

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
31 July 2003 (31.07.2003)

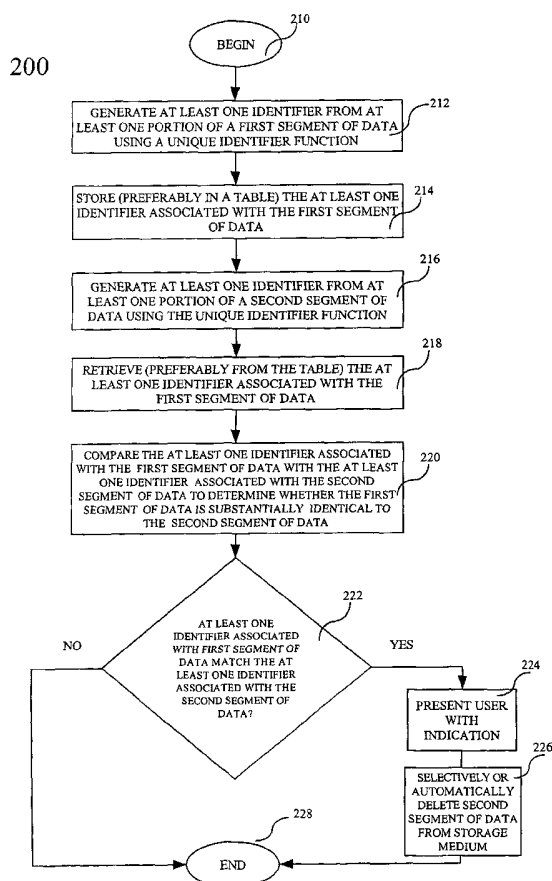
PCT

(10) International Publication Number
WO 03/062996 A1

- (51) International Patent Classification⁷: **G06F 12/00** Apt. D, 9339 Notre Dame Drive, Indianapolis, IN 46240 (US).
- (21) International Application Number: PCT/US03/01194
- (22) International Filing Date: 15 January 2003 (15.01.2003)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data: 10/051,999 17 January 2002 (17.01.2002) US
- (71) Applicant (for all designated States except US): **THOMSON LICENSING S.A.** [FR/FR]; 46, quai A. Le Gallo, F-92648 Boulogne (FR).
- (71) Applicants and
- (72) Inventors: **SCHULTZ, Mark, Alan** [US/US]; 4437 Somerset Way S., Carmel, IN 46033 (US). **LIN, Shu** [CN/US];
- (72) Inventor; and
- (75) Inventor/Applicant (for US only): **KELLY, Michael, Gene** [US/US]; 12115 Wellingham Court, Fort Wayne, IN 46845 (US).
- (74) Agents: **TRIPOLI, Joseph, S.** et al.; c/o Thomson Multimedia Licensing Inc., Suite 2, 2 Independence Way, Princeton, NJ 08540 (US).
- (81) Designated States (national): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, OM, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

[Continued on next page]

(54) Title: SYSTEM AND METHOD FOR SEARCHING FOR DUPLICATE DATA



(57) Abstract: The invention concerns a method (200) and system (100) for searching for duplicate data. The method (200) includes the steps of: generating (212) at least one identifier from at least one portion of a first segment of data using a unique identifier function; generating (216) at least one identifier from at least one portion of a second segment of data using the unique identifier function; and comparing (220) at least one identifier associated with the first segment of data with at least one identifier associated with the second segment of data to determine whether the first segment of data is substantially identical to the second segment of data.

WO 03/062996 A1



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

— *before the expiration of the time limit for amending the claims and to be republished in the event of receipt of amendments*

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.

Published:

— *with international search report*

SYSTEM AND METHOD FOR SEARCHING FOR DUPLICATE DATA

BACKGROUND OF THE INVENTION

1. Technical Field

5 The inventive arrangements relate generally to recording systems and more particularly to multimedia recording systems that record digitally encoded signals onto disc media such as hard drives and recordable optical discs.

2. Description of Related Art

 Currently, many forms of data can be recorded onto many different types of
10 storage media. As an example, many consumers record television programs or music onto an optical disc medium or a hard disc drive (HDD). As technology has improved, the storage capacity of optical disc media and HDD has significantly increased. In fact, some HDDs can store well over 50 gigabytes of data. As such, a consumer can record a large number of programs or songs on this type of storage
15 medium.

 When data is recorded onto a recordable storage medium, the recordable storage medium device typically permits the user to enter a title for purposes of identifying the recorded work. These titles may be useful when the user wishes to locate a particular piece of recorded data to determine whether the user has
20 previously recorded such data. Significantly, however, this process of searching may be laborious, inefficient and prone to errors, as the storage medium may contain hundreds or even thousands of titles. This problem may be particularly acute if the

storage medium is a large HDD or if the titles of certain data segments were given default titles.

Even assuming a data segment on a storage medium could be located relatively easily by searching for the title, a particular title may be the same for
5 different data segments. For example, if a song is recorded on a storage medium and is given a title based on the name of the song, a second song can be recorded later that has a name that is identical to the name of the first song. This confusion may occur, for example, if two separate artists record different versions of the same song. When recording the second song, the user may check the titles of the songs
10 previously recorded and may mistakenly assume that the second song has already been recorded. Thus, a need exists for a system and method for searching for duplicate data without increasing system costs or complexity and further reducing the possibility for errors when searching and considering deletion of duplicate data.

Summary of the Invention

15 The present invention concerns a method of searching for duplicate data. The method includes the steps of: generating at least one identifier from at least one portion of a first segment of data using a unique identifier function; generating at least one identifier from at least one corresponding portion of a second segment of data using the unique identifier function; and comparing the at least one identifier
20 associated with the first segment of data with the at least one identifier associated with the second segment of data to determine whether the first segment of data is substantially identical to the second segment of data.

In one arrangement, the step of generating at least one identifier from at least one portion of a first segment of data can include the step of generating at least one

identifier from the at least one portion of the first segment data using a unique identifier function as the first segment of data is recorded onto a storage medium or after the first segment of data is recorded onto the storage medium. In addition, the step of generating at least one identifier from at least one portion of a second
5 segment of data can include the step of generating the at least one identifier from the at least one portion of a second segment of data using the unique identifier function as the second segment of data is recorded onto the storage medium. Moreover, the step of generating at least one identifier from at least one portion of a second segment of data can occur as the second segment of data is recorded onto a
10 different storage medium.

In one aspect, the first segment of data and the second segment of data can be segments of multimedia data. The method can also include the steps of: storing in a table the at least one identifier associated with the first segment of data; and retrieving from the table the at least one identifier associated with the first
15 segment of data prior to the comparing step. In addition, the method can include the step of presenting an indication that the first segment of data is substantially identical to the second segment of data when at least one identifier associated with the first segment of data matches the at least one identifier associated with the second segment of data.

20 In another arrangement, the size of the at least one portion of the first segment of data and the at least one portion of the second segment of data can be based on a temporal measurement or a bit measurement. The at least one portion of the first segment of data can correspond temporally or correspond bit by bit with the at least one portion of the second segment of data. In another aspect, the at

least one identifier associated with the first segment of data and the at least one identifier associated with the second segment of data can be hash values, and the unique identifier function can be a hash function in which the hash value associated with the first segment of data will equal a hash value associated with the second
5 segment of data when the first segment of data and the second segment of data are identical.

Also, the comparing step can include the step of comparing a plurality of identifiers associated with the first segment of data with a plurality of identifiers associated with the second segment of data to determine whether the first segment
10 of data is substantially identical to the second segment of data. Further, the comparing step can include the step of comparing a plurality of identifiers associated with a first set of segments of data with a plurality of identifiers associated with a second set of segments of data to determine whether the first set of segments of data is substantially identical to the second set of segments of data.

15 The present invention also concerns a system for searching for duplicate data. The system includes: a controller for reading data from and writing data to a storage medium; and a processor, wherein the processor is programmed to: generate at least one identifier from at least one portion of a first segment of data using a unique identifier function; generate at least one identifier from at least one
20 corresponding portion of a second segment of data using the unique identifier function; and compare the at least one identifier associated with the first segment of data with the at least one identifier associated with the second segment of data to determine whether the first segment of data is substantially identical to the second

segment of data. The system also includes suitable software and circuitry to implement the methods as described above.

Brief Description of the Drawings

FIG. 1 is a block diagram of a system that can search for duplicate data in accordance with the inventive arrangements herein.

FIG. 2 is a flow chart that illustrates an operation of searching for duplicate data in accordance with the inventive arrangements.

Detailed Description of the Preferred Embodiments

A system 100 for implementing the various advanced operating features in accordance with the inventive arrangements is shown in block diagram form in FIG. 1. The invention, however, is not limited to the particular system illustrated in FIG. 1, as the invention can be practiced with any other system capable of receiving a digitally encoded signal. In addition, the system 100 is not limited to reading data from or writing data to any particular type of storage medium, as any storage medium capable of storing digitally encoded data can be used with the system 100.

The system 100 can include a controller 110 for reading data from and writing data to a storage medium 112. The controller can also read data from and write data to a different storage medium or memory 120. The system 100 can also have a microprocessor 114, a table or memory 116 and a display 118. Control and data interfaces can also be provided for permitting the microprocessor 114 to control the operation of the controller 110 and the display 118 and to retrieve information stored in the table 116. Suitable software or firmware can be provided in memory for the conventional operations performed by the microprocessor 114. Further, program routines can be provided for the microprocessor 114 in accordance with the inventive

arrangements. Moreover, any other suitable software or circuitry can be used in place of the microprocessor 114.

In operation, the controller 110 can write a first segment of data to the storage medium 112. In one arrangement, as the first segment of data is recorded to the storage medium 112, the microprocessor 114 can generate at least one identifier from at least one portion of the first segment of data using a unique identifier function. Once the microprocessor 114 generates at least one identifier from the at least one portion of the first segment of data, the microprocessor 114 can transfer this at least one identifier to the table 116. In another arrangement, the at least one identifier associated with the first segment of data can be generated anytime after the first segment of data has been recorded to the storage medium 112.

The microprocessor 114 can also generate at least one identifier for at least one corresponding portion of a second segment of data using the unique identifier function. The microprocessor 114 can generate the at least one identifier associated with the second segment of data as the second segment of data is recorded onto the storage medium 112 or, alternatively, as the second segment of data is recorded in the memory 120. It is understood that the memory 120 can be any suitable form of memory for storing digitally encoded data.

Once generated, the microprocessor 114 can retrieve from the table 116 the at least one identifier associated with the first segment of data. The microprocessor 114 can then compare the at least one identifier associated with the first segment of data to the at least one identifier associated with the second segment of data to determine whether the first segment of data is substantially identical to the second segment of data. If the at least one identifier associated with the first segment of

data matches the at least one identifier for the second segment of data, then the first segment of data is substantially identical if not completely identical to the second segment of data. The microprocessor 114 can then present an indication to a user through the display 118 that the two segments of data are identical. The overall operation of the invention will be discussed in greater detail below.

SEARCHING FOR DUPLICATE DATA

FIG. 2 illustrates a flowchart 200 that demonstrates an operation for searching for duplicate or identical data. At step 210, the process can begin. As shown at step 212, at least one identifier from at least one portion of a first segment of data can be generated using a unique identifier function. The identifier can be generated as the first segment of data is recorded onto a storage medium. Conversely, the identifier can be generated anytime after the first segment of data has been recorded onto the storage medium.

The first segment of data can be any suitable type of data including text-based data, audio, video or any combination thereof or any other suitable form of data. The first segment of data can also be an encrypted or non-encrypted segment of data. Further, an identifier can be generated from any portion of the first segment of data, including non-consecutive or non-successive portions of the first segment of data. Moreover, more than one identifier can be generated from any portion of data contained in the first segment of data. The size of the portion of data from which an identifier is to be generated can be based on a temporal measurement or a bit measurement.

As an example, if the first segment of data is a song, the identifier can be generated from the entire song such that the at least one portion includes the entire

first segment of data. As another example, the song can be divided into two discrete portions: a beginning portion and an ending portion. If the size of both of these portions of the song is to be based on a temporal measurement, then the beginning portion can include the first 30 seconds of the song and the ending portion can include the last 30 seconds of the song. According to the inventive arrangements, these two portions of the song can be combined and at least one identifier can be generated from this combination. Thus, one or more identifiers per segment of data can be used for comparison with a corresponding number of identifiers associated with another segment of data.

Continuing with the example, an identifier can be generated from both portions of the song such that two separate identifiers are generated from the same song. Alternatively, an identifier can be generated from the temporal measurement between the beginning and ending portions. Moreover, if the size of the at least one portion of data is based on a bit measurement, then an identifier can be generated from, for example, the first 1 Mb of data in the song. It is noted, however, that the invention is not limited to the foregoing examples, as any number of identifiers can be generated from any number of portions of a first segment (including when the at least one portion includes the entire first segment of data) of any suitable type of data.

Referring back to the flowchart 200, the at least one identifier associated with the first segment of data can preferably be stored in table, as shown at step 214. At step 216, at least one identifier from at least one portion of a second segment of data can be generated using the unique identifier function. The generation of the at least one identifier associated with the second segment of data can be in accordance with

the process discussed in relation to the first segment of data (see step 212 discussion). To increase accuracy, however, the portion(s) of the second segment of data from which the at least one identifier is generated can correspond to the at least one portion of the first segment of data. This correspondence can be

5 temporally based or a bit by bit basis.

For example, if the first segment of data is a song and the at least one identifier associated with the first segment of data is generated from the entire song (the at least one portion includes the entire first segment of data), then to maximize accuracy, the at least one identifier associated with the second segment of data can

10 be generated from the entire song (assuming that the second segment of data is indeed a song). As another example, if the at least one portion of the first segment of data includes the first 1Mb of data and the at least one identifier associated with the first segment is generated from that portion, then it is preferred to generate the at least one identifier associated with the second segment of data from the first 1Mb of

15 data in the second segment of data.

In one arrangement, the at least one identifier associated with the second segment of data can be generated as the second segment of data is recorded onto the same storage medium onto which the first segment of data was recorded. Conversely, the at least one identifier associated with the second segment of data

20 can be generated as the second segment of data is recorded onto a different storage medium.

At step 218, once the appropriate identifier or identifiers have been generated from the second segment of data, the at least one identifier associated with the first segment of data can be retrieved from memory, preferably the table. At step 220,

the at least one identifier associated with the first segment of data can be compared to the at least one identifier associated with the second segment of data to determine whether the first segment of data is substantially identical to the second segment of data. If the identifiers are equal, then the first segment of data is virtually always
5 identical to the second segment of data. At decision block 222, when at least one identifier associated with the first segment of data matches the at least one identifier associated with the second segment, then a user can be presented with an indication that the first segment of data is substantially identical to the second segment of data, as shown at step 224. Further, at step 226, any portion of the
10 second segment of data that was recorded onto a storage medium for purposes of the comparison step can be deleted either selectively by a user or automatically. The process can end at step 228.

In one arrangement, the at least one identifier can be a hash value. In addition, the unique identifier function can be a hash function. A hash value
15 associated with the first segment of data can equal a hash value associated with the second segment of data when the first segment of data and the second segment of data are identical or substantially identical. An example of several hash functions that can be used to practice the invention is an exclusive-or function. It is understood, however, that the invention is not limited to this particular hash function,
20 as any other suitable hash function can be used.

Although the present invention has been described in conjunction with the embodiments disclosed herein, it should be understood that the foregoing description is intended to illustrate and not limit the scope of the invention as defined by the claims.

CLAIMS

1 1. A method of searching for duplicate data, comprising the steps of:
2 generating at least one identifier from at least one portion of a first
3 segment of data using a unique identifier function;
4 generating at least one identifier from at least one portion of a second
5 segment of data using the unique identifier function; and
6 comparing the at least one identifier associated with the first segment
7 of data with the at least one identifier associated with the second segment of data to
8 determine whether the first segment of data is substantially identical to the second
9 segment of data.

1 2. The method according to claim 1, wherein said step of generating at
2 least one identifier from at least one portion of the first segment of data comprises
3 the step of generating at least one identifier from at least one portion of the first
4 segment data using the unique identifier function as the first segment of data is
5 recorded onto a storage medium.

1 3. The method according to claim 2, wherein said step of generating at
2 least one identifier from at least one portion of the second segment of data
3 comprises the step of generating the at least one identifier from at least one portion

4 of the second segment of data using the unique identifier function as the second
5 segment of data is recorded onto the storage medium.

1 4. The method according to claim 2, wherein said step of generating at
2 least one identifier from at least one portion of the second segment of data
3 comprises the step of generating at least one identifier from at least one portion of
4 the second segment of data using the unique identifier function as the second
5 segment of data is recorded onto a different storage medium.

1 5. The method according to claim 1, wherein said step of generating at
2 least one identifier from at least one portion of the first segment of data occurs after
3 the first segment of data is recorded onto the storage medium.

1 6. The method according to claim 5, wherein said step of generating at
2 least one identifier from at least one portion of a second segment of data comprises
3 the step of generating at least one identifier from at least one corresponding portion
4 of the second segment of data using the unique identifier function as the second
5 segment of data is recorded onto the storage medium.

1 7. The method according to claim 5, wherein said step of generating at
2 least one identifier from at least one portion of the second segment of data
3 comprises the step of generating the at least one identifier from the at least one

4 portion of the second segment of data using the unique identifier function as the
5 second segment of data is recorded onto a different storage medium.

1 8. The method according to claim 1, wherein the first segment of data and
2 the second segment of data are segments of multimedia data.

1 9. The method according to claim 1, further comprising the steps of:
2 storing in a table the at least one identifier associated with the first
3 segment of data; and
4 retrieving from the table the at least one identifier associated with the
5 first segment of data prior to said comparing step.

1 10. The method according to claim 1, further comprising the step of
2 presenting an indication that the first segment of data is substantially identical to the
3 second segment of data when at least one identifier associated with the first
4 segment of data matches the at least one identifier associated with the second
5 segment of data.

1 11. The method according to claim 1, wherein a size of the at least one
2 portion of the first segment of data and the at least one portion of the second
3 segment of data is based on a temporal measurement, wherein the at least one

4 portion of the first segment of data corresponds temporally with the at least one
5 portion of the second segment.

1 12. The method according to claim 1, wherein a size of the at least one
2 portion of the first segment of data and the at least one portion of the second
3 segment of data is based on a bit measurement, wherein the at least one portion of
4 the first segment of data corresponds bit by bit with the at least one portion of the
5 second segment.

1 13. The method according to claim 1, wherein the at least one identifier
2 associated with the first segment of data and the at least one identifier associated
3 with the second segment of data are hash values and the unique identifier function is
4 a hash function wherein a hash value associated with the first segment of data will
5 equal a hash value associated with the second segment of data when the first
6 segment of data and the second segment of data are identical.

1 14. The method according to claim 1, wherein said comparing step
2 comprises the step of comparing a plurality of identifiers associated with the first
3 segment of data with a plurality of identifiers associated with the second segment of
4 data to determine whether the first segment of data is substantially identical to the

5 second segment of data.

1 15. The method according to claim 1, wherein said comparing step
2 comprises the step of comparing a plurality of identifiers associated with a first set of
3 segments of data with a plurality of identifiers associated with a second set of
4 segments of data to determine whether the first set of segments of data is
5 substantially identical to the second set of segments of data.

1 16. A system for searching for duplicate data, comprising:
2 a controller for reading data from and writing data to a storage medium;
3 and
4 a processor coupled to the controller, wherein the processor is
5 programmed to:
6 generate at least one identifier from at least one portion of a first
7 segment of data using a unique identifier function;
8 generate at least one identifier from at least one portion of a
9 second segment of data using the unique identifier function; and
10 compare the at least one identifier associated with the first
11 segment of data with the at least one identifier associated with the second segment
12 of data to determine whether the first segment of data is substantially identical to the

13 second segment of data.

1 17. The system according to claim 16, wherein the processor is further
2 programmed to generate at least one identifier from the at least one portion of the
3 first segment of data using a unique identifier function as the first segment of data is
4 recorded onto the storage medium.

1 18. The system according to claim 17, wherein the processor is further
2 programmed to generate the at least one identifier from the at least one portion of
3 the second segment of data using the unique identifier function as the second
4 segment of data is recorded onto the storage medium.

1 19. The system according to claim 17, wherein the processor is further
2 programmed to generate the at least one identifier from the at least one portion of
3 the second segment of data using the unique identifier function as the second
4 segment of data is recorded onto a different storage medium.

1 20. The system according to claim 16, wherein the processor is further
2 programmed to generate the at least one identifier from the at least one portion of
3 the first segment of data using the unique identifier function after the first segment of

4 data is recorded onto the storage medium.

1 21. The system according to claim 20, wherein the processor is further
2 programmed to generate the at least one identifier from the at least one portion of
3 the second segment of data using the unique identifier function as the second
4 segment of data is recorded onto the storage medium.

1 22. The system according to claim 20, wherein the processor is further
2 programmed to generate the at least one unique identifier from at least one
3 corresponding portion of the second segment of data using the unique identifier
4 function as the second segment of data is recorded onto a different storage medium.

1 23. The system according to claim 16, wherein the first segment of data
2 and the second segment of data are segments of multimedia data.

1 24. The system according to claim 16, further comprising a table, wherein
2 the processor is further programmed to:

3 store in the table at least one identifier associated with the first
4 segment of data; and

5 retrieve from the table at least one identifier associated with the first

6 segment of data prior to said comparing step.

1 25. The system according to claim 16, wherein the processor is further
2 programmed to present an indication that the first segment of data is substantially
3 identical to the second segment of data when at least one identifier associated with
4 the first segment of data matches the at least one identifier associated with the
5 second segment.

26. The system according to claim 16, wherein the at least one identifier associated with the first segment of data and the at least one identifier associated with the second segment of data are hash values and the unique identifier function is a hash function, wherein the processor determines if a hash value associated with the first segment of data equals a hash value associated with the second segment of data indicative of the first segment of data and the second segment of data being substantially identical.

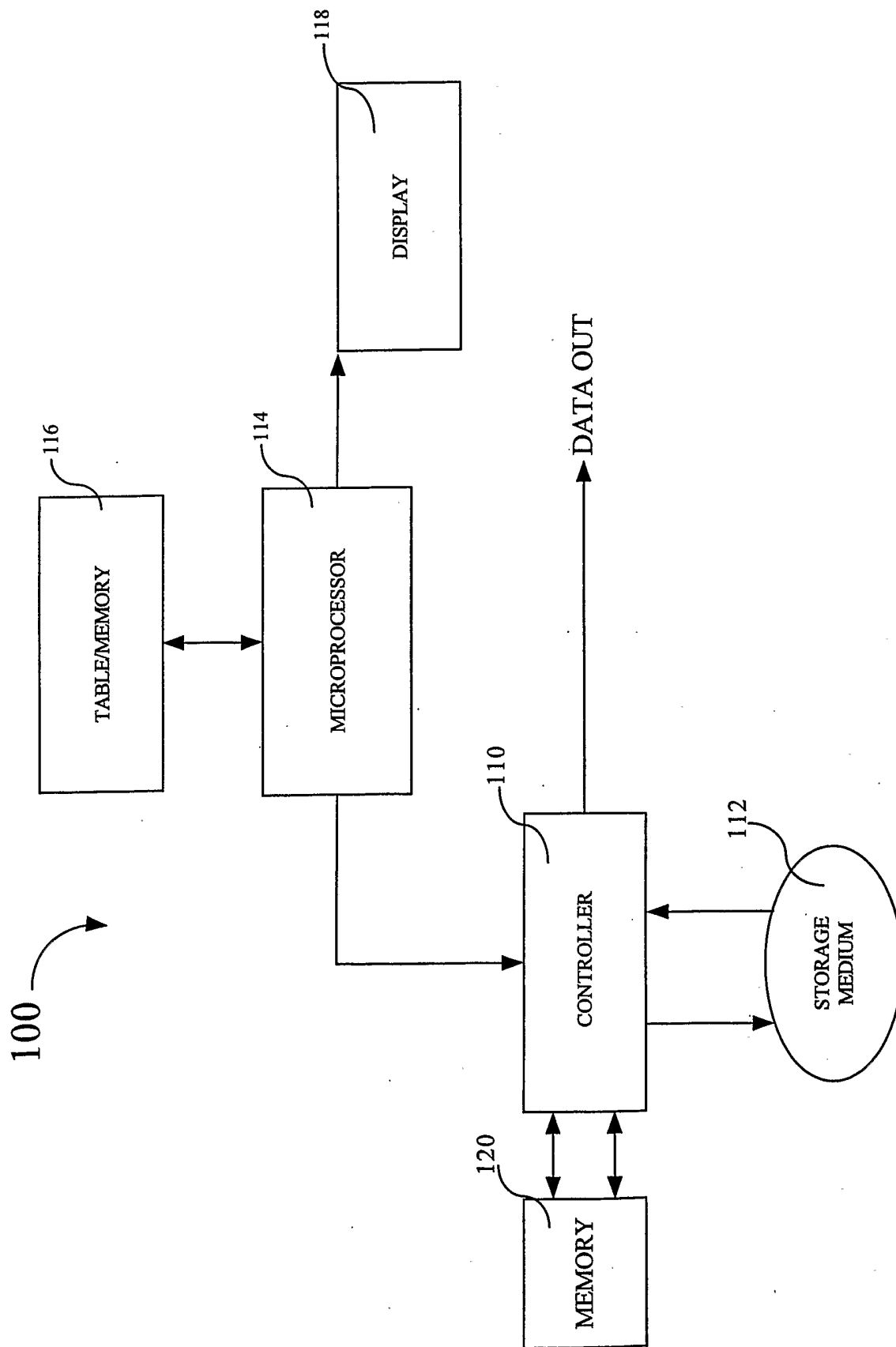


FIG. 1

2/2

200

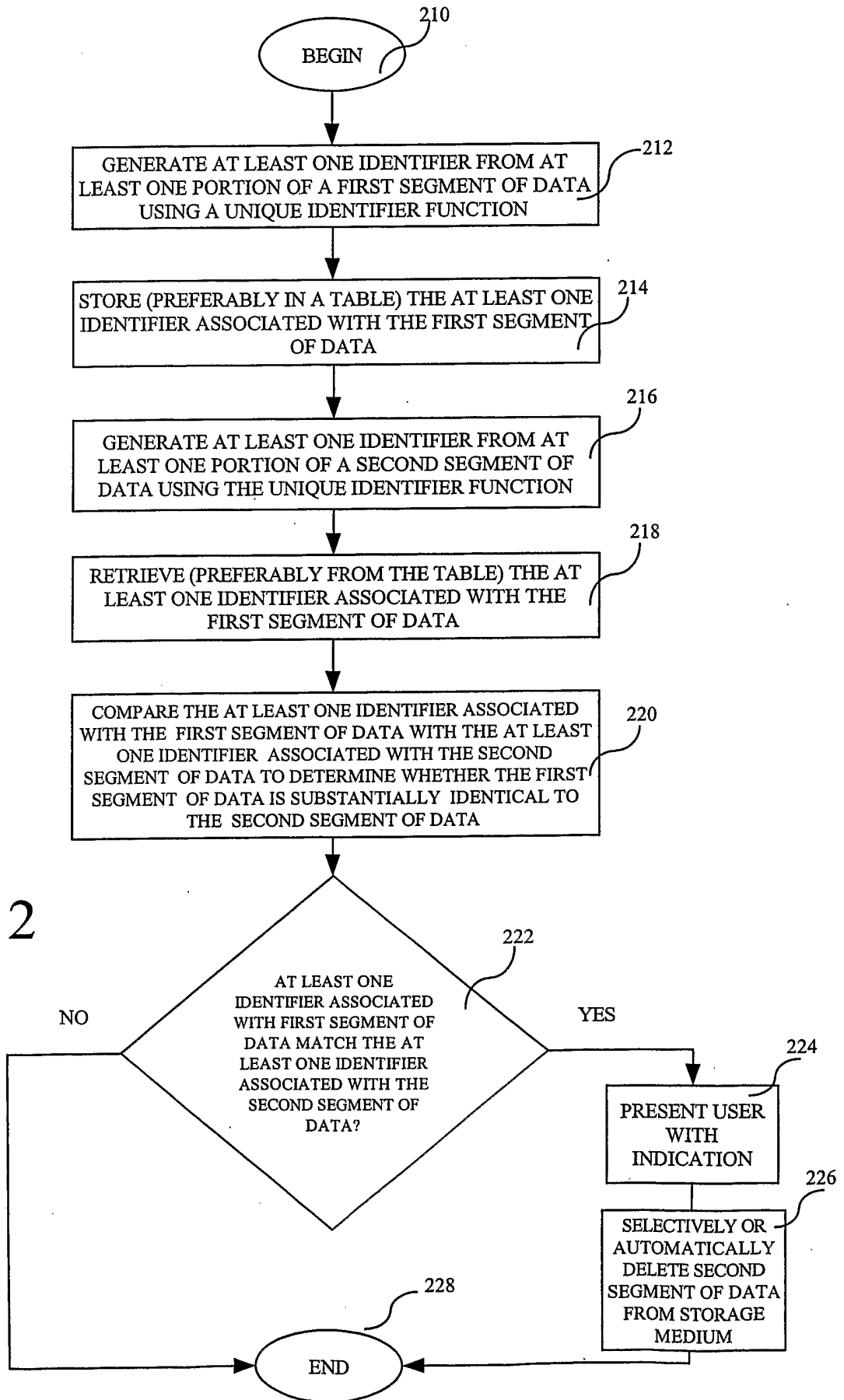


FIG. 2

INTERNATIONAL SEARCH REPORT

International application No.

PCT/US03/01194

A. CLASSIFICATION OF SUBJECT MATTER

IPC(7) : G06F 12/00

US CL : 711/154

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)

U.S. : 711/154

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 practicable, search terms used)

IEEE, ACM

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X --- Y	US A1 2001/0037323 (MOULTON et al.) 01 November 2001. Figure 3, pg 2 paragraph 11, pg 1 par 5,	1-4,8-9,12-19,23-26 ----- 7
X,P --- Y,P	US 6,389,433 B1 (BOLOSKY et al) 14 May 2002, Fig 4,9	1-6,9,14-22,24-25 ----- 7
A,P	US A1 2002/0194197 (FLANK) 19 December 2002. Paragraph 57, claims 27-29	11,12
A, P	US A1 2002/0194198 (FLANK et al.) 19 December 2002 Paragraph 55, claims 7-9	11,12

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	"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
"E" earlier application or patent published on or after the international filing date	"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
"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)	"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
"O" document referring to an oral disclosure, use, exhibition or other means	"&" document member of the same patent family
"P" document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search

11 April 2003 (11.04.2003)

Date of mailing of the international search report

21 MAY 2003

Name and mailing address of the ISA/US
Commissioner of Patents and Trademarks
Box PCT
Washington, D.C. 20231

Facsimile No. (703)305-3230

Authorized officer

Do Yoo

Peggy Hanrod

Telephone No. (703)305-3900