



- (51) International Patent Classification:
G06F 11/00 (2006.01)
- (21) International Application Number:
PCT/CN2015/095840
- (22) International Filing Date:
27 November 2015 (27.11.2015)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data:
EP15162563.9 7 April 2015 (07.04.2015) EP
- (71) Applicant: HUAWEI TECHNOLOGIES CO., LTD.
[CN/CN]; Huawei Administration Building, Bantian, Long-gang District, Shenzhen, Guangdong 518129 (CN).
- (72) Inventors: LEVY, Eliezer; c/o Huawei Technologies Duesseldorf GmbH, Riesstr. 25, Munich, 80992 (DE). CHEN, Zhibiao; c/o Huawei Technologies Duesseldorf GmbH, Riesstr. 25, Munich, 80992 (DE). DAR, Usama; c/o Huawei Technologies Duesseldorf GmbH, Riesstr. 25, Munich, 80992 (DE). AVITZUR, Aharon; c/o Huawei Technologies Duesseldorf GmbH, Riesstr. 25, Munich, 80992 (DE). GOIKHMAN, Shay; c/o Huawei Technologies Duesseldorf GmbH, Riesstr. 25, Munich, 80992 (DE). WOLSKI, Antoni; c/o Huawei Technologies Duesseldorf GmbH, Riesstr. 25, Munich, 80992 (DE).

- (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, BN, 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, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, 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, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, 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, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

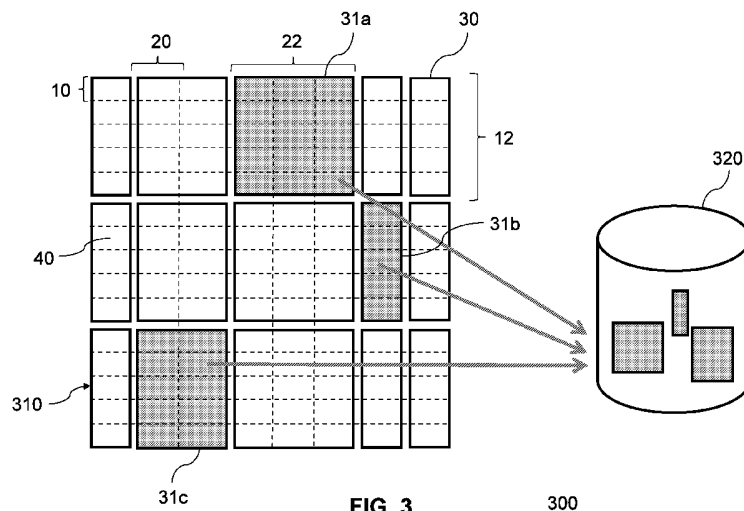
Declarations under Rule 4.17:

— as to applicant's entitlement to apply for and be granted a patent (Rule 4.17(ii))

Published:

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

(54) Title: METHOD, APPARATUS AND DATA STRUCTURE FOR COPYING VALUES OF TABLE OF DATABASE



(57) Abstract: A method for copying values of a table of a database between a primary memory and a secondary memory is provided, wherein the table is organized in a plurality of stripes and a plurality of vertical partitions, wherein a stripe comprises at least two rows of the table and a vertical partition comprises one or more columns of the table, wherein the table is stored as a plurality of segments, wherein a segment comprises values at a cross-section of a stripe and a vertical partition, and wherein a segment stores adjacent column values in adjacent locations of the primary or the secondary memory, the method comprising a step of selecting one or more segments and copying the one or more selected segments between the primary memory and the secondary memory.

METHOD, APPARATUS AND DATA STRUCTURE FOR COPYING VALUES OF TABLE OF DATABASE

TECHNICAL FIELD

[0001] The present invention relates to a method and an apparatus for copying values of a table of a database. The present invention also relates to a computer-readable storage medium and a data structure.

BACKGROUND

[0002] The traditional database storage layout of disk-based databases involves fixed-sized data pages composed of table rows. The data pages are fetched from disk and maintained in a shared page buffer pool, for efficient usage by concurrent applications. The page size and structure is optimized for efficient disk I/O. The solution is not suitable for in-memory databases because the overhead of the shared page buffer pool is too large. Especially, several steps of indirect addressing are applied to access a single data item (e.g. a column value). Additionally, in in-memory databases, there is no need to buffer the data on the way from, or to disk. The solution is also called a row store, or a row-oriented database because of the row orientation of the page structure. Row store is beneficial in on-line transaction processing (OLTP) databases. This is because the number of columns is small, and several column values are likely to be processed in a database operation like insert or select. On the other hand, in a column store, or a columnar database, the data is organized by placing the column values close to each other. Columnar databases are suited for analytical processing, because the number of columns can be very large (even in the range of hundreds or thousands) and most query operations are column-wise. By placing the column values adjacently, the efficiency of disk I/O and memory access is improved. Especially, in modern hardware platforms characterized by asymmetric memory (exemplified by the Non-Uniform Memory Access—NUMA—architecture), multi-level memory caches and vector processing (SIMD—single instruction multiple data) units, the advantage of processing the memory contents sequentially is significant. Thus, all current implementations of in-memory analytics databases employ large column objects allocated in memory. The disadvantage of that solution is that it does not allow for data migration between the memory and disk. The data migration is required because there is a need to free memory from the data that has not been used recently, and restore it when it is needed. Specifically, in columnar databases, there is also a need to move cold (not used recently) columns from memory to disk.

[0003] Most in-memory databases are not designed for data migration. Typically, an in-memory table has to reside in memory in its totality. That also applies to columnar databases where column data is stored in fixed-size column vectors. One prior art solution, called anti-caching, applies moving of the data to disk on a row-by-row basis. That is not suitable to columnar databases when a need to migrate data on a column basis emerges.

[0004] Regarding the data organization in memory, various solutions are proposed, tending to balance the needs of the row-oriented and column oriented processing. However, no solution has been offered to the problem of efficient copying of values, in particular data migration.

SUMMARY OF THE INVENTION

[0005] In view of the above, one object of the present invention is to provide a method, an apparatus, a computer-readable storage medium and a data structure that solve at least one of the above-mentioned problems of the prior art. The foregoing and other objects are achieved by the features of the independent claims. Further implementation forms are apparent from the dependent claims, the description and the Figures.

[0006] A first aspect of the invention provides a method for copying values of a table of a database between a primary memory and a secondary memory, wherein the table is organized in a plurality of stripes and a plurality of vertical partitions, wherein a stripe comprises at least two rows of the table and a vertical partition comprises one or more columns of the table, wherein the table is stored as a plurality of segments, wherein a segment comprises values at a cross-section of a stripe and a vertical partition, and wherein a segment stores adjacent column values in adjacent locations of the primary or the secondary memory, the method comprising a step of selecting one or more segments and copying the one or more selected segments between the primary memory and the secondary memory.

[0007] In the following copying also comprises evicting, i.e., copying the data from source memory to destination memory and deleting the data from the source memory.

[0008] The data of the table of database is organized and/or copied in segments. Segments can be units of database checkpointing and data migration. The segments can adhere to the principles of a columnar database in that they can contain sequences of column values of one or more columns of a table.

[0009] The organization of the values of the table in segments, which each store adjacent column values in adjacent memory locations, allows for efficient parallel column-wise processing of data in memory, while still being able to evict or restore data from memory efficiently. In embodiments of the invention, a consistent persistent image of the database can be maintained, and "cold" data can be kept on disk.

[0010] A segment contains values at a cross-section of a set of rows and one or more columns of a table. It contains one or more column value vectors and is thus suitable for column-wise computing. In a segment, column values are organized in adjacent memory locations, suitable for efficient scanning and vector processing (single instruction multiple data, SIMD). Column values are aligned, in segments, in such a way that a set of horizontal segments, called a stripe, represents a set of rows. Preferably, a segment can be stored on the secondary memory in binary form, without a need of reformatting. Upon restore of a segment, the column vectors can be available at pre-defined offsets of the segment memory

block. In other words, in embodiments of the invention, a segment is a data unit that is shared between the checkpoint activity and data migration activity.

[0011] In an embodiment of the invention, a segment comprises values of a cross-section of a connected set of rows and a connected set of one or more columns of a table of the database, but comprises no values from at least one column of the table. In other words, segments can comprise values of a rectangular region of the table of the database. In particular, it can be foreseen that a segment only consists of the values of a rectangular region of the table of the database. In other embodiments, the segment might consist of the values of the rectangular region of the table of the database and some additional metadata, but no values from cells in the table that are outside the mentioned rectangular region.

[0012] In the following, values of a table can refer to any kind of data that can be stored in a table, in particular any kind of data that can be stored in a cell of a table. Among others, a value can be an integer value, a floating point value, or any more complex variable or data structure.

[0013] According to a first possible implementation of the method according to the first aspect, the method is implemented in a database management system and/or the primary memory is a volatile memory and the secondary memory is a persistent memory, in particular a hard disk. Implementing the method in a database management system has the advantage that the database management system can efficiently copy values of a table of the database between primary and secondary memory. If the secondary memory is a persistent memory, the method provides an efficient way of backing up data to a persistent storage.

[0014] According to a second possible implementation of the method according to the first aspect, the method is a checkpointing method for copying one or more changed segments from the primary memory of the database to the secondary memory. Storing some segments in a volatile memory and other segments in a persistent memory can be advantageous e.g. if the volatile memory has a faster access time than the persistent memory.

[0015] In embodiments of the invention, changed segments, i.e., segments that have changed in the memory, are checkpointed (copied) to disk at checkpoint intervals. The checkpoint intervals can be regular intervals determined by a clock unit, or they can be irregular intervals that are determined based on a rule, e.g. based on when a sufficient number of changes that need checkpointing has accumulated.

[0016] According to a third possible implementation of the method according to the first aspect, the method of the first aspect comprises the steps:

- selecting one or more segments in the primary memory that have changed,
- freezing the one or more selected segments, such that a state of the one or more selected segments is preserved,
- copying them one or more frozen segments to the secondary memory, and
- releasing the frozen segments.

[0017] Freezing the one or more selected segments can be implemented e.g. by preventing write access to the segment. This allows that a consistent state of the one or more selected segments is copied to the secondary memory. Preferably, after copying, the frozen segments are released in the primary memory.

[0018] According to a fourth possible implementation of the method according to the first aspect, freezing the one or more selected segments comprises shadowing, copying on write, or locking of the one or more selected segments. With a copy-on-write freezing method, the frozen segments are "forgotten" (i.e. removed from memory) and replaced with changed copies. Freezing the selected segments can also be performed with further freezing methods that are known to the person skilled in the art.

[0019] Releasing the frozen segments finalizes the checkpointing operation and makes the selected segments available for new changes.

[0020] According to a fifth possible implementation, the method of the first aspect is a method for evicting data from the primary memory to the secondary memory.

[0021] According to a sixth possible implementation of the method according to the first aspect, the method comprises the steps:

- detecting a need for data eviction,
- determining an extent of the needed data eviction,
- selecting one or more segments to be evicted,
- for each selected segment, determining whether the segment has been checkpointed,
- for each selected segment that has not been checkpointed, copying the segment from the primary memory to the secondary memory, and
- deleting the selected segments from the primary memory.

[0022] This implementation provides an efficient way of segment eviction. A special case of a segment eviction is an eviction of a full column, wherein the full column can be evicted by evicting a plurality of segments that make up the full column. For example, if the table of the database is organized into four stripes, a full column can be evicted by evicting the corresponding four segments.

[0023] In particular, the method according to the fifth possible implementation can further comprise a step of, for each identified segment that has been checkpointed, marking the segment as evicted.

[0024] When an eviction request is issued, for example the one or more least recently used segments can be moved to disk. Because segments that are resident on the secondary memory are shared between the checkpoint and data migration, the method can include an optimization such that, upon eviction, the one or more segments need not be written to disk if they have been already checkpointed before. In that case, the one or more segments are

only removed from primary memory and marked as evicted. Most of the segments being evicted satisfy the condition because, in most cases, the evicted segments have not been recently used (e.g. changed).

[0025] According to a seventh possible implementation, in the method of the sixth implementation the step of selecting the segments to be evicted comprises:

- selecting least recently used segments,
- selecting one or more full columns that are not likely to be used, and/or
- selecting a segment based on a selection criterion that is based on an age of the data.

[0026] Thanks to the grid-like layout of the segments, various collections of segments can be selected for eviction. Possible choices of the selecting method also include selecting spans of rows, with segment granularity, based on some criterion like the age of data. On segment restore, demand paging can be used whereby the segments are actually loaded to memory when they are needed.

[0027] According to an eighth possible implementation, the one or more selected segments are encrypted and/or compressed before they are copied to the secondary memory. Encrypting the one or more selected segments has the advantage that unauthorized access to the data of the database is prevented. For example, if an intruder achieves unauthorized access to the secondary memory, he could not retrieve the information content of the database table if the segments are encrypted before copying to the secondary memory.

[0028] According to a ninth possible implementation, the copying is performed as part of restoring data from the secondary memory back to the primary memory.

[0029] According to a tenth possible implementation, the method according to the first aspect further comprises the steps:

- detecting a need for restoring of data,
- determining which data need restoring,
- selecting one or more segments to be restored,
- determining whether the primary memory comprises sufficient free space for restoring the selected segments,
- if there is not sufficient free space, freeing space in the primary memory by evicting data, in particular one or more segments, from the primary memory,
- restoring the selected segments from the secondary memory back to the primary memory.

[0030] According to an eleventh possible implementation, the one or more selected segments are decrypted and/or decompressed before they are restored from the secondary memory back to the primary memory. Preferably, the decryption or decompression is performed

“on the fly”, such that e.g. a database management system is not even aware that segments of its database were temporarily encrypted and/or compressed.

- [0031] If the segments are stored on the secondary memory in an encrypted or compressed form, adjacent column values are stored in adjacent memory locations only in the sense that their decrypted or decompressed version corresponds to storing adjacent column values in adjacent (uncompressed / unencrypted) memory locations.
- [0032] A second aspect of the present invention provides apparatus for copying values of a table of a database between a primary memory and a secondary memory, wherein the table is organized in stripes and vertical partitions and a stripe comprises at least two rows of the table and a vertical partition comprises one or more columns of the table, wherein the table is stored as a plurality of segments, wherein a segment comprises values at a cross-section of a stripe and a vertical partition, and wherein a segment stores adjacent column values in adjacent locations of the primary memory, wherein the apparatus is configured to select one or more segments and copy the one or more selected segments between the primary memory and the secondary memory.
- [0033] In particular the second aspect of the invention provides an apparatus that is configured to carry out the method of the first aspect of the invention and/or one or more of the implementations of the first aspect of the invention.
- [0034] A third aspect of the invention provides a computer-readable storage medium, comprising program code, the program code comprising instructions for carrying out the method of the first aspect of the invention and/or one or more of the implementations of the first aspect of the invention.
- [0035] A fourth aspect of the invention provides a data structure of a database, comprising a plurality of segments of a table of the database, wherein the table is organized in stripes and vertical partitions and a stripe comprises at least two rows of the table and a vertical partition comprises one or more columns of the table, wherein the table is stored in the data structure as a plurality of segments, wherein a segment comprises values at a cross-section of a stripe and a vertical partition, and wherein a segment stores adjacent column values in adjacent locations of a primary memory.

BRIEF DESCRIPTION OF THE DRAWINGS

- [0036] To illustrate the technical features of embodiments of the present invention more clearly, the accompanying drawings provided for describing the embodiments are introduced briefly in the following. The accompanying drawings in the following description are merely some embodiments of the present invention, but modifications on these embodiments are possible without departing from the scope of the present invention as defined in the claims.

[0037] FIG. 1 shows a high level diagram of a system for copying values of a table of a database that is not in accordance with the present invention,

[0038] FIG. 2 shows a schematic illustration of a table that is organized according to embodiments of the invention,

[0039] FIG. 3 shows a schematic illustration of a table of a database and a method for checkpointing values of the table in an embodiment of the invention that is related to checkpointing,

[0040] FIG. 4 shows a flow chart of a method for checkpointing a database according to an embodiment of the invention,

[0041] FIG. 5 shows a schematic illustration of a table of a database and a method for evicting segments of that table to a secondary memory in accordance with the present invention,

[0042] FIG. 6 shows a flow chart of a method for evicting segments of a table according to an embodiment of the invention,

[0043] FIG. 7 shows a schematic illustration of a table of a database and a method for restoring segments from the secondary memory in accordance with the present invention, and

[0044] FIG. 8 shows a flow chart of a method for restoring segments of a table according to an embodiment of the invention.

DETAILED DESCRIPTION OF THE EMBODIMENTS

[0045] FIG. 1 is a high level diagram of a system for copying data that is not part of the present invention. The system 100 shown in FIG. 1 comprises a primary memory 110 and a secondary memory 120. The primary memory 110 is organized in a plurality of rows 10a that are indexed by an index 112. On the secondary memory 120, rows 10b are stored that were evicted from the primary memory 110. On a row miss, i.e., when access is required to a row that is missing in primary memory 110, the corresponding evicted row 10b is restored to primary memory 110. No column-wise eviction and restore is possible.

[0046] FIG. 2 shows a schematic representation of a table 200 that is organized according to embodiments of the invention. This schematic representation also indicates the layout of a data structure according to an aspect of the present invention. The table 200 comprises rows 10 and columns 20. The table 200 is organized in stripes 12 that comprise a plurality of rows 10. In the example shown in FIG. 2 the first, second and third stripe 12 each comprises five rows. The table 200 is furthermore organized in vertical partitions 22 that comprise one or more columns 20. In the example shown in FIG. 2, the first vertical partition comprises one column, the second vertical partition comprises two columns, the third vertical partition 22 comprises three columns and the fourth and fifth vertical partition each comprise one column. Each cell of the table stores one column value 40.

[0047] In embodiments of the invention, the organization of the table is reflected in the physical storage of the table such that the segments of the table are stored subsequently on a physical storage. Furthermore, the segments are units of database checkpointing and data migration. The segments adhere to the principles of a columnar database in that they contain sequences of column values of one or more columns of the table.

[0048] FIG. 3 illustrates an apparatus and a method for checkpointing data of an in-memory database in accordance with the present invention. A database checkpoint is a persistent storage representing a consistent state of a database. Database checkpoints are used e.g. upon recovery from computer system failures. The operation of creating a checkpoint is called checkpointing. A preferable way to do checkpointing is to store only the data that has changed since the previous checkpointing. This is called incremental checkpointing. An execution of an incremental checkpointing using the invention method is illustrated in FIG. 3.

[0049] The system 300 shown in FIG. 3 comprises a table 310 that comprises rows 10 and columns 20 and that is organized in stripes 12 and vertical partitions 22. The plurality of stripes 12 and vertical partitions 22 define segments 30, wherein each segment 30 comprises the column values 40 that are contained in a cross section of a stripe 12 and a vertical partition 22.

[0050] Whenever a change to a segment is detected, the one or more changed segments 31a, 31b, 31c are copied to the secondary memory 320. In that way, the secondary memory 320 functions as a checkpoint storage for changed segments. On a subsequent segment restore, demand paging can be used whereby segments are actually loaded to primary memory 310 when they are needed. For example, the primary memory 310 can be the working memory of a computing device which is hosting the database and a database management system managing the database.

[0051] FIG. 4 shows a flow chart of a checkpointing method in accordance with the present invention. In a first step S410, a new checkpoint is begun (e.g. by starting the time interval corresponding to a new checkpoint). Then, in step S420, one or more changed segments are identified and frozen, i.e. a further modification of the changed segments is prevented. There are various methods to implement freezing, like shadowing, copy on write, or locking. With a copy-on-write freezing method, the frozen segments are "forgotten" (i.e. removed from memory) and replaced with changed copies.

[0052] In step S430, the changed segments are copied to a checkpoint storage, which can be a persistent storage. In step S440, the frozen segments are released, i.e. the frozen segments are made available for future modifications. In step S450, the method ends. In other embodiments, instead of the method ending in step S450, the method can be executed iteratively, i.e. a new checkpoint is begun.

[0053] A further embodiment of the present invention is related to an implementation of the operation of data migration from primary memory to persistent storage by way of segment eviction. A special case of a segment eviction is an eviction of a full column, as illustrated in FIG. 5.

[0054] The system 500 shown in FIG. 5 comprises a table 510 of an in-memory database that is primarily stored in a primary memory. Furthermore, the system 500 comprises a secondary memory 520 that acts as a segment storage for one or more segments 32a, 32b, 32c that are evicted from the primary memory.

[0055] The steps of carrying out the segment eviction are illustrated in the flow chart in FIG. 6. When there is a need for data eviction (for example, when there is no room in the primary memory for new data), the eviction method is invoked. In step S610, the method is invoked. In step S620, a need for eviction is detected, and an extent (in terms of the amount of free space needed) is evaluated. In this step, a selection method can be chosen, for example by selecting full columns for eviction. In step S630, the segments to be evicted are selected. As a special case, the whole column can be indicated for eviction. Once the segments are identified, in step S640, for each evicted segment, a check is done whether that segment has been checkpointed in its current state. If so, no storing to persistent storage is needed. The segment is removed from the primary memory. If the segment has not been checkpointed, in step S650 it is moved to the secondary memory, e.g. a persistent storage such as a hard disk. In step S660, the method ends. In other embodiments of the invention, the previous steps can be carried out iteratively, i.e., instead of ending the method in step S660, execution of the method is continued at step S620.

[0056] In embodiments of the invention, as indicated in FIG. 6, steps S640 and S650 are performed for each segment of the table.

[0057] A third embodiment of the invention is an implementation of an operation of segment restore, after the segments have been evicted before. The need for restore can result from a query that accesses the evicted data, for example an evicted column. FIG. 7 illustrates a system 700 where a plurality of segments 33a, 33b, 33c are restored from a segment storage 720 to a table 710 that is stored in a primary memory.

[0058] The corresponding method steps are illustrated in the flow chart shown in FIG. 8. Execution of the method begins in step S810. The need for restore (e.g. caused by a query requesting the evicted data) is evaluated in step S820. The segments to be restored are identified in step S830. In step S840, a check is performed whether there is enough memory space for restoring the segments. If not, in step S850, the necessary memory space is freed by evicting segments from the primary memory. In step S860, the identified segments are restored from the secondary memory to the primary memory. In step S870, the method terminates.

[0059] The foregoing descriptions are only implementation manners of the present invention, but the protection of the scope of the present invention is not limited to this. Any variations or replacements can be easily made through the person skilled in the art. Therefore, the protection scope of the present invention should be subject to the protection scope of the attached claims.

CLAIMS

1. A method for copying values of a table of a database between a primary memory and a secondary memory, wherein the table is organized in a plurality of stripes and a plurality of vertical partitions, wherein a stripe comprises at least two rows of the table and a vertical partition comprises one or more columns of the table, wherein the table is stored as a plurality of segments, wherein a segment comprises values at a cross-section of a stripe and a vertical partition, and wherein a segment stores adjacent column values in adjacent locations of the primary or the secondary memory, the method comprising a step of selecting one or more segments and copying the one or more selected segments between the primary memory and the secondary memory.
2. The method of claim 1, wherein the method is implemented in a database management system and/or wherein the primary memory is a volatile memory and the secondary memory is a persistent memory, in particular a hard disk.
3. The method of one of the previous claims, wherein the method is a checkpointing method for copying one or more changed segments from the primary memory of the database to the secondary memory.
4. The method of claim 3, comprising the steps:
 - selecting one or more segments in the primary memory that have changed,
 - freezing the one or more selected segments such that a state of the one or more selected segments is preserved,
 - copying the one or more frozen segments to the secondary memory, and
 - releasing the frozen segments.
5. The method of claim 4, wherein freezing the one or more selected segments comprises shadowing, copying on write, or locking of the one or more selected segments.
6. The method of one of the previous claims, wherein the method is a method for evicting data from the primary memory to the secondary memory.
7. The method of claim 6, comprising the steps:

- detecting a need for data eviction,
 - determining an extent of the needed data eviction,
 - selecting one or more segments to be evicted,
 - for each selected segment, determining whether the segment has been checkpointed,
 - for each selected segment that has not been checkpointed, copying the segment from the primary memory to the secondary memory, and
 - deleting the selected segments from the primary memory.
8. The method of claim 7, wherein selecting the one or more segments to be evicted comprises:
- selecting one or more least recently used segments,
 - selecting one or more full columns that are not likely to be used, and/or
 - selecting one or more segments based on a selection criterion that is based on an age of the data.
9. The method of one of the previous claims, wherein the one or more selected segments are encrypted and/or compressed before they are copied to the secondary memory.
10. The method of one of the previous claims, wherein the copying is performed as part of restoring data from the secondary memory back to the primary memory.
11. The method of one of the previous claims, comprising the steps of:
- detecting a need for restoring of data,
 - determining which data need restoring,
 - selecting one or more segments to be restored,
 - determining whether the primary memory comprises sufficient free space for restoring the selected segments,
 - if there is not sufficient free space, freeing space in the primary memory by evicting data, in particular one or more segments, from the primary memory,
 - restoring the selected segments from the secondary memory back to the primary memory.

12. The method of claim 11, wherein the one or more selected segments are decrypted and/or decompressed before they are restored from the secondary memory back to the primary memory.
13. An apparatus for copying values of a table of a database between a primary memory and a secondary memory, wherein the table is organized in a plurality of stripes and a plurality of vertical partitions and a stripe comprises at least two rows of the table and a vertical partition comprises one or more columns of the table, wherein the table is stored as a plurality of segments, wherein a segment comprises values at a cross-section of a stripe and a vertical partition, and wherein a segment stores adjacent column values in adjacent locations of the primary memory or the secondary memory, wherein the apparatus is configured to select one or more segments and copy the one or more segments between the primary memory and the secondary memory, wherein in particular the apparatus is configured to carry out the method of one of the previous claims.
14. A computer-readable storage medium, comprising program code, the program code comprising instructions for carrying out the method of one of claims 1 to 12.
15. Data structure of a database, comprising a plurality of segments of a table of the database, wherein the table is organized in a plurality of stripes and a plurality of vertical partitions and a stripe comprises at least two rows of the table and a vertical partition comprises one or more columns of the table, wherein the table is stored in the data structure as a plurality of segments, wherein a segment comprises values at a cross-section of a stripe and a vertical partition, and wherein a segment stores adjacent column values in adjacent memory locations.

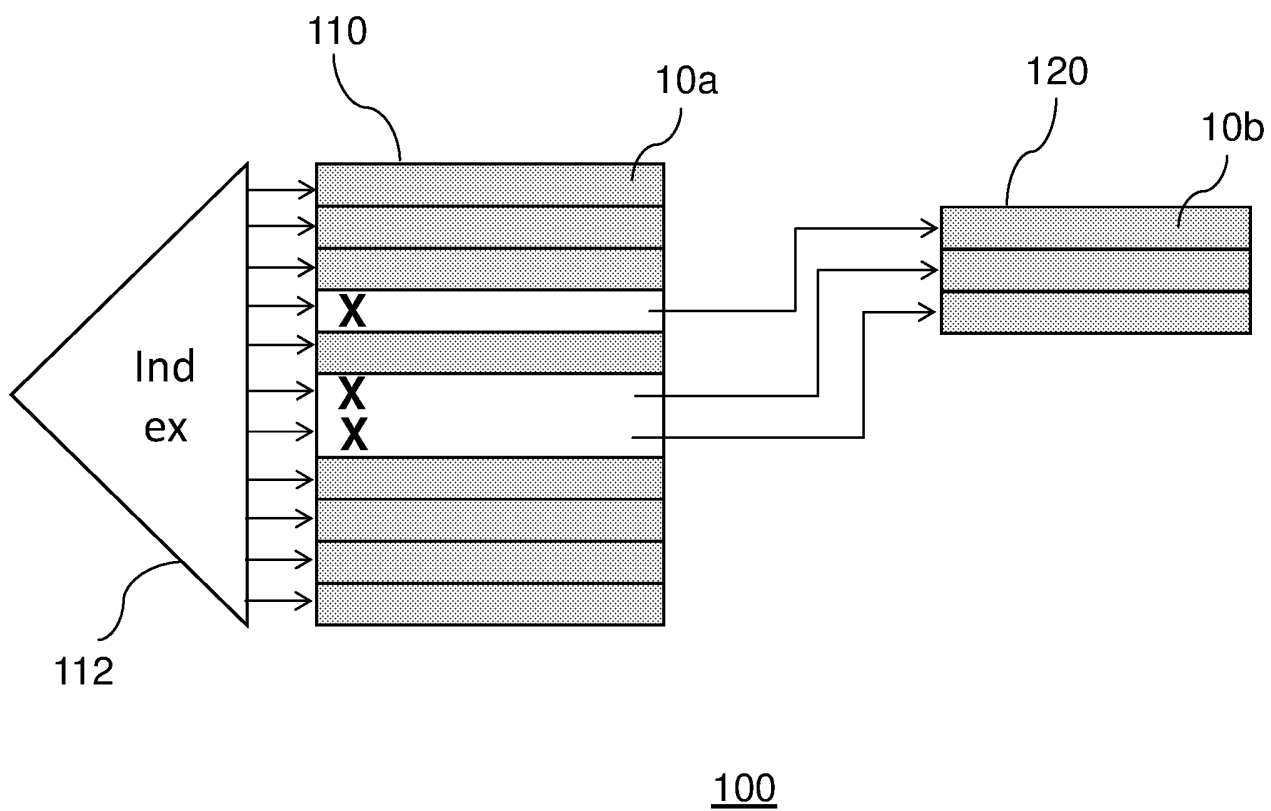
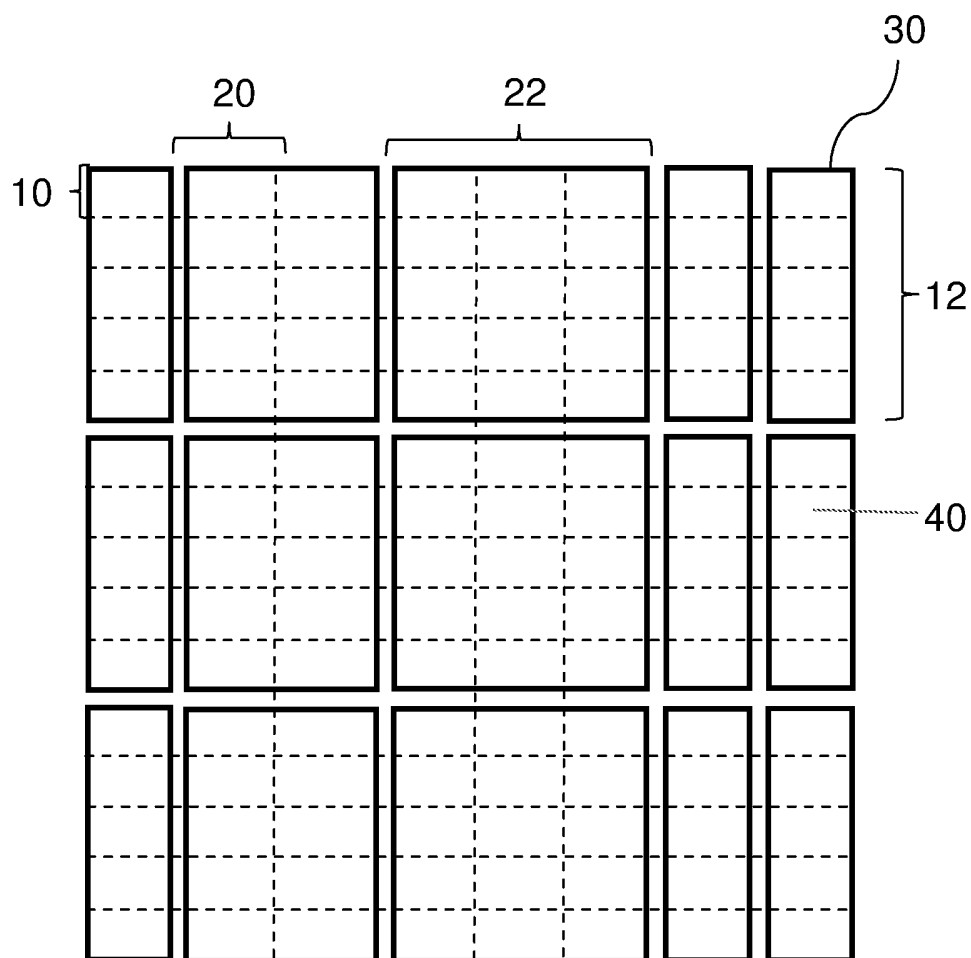
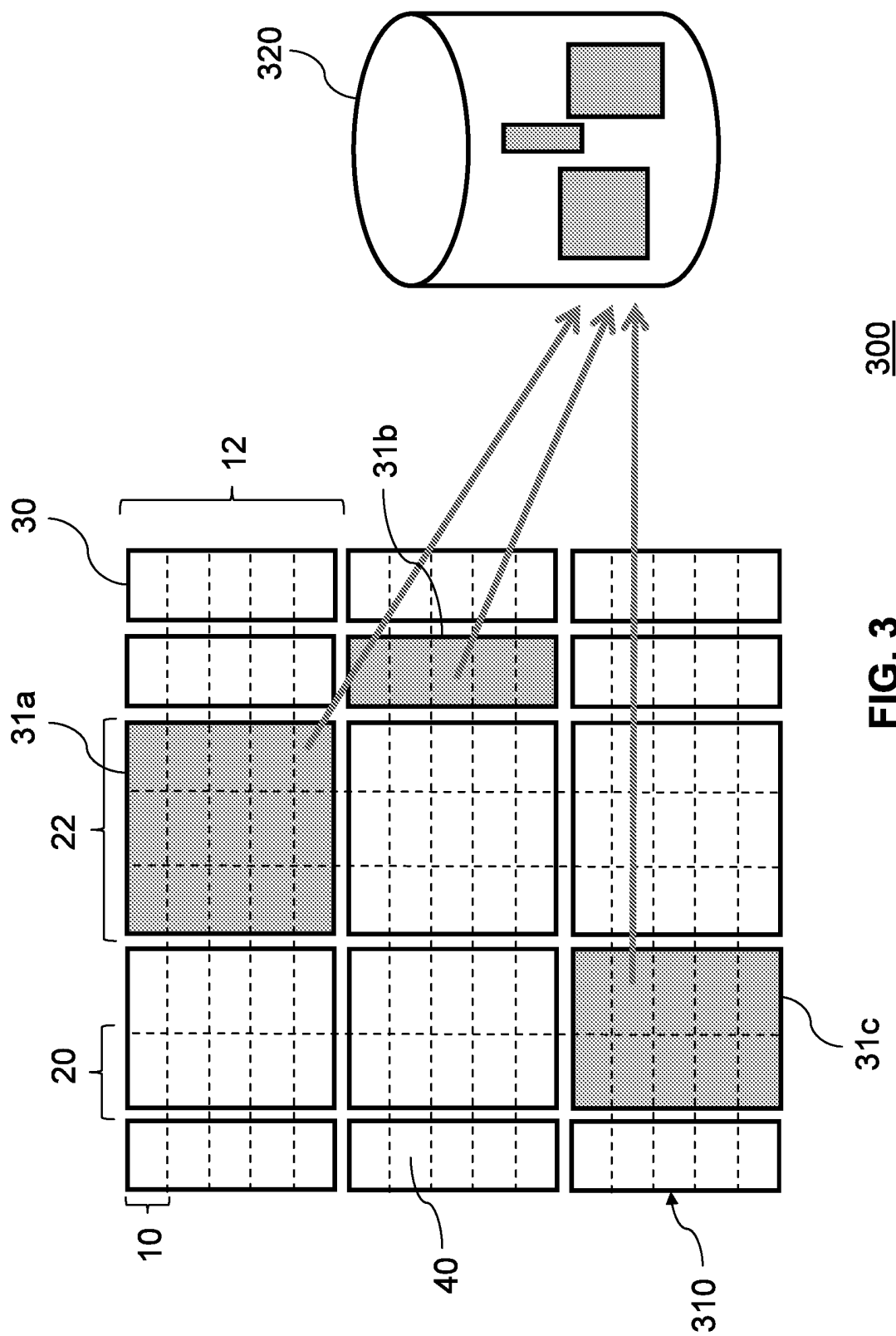


FIG. 1



200

FIG. 2



300

FIG. 3

31c

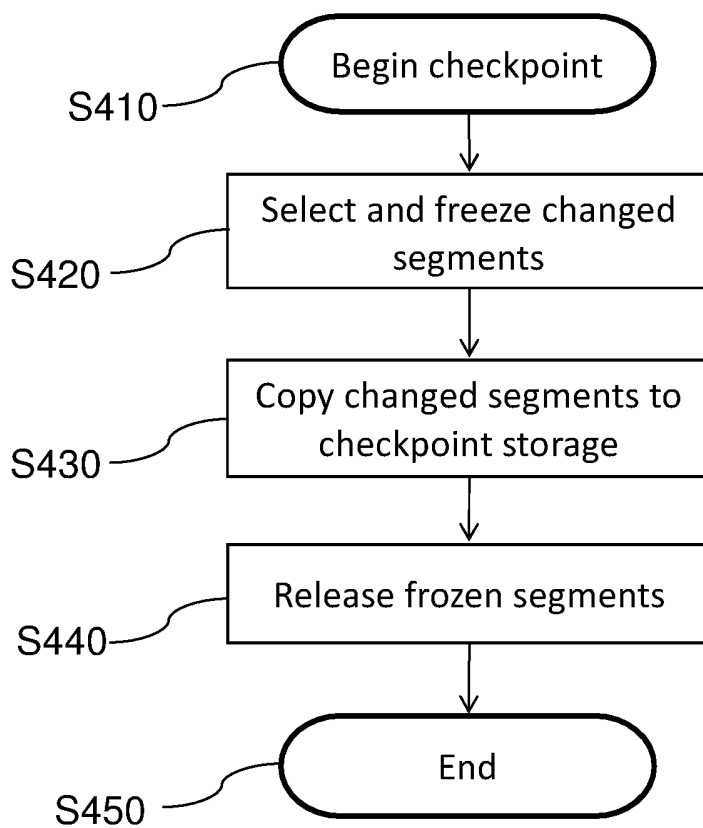
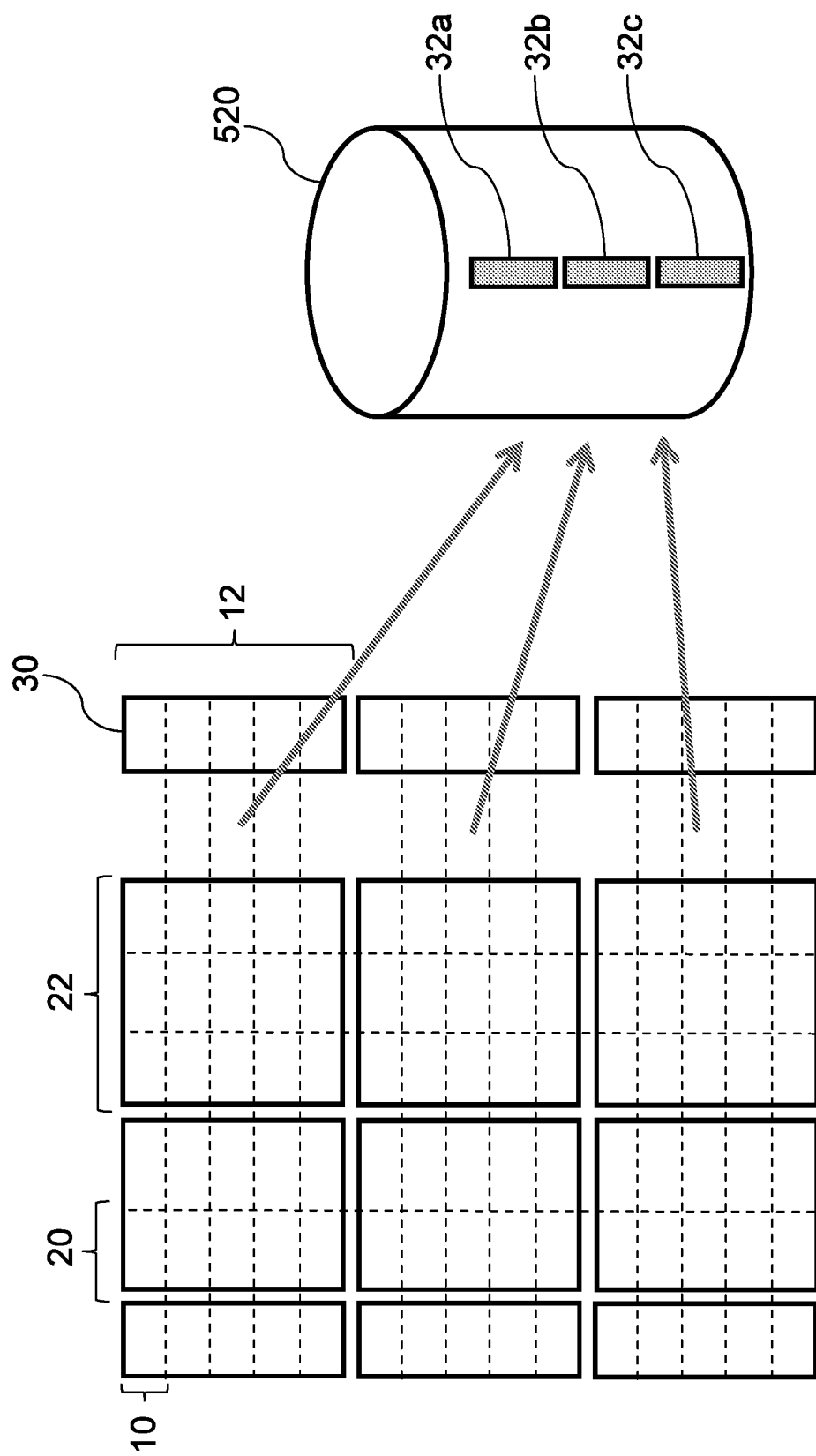


FIG. 4



500

FIG. 5

510

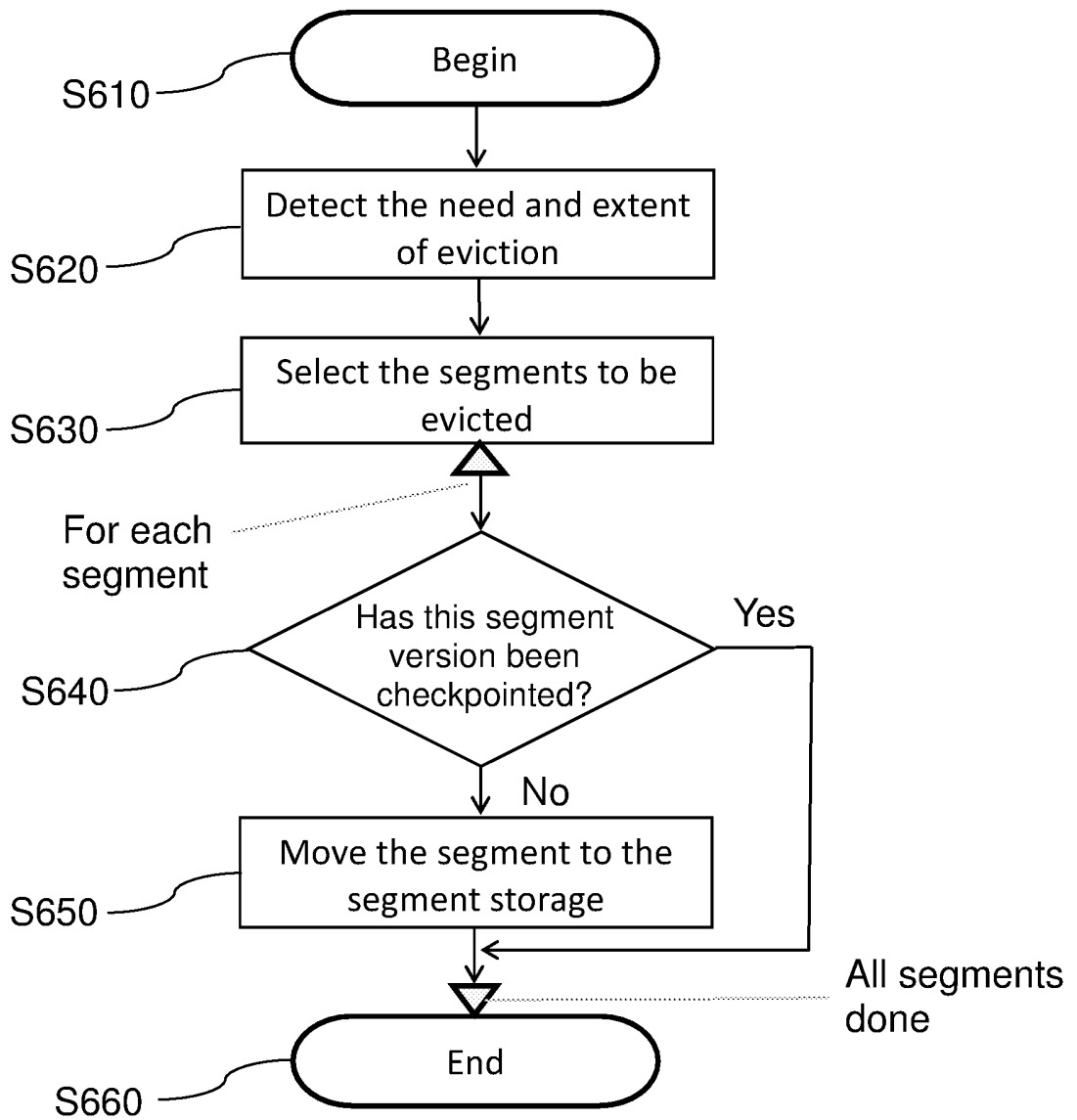
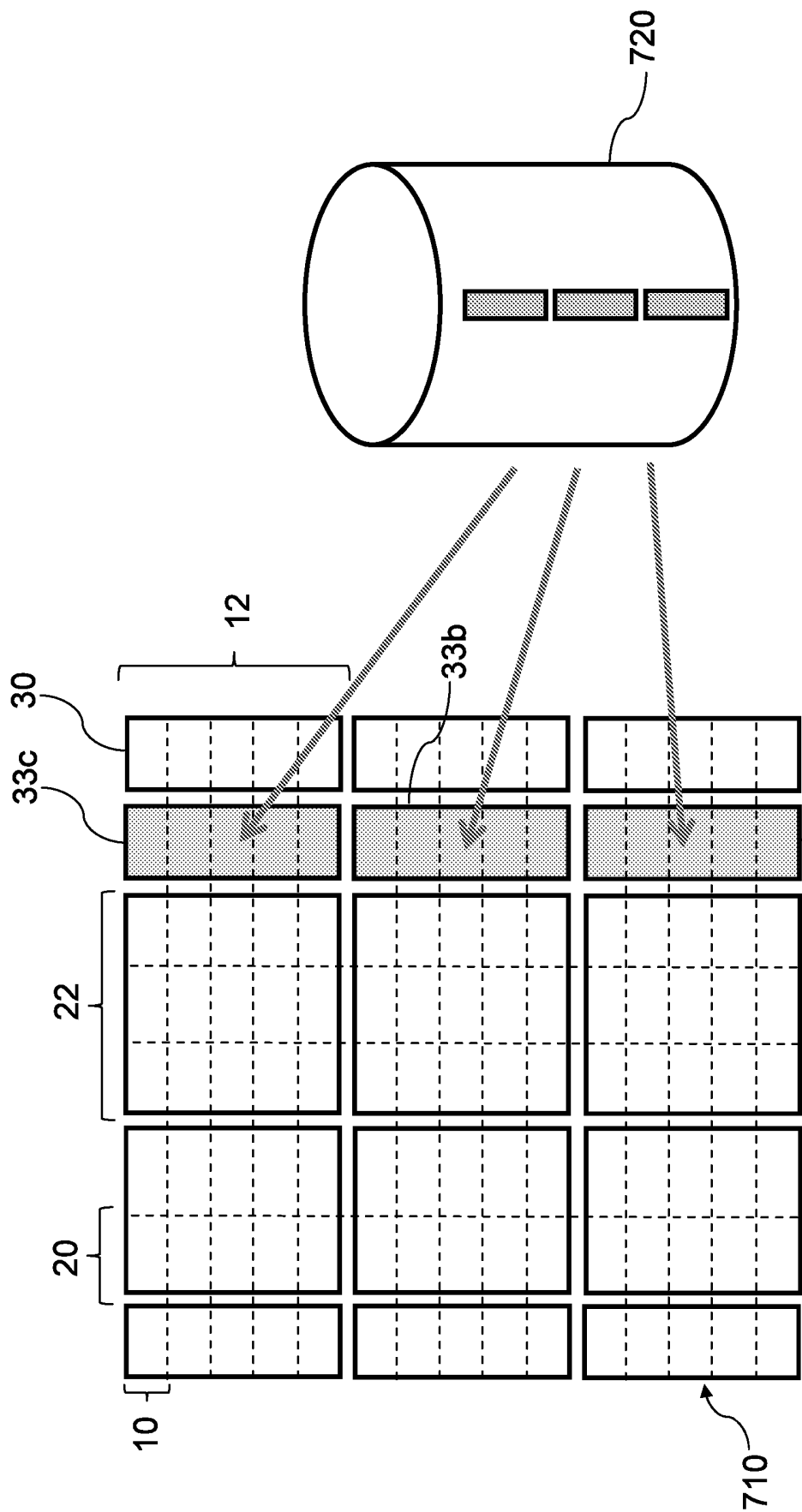


FIG. 6



33a
33b
33c
22
20
10
12
710
720
FIG. 7
700

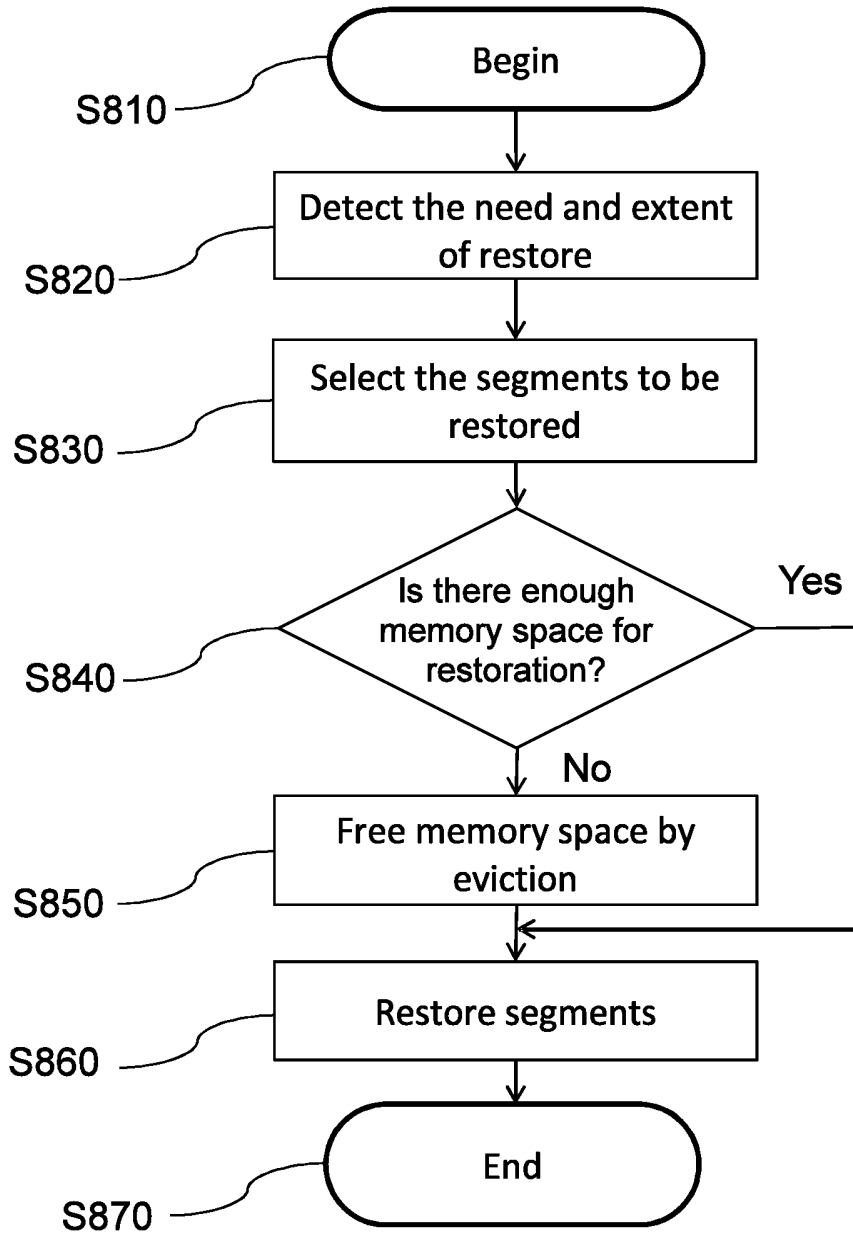


FIG. 8

INTERNATIONAL SEARCH REPORT

International application No.

PCT/CN2015/095840

A. CLASSIFICATION OF SUBJECT MATTER		
G06F 11/00(2006.01)i		
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 practicable, search terms used)		
CNPAT,WPI,EPODOC,CNKI:memory,storage,table,database,primary,second,stripe?,vertical,column?,row?,cross,value?,segment,copy+		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 8448021 B1 (DSSD, INC.) 21 May 2013 (2013-05-21) description, column 1, line 24 to column 2, line 23, and figures 1-9D	1-15
A	CN 1348135 A (WUXI YONGZHONG SCIENCE & TECHNOLOGY CO., LTD.) 08 May 2002 (2002-05-08) the whole document	1-15
A	CN 1299096 A (CHANG, CU YU) 13 June 2001 (2001-06-13) the whole document	1-15
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> 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”
“E”	earlier application or patent but published on or after the international filing date	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
“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)	“X”
“O”	document referring to an oral disclosure, use, exhibition or other means	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
“P”	document published prior to the international filing date but later than the priority date claimed	“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	
04 February 2016	01 March 2016	
Name and mailing address of the ISA/CN	Authorized officer	
STATE INTELLECTUAL PROPERTY OFFICE OF THE P.R.CHINA 6, Xitucheng Rd., Jimen Bridge, Haidian District, Beijing 100088, China	WANG,Xiaofei	
Facsimile No. (86-10)62019451	Telephone No. (86-10)62413918	

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.
PCT/CN2015/095840

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)			Publication date (day/month/year)
US	8448021	B1	21 May 2013	US	2015095697	A1	02 April 2015
				US	8316260	B1	20 November 2012
				EP	2828754	A1	28 January 2015
				WO	2013142646	A1	26 September 2013
				JP	2015516630	A	11 June 2015
				CN	104272261	A	07 January 2015
				US	8327185	B1	04 December 2012
.....							
CN	1348135	A	08 May 2002	None			
.....							
CN	1299096	A	13 June 2001	None			
.....							