

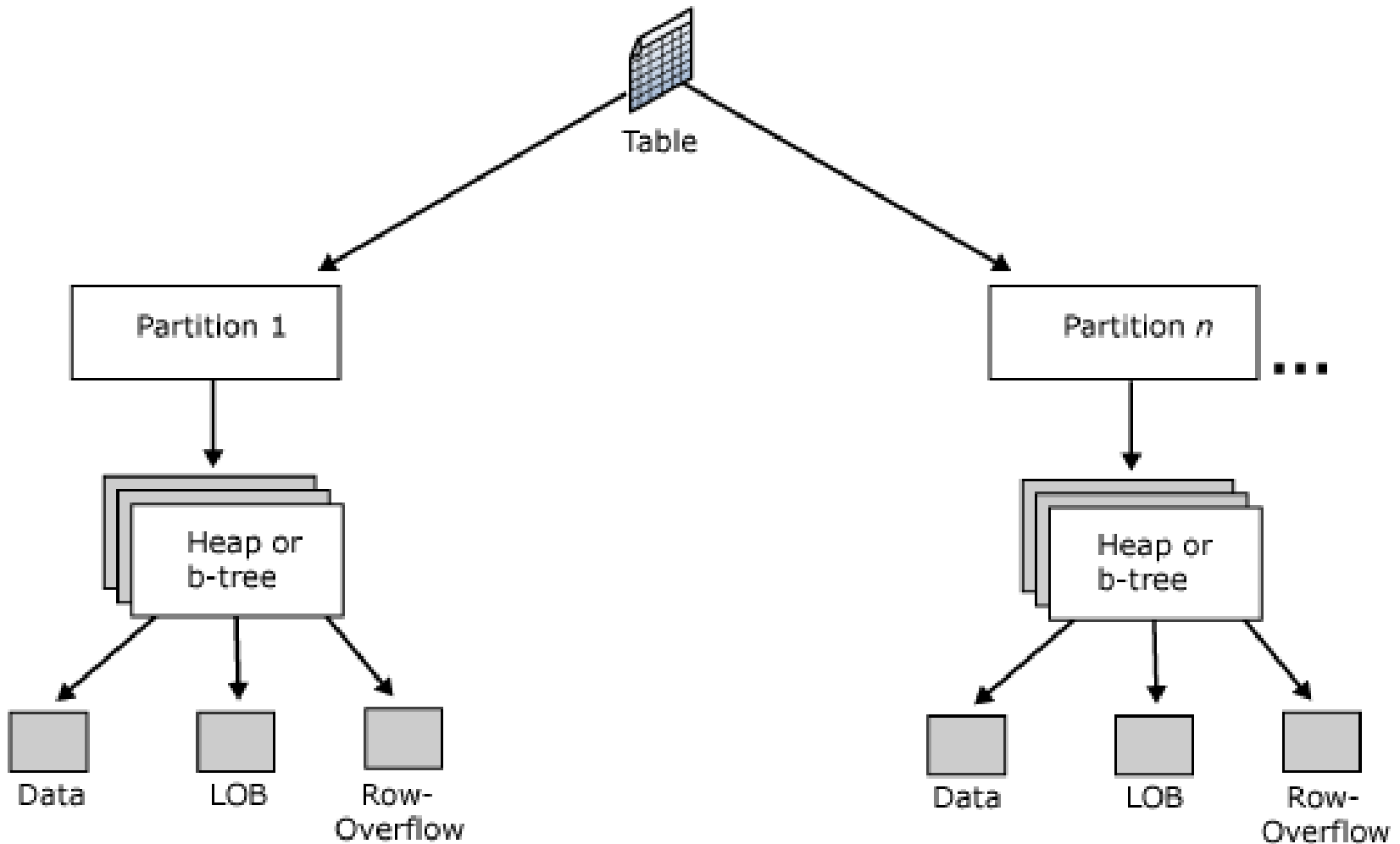
Physical Organization: SQL Server

Leggere Cap 7 Riguzzi et al. Sistemi Informativi

Tables

- Tables and indexes are stored as a collection of 8 KB pages
- A table is divided in one or more **partitions**
- Each partition contains data rows in either a **heap** or a **clustered table**.
- The pages of the heap or clustered index are managed in one or more **allocation units**, depending on the column types in the data rows.

Tables



Partitions

- Table and index pages are divided in one or more partitions.
- By default, a table or index has only one partition that contains all the table or index pages. The partition resides in a single filegroup.
- When a table or index uses multiple partitions, the data is partitioned horizontally so that groups of rows are mapped into individual partitions, based on a specified column.

Partitions

- The partitions can be put on one or more filegroups in the database. The table or index is treated as a single logical entity when queries or updates are performed on the data.
- To view the partitions used by a table or index, use the **sys.partitions** catalog view.

Organization of a Partition

- SQL Server tables use one of two methods to organize their data pages within a partition:
- **Clustered tables:** tables that have a clustered index. The data rows are stored in order based on the clustered index key. The clustered index is implemented as a B+tree index.
 - The pages in each level of the index, including the data pages in the leaf level, are linked in a doubly-linked list.
- **Heaps:** tables that have no clustered index. The data rows are not stored in any particular order, and there is no particular order to the sequence of the data pages. The data pages are not linked in a linked list.

Allocation Units

- An allocation unit is a collection of pages within a heap or B+tree used to manage data based on their page type.

Allocation unit type	Is used to manage
IN_ROW_DATA	Data or index rows that contain all data, except large object (LOB) data. Pages are of type Data or Index.
LOB_DATA	Large object data stored in one or more of these data types: text , ntext , image , xml , varchar(max) , nvarchar(max) , varbinary(max) , or CLR user-defined types (CLR UDT). Pages are of type Text/Image.
ROW_OVERFLOW_DATA	Variable length data stored in varchar , nvarchar , varbinary , or sql_variant columns that exceed the 8,060 byte row size limit. Pages are of type Data.

IN_ROW_DATA Allocation Unit

- For every partition used by a table (heap or clustered table), index, or indexed view, there is one IN_ROW_DATA allocation unit that is made up of a collection of data pages.
- This allocation unit also contains additional collections of pages to implement each nonclustered and XML index defined for the table or view.

ROW_OVERFLOW_DATA Allocation Unit

- For every partition used by a table (heap or clustered table), index, or indexed view, there can be one ROW_OVERFLOW_DATA allocation unit.
- This allocation unit contains zero (0) pages until a data row with variable length columns (**varchar**, **nvarchar**, **varbinary**, or **sql_variant**) in the IN_ROW_DATA allocation unit exceeds the 8 KB row size limit.
- When the size limitation is reached, SQL Server moves the column with the largest width from that row to a page in the ROW_OVERFLOW_DATA allocation unit. A 24-byte pointer to this off-row data is maintained on the original page.

LOB_DATA Allocation Unit

- When a table or index has one or more LOB data types, one LOB_DATA allocation unit per partition is allocated to manage the storage of that data.
- The LOB data types include **text**, **ntext**, **image**, **xml**, **varchar(max)**, **nvarchar(max)**, **varbinary(max)**, and CLR user-defined types.

Heaps

- A heap stores a table without a clustered index.
- When a heap has multiple partitions, each partition has a heap structure that contains the data for that specific partition.
- At a minimum, each heap will have one `IN_ROW_DATA` allocation unit per partition. The heap may also have one `LOB_DATA` allocation unit per partition and one `ROW_OVERFLOW_DATA` allocation unit per partition.

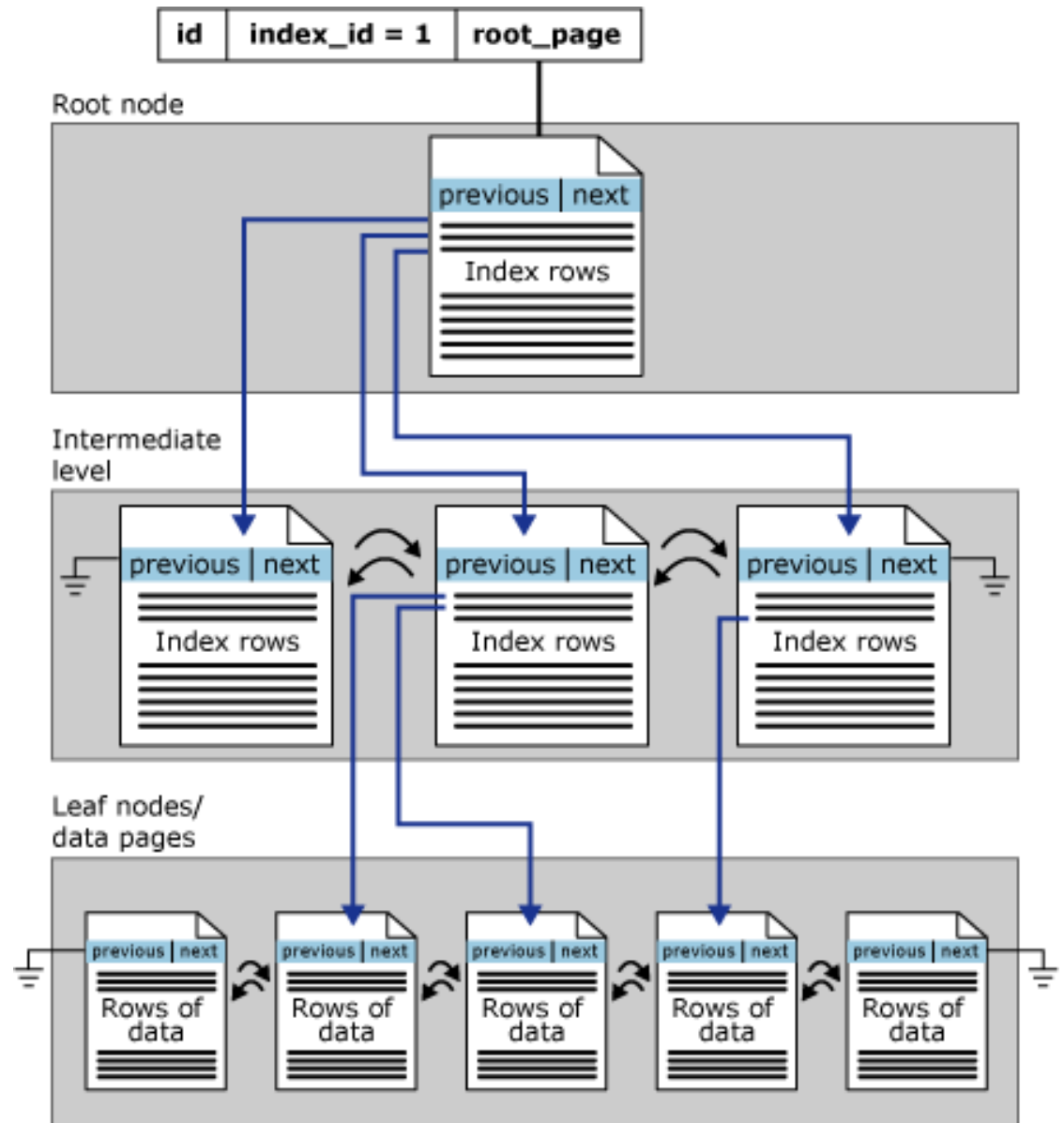
Clustered Tables

- They are organized as B+-trees
- One page per node
- The leaf nodes contain the data pages of the underlying table.
- Thus the data is stored inside the clustered index
- The pages in each level of the index are linked in a doubly-linked list.
- When a clustered index has multiple partitions, each partition has a B+-tree structure that contains the data for that specific partition.

Clustered Tables

- The pages in the data chain and the rows in them are ordered on the value of the clustered index key.
- All inserts are made at the point where the key value in the inserted row fits in the ordering sequence among existing rows.
- At a minimum, each clustered index will have one `IN_ROW_DATA` allocation unit per partition. The clustered index may also have one `LOB_DATA` allocation unit per partition and one `ROW_OVERFLOW_DATA` allocation unit per partition.

A Clustered Table in a Single Partition



Nonclustered Indexes

- Same B+-tree structure as clustered tables but
 - The leaf layer of a nonclustered index is made up of index pages instead of data pages.
 - The data rows of the underlying table are not sorted and stored in order based on their nonclustered keys
- When a nonclustered index has multiple partitions, each partition has a B+-tree structure that contains the data for that specific partition.

Nonclustered Indexes

- Nonclustered indexes can be defined on a table or view with a clustered index or a heap.
- Each index row in the leaves of the nonclustered index contains the key value and a row locator. This locator points to the data row in the clustered index or heap having the key value.
- At a minimum, each nonclustered index will have one `IN_ROW_DATA` allocation unit per partition. The unclustered index may also have one `LOB_DATA` allocation unit per partition and one `ROW_OVERFLOW_DATA` allocation unit per partition.

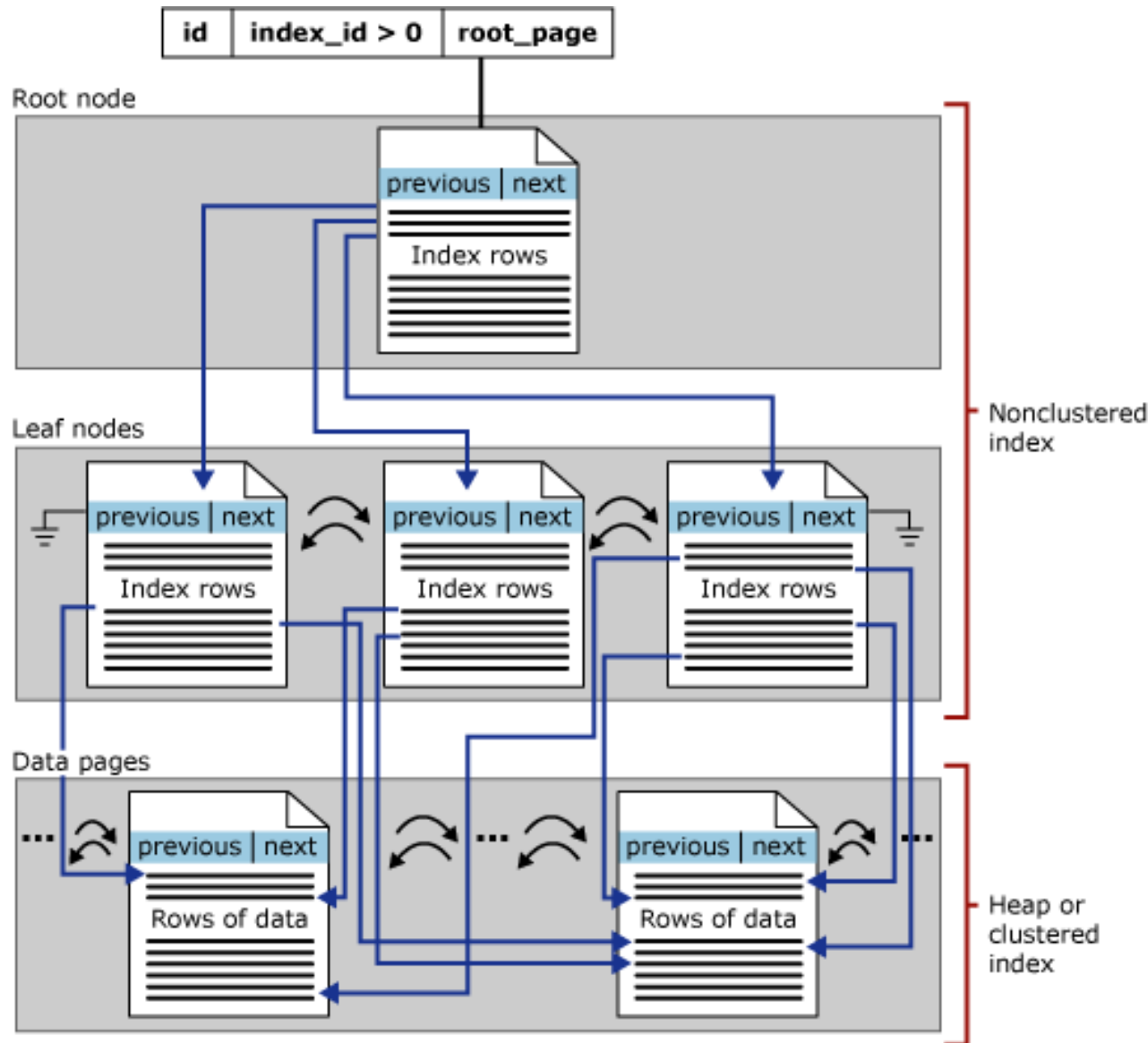
Row Locators

- If the table is a heap the row locator is a pointer to the row. The pointer is built from the file identifier (ID), page number, and number of the row on the page. The whole pointer is known as a Row ID (RID).
- If the table has a clustered index, the row locator is the clustered index key for the row.
 - If the clustered index is not a unique index, SQL Server 2005 makes any duplicate keys unique by adding an internally generated value called a **uniqueifier**. This four-byte value is not visible to users.
 - SQL Server retrieves the data row by searching the clustered index using the clustered index key stored in the leaf row of the nonclustered index.

Nonclustered Indexes

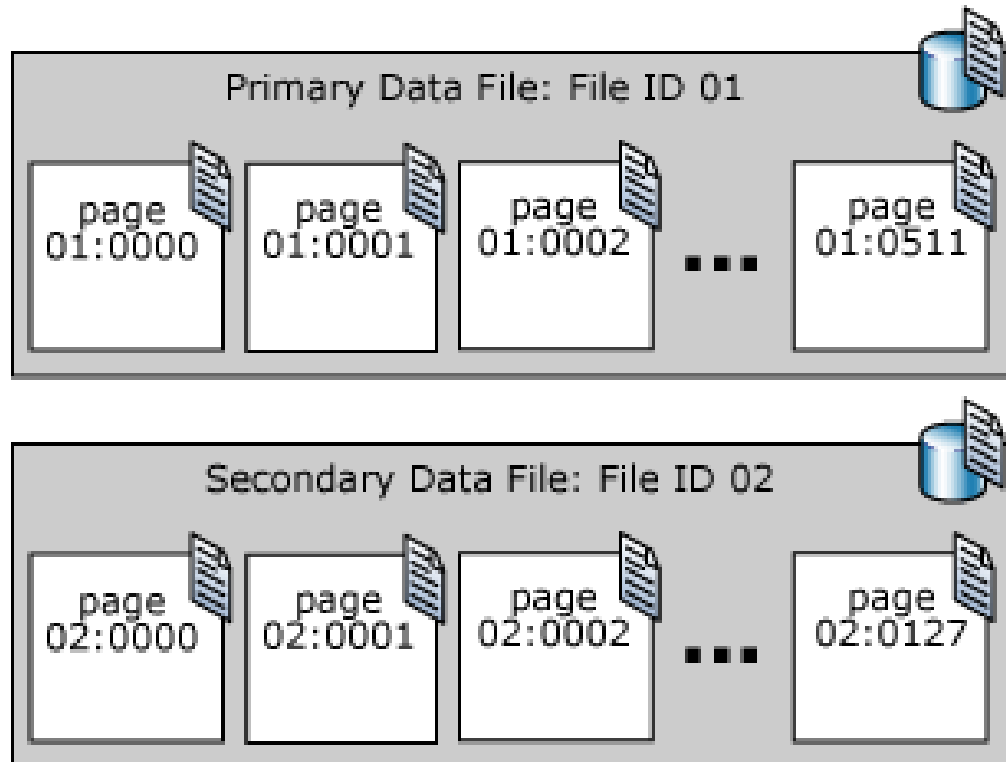
- When a nonclustered index has multiple partitions, each partition has a B+tree structure that contains the index rows for that specific partition.

A Nonclustered Index in a Single Partition



Pages

- Pages in file are numbered sequentially, starting from 0



Pages

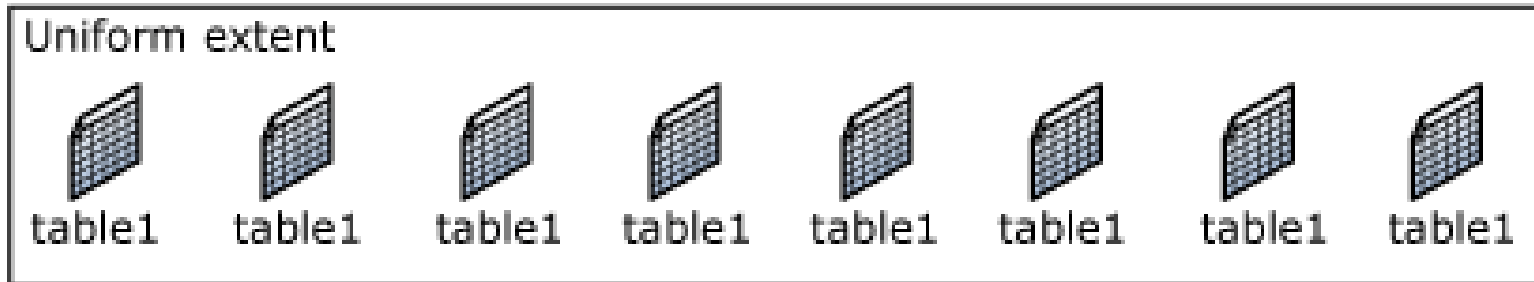
- The first page of a file contains information about the attributes of the file
- Other pages at the beginning of the file can be used for containing system information, such as allocation maps

Management of Space on Disks

- Space on disks is managed in **extents**
- An extent is eight physically contiguous pages, or 64 KB. This means SQL Server databases have 16 extents per megabyte.
- SQL Server has two types of extents:
 - Uniform extents are owned by a single object; all eight pages in the extent can only be used by the owning object.
 - Mixed extents are shared by up to eight objects. Each of the eight pages in the extent can be owned by a different object.

Extents

- A new table or index is generally allocated pages from mixed extents. When the table or index grows to the point that it has eight pages, it then switches to use uniform extents for subsequent allocations.



Space Allocation

- The data structures that manage extent allocations and track free space have a relatively simple structure. Benefits:
 - The free space information is densely packed, so relatively few pages contain this information.
 - Most of the allocation information is not chained together. This simplifies the maintenance of the allocation information.

Space Allocation

- SQL Server uses two types of allocation maps to record the allocation of extents:
- Global Allocation Map (GAM): GAM pages record what extents have been allocated. Each GAM covers 64,000 extents, or almost 4 GB of data. The GAM has one bit for each extent. If the bit is 1, the extent is free; if the bit is 0, the extent is allocated.
- Shared Global Allocation Map (SGAM): SGAM pages record which extents are currently being used as mixed extents and also have at least one unused page. Each SGAM covers 64,000 extents, or almost 4 GB of data. The SGAM has one bit for each extent. If the bit is 1, the extent is being used as a mixed extent and has a free page. If the bit is 0, the extent is not used as a mixed extent, or it is a mixed extent and all its pages are being used.

GAM and SGAM Bits

Current use of extent	GAM bit setting	SGAM bit setting
Free, not being used	1	0
Uniform extent, or full mixed extent	0	0
Mixed extent with free pages	0	1

Space Allocation Algorithm

- To allocate a uniform extent, the Database Engine searches the GAM for a 1 bit and sets it to 0.
- To find a mixed extent with free pages, the Database Engine searches the SGAM for a 1 bit.
- To allocate a mixed extent, the Database Engine searches the GAM for a 1 bit, sets it to 0, and then also sets the corresponding bit in the SGAM to 1.
- To deallocate an extent, the Database Engine makes sure that the GAM bit is set to 1 and the SGAM bit is set to 0.
- The algorithms that are actually used internally by the Database Engine are more sophisticated but they do not have to manage chains of extent allocation information.

Tracking Free Space

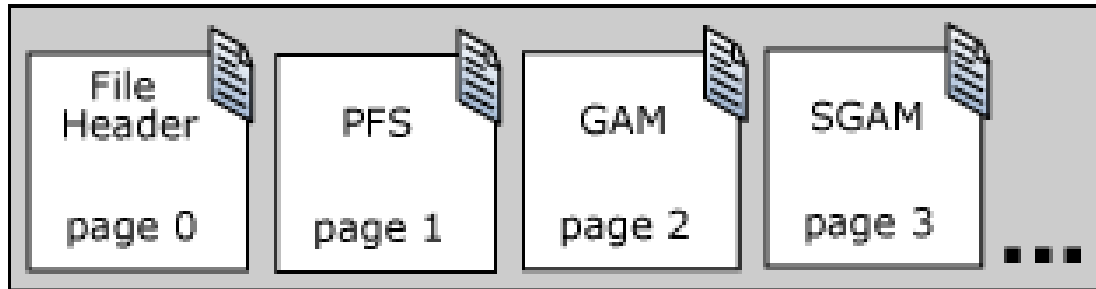
- Page Free Space (PFS) pages record the allocation status of each page, whether an individual page has been allocated, and the amount of free space on each page.
- The PFS has one byte for each page, recording whether the page is allocated, and if so, whether it is empty, 1 to 50 percent full, 51 to 80 percent full, 81 to 95 percent full, or 96 to 100 percent full.
- After an extent has been allocated to an object, the Database Engine uses the PFS pages to record which pages in the extent are allocated or free. This information is used when the Database Engine has to allocate a new page.

Tracking Free Space

- The amount of free space in a page is only maintained for heap and Text/Image pages. It is used when the Database Engine has to find a page with free space available to hold a newly inserted row.
- Indexes do not require that the page free space be tracked, because the point at which to insert a new row is set by the index key values

File structure

- A PFS page is the first page after the file header page in a data file (page number 1). This is followed by a GAM page (page number 2), and then an SGAM page (page 3). There is a PFS page approximately 8,000 pages in size after the first PFS page. There is another GAM page 64,000 extents after the first GAM page on page 2, and another SGAM page 64,000 extents after the first SGAM page on page 3.



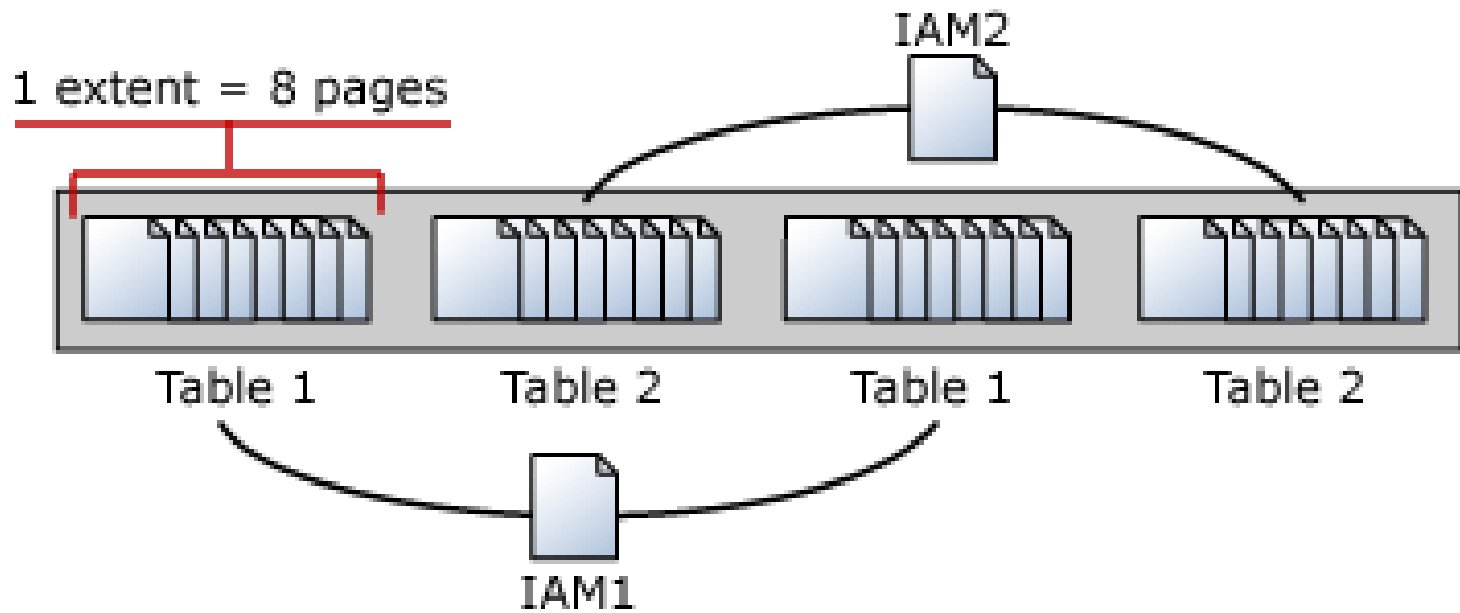
Space Used by Allocation Units

- An Index Allocation Map (IAM) page maps the extents in a 4-gigabyte (GB) part of a database file used by an allocation unit
- An IAM page has the same coverage as a GAM or SGAM page.
- If the allocation unit contains extents from more than one file, or more than one 4-GB range of a file, there will be multiple IAM pages linked in an IAM chain.
- Therefore, each allocation unit has at least one IAM page for each file on which it has extents.

IAM Pages

- An IAM page has a header that indicates the starting extent of the range of extents mapped by the IAM page.
- The IAM page also has a large bitmap in which each bit represents one extent. The first bit in the map represents the first extent in the range, the second bit represents the second extent, and so on.
- If a bit is 0, the extent it represents is not allocated to the allocation unit owning the IAM. If the bit is 1, the extent it represents is allocated to the allocation unit owning the IAM page.

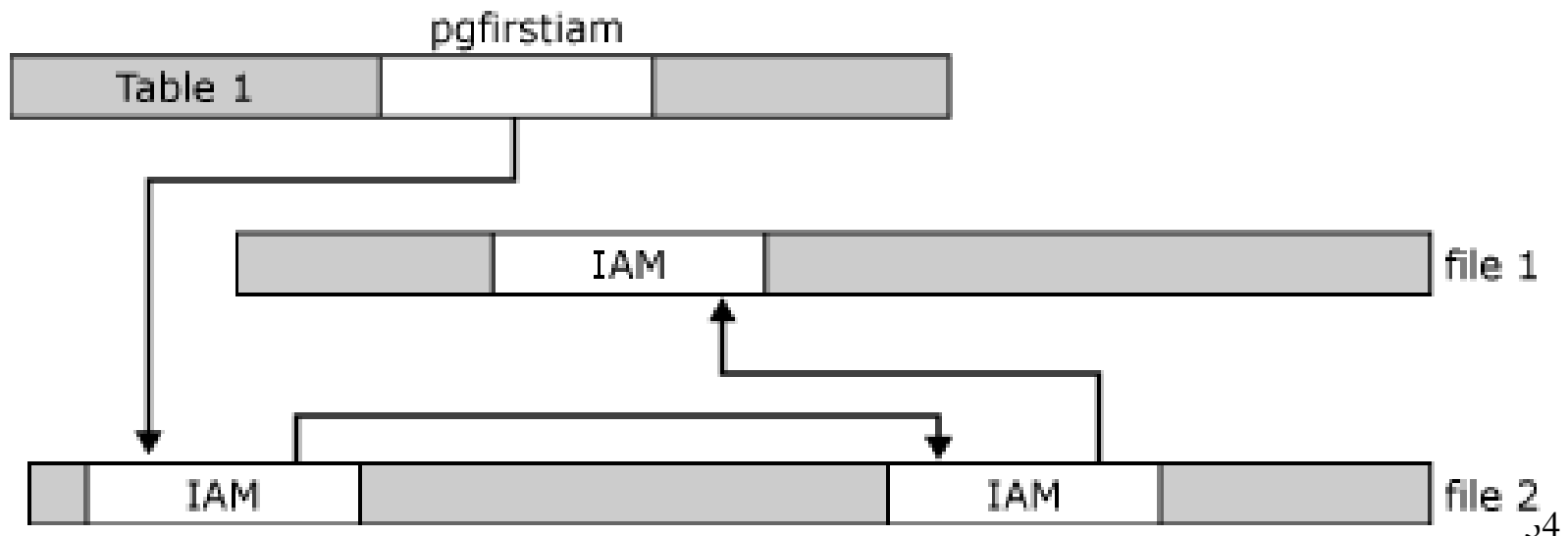
IAM Pages



IAM Pages

- IAM pages are allocated as required for each allocation unit and are located randomly in the file.
- The catalog view, **sys.system_internals_allocation_units** points to the first IAM page for an allocation unit.

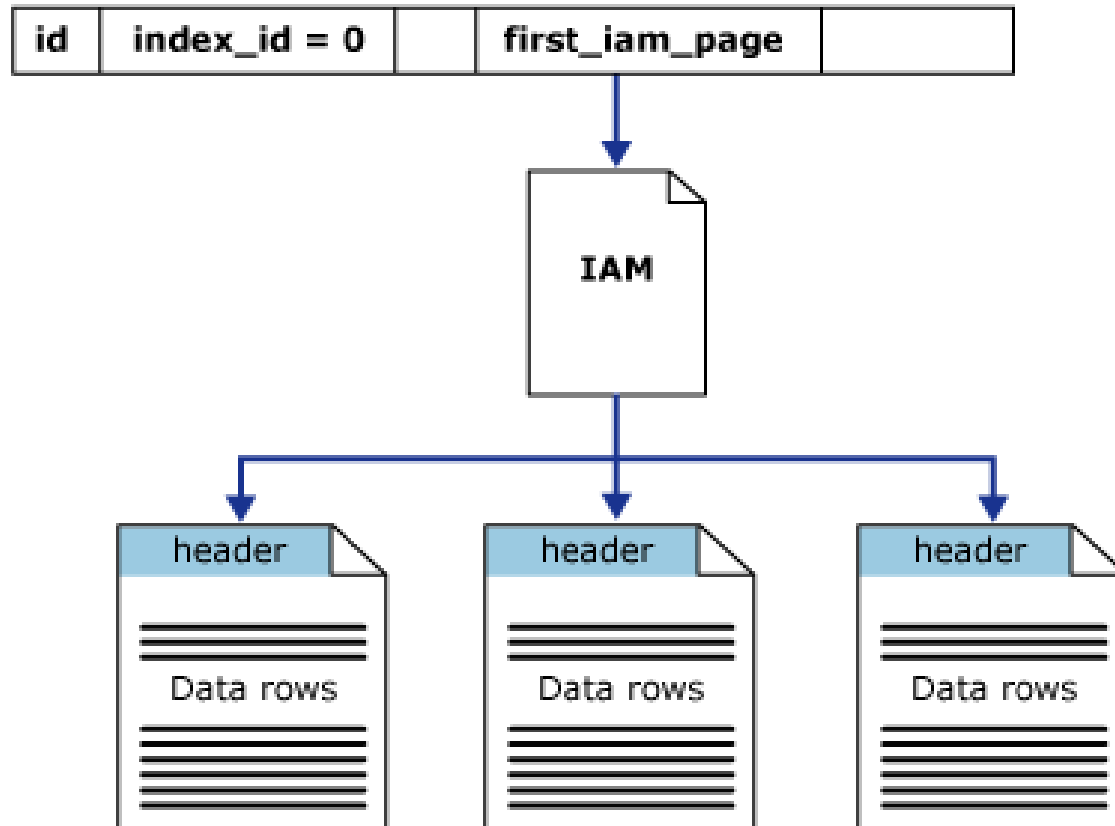
sys.system_internals_allocation_units



Heap Organization

- The IAM pages are used to move through the heap. The only logical connection between data pages is the information recorded in the IAM pages.
- The column **first_iam_page** in the **sys.system_internals_allocation_units** catalog view points to the first IAM page in the chain of IAM pages that manage the space allocated to the heap in a specific partition.
- Table scans or serial reads of a heap can be performed by scanning the IAM pages to find the extents that are holding pages for the heap. Because the IAM represents extents in the same order that they exist in the data files, this means that serial heap scans progress sequentially through each file.

Heap Organization



Space Allocation Algorithms

- When the SQL Server Database Engine has to insert a new row in a heap or Text/Image page and no space is available in the current page, it uses the IAM and PFS pages to find a page with sufficient space to hold the row.
- The Database Engine uses the IAM pages to find the extents allocated to the allocation unit. For each extent, the Database Engine searches the PFS pages to see if there is a page that can be used.
- Each IAM and PFS page covers lots of data pages, so there are few IAM and PFS pages in a database. This means that the IAM and PFS pages are generally in memory in the SQL Server buffer pool, so they can be searched quickly.

Space Allocation Algorithms

- For clustered tables, the insertion point of a new row is set by the index key. In this case, the search process previously described does not occur.
- The Database Engine allocates a new extent to an allocation unit only when it cannot quickly find a page in an existing extent with sufficient space to hold the row being inserted.

System Databases

- System databases store metadata regarding the databases available in the instance

System database	Description
master Database	Records all the system-level information for an instance of SQL Server.
msdb Database	Is used by SQL Server Agent for scheduling alerts and jobs.
model Database	Is used as the template for all databases created on the instance of SQL Server. Modifications made to the model database, such as database size, collation, recovery model, and other database options, are applied to any databases created afterward.
Resource Database	Is a read-only database that contains system objects that are included with SQL Server. System objects are physically persisted in the Resource database, but they logically appear in the sys schema of every database.
tempdb Database	Is a workspace for holding temporary objects or intermediate result sets.