drawings, in which components designated by same or similar reference numerals indicate components having same or similar functionality and construction and wherein:

Fig. 1 shows a block diagram of a system according to present invention;

Fig. 2 shows a general architecture of modules implemented in the system of Fig. 1;

Fig. 3 shows the architecture of Fig. 2 in more detail;

Fig. 4 shows an architecture of a block of Fig. 3 according to a first embodiment; and

Fig. 5 shows an architecture of a block of Fig. 3 according to a second embodiment; and

Fig. 6 shows one step of an internal process of Fig. 4 or 5.

**Best mode for Carrying Out the Invention**

With reference to Fig. 1 a system for content identification (system) 10 comprises one (or more) computer servers 12, that host software elements arranged for exposing, for instance, Web applications to a plurality of users or end users connected to the servers 12 through a network 16 by means of computer terminals 14, as for instance personal computers.

The invocation of such applications, for instance the search and navigation of resources, is triggered by the servers 12 that access to services provided by at least one classification server 20a for automatic recognition of the correct information to be used for enriching searching and navigation of the end users.

According to the preferred embodiment of present invention the system 10 is arranged for supplying a reliable and semantically validated information about a multimedia resource and is used in real-time.

Nothing prevents parts of system 10, as for instance servers 12 and classification server 20a, however, from being used in a 'batch' mode, where an entire collection, or portion thereof, is indexed at a time, and the resulting information is stored for later usage by user interfaces or other applications.

According to the preferred embodiment server 20a comprises a set of computer modules or a package 20 (Fig. 1, Fig. 2) having an architecture arranged for receiving:

- in a first input 21a metadata 21 of a single resource d; such metadata comprises, for instance, a plurality of atomic units of content 31a, 31b, ..., 31n (Fig. 2, Fig. 3, Fig. 4) as for instance: title, tags, description, comments, or a subset thereof (as already reported); and

- in a second input or auxiliary input 23 a constant information not depending on specific inputs; such a constant information comprises, for example, one or more of the following inputs:

- a list of target categories 23a;

- a semantic lexical network 23b, for instance the semantic lexical network WordNet as described in "Introduction to WordNet: an on-line lexical database" by G.A. Miller, R. Beckwith, C. Fellbaum, D. Gross, K. J. Miller (Int. J Lexicography, vol. 3, pp. 235-244, January 1990;

- a set of mappings (Category senses) 23c of each category onto a lexical network;

- a set of additional pre-processing tables (Dictionary) 23d.

The package is further arranged for generating in output 25a:

- a plurality of numbers 25 estimating relevance of a resource with respect to target categories; each of such numbers represents, for instance, a number or category weight CW(i) 25 that, for each of the target categories, estimates the relevance of the resource to the corresponding category.

Numbers 25, according to present embodiment, are in a range from 0 to 1 and represent the automatic recognised correct information to be used for enriching searching and navigation of the end users.

In particular, the output 25a of the system 10 comprises a category weight CW(i) for each target category C(i). Each of such weights 25 is a real number in the range from 0 to 1 and a higher value of the category weight means that resource d is more relevant to that specific target category.

Once target categories have been defined as a constant information, the output of the overall classification system simply consists of a different weight associated to each category. An example of the output is shown in Table 1, where each category is given a weight by means of a real number, from 0.000 to 1.000, that estimates the relevance of the resource being classified with each possible target category. The weights are such that 0.000 means "totally not relevant", 1.000 means "with maximum relevance" and intermediate values represent intermediate degrees of relevance.

| Category ID=i | C(i) | CW(i) |
|---|---|---|
| 1 | anime and AMV | 0.644 |
| 2 | action and adventure | 0.122 |
| 3 | movies and trailers | 0.832 |
| 4 | science and technology | 0.000 |
| 5 | web | 0.123 |
| ... | ... | ... |
| N | politics | 0.388 |

Table 1

According to the preferred embodiment of the system 10 the adopted strategy implemented with the package (classification package) 20 is to process a plurality of atomic units of content (metadata) 31a, 31b, ..., 31n (Fig. 3) independently in a set of

parallel "Estimator" blocks 33a, 33b, ..., 33n, and to compute category weights 35a, 35b, ..., 35n independently for each metadata (title or tags or description etc).

These category weights 35a, 35b, ..., 35n are combined by a final "Merge" module 37 so as to classify the multimedia content d into at least one category or to infer that the multimedia content is not relevant to any category.

Each of the parallel Estimator blocks (Estimators) 33a, 33b, ..., 33n works in the same way, by analyzing one field of input resource (metadata field) 31a, 31b, ..., 31n, and by producing a complete list of category weights for each metadata field. Each Estimator 33a, 33b, ..., 33n will give the measure of relevance stemming from the input information 31a, 31b, ..., 31n it received. If the various metadata fields 31a, 31b, ..., 31n are consistent with each other, then it is likely that the estimated target category weights 35a, 35b, ..., 35n would be in agreement, too. Otherwise, for badly assorted or inconsistent metadata, the category weights would be inconsistent as well.

The Merge module 37 measures the similarity among the category weights and outputs a final relevance weight 25 and category or relevance weights and corresponding categories, only in those cases where a sufficient agreement is found.

Each of the Estimators 33a, 33b, ..., 33n respectively, has, preferably, the same external interface and does operate in the same way but on a different metadata field.

In the following, for the sake of simplicity, operation of one Estimator is disclosed by using as reference the Estimator 33a, and by assuming that other estimators are substantially equivalent.

Estimator 33a (Fig. 4) is arranged for operating as follows:

- receiving in input one metadata field 31a of the resource metadata 21. Such a field is represented as a text string, interpreted as a set of words. The Estimator extracts the individual words from the text string (by breaking it at separator characters, such as spaces, commas, etc);

- generating in output a set of category weights CW(i) 35a, for all target categories C(i) reported in the list of target categories 23a.

The estimator 33a might also refuse to make a classification, whenever the characteristics of the input string are not sufficient for unambiguously identifying the relevant categories, as will be explained in the following; in such a case, no category weights are produced, but just an unknown token or an information corresponding to an unknown token.

The Estimator 33a, according to the preferred embodiment, is arranged for using the list of target categories 23a, at least the semantic lexical network 23b and the associated set of mappings 23c of each category onto the lexical network (Category senses).

More preferably the Estimator 33a uses also the set of additional pre-processing tables (Dictionary) 23d.

According to present exemplified embodiment the internal working of the Estimator 33a is based on a semantic representation of information. For this purpose, the semantic lexical network 23b is used to represent the semantic senses related to the words of a certain resource description in the input string.

In the present example, the semantic lexical network WordNet is used, but, as easily comprehensible to a technician in the field, other semantic lexical networks might also be used without departing from the invention as claimed.

As known a semantic lexical network (lexical network or WordNet) 23b consists of a long list of senses. Each sense is associated to a list of lemmas, i.e., words that might signify that sense in English or in a language used by the semantic lexical network; such association is called, as known, semantic mapping.

Under the semantic mapping, each lemma or word may be associated to multiple senses (polysemy, i.e., multiple meanings for a single word), and each sense may be associated to multiple lemmas (synonymy, i.e., multiple words with the same meaning).

In WordNet 23b each different sense is identified by a different numerical value, called offset.

In the lexical network 23b, preferably, senses are connected to each other by means of semantic (binary) relationships r. Such relationships r specify the various kinds of similarities, affinities, and connections among the various senses.

Examples of semantic relationships are: hypernym (a sense is more general than another) or its inverse hyponym (a sense is more specific), part meronym (a sense describes a sub-part, or component, of the other sense) and its inverse part holonym (a sense describes the whole of another part), etc. Each relationship r may be represented as a set of triples: $o_i \rightarrow r \rightarrow o_j$, where $o_i$ and $o_j$ are numerical offsets of two senses, and r is one of the semantic relationships supported by WordNet 23b.

As clearly understandable by a technician in the field, other semantic networks possess similar sets of semantic relationships and may be used by the method disclosed in present invention.

Estimator block 33a comprises a plurality of semantic processes, 41, 43, 45 and 63 that rely on a common representation of a given entity.

This common representation of the semantic information is encoded in a common form, i.e., a weighted set of relevant senses or SenseSet.

From a practical point of view, a SenseSet is a table listing a set of relevant senses, each mapped to a weight value (a real positive number) that represents the "importance" of that sense to influence the meaning of the entity.

According to present invention, from a formal point of view, a SenseSet SS is represented as an incomplete mapping (offset → weight):

$$SS = \{o_1 \rightarrow w_1, o_2 \rightarrow w_2, ...\}.$$

Such a mapping may be represented by means of a table as synthetically reported below in Table 2:

| Sense Offset | Sense Weight |
|---|---|
| 132032 | 3.452 |
| 212031 | 0.122 |
| 100102 | 12.321 |
| 89544 | 23.313 |
| 212346 | 0.843 |

Table 2

Wherein:

- Sense Offset is a constant value read from WordNet;

- Sense Weight is the result of a semantic process as disclosed hereinafter below.

The SenseSet may also be interpreted as a vector on a multi-dimensional space, whose components are the various sense offsets, and whose projections along each component are represented by the sense weights. In such interpretation, all Sense Offsets that do not appear in the SenseSet are implicitly assumed to have a null (zero) projection on the corresponding component. With this vector interpretation, the usual vector operations (such as sum, multiplication by a constant, scalar product, etc) may be trivially extended to operate on SenseSets, too.

Each semantic process inside the estimator blocks receives and generates in input and output, respectively, one or more SenseSet.

The Estimator block 33a elaborates the possible meanings of the input words by means of semantic processes, 41, 43, 45 63, and tries to determine the most relevant senses associated with the input and then relates them to the target categories.

The Estimator 33a proceeds in three subsequent steps:

1.      Pre-Processor block 41: such a block is arranged to find the widest possible SenseSet $SS_F$ representing all the possible meanings of the input text.

2.      Expander block 43: such a block is arranged to navigate the semantic network 23b to identify senses that are recurring in the input SenseSet $SS_F$, and to isolate and delete senses that are marginal (i.e., possible interpretations of the text that are not supported by the context in which that text appears), and compute an updated "expanded" SenseSet $SS_E$ encoding this new (disambiguated) knowledge.

3.      Matching block 45: such a block is arranged to compare the expanded SenseSet $SS_E$ computed for this input text to the SenseSets $SS_{EC}$ representing SenseSets corresponding to target categories C(i); the matching block 45 computes a similarity measure between $SS_E$ and each of the $SS_{EC}$. Such computation comprises the category weights CW(i).

In the following description a detailed disclosure is made of embodiments of the implemented method according to present invention.

The Pre-Processor block (Pre-Processor) 41 provides the following architecture:

- in input the same input of the Estimator block, i.e., a resource metadata field 31a comprising a set of words;

- in Output a SenseSet $SS_F$ representing all the possible interpretations of the input words.

The Pre-Processor is arranged to use the Word-Net semantic network (used for mapping lemmas to senses) 23b and, preferably, at least one of the other auxiliary inputs as for instance the user-generated set of pre-processing instructions (Dictionary) 23d or the SenseSets corresponding to each of the target categories 23c.

The goal of the Pre-Processor 41 is to extract the possible (semantic) interpretations from the input words. The problems solved by the Pre-Processor are of dual nature:

- many words, frequently found in social networks, are not listed in English dictionaries (example: "anime"), or are used with a different meaning than the English standard one (example: "cool" is rarely used to represent fresh temperatures);

- very often, a clear meaning can be attributed to a pair of words, but not to a single word. Example: "hop" may be a jump, with some ambiguity, but "hip hop" is a music and dance genre, without ambiguity.

According to a preferred embodiment of present invention, the Pre-Processor 41 works on (unordered) word pairs, instead of single words.

In particular, for generality and ease of implementation, individual words are encoded as the pair of a word with itself.

As an example, the input string "live madonna concert" is transformed to the following set of word pairs: "live live", "live madonna", "live concert", "madonna madonna", "madonna concert", "concert concert".

Each word pair is then, preferably, looked up into the Dictionary of pre-defined actions (Dictionary) 23d. According to the action stored in the Dictionary 23d for the current word pair, a new SenseSet may be generated, and used for computing the final output SenseSet $SS_F$.

According to the preferred embodiment, the list of possible Actions that may be specified in the Dictionary is shown in Table 3, where the effect of each action is described, and an ActionInfo additional information is specified.

As clearly apparent to a technician in the field, according to further embodiments the list and effect of possible actions may be differently shaped.

| Group | Action | Description | ActionInfo |
|---|---|---|---|
| 1° | Discard | Discard this word pair without any further action | |
| | Remap | Replace this word pair with a specified word pair, if it is not yet in the list of word pairs to be processed, and restart processing with the new pair | a word pair |
| 2° | Category | This word pair generates an output SenseSet equal to the weighted sum of the SenseSet(s) of some specified target categories | weighted list of Category IDs: list of (C(i),AWC(i)) |
| | Semantic | This word pair generates the output SenseSet obtained by looking up both words in WordNet, and by taking all senses relevant to either word | |
| | Senses | This word pair generates the output SenseSet stored in the Dictionary | a SenseSet SS |
| 3° | Default | This word pair was not found in the | |

| | | Dictionary. Proceed as in the case Semantic | |
|---|---|---|---|

Table 3

According to the preferred embodiment, actions of the pre-processor block are divided in three groups:

**1° -       actions that generate no output (Discard and Remap):**

The first action is used to discard frequently used but semantically useless words (often called "stopwords"), such as "the", "for", "of", etc. Remapping is used to normalize a given word pair to an equivalent way expressing the same information (e.g., remapping the word "hiphop", that appears as the word pair "hiphop hiphop", to the pair "hip hop").

**2° -       actions that generate a SenseSet as output (Category, Semantic and Senses) wherein:**

–        Category identifies word pairs that may be mapped unambiguously to a category or a set of target categories (possibly with different weights); in this case the SenseSet(s) of the selected Category/ies are returned.

–        Semantic represents a word pair that is used in the same sense(s) of the English word, therefore the correct SenseSet corresponds to the meaning(s) that WordNet 23b already assigns to those words; in this case, the returned SenseSet is computed by looking up the two words in WordNet.

–        Senses covers words whose meaning is not exactly matching a category, nor exactly corresponding to the WordNet senses; in this case, a custom SenseSet is stored in the dictionary, and is returned as a result.

**3° -     Default action:**

This is an action that is taken whenever the word pair is not found in the Dictionary. In this case it behaves like in the Semantic case, i.e., looking up the words in WordNet 23b.

The default action is also used, for instance, in embodiments where Dictionary is not used.

The different types of actions of the Pre-Processor block are of different semantic relevance, since they represent different degrees of probability that the "right" meaning of the pair of words has been identified, due to the different quality of source information.

Therefore, the four types of actions that may return one or more SenseSets (Category, Semantic, Senses, Default) are associated to an Action Weight (AW) used to influence the contribution of the considered action onto the final returned SenseSet $SS_F$.

For instance, Action Weight may be a table of coefficients as listed in the following Table 4.

| Action type | Preferred range for AW | Example value for AW |
|-------------|------------------------|----------------------|
| Senses | 2.0 – 6.0 | 3.0 |
| Category | 2.0 – 4.0 | 2.5 |
| Semantics | 1.0 – 2.0 | 1.5 |
| Default | 1.0 | 1.0 |

Table 4

Wherein:

(1)

$$AW(Senses) \geq AW(Category) \geq AW(Semantic) > AW(Default) > 0$$

In particular, the final SenseSet $SS_F$ returned by the Pre-Processor block is computed, for instance according to the expression (2), as the weighted sum of all SenseSets returned by the Actions triggered by the Dictionary or by applying the default action, where weights, as listed in Table 4, depend on the type of the i-th Action:

(2)

$$SS_F = \sum_{i \in \{\text{word pairs}\}} AW(Action_i) \cdot SS_i^{Action_i}$$

In the above equation, multiplication of a SenseSet by a real value relies on the vector interpretation of SenseSets, and therefore corresponds to multiplying each of the SenseSet weights by a certain real value (AW).

In the case of "Senses", the returned SenseSet $SS_{i,Senses}$ is taken with weights already stored in the Dictionary.

In the case of "Category", according to the preferred embodiment, one or more SenseSets are returned, depending on how many categories are listed for the word pair in the Dictionary.

More preferably, for each listed category, the SenseSet corresponding to the Category ID (taken from the "Category senses" auxiliary input), may be multiplied by an "Action-Category Weight" (ACW(j)) coefficient specified in the Dictionary for the Category C(j), for instance according to the expression (3).

(3)

$$SS_i^{Category} = \sum_{C(j) \in \{\text{Category list}\}} AWC(j) \cdot SS_{C(j)}$$

In the case of "Semantic" and "Default", a SenseSet is built by taking all the relevant senses in WordNet associated with the two words of the pair, and a constant conventional weight 1.0 is assigned to each sense in the SenseSet $SS_i$ (a sort of "diagonal SenseSet," since there is no additional information for ranking the relevance of various possible senses). Therefore these two action types behave in the same way, except that, later on, the resulting SenseSet will be multiplied by a different Action Weight AW according to (1) and (2).

If WordNet lookup fails for some word, then no SenseSet is returned for that word.

At the end of the procedure a SenseSet $SS_F$, computed according to expression (2), is returned by the Pre-Processor and is given to the Sense Expander block 43.

The Sense Expander block (Expander block) 43 provides the following architecture:

- in input the SenseSet $SS_F$ returned by the Pre-Processor block 41, representing all the possible interpretations of the input words;

- in output an "expanded" SenseSet $SS_E$ representing the most likely and relevant interpretations of the input words;

The Expander block 43, according to the preferred embodiment, uses the semantic lexical network 23b as auxiliary input, as for instance the WordNet semantic network, for navigating semantic relationships among senses.

The goal of the Sense Expander is to extract the most likely coherent interpretation of the input words, as interpreted by the Pre-Processor, and tackling the following problems:

- The input SenseSet $SS_F$ gives a "wide" breadth of interpretation since it contains any sense that might be associated with the input words. This means that, if a word has multiple meanings, all these meanings will be present in $SS_F$, even if only one of these meanings will be really relevant to the resource being classified. Therefore, senses that are out of the "dominant context" must be 'penalised'. In other words, interpretations that would be plausible if a word is taken out of context, but that become "isolated" in the current context defined by the other words, must be penalized.

- Two or more words with extremely similar meanings might have generated a SenseSet $SS_F$ that does not contain identical senses, but only very similar senses. This is due to the very high number of senses in WordNet 23b and to the fact that even small nuances in meaning are represented as different sense offsets.

It is, therefore, important to recognize which senses in $SS_F$ are "near" enough to be considered having essentially the same meaning, and 'strengthen' them, as well as including in the SenseSet $SS_F$ new senses that are strongly connected with existing relevant senses.

According to the preferred embodiment of present invention, the Expander block 43 tries to determine the best interpretation, by navigating the semantic relationships defined in WordNet 23b, and by exploiting the following assumption:

- if the interpretation of the input words is coherent, then it is expected that a large portion of senses in $SS_F$ will be clustered near to one (or a few) group(s) of relevant senses, while other non relevant senses would be isolated and poorly connected with the relevant ones. The notion of "distance" here is defined by considering the traversal of the semantic relationships.

The internal process in the Expander block 43 computes a new "expanded" SenseSet $SS_E$ in which the weights of the relevant senses are greatly increased, and the weights of the non relevant senses have a much lower value. To do this, the weight of each sense is recursively "propagated" to neighbouring senses, as explained below, by adding to each of these neighbours a fraction of the weight of the considered sense.

The notion of neighbouring sense is assumed, in general, by:

1. selecting a subset of WordNet relationships r that should be followed to find neighbouring nodes;

2. defining a fractional relationship weight RW(r), in the range between 0.0 - 1.0, that defines which fraction of the previous node's weight is carried over to the neighbouring nodes, when they are connected through r.

In particular, RW is chosen to favour the generalization of the current senses, in order to find (and weight more) general senses that are common ancestors to most senses in $SS_F$.

According to present exemplified embodiment RW is called "Generalizing Relationship Weights" or GRW and the process is illustrated through the following expression (4):

(4)

$$SS_E = \text{Expand}(SS_F, GRW)$$

Wherein:

GRW is a table of values as shown, for instance, in the following Table 5 wherein L means "low weight" and H means "higher weight" than L.

Preferably, L may be in the range of 0.0 - 0.5.

More preferably L has a value of 0.2 and H a value of 0.5.

| Relationship r | Preferred weight GRW(r) |
|---|---|
| ANTONYM | 0 |
| ATTRIBUTE | 0 |
| CATEGORY | H |
| CATEGORY_MEMBER | 0 |
| CAUSE | 0 |
| DERIVED | L |
| ENTAILED_BY | L |
| ENTAILMENT | 0 |
| HYPERNYM | H |
| HYPONYM | 0 |
| INSTANCE_HYPERNYM | H |
| INSTANCES_HYPONYM | 0 |
| MEMBER_HOLONYM | H |
| MEMBER_MERONYM | 0 |
| NOMINALIZATION | 0 |
| PART_HOLONYM | H |
| PART_MERONYM | 0 |
| PARTICIPLE_OF | 0 |
| PERTAINYM | L |
| REGION | L |
| REGION_MEMBER | 0 |
| SEE_ALSO | 0 |
| SIMILAR_TO | L |
| SUBSTANCE_HOLONYM | L |
| SUBSTANCE_MERONYM | 0 |
| USAGE | L |
| USAGE_MEMBER | 0 |

TABLE 5

The expanded SenseSet $SS_E$ is determined by computing the fixed-point of an equation (5), that analyzes the semantic network as illustrated in Figure 6.

According to present disclosure such equation is defined a "recursive process" since its goal is the computation of all the weights $w_E()$, wherein the computation of each value $w_E(o)$ depends on the value of other values $w_E(o')$, which on turn depends on other $w_E(o'')$, and so on. This definition ("recursive process" or function) is in accordance with similar definition of recursive processes or functions as disclosed, for example, by Lawrence Page in "Method for node ranking in a linked database", U.S. Patent number: 6285999, September 4, 2001.

(5)

$$w_E(o) = \sum_{o':o' \xrightarrow{r} o} (w_E(o') \cdot RW(r)) + w_F(o)$$

Wherein:

- w_E(o) is the weight of the sense with offset o in the SenseSet SS_E;
- w_F (o) is the weight of the sense with offset o in the SenseSet SS_F;
- w_E(o') is the weight of the sense with offset o' in the SenseSet SS_E; and
- o' is the offset of a sense in WordNet from which it is possible to reach the sense with offset o by traversing a relationship r in the selected subset.

The expression (5) discloses a recursive process arranged for stating the weight w_E (o) that each sense with offset (o) should have in the expanded SenseSet SS_E, by considering all selected WordNet relationships r pointing to the node corresponding to the sense with offset o.

For each such relationship, o' is the offset of the node at the other end of the relationship and the weight w_E(o') is summed up, weighted according to the weight assigned to the involved relationship RW(r).

Weights w_F(o) of the initial SenseSet SS_F are the external starting points of the recursive process and are added as an additional term in expression (5).

Whenever a sense offset o changes its weight w_E(o) from a previous value ($w_E^{old}(o)$) to a new value ($w_E^{new}(o)$) due to the evaluation of the recursive expression (5), all weights for other sense offsets o" dependent from o, i.e., senses for which there is at least one relationship o→r→o", need to be re-evaluated (Fig. 6).

Equation (5) is computed by repeatedly computing all w_E(o), and re-computing each of them when at least one of the related weights w_E(o') changes. The process is therefore repeatedly executed until no further changes occur (thus reaching an exact solution), or by stopping the computation when an approximate solution is reached. The approximate solution is defined by a predefined threshold $\tau$, and the process is stopped when all weight variations are below the threshold: $w_E^{new}(o) - w_E^{old}(o) < \tau$ wherein the threshold $\tau$ may be, for instance, in a range of 0.001 - 0.1, more preferably a value of 0.01.

The resulting values w_E(o) are used as components (o → w_E(o)) of the resulting expanded SenseSet SS_E, and are returned as the output of the Sense Expander block 43 (Fig. 4). Such output is used by the Matching block 45.

The Matching block 45 is the last step of the Evaluator block 33a and comprises the process of measuring the similarity between the expanded SenseSet $SS_E$ and the pre-defined target categories.

Each category C(i) is semantically described by a suitable Category SenseSet $SS_{C(i)}$ (Category Senses 23c in Figure 4).

The Matching block 45 provides the following architecture:

- in a first input the expanded SenseSet $SS_E$ returned by the Sense Expander block 43, representing the most plausible coherent interpretations of the input words;

- in a second input the SenseSets $SS_{EC(i)}$ describing all the target categories as described below;

- in output a "relevance weight" CW(i) for each of the target categories C(i), as exemplified in Table 1.

According to the preferred embodiment, the Category SenseSets $SS_{C(i)}$ are manually defined, by choosing in WordNet 23b the most relevant general terms that subsume the actual meaning of the category C(i). The definition of such Category SenseSets is not repeated for each classified resource, but it needs to be defined only once in a configuration phase, when the system is personalized to a specific target domain and target application.

Such SenseSets are, in general more abstract and much more concise (only a handful of sense offsets listed) than $SS_E$. This means that a direct comparison is, in general, not feasible, unless the Category SenseSets $SS_{C(i)}$ are processed with an Expansion process similar to the one already disclosed.

In particular, before actual comparison, each Category SenseSet $SS_{C(i)}$ is first "expanded" through a Category Sense Expander block 63, yielding an "Expanded Category SenseSet" $SS_{EC(i)}$.

This process adopts an expression identical to that (5) described in relation to the Sense Expander block 43, except for the relationship weights. As a matter of fact, while the expansion in the Sense Expander block 43 is aimed at "generalizing" the existing senses, in the Category Sense Expander block 63 it is necessary to "spread" the (few) general senses describing each category into a larger set of more concrete senses.

In this phase, it is therefore used as Relationship Weights a set of "Analyzing Relationship Weights" RW = ARW as shown, for instance in the following Table 6:

| Relationship r | Preferred weight ARW(r) |
|---|---|
| ANTONYM | 0 |

| ATTRIBUTE | 0 |
|---|---|
| CATEGORY | 0 |
| CATEGORY_MEMBER | H |
| CAUSE | 0 |
| DERIVED | 0 |
| ENTAILED_BY | 0 |
| ENTAILMENT | L |
| HYPERNYM | L |
| HYPONYM | H |
| INSTANCE_HYPERNYM | 0 |
| INSTANCES_HYPONYM | H |
| MEMBER_HOLONYM | 0 |
| MEMBER_MERONYM | H |
| NOMINALIZATION | 0 |
| PART_HOLONYM | 0 |
| PART_MERONYM | H |
| PARTICIPLE_OF | 0 |
| PERTAINYM | 0 |
| REGION | 0 |
| REGION_MEMBER | L |
| SEE_ALSO | 0 |
| SIMILAR_TO | 0 |
| SUBSTANCE_HOLONYM | 0 |
| SUBSTANCE_MERONYM | L |
| USAGE | 0 |
| USAGE_MEMBER | L |

TABLE 6


And the following expression:

(6)

$$SS_{EC(i)} = \text{Expand}(SS_{C(i)}, ARW)$$


In the matching block 45 the actual comparison of the Expanded SenseSet $SS_E$ and the various Expanded Category SenseSets $SS_{EC(i)}$ for the various categories C(i) is done by interpreting each SenseSet as a vector in a high-dimensional space (with as many dimensions as sense offsets in WordNet, i.e., over 100,000 dimensions), and by computing the "cosine of the angle" between the vector corresponding to $SS_E$ and each of the vectors corresponding to the various $SS_{EC(i)}$.

Each Category C(i) is therefore assigned to a Category Weight CW(i) computed according to cosine-similarity.

Computation is done, for instance, according to equation (7), where the symbol x stands for the scalar product of two vectors. The scalar product, and the modulus operator, are computed according to the vector interpretation of SenseSets.

(7)

$$CW(i) = \cos(SS_E \angle SS_{EC(i)}) = \frac{SS_E \times SS_{EC(i)}}{|SS_E| \cdot |SS_{EC(i)}|}$$

The final category weights CW(i) (one positive real number, in the range [0,1], per each category) are returned as the final result of the Matching block 45, and indirectly, as the result of the whole Estimator block 33a.

As above disclosed the system 10 is able to automatically extract the most reasonable category weights, starting from the information implicit in the input set of words.

However, the computed category weights might not be correct, due to at least one of the following reasons:

- the input words are insufficient to create a coherent representation of the intended meaning. For examples, in case of less that 4-5 words, it's extremely unlikely that the process could identify a semantic cluster;

- the input words are too semantically dispersed, i.e., they are not related to each other in any way. This means that the input lacks internal coherence. No useful information can be derived in such a situation;

- too many input words are not found in WordNet 23b nor in the Dictionary 23d. This means that very infrequent words have been used, or infrequent acronyms, or proper nouns, etc. In these cases, the words are necessarily discarded.

In these cases each Evaluator block, 33a, 33b, ..., 33n, is forced anyway to make a prediction about the relevance of categories, but such prediction is very likely to be wrong, due to the lack of meaningful and/or usable and/or coherent information.

To overcome this issue, according to a second embodiment of present invention it is provided an alternative set of Evaluator blocks, of which it is shown only the first one 133a (Fig. 5), the others having a structure equivalent to the first one.

The Evaluator blocks according to the second embodiment are alternative to blocks 33a, 33b, ..., 33n.

Each of said alternative blocks further comprises a control or inhibition process or block 48, that monitors the progress of the processes within the Evaluator block, as for instance 133a, and estimates whether the computed classification would be correct.

If the inhibition block 48 suspects that the classification would be wrong, then such a process 48 "inhibits" the output of the matching block 45, so that no classification is made (which is better than a wrong classification).

In the following, as already mentioned, for the sake of simplicity, operation of one Estimator is disclosed by using as reference Estimator 133a, assuming that other estimators, according to the second embodiment, are substantially equivalent.

The inhibition process or logic (inhibitor) 48 relies on purely statistical information about the inputs and outputs of the pre-processing block 41, and the Expander block 43 as reported in Table 7 as reported below, and never considers the actual value of the input words nor of the various SenseSets.

The inhibitor block 48 works by comparing, from a statistical point of view, the statistical indicators of Table 7, with a model representing 'right' and 'wrong' predictions, trained on a sufficiently large set of manually classified documents.

More preferably, the inhibitor block 48 is trained on a statistically significant sample of resources and is able to predict (with a margin of error) whether the identified Categories would be correct, or not.

The inhibitor block 48 receives in input the set of fields listed in Table 7, and produces in output a Boolean value: inhibit/don't inhibit.

| ID Description |
| --- |
| 1 Number of words in the initial text string |
| 2 Number of words composed of only alphabetic characters |
| 3 Number of word pairs found in the Dictionary |
| 4 Number of word pairs with Action=Semantic |
| 5 Number of word pairs with Action=Category |
| 6 Number of distinct senses generated by Action=Semantic |
| 7 Number of distinct senses generated by Action=Category |
| 8 Number of senses after Sense Expansion (in $SS_E$) |
| 9 Category Weight for the highest-ranked Category |
| 10 Category Weight for the second highest-ranked Category |

Table 7

The classification process adopted here is based on a Support Vector Machine known per sè in the field of Data Mining.

Every time that a "bad" classification is inhibited, the precision of the result is increased. On the other hand, if a "good" classification is inhibited, the recall might be reduced. Therefore the adoption of the Inhibitor improves the precision of prediction, at the cost of a possible reduction in the recall.

The Merge block 37 (Fig. 3) has the goal of comparing the information generated by the various Estimator blocks, for instance 33a or 133a, and of determining whether such information is consistent.

When Estimator blocks, for instance 33a or 133a (Fig. 3, Fig. 4, Fig. 5), agree on a common interpretation of the resource metadata (meaning that the Target Categories that receive a high ranking are nearly the same, and the Target Categories receiving a low or null ranking are nearly the same), then the Merge block 37 averages the relevant category weights and returns the overall result.

In the contrary case, when the category rankings disagree, no category weights are returned at all, and in this case the whole system remains "silent" and no classification is made, because inconsistent information has been detected in the input information.

Processes used in the Merge block may be different. For example, the process defined by D. Wolfran and H.A. Olson in "A method for comparing large scale inter-indexer consistency using ir modelling" in Canadian Association for Information Science, Canadian Association for Information Science, 2007 may be used.

The processes used in the Merge block fall into the general problem known in the literature as "Combining multiple classifiers", for which several solutions exist. For example the Merge block could adopt the methods described by L. Xu, A. Krzyzak, and C.Y.Suen in "Methods of Combining Multiple Classifiers and Their Applications to Handwriting Recognition" in IEEE Transactions on Systems, Man and Cybernetics, vol. 22, no. 3, May/June 1992, by computing a resulting classifier based on the average value (equation (4) in the cited paper) or on the median value (equation (7)).

The known combination of multiple classifiers can be considered as corresponding to the merge block as disclosed according to present invention.

As a matter of fact the category weights can be considered as type three classifiers as defined in the known document.

It is also apparent that, in the case in which there is only one Evaluator block (because there is only one kind of metadata input, or all metadata fields have been concatenated into one string), the Merge block is not needed.

Operation of the system as disclosed is the following.

As soon as an end user or a Web service provider desires to file, search, visualise, ... a resource d and an associated resource description 21 into a computer server 12, such server 12 accesses the classification server 20a and requests an on-line automatic classification of the resource d and of the associated resource description 21.

The classification server 20a, by using the classification package 20 and auxiliary input 23 is arranged to automatically generate in output an unbiased classification associated to the evaluated resource.

As a result of the above process, resources may be searched by end users of the system 10 without any loss of time. As a matter of fact, wrongly classified resources are avoided.

The same result is granted to Web applications provided by service providers on servers 12 as for instance for advertisement services, research services, etc.

The above summarised process may be exemplified as follows through a real example with real numbers and results.

Let us consider a multimedia resource, for example a YouTube video, with the following metadata:

| Video ID | _bsniYwSaWg |
| Title | Britney Spears - Baby One More Time Pop Music Video |
| Tags | Britney Spears - Baby One More Time |
| Description | Britney Spears - Baby One More Time from the album - Baby One More Time<br>(C) 2003 Zomba Records |
| Author | BritneyTV |
| Upload date | 11 march 2007 |
| Category | Music |

Table 8

Assuming to describe the Estimator block 33a arranged to process metadata "Title" in Table 8, the input string considered by the Estimator block 33a is "Britney Spears - Baby One More Time Pop Music Video".

The pre-processor block 41 identifies the following words:

1.      Britney

2.      Spears

3.      Baby

4.      One

5    5.      More

6.      Time

7.      Pop

8.      Music

9.      Video

10

The elaboration of the pre-processor block 41 of the word pairs gives the following results:

| Action type | Word pairs | Notes |
|---|---|---|
| DISCARD | 1 | - The pair (britney, britney) is marked as DISCARD in the dictionary. |
| DEFAULT | 34 | Word pairs that are not found in the dictionary. |
| REMAP | 7 | - (One, One) is remapped to (one, one); the new pair (one, one) then generates a DISCARD action.<br>- (Time, Time) is remapped to (time, time); the new pair (time, time) then generates a SEMANTICS action.<br>- (Pop, Pop) is remapped to (pop, pop); the new pair (pop, pop) then generates a SEMANTICS action.<br>(Music, Music) is remapped to (music, music); the new pair (music, music) then generates a SEMANTICS action.<br>…etc…<br>In total, for instance, from the REMAP actions we have<br>-   1 DISCARD<br>-   5 SEMANTICS<br>-   1 CATEGORY |
| CATEGORY | 2 | - (britney, spears) generates a CATEGORY action for category number 3.<br>- (music, music) generates a CATEGORY action for category number 3. |
| SEMANTICS | 41 | - Number of Word pairs that are found in the dictionary with Action = Semantics |
| TOTAL | 85 | Total number of considered word pairs |

Table 9

15    All the actions marked SEMANTICS, CATEGORY, DEFAULT in Table 9 generate a SenseSet to be added to the result SS$_F$ of the pre-processor block 41. Some of these

SenseSets may be identical (in particular when multiple CATEGORY actions point to the same Category).

In this example, $SS_F$ is composed of 14 distinct SenseSets $SS_i$.

The weighted sum of all $SS_i$ gives the $SS_F$, that in this case is composed of 69 different sense offsets.

The following Table 10 shows a portion (less than half) of $SS_F$, from which we may already appreciate that all the possible meanings of the words in the title are taken into account.

We may also appreciate that several lemmas (words) are associated to several senses (all inflections of the meaning of the word).

| Sense offset | Weight | Lemma(s) corresponding to the sense offset |
|---|---|---|
| 99341 | 1.0 | ADVERB - more, to_a_greater_extent, |
| 99712 | 1.0 | ADVERB - more, |
| 505410 | 1.0 | ADJECTIVE - matchless, nonpareil, one(a), one_and_only(a), peerless, unmatched, unmatchable, unrivaled, unrivalled, |
| 702642 | 1.0 | ADJECTIVE - one(a), |
| 796767 | 1.0 | NOUN - baby, |
| 1322221 | 1.0 | NOUN - baby, |
| 1330506 | 1.0 | ADJECTIVE - one, |
| 1444887 | 1.0 | VERB - spear, |
| 1555133 | 1.0 | ADJECTIVE - more(a), more_than, |
| 1556355 | 1.0 | ADJECTIVE - more(a), |
| 1677623 | 1.0 | ADJECTIVE - one(a), |
| 2064427 | 1.0 | ADJECTIVE - one(a), |
| 2186338 | 1.0 | ADJECTIVE - one, 1, i, ane, |
| 2477885 | 1.0 | ADJECTIVE - one(a), unitary, |
| 2570267 | 1.0 | VERB - pamper, featherbed, cosset, cocker, baby, coddle, mollycoddle, spoil, indulge, |
| 2714200 | 1.0 | VERB - spear, spear_up, |
| 4270891 | 1.0 | NOUN - spear, lance, shaft, |
| 4271148 | 1.0 | NOUN - spear, gig, fizgig, fishgig, lance, |
| 5870055 | 1.0 | NOUN - one, |
| 9827363 | 1.0 | NOUN - baby, babe, sister, |
| 9827519 | 1.0 | NOUN - baby, |
| 9827683 | 1.0 | NOUN - baby, babe, infant, |
| 9828216 | 1.0 | NOUN - baby, |
| 9918554 | 1.0 | NOUN - child, baby, |
| 11190183 | 1.0 | NOUN - More, Thomas_More, Sir_Thomas_More, |
| 13742573 | 1.0 | NOUN - one, 1, I, ace, single, unity, |
| 28270 | 2.0 | NOUN - time, |
| 296973 | 2.0 | VERB - time, |
| 297906 | 2.0 | VERB - time, |

| 309582 | 2.0 | VERB - pop, |
|--------|-----|-------------|
| 309792 | 2.0 | VERB - pop, |
| ... | ... | ... etc ... |

Table 10

This is the output of the Pre-processor block 41, that is then passed on to the Sense Expander block 43.

The Sense Expander computes the Expanded SenseSet $SS_E$, that is much larger (217 distinct sense offsets, in this case), because new relevant and related senses are now included, as above disclosed (Table 11).

| Sense offset | Weight | Lemma(s) corresponding to the sense offset |
|--------------|--------|---------------------------------------------|
| 2157519 | 0.59 | VERB - crop_up, pop_up, pop, |
| 2185988 | 0.59 | VERB - pop, |
| 2186192 | 0.59 | VERB - pop, |
| 4534127 | 0.59 | NOUN - video_recording, video, |
| 4991738 | 0.59 | NOUN - meter, metre, time, |
| 5718556 | 0.59 | NOUN - music, euphony, |
| 5718935 | 0.59 | NOUN - music, |
| 6277803 | 0.59 | NOUN - video, picture, |
| 6277992 | 0.59 | NOUN - video, |
| 7059962 | 0.59 | NOUN - pop_music, pop, |
| 7288215 | 0.59 | NOUN - time, |
| 7309599 | 0.59 | NOUN - time, clip, |
| 7390400 | 0.59 | NOUN - pop, popping, |
| 7927512 | 0.59 | NOUN - pop, soda, soda_pop, soda_water, tonic, |
| 9988063 | 0.59 | NOUN - dad, dada, daddy, pa, papa, pappa, pop, |
| 15122231 | 0.59 | NOUN - time, |
| 15129927 | 0.59 | NOUN - clock_time, time, |
| 15135822 | 0.59 | NOUN - fourth_dimension, time, |
| 15224692 | 0.59 | NOUN - prison_term, sentence, time, |
| 15245515 | 0.59 | NOUN - time, |
| 15270431 | 0.59 | NOUN - time, |
| 414518 | 0.59 | ADJECTIVE - popular, pop, |
| 543233 | 0.60 | NOUN - music, |
| 6277280 | 0.64 | NOUN - television, telecasting, TV, video, |
| 2186338 | 0.64 | ADJECTIVE - one, 1, i, ane, |
| 7020895 | 0.74 | NOUN - music, |
| 2183611 | 0.92 | ADJECTIVE - cardinal, |
| ... | ... | ... etc ... |

Table 11

This Expanded SenseSet $SS_E$ just computed by the Sense Expander Block 43 is then given as input to the Matching Block 45, that computes the cosine-similarity of $SS_E$

with all the Categories Expanded Senses $SS_{EC(i)}$ returned by the Category Sense Expander block 63. Such similarities are the final Category Weights CW(i) returned by the Matching Block 45. A subset resulting from these comparisons are shown in the following Table 12:

| Category ID $i$ | Category name | Cosine similarity $CW(i)$ |
|---|---|---|
| … | … | … |
| 4 | science and technology | 0.0010 |
| 5 | web | 0.0 |
| … | … | … |
| 12 | music videos | 0.058 |
| … | … | … |
| 15 | cars and vehicles | 0.0 |
| … | … | … |
| 26 | Tv programs and shows | 0.027 |
| … | … | … |
| 29 | music | 0.024 |
| … | … | … |
| 44 | video game and software | 0.0070 |
| 45 | sports | 0.0060 |
| … | … | … |

Table 12

We may observe that for many categories, the Category Weight CW(i) is null (0.0), i.e., they are totally non-relevant categories. Other categories are assigned a higher or lower level, depending on the stronger or weaker similarity with the resource expanded SenseSet.

The final result of title classification process may therefore be presented as the "ranked" list of top-relevant categories for the resource analysed, as in the present example a video:

1.     celebrities (0.228)

2.     music videos (0.058)

3.     anime and AMV (0.042)

4.     Tv programs and shows (0.027)

5.     music (0.024)

The same process, obviously, will be executed for all other resource descriptions so as to obtain, at the end, by means of the "Merge" module 37 an automatic unbiased classification the resource d into at least one category.

Of course, obvious changes and/or variations to the above disclosure are possible, as regards devices and connections, as well as details of the described construction and operation method without departing from the scope of the invention as defined by the claims that follow.

CLAIMS

1. System for automatic classification of resources, in particular multimedia resources, comprising

- at least one computer system (12, 20a) having stored therein resources (d) and metadata (21) associated to said resources, said metadata including a plurality of elements, each of said elements comprising a different textual metadata field;

- at least one set of semantic processes (41, 43, 45, 63) arranged for managing one metadata field and for generating in output, by using as a reference a constant information including at least a semantic lexical network (23b), a plurality of category weights (CW(i), 25) that are representative of the classification of the metadata field, wherein each category weight is associated to a target category reported in a predetermined list of target categories (23a).

2. System according to Claim 1 characterised in that said different textual metadata fields comprise at least one element selected in the group comprising

- a title;

- a description;

- a site category;

- tags;

- user information;

- comments;

- bookmarks.

3. System according to Claim 1 or 2 characterised in that said constant information further comprises at least one constant information selected in the group comprising

- a set of mappings (23c) of each target category onto a semantic lexical network (23b);

- a set of pre-processing tables 23d.

4. System according to any one of claims 1 to 3 characterised in that said set of semantic processes (41, 43, 45, 63) comprises

- at least one pre-processor block (41) arranged to find the widest possible senses ($SS_F$) representing all the possible meanings of the metadata type;

- at least one expander block (43) arranged to navigate the semantic lexical network (23b) to identify senses that are recurring in the widest possible senses ($SS_F$), and to isolate and delete senses that are marginal, and compute expanded senses ($SS_E$) encoding a new disambiguated knowledge;

- at least one matching block (45) arranged to compare the expanded senses ($SS_E$) to constant senses ($SS_{EC(i)}$) corresponding to target categories ($C(i)$), and to assign said plurality of category weights ($CW(i)$).

5. System according to claim 4 characterised by

- at least an inhibition block (48) that inhibits the output of the matching block (45), on the basis of statistical information, so that no classification is made of the metadata field.

6. Method for automatically assigning a classification to a resource accessible in a computer network, the method being arranged to automatically execute the following steps

- receiving in input to a computer system (12, 20a) metadata (21) associated to the resource (d), said metadata including a plurality of elements, each of said elements comprising a different textual metadata field;

- managing one metadata field by means of one set of semantic processes (41, 43, 45, 63), said set of semantic processes being arranged, by using as a reference a constant information including at least a semantic lexical network (23b), to:

       - associate to said one metadata field one predetermined list of target categories ($C(i)$);

       - estimating for each target category a category weight ($CW(i)$) to be associated to said one metadata field so that each metadata field is weighted by category weights of the list of categories, each weight being representative of the relevance of that metadata field to each target category, wherein higher weights represent higher relevance of said metadata field to the category.

7. Method according to claim 6 wherein said step of managing one metadata field comprises the steps of:

       - finding the widest possible senses ($SS_F$) representing all the possible meanings of the metadata field;

       - navigating the semantic lexical network (23b) to identify senses that are recurring in the widest possible senses ($SS_F$) by isolating and deleting senses that are marginal, and by computing expanded senses ($SS_E$) encoding a new disambiguated knowledge

       - comparing the expanded senses ($SS_E$) to constant senses ($SS_{EC}$) corresponding to target categories ($C(i)$), and assigning said plurality of category weights ($CW(i)$) to said metadata field.

8. Method according to claim 7 wherein said step of finding the widest possible senses ($SS_F$) representing all the possible meanings of the metadata field comprises the step of

       - working on word pairs of said metadata field.

9. Method according to any one of claims 7 to 8 wherein the step of comparing the expanded senses ($SS_E$) to constant senses ($SS_{EC}$) is followed by the further step of:
- inhibiting the comparing result on the basis of statistical information, so that no classification is made of the metadata field.

10. Computer program product or set of computer program products loadable in the memory of at least one computer and including software code portions arranged to perform, when the product is run on at least one computer, the method according to any one of claims 6 to 9.
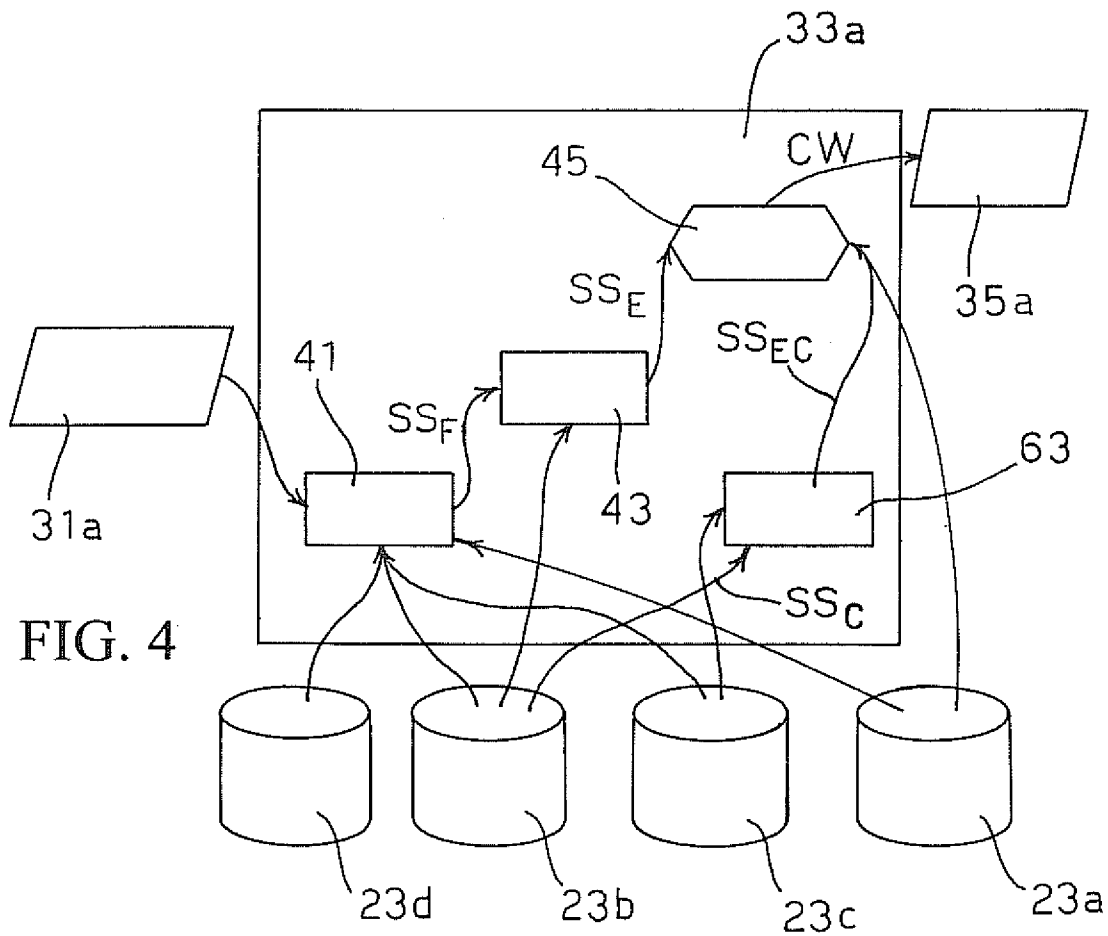
## 1/3



FIG. 1



FIG. 2

FIG. 3



FIG. 4

FIG. 5



FIG. 6

# INTERNATIONAL SEARCH REPORT

## A. CLASSIFICATION OF SUBJECT MATTER
INV.  G06F17/30
ADD.

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, INSPEC

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category* | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|---|---|---|
| X | GHARIB T F ET AL:  "Web document clustering approach using wordnet lexical categories and fuzzy clustering" COMPUTER AND INFORMATION TECHNOLOGY, 2008. ICCIT 2008. 11TH INTERNATIONAL CONFERENCE ON, IEEE, PISCATAWAY, NJ, USA, 24 December 2008 (2008-12-24), pages 48-55, XP031443129 ISBN: 978-1-4244-2135-0 * abstract page 49 - page 51 -----  -/-- | 1-10 |

☒ Further documents are listed in the continuation of Box C.                ☐ See patent family annex.

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

| Date of the actual completion of the international search | Date of mailing of the international search report |
|---|---|
| 8 November 2010 | 15/11/2010 |

| Name and mailing address of the ISA/ | Authorized officer |
|---|---|
| European Patent Office, P.B. 5818 Patentlaan 2 NL – 2280 HV Rijswijk Tel. (+31–70) 340–2040, Fax: (+31–70) 340–3016 | Bernardi, Luca |

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

| Category* | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|---|---|---|
| X | PANDYA A ET AL: "Text similarity measurement using concept representation of texts" PATTERN RECOGNITION AND MACHINE INTELLIGENCE. FIRST INTERNATIONAL CONFERENCE, PREMI 2005. PROCEEDINGS (LECTURE NOTES IN COMPUTER SCIENCE VOL.3776) SPRINGER-VERLAG BERLIN, GERMANY, 2005, pages 678-683, XP002582240 ISBN: 3-540-30506-8 the whole document | 1-10 |
| A | CHIN O S ET AL: "Automatic discovery of concepts from text" PROCEEDINGS OF THE IEEE/WIC/ACM INTERNATIONAL CONFERENCE ON WEB INTELLIGENCE IEEE COMPUT. SOC LOS ALAMITOS, CA, USA, 2006, page 4 PP., XP002582242 ISBN: 0-7695-2747-7 the whole document | 1-10 |
| A | YING LIU ET AL: "Using WordNet to disambiguate word senses for text classification" COMPUTATIONAL SCIENCE-ICCS 2007. 7TH INTERNATIONAL CONFERENCE. PROCEEDINGS, PART III (LECTURE NOTES IN COMPUTER SCIENCE VOL.4489) SPRINGER BERLIN, GERMANY, 2007, pages 781-789, XP002582243 ISBN: 3-540-72587-3 the whole document | 1-10 |
| A | Wolfram D et al.: "A Method for Comparing large Scale Inter-indexer Consistency Using IR Modeling" Canadian Association for Information Science 2007, pages 1-8, XP002582270 Retrieved from the Internet: URL:http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.83.8621 [retrieved on 2010-05-11] cited in the application the whole document | 1-10 |