

MaxDB

Introduction to MaxDB Internals
Version 7.6

Heike Gursch
Werner Thesing

THE BEST-RUN BUSINESSES RUN SAP™



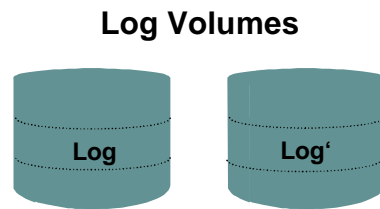
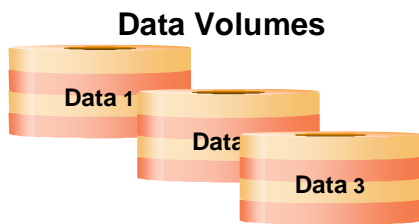


Database Kernel

Process Layer

Memory Level

Data Storage





Buffer Interface

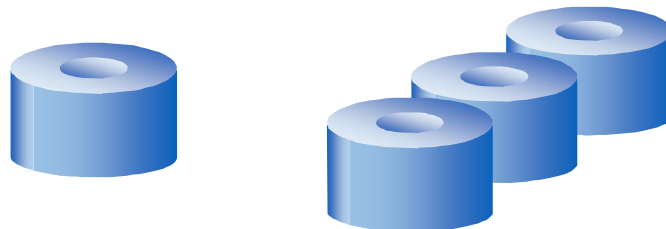
SQL Manager (AK)

Buffer Interface

Data Access Manager
communication base layer (KB)

Procedural Interface

Data Access Manager
base data layer (BD)



The MaxDB kernel consists of three levels:

SQL statements sent to the database system are received by the highest level, the SQL Manager (previously known as the application communication layer).

They are then sent in compressed form to the Data Access Manager. The Data Access Manager is divided into the communication base layer (KB) and the base data (BD) layer.

The necessary individual queries are determined on the communication base layer.

The required data is then procured on the lower-level base data layer.

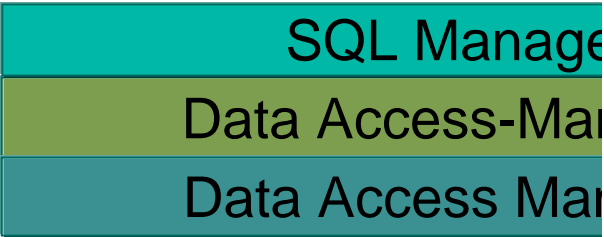
The data records are then sent through the layers to the querying application.

In recent years, the KB and BD layers have increasingly merged.

Since version 7.4 there has no longer been a sharp division between the two layers.



```
SELECT name1, name2, tel1
FROM kna1
WHERE name1 = 'GARDENER'
or name1 = 'GARDENERS'
```



Catalog entries for table kna1

Table kna1: ID 3421; Node p15714

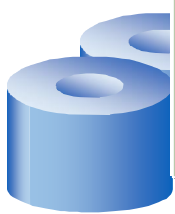
Column	Type	Length
mandt	NUMCHAR	3
kunnr	NUMCHAR	8
name1	CHAR	20
tel1	CHAR	20

Index kna1__1: ID 3422; Node p15714

Column	Type	Length
Name1	CHAR	20
Name2	CHAR	20
Prim.Key	NUMCHAR	13

Access strategy:

- 1) Perform Range Scan on Index with ID 3422 located on node p15714
- 2) Select Rows by Primary Key on table with ID 3421 located on node p15714



Let's have a look at a relatively simple SELECT statement:

The SQL Manager receives an SQL statement that was sent as an SQL packet and first checks the syntax for correctness.

It converts the incoming SQL packets into stack code with various entries, each with a length of 8 bytes.

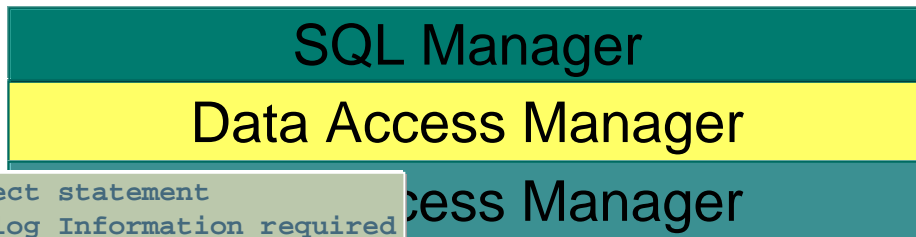
Catalog management, that is, the determination of information about the existing tables and their columns, also takes place on this layer.

Finally, the Optimizer in the SQL Manager decides if the data should be accessed via a secondary index.

In our simple example, we assume that the Optimizer has opted for an access strategy via the secondary index kna1_1 .

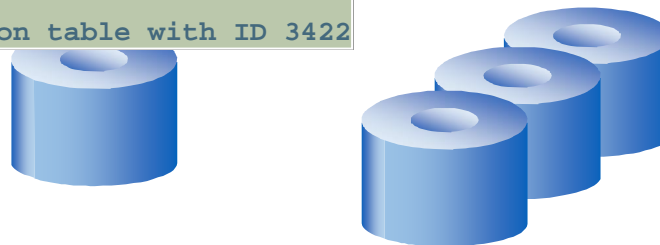


```
SELECT name1, name2, tel1  
FROM kna1  
WHERE name1 = 'GARDENER'  
or name1 = 'GARDENERS'
```



Select statement
No log information required

Locking
No locks held on table with ID 3422



Transaction management and lock management are executed on the communication base layer of the Data Access Manager.

Transaction management manages information about active transactions in the database as well as the changes they make to data.

Lock management provides memory lists containing all objects of an instance that are locked at a certain point in time as well the transactions that correspond to them.

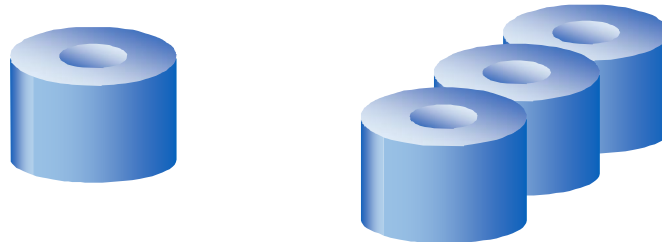
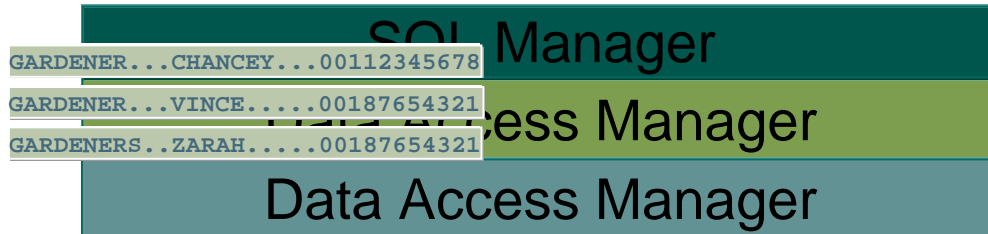
In the case of complex SQL statements (such as joins), access to individual tables and indexes is coordinated.

In our simple example, the KB layer will first access the index table kna1_1 and then, using the determined index entries, directly access the appropriate rows of table kna1.

Data Access Manager (BD layer)



```
select NAME1, NAME2, TEL1  
from KNA1  
where NAME1 = 'GARDENER' or NAME1 = 'GARDENERS'
```

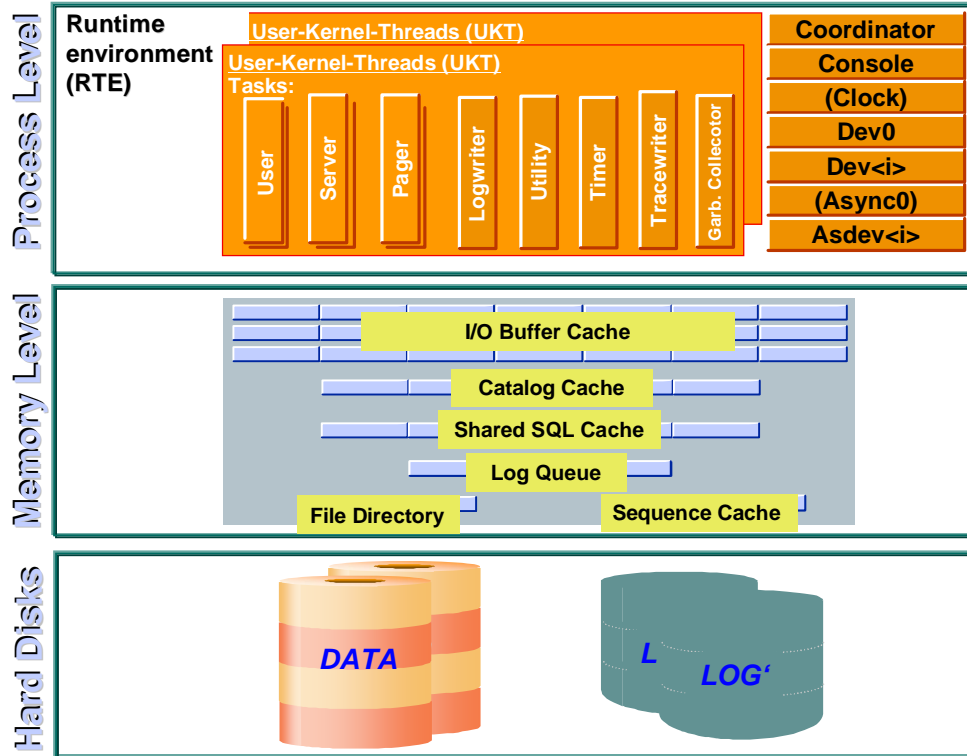


© SAP 2007 / MaxDB 76 Internals – Introduction /Page 7

The base data layer of the Data Access Manager executes accesses to data that is either still located in the data volumes or is already in the cache.

The BD layer returns the results records to the KB layer.

Database Instance



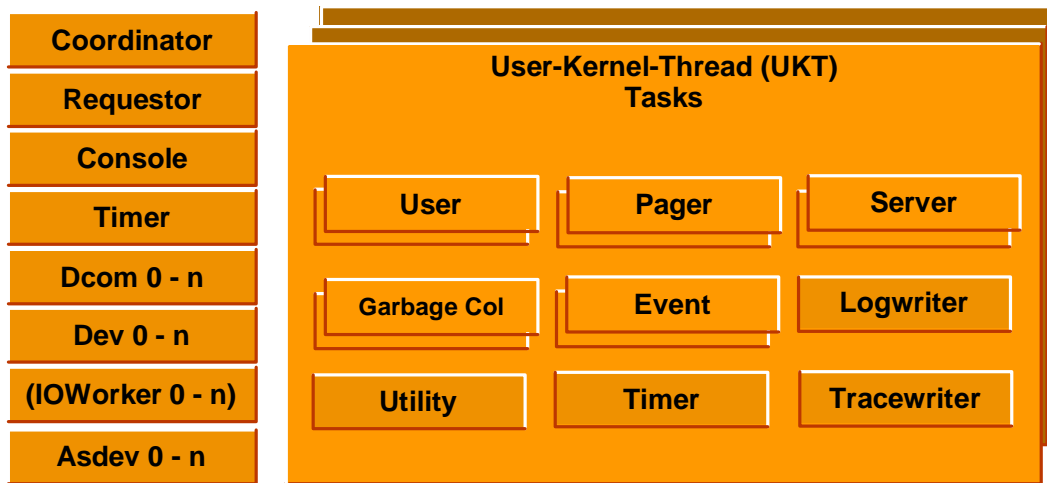
© SAP 2007 / MaxDB 76 Internals – Introduction /Page 8

A database instance consists of the following three parts:

Database kernel (process level)

Caches (memory level)

Volumes (hard disk level)

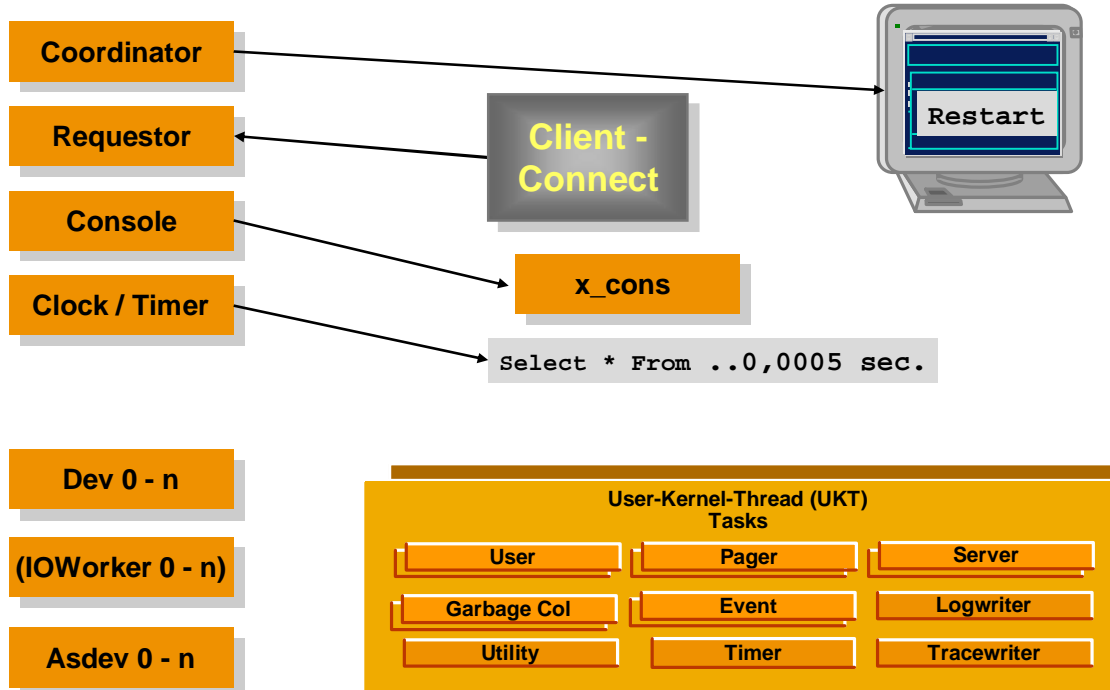


A user kernel thread forms a subset of all tasks (internal tasking).

The database kernel runs as one process divided into threads. Threads can be active in parallel on several processors within the operating system. Threads perform various tasks.

User kernel threads (UKT) consist of several tasks that perform various tasks. This tasking enables more efficient coordination of tasks than operating system tasking that employs individual threads.

The runtime environment (RTE) defines the structure of the process and the user kernel thread.



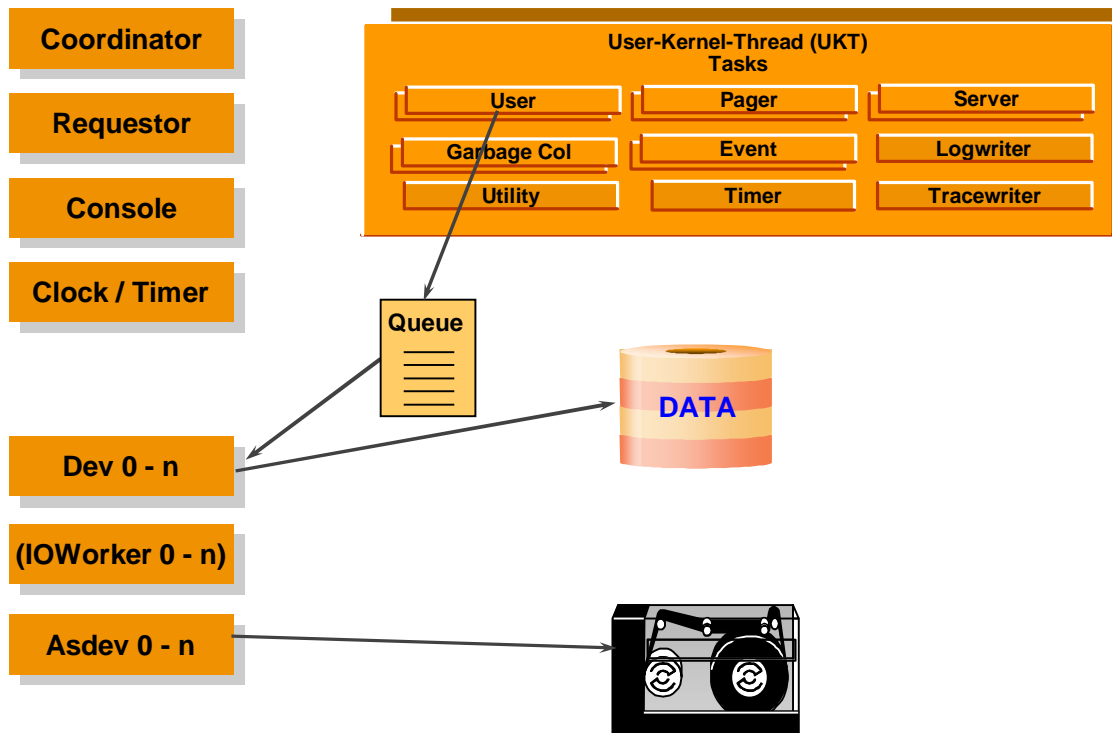
When the runtime environment is started, that is, when the database instance is started in the Admin state, first the **coordinator thread** is generated. This thread is of particular importance.

- When started, the coordinator thread uses database parameters to determine the configuration of the memory and process of the database instance. For this reason, changes to database parameters take effect only after you have restarted the database instance.
- The coordinator thread also coordinates the start procedures of the other threads and monitors them while the database instance is in operation.
- If operating errors occur, the coordinator thread can stop other threads.

The **requestor thread** receives logon requests from the user processes to the database. The logon is assigned to a task within the user kernel thread.

The **console thread** collects information about the internal states of the database kernel when x_cons is being used.

The **clock thread** and the **timer thread** calculate internal times, for example to determine how much time was required to execute an SQL statement.



© SAP 2007 / MaxDB 76 Internals – Introduction /Page 11

Dev threads are responsible for processing the write and read requests to and from data and log volumes that are requested by the corresponding tasks. MaxDB supports asynchronous I/O requests. On Windows, the asynchronous I/O of the operating system is used.

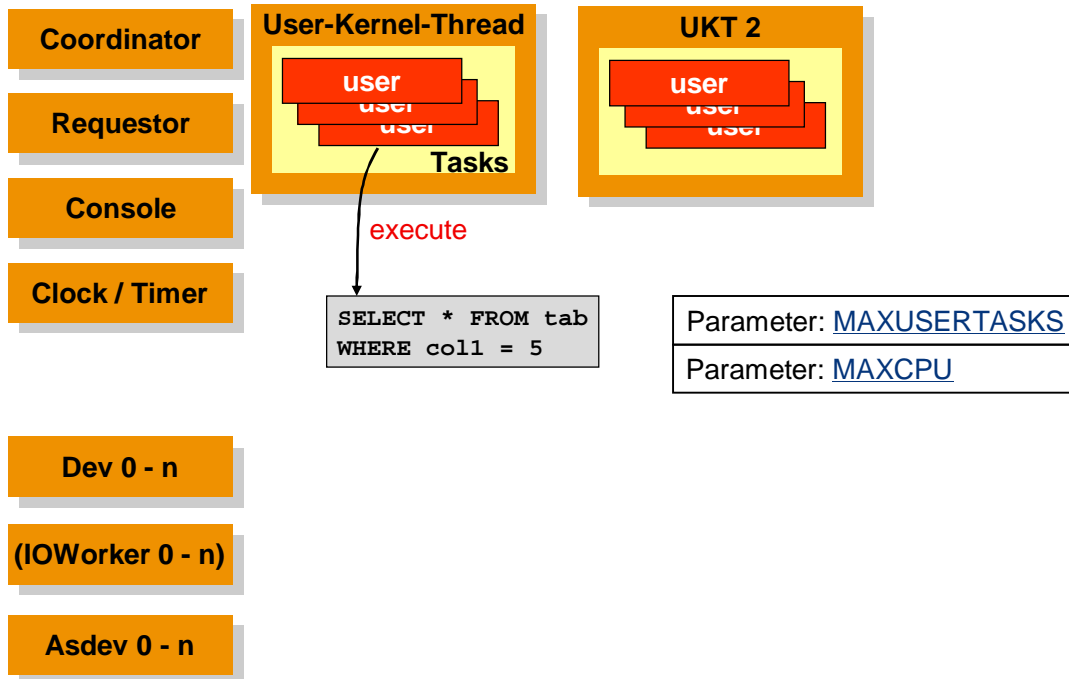
The number of dev threads is primarily dependent on the number of volumes in the database instance. As a rule, two dev threads are activated for each data and log volume and one for the writing of the database instance. When using the asynchronous I/O of the operating system, only one dev thread is started in each case as the I/O calls of the operating system are used.

The thread **dev0** plays a special role; dev0 coordinates and monitors the dev threads. For example, if a volume fails during operation (*Bad Volume*), dev0 ensures that the corresponding dev threads are terminated. If the database is enlarged while in operation through the addition of another data volume, dev0 ensures that new dev threads are generated.

If the asynchronous I/O calls of the operating system are being used, the user kernel threads send their requests directly to the operating system. The reply from the operating system is put in a queue. The queue is read by the **IOWorker threads**.

Threads are also temporarily activated to read and write data for data backups. These threads are called **asdev<i>**. Their number depends on the number of data volumes and of the number of backup devices.

User Kernel Threads and Tasks



© SAP 2007 / MaxDB 76 Internals – Introduction / Page 12

Each user session is assigned exactly one **user task** at logon.

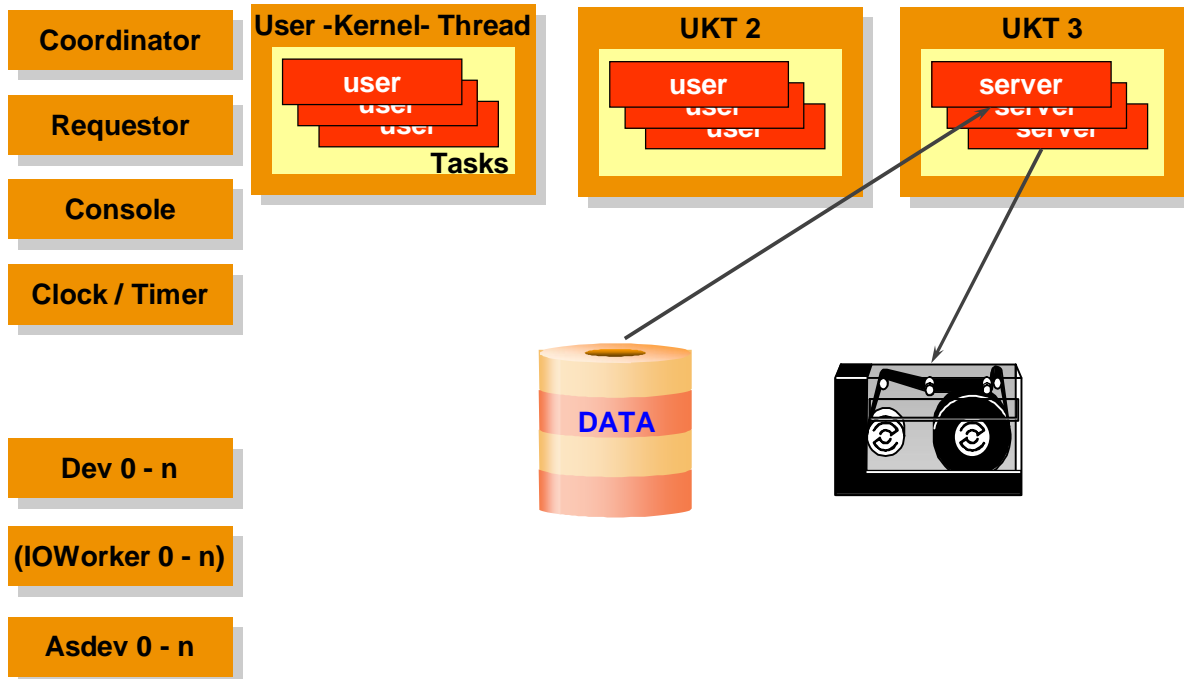
The maximum number of available user tasks is determined by the database parameter `MAXUSERTASKS`. This parameter also restricts the number of user sessions that can be logged on to the database system simultaneously.

The database parameter `_MAXTASK_STACK` determines the memory usage of the user tasks.

The general database parameter `MAXCPU` specifies the number of user kernel threads among which the user tasks are distributed. Other tasks and global threads use very little CPU time. The parameter `MAXCPU` allows you to specify how many processors the database should use in parallel.

As of version 7.4.03, user tasks can switch from one UKT to another if the previously-responsible UKT is overburdened. This results in better scaling for multiprocessor servers (SMP). To use this function, set the parameter `LOAD_BALANCING_CHK` to a value greater than 0.

Server Tasks



© SAP 2007 / MaxDB 76 Internals – Introduction /Page 13

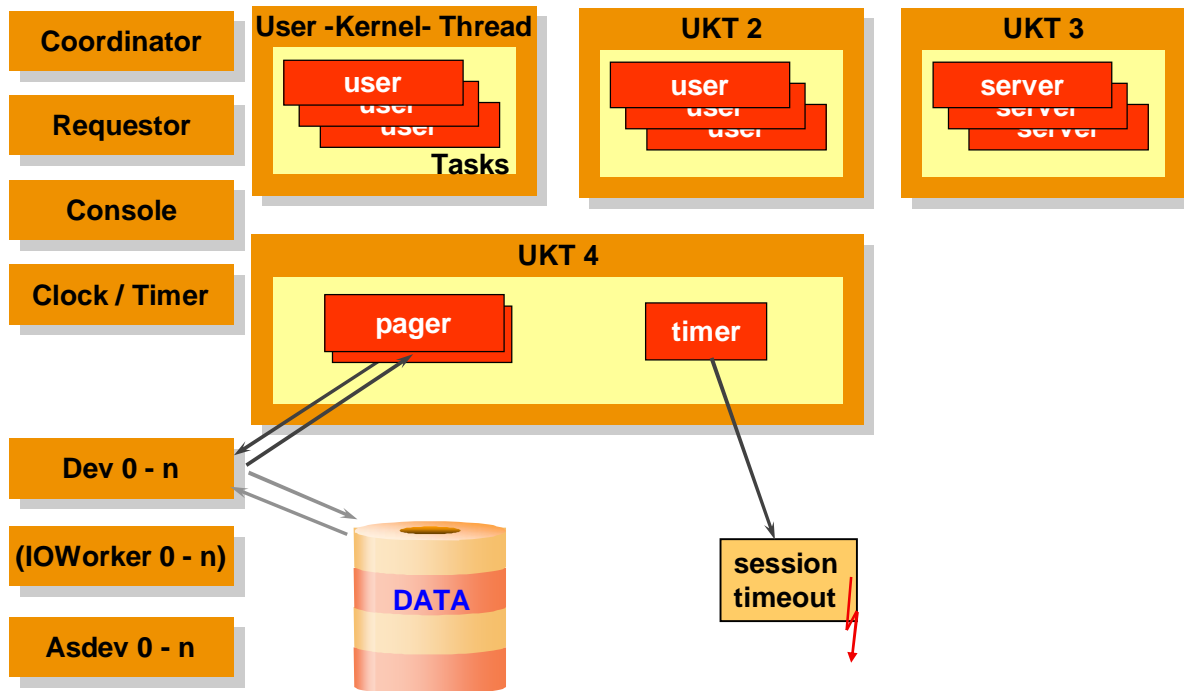
Server tasks are primarily used to back up data. Some server tasks read from the data volumes; others write to the backup medium.

The CREATE INDEX statement instructs the server tasks to read the table data in parallel from the data volumes.

With DROP TABLE statement, server tasks delete the data in the tables asynchronously. Users do not have to wait for all data to be deleted.

The system automatically calculates the number of server tasks needed in the configuration of the database instance from the number of data volumes and the number of backup devices.

Pager and Timer Tasks



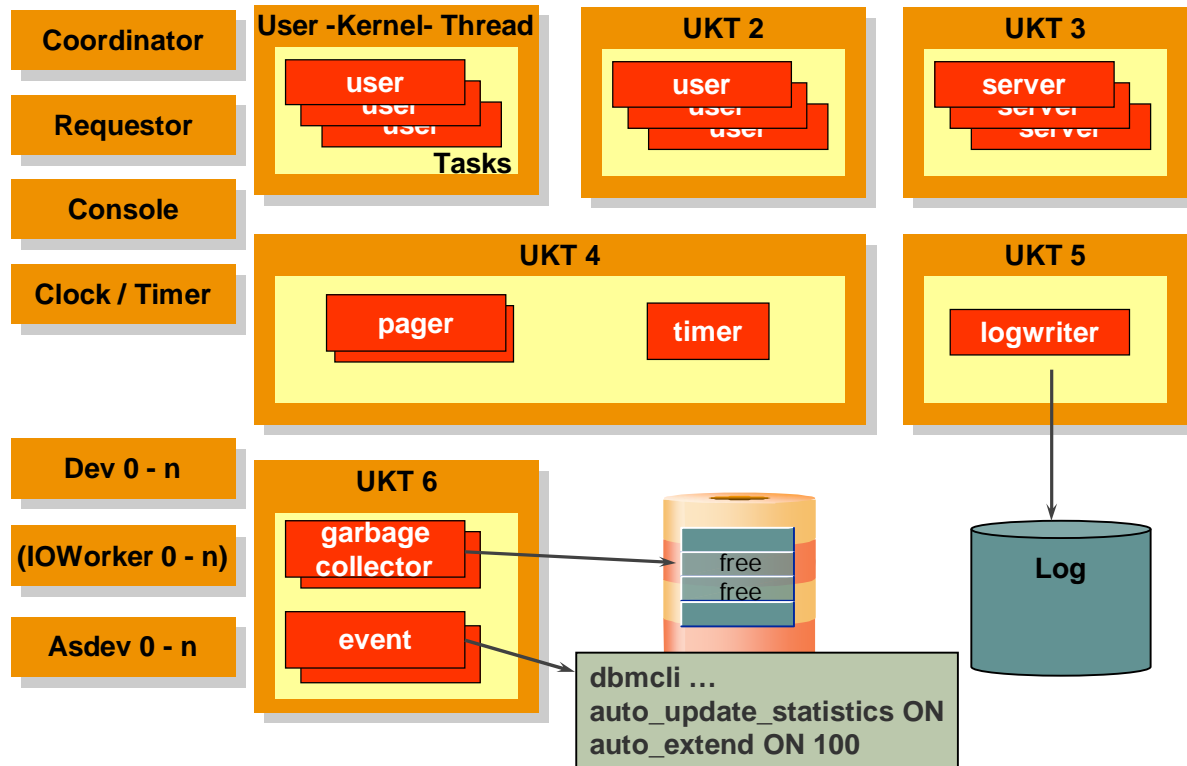
© SAP 2007 / MaxDB 76 Internals – Introduction /Page 14

Pager tasks are responsible for writing data from the data cache to the data volumes. They become active when a savepoint is being executed.

The system calculates the number of pagers. It depends primarily on the data cache size and the number of data volumes.

The **timer task** is used for handling all types of timeout situations (such as session timeouts and lock request timeouts).

Garbage Collectors, Logwriter and Event Tasks



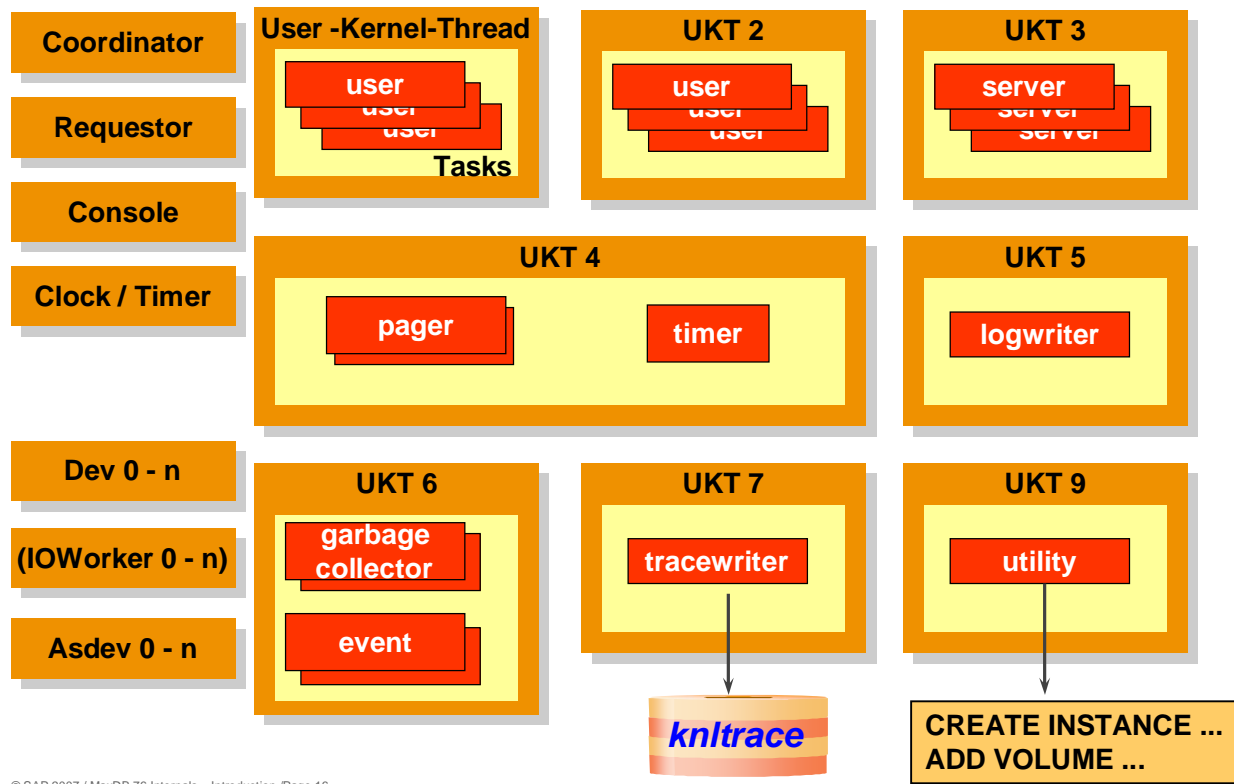
© SAP 2007 / MaxDB 76 Internals – Introduction /Page 15

The **Logwriter task** is responsible for writing redo log entries in the log volumes.

Garbage collectors release undo log files to free space management.

Event tasks send messages about database events to the Database Manager (e.g. DBMGUI). You can use the Event Dispatcher to define reactions. For example, you can have the database enlarge automatically when it gets full (`auto_extend`). You can also have the Event Dispatcher automatically update statistics when certain events occur (`auto_update_statistics`).

Utility-, Tracewriter Task



© SAP 2007 / MaxDB 76 Internals – Introduction /Page 16

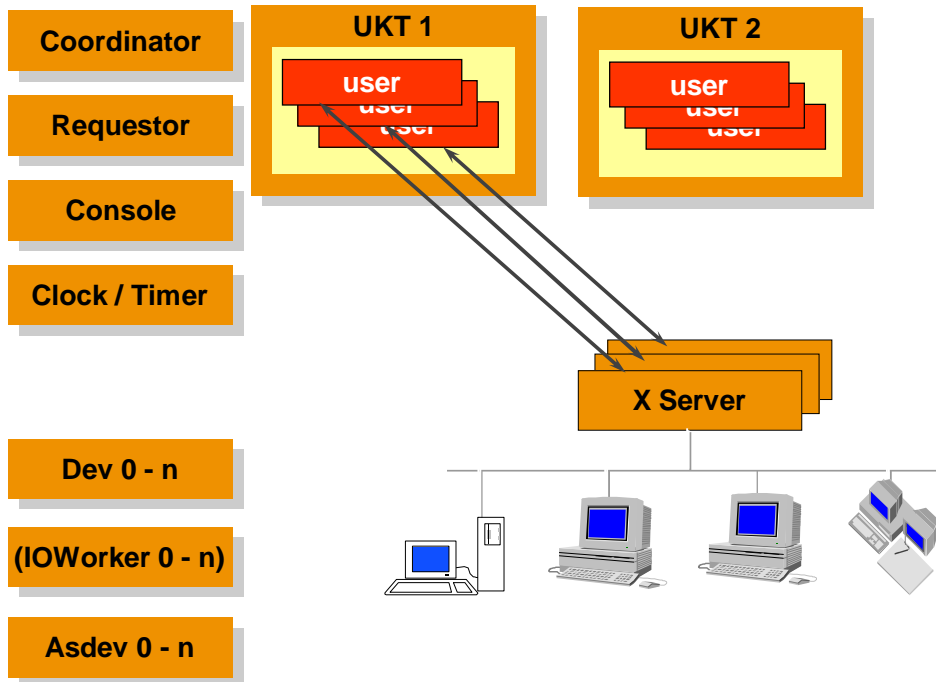
MaxDB offers the possibility of writing a special log, the database trace. The active tasks write the trace data to a buffer. If requested, the **trace writer task** writes the data from the buffer to the file knltrace.

The **utility task** is reserved exclusively for the administration of the database instance.

Automatic log backup can be executed in parallel with other administration tasks as it does not occupy a utility session after it is activated.

As of version 7.6 all administration actions are executed via user tasks. For compatibility reasons, the utility task will be retained for the foreseeable future although it will not be used by the DBM server in future versions.

Remote SQL Server

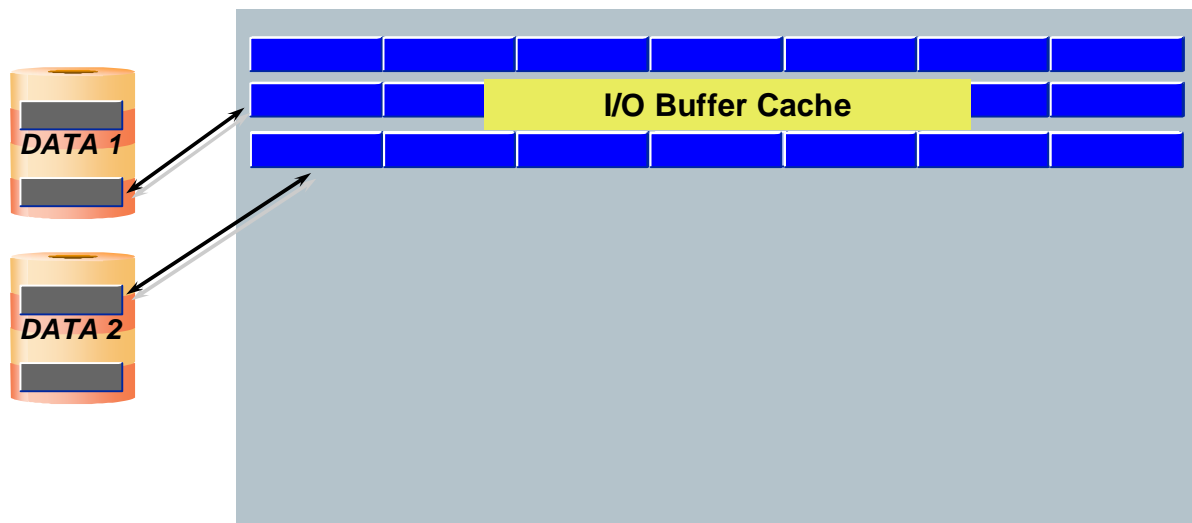


© SAP 2007 / MaxDB 76 Internals – Introduction /Page 17

The **X-Server** is the communication server of the database system. It listens at a service port for connection requests from database applications and database tools. In the process list, this process is called **vserver**. A new vserver process is generated for every user process that logs on to the database remotely. The generating process serves the user; the new process waits for the next user logon. On Windows, an additional thread is started for the user logon.

On Windows, the X-Server runs as a service.

Local user sessions communicate with the database instance via a shared memory.



© SAP 2007 / MaxDB 76 Internals – Introduction /Page 18

Read and write operations to the volumes are buffered in order to save time-critical disk accesses. The corresponding main memory structures are called caches. They can be dimensioned according to the user profile. The **I/O buffer cache** contains the last read- or write-accessed pages of the data volumes. It is shared by all simultaneously active users.

The hit rate, that is, the ratio of successful accesses to the total number of accesses to the **I/O buffer cache**, is a crucial measure of performance. It should be greater than 98%. Successful access means that the required data was already available in the data cache.

In addition to data pages, the I/O buffer cache also contains **converter pages**. Converter pages, like data pages, are stored in the data volumes. They store the assignment of the logical data page numbers to their physical position in the data volumes.

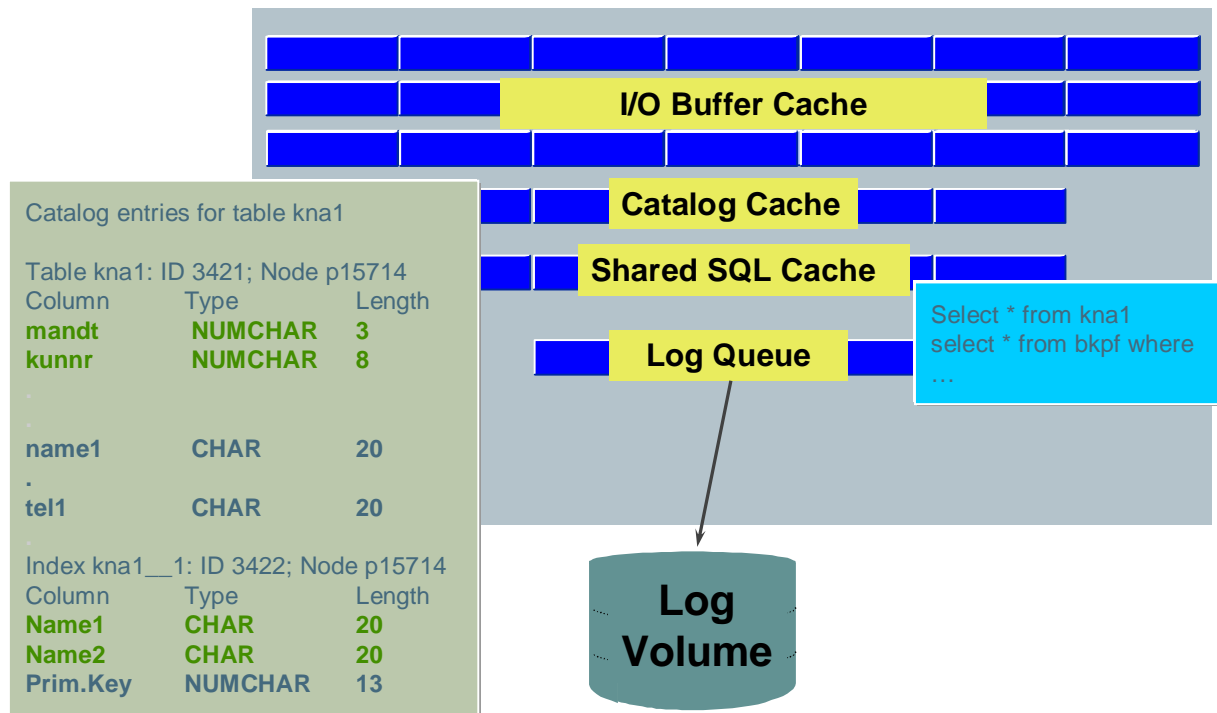
The number of converter pages is calculated automatically. It can increase with data growth. When deletions are performed, converter pages are released.

All converter pages are held in the cache. Each converter page contains 1861 entries that are managed for data pages.

The database parameter used for setting the size of the I/O buffer cache is called `CACHE_SIZE`.

In versions < 7.4, data pages and converter pages were managed in separate caches, the data cache and the converter cache.

Catalog Cache, Log Queue



© SAP 2007 / MaxDB 76 Internals – Introduction /Page 19

In the **catalog cache**, the database system stores user-specific data and global data from the database catalog. Data that has been displaced from the catalog cache is moved to the data cache. A catalog cache is assigned to each user task at the start of the database session; this catalog cache is then released after the database session has ended. The total of all catalog caches can increase up to the value that has been configured with the database parameter CAT_CACHE_SUPPLY.

The hit rate for the catalog cache should be over 85 percent. The lower limit for the hit rate depends largely on the application.

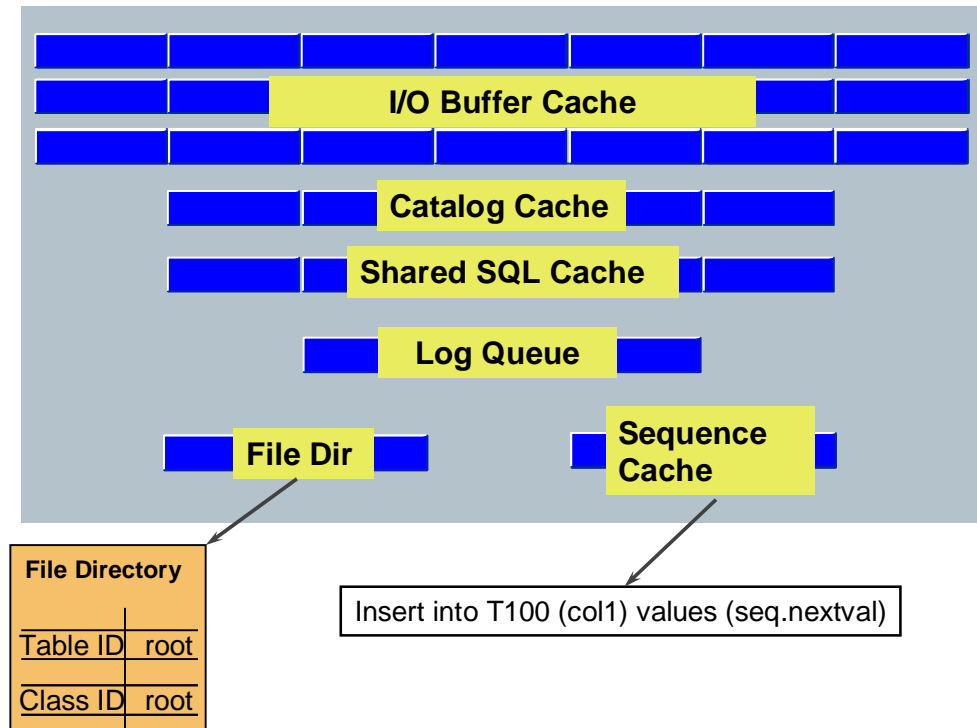
All statements to be executed will be held in the **shared SQL cache** along with their execution plans. The shared SQL cache is shared by all users, that is, a statement is only accepted once. When shared SQL is not being used, the statements of each user are kept in the catalog cache. The setting NO for database parameter SHAREDSQL deactivates use of this cache.

A statement, together with the statement text, is stored only once in the shared SQL cache. This allows you to see which statements are active in the database at any time. The shared SQL manager also collects monitor data that without the shared SQL are only visible with an active command or resource monitor. Shared SQL offers the following advantages over versions prior to 7.6:

- A statement only has to be prepared once rather than for each user.
- Statements are stored only once. That saves space in the main memory.
- Storing the data in the main memory offers better monitoring possibilities.

The **log queue** allows log entries to be written to the log volumes asynchronously and increases the likelihood of group commits. In a group commit, several write transactions are completed in the log area with an I/O.

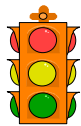
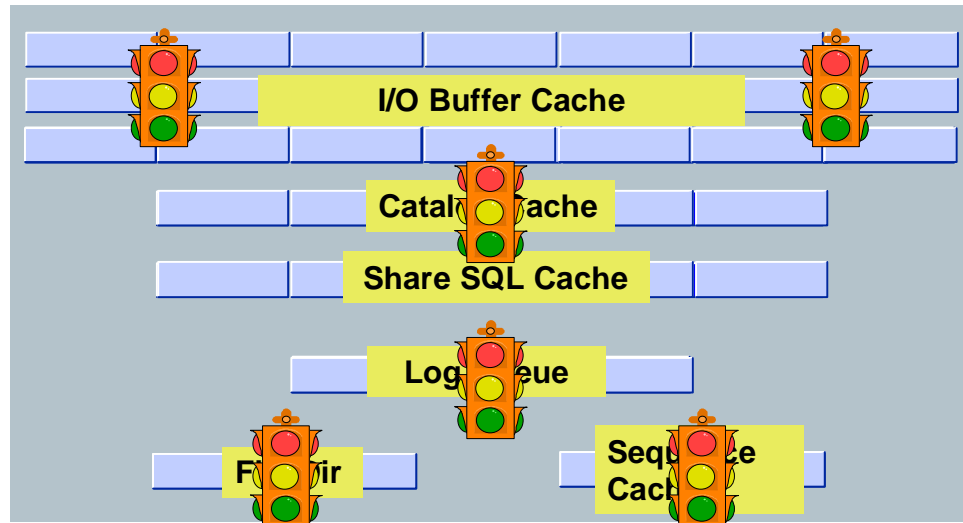
File Directory, Sequence Cache



© SAP 2007 / MaxDB 76 Internals – Introduction /Page 20

The **file directory** is required for the internal organization of the database instance. The assignments of the root pages of the B*-trees to the table IDs and a type flag are stored here. The type flag specifies whether the table contains primary data, secondary key or LONG data. Mit Version 7.6 wurde das File-Directory neu implementiert. Es enthält nun zusätzlich Angaben über die Anzahl der Sätze und die Größen der Tabellen. Die Ermittlung der Anzahl der Sätze einer großen Tabelle mit "SELECT COUNT(*) FROM <Tabelle>" kann damit in sehr kurzer Zeit erfolgen.

The **sequence cache** stores current data on number generators.



Critical sections in memory are protected by synchronization mechanisms (Regions)

Accesses to caches can be synchronized over one or more regions. Depending on its size, the data cache is comprised of 8 to 64 segments of the same size, each of which is protected by exactly one region.

If a task or thread accesses a critical section, the region locks this session for all other tasks or processes.

Other main memory structures are also managed via synchronization mechanisms provided by the database.

Reader-writer locks are used to synchronize the shared SQL cache. Reader-writer locks are used in version 7.5 and up. In contrast to regions, reader-writer locks make it possible to distinguish between shared and exclusive locks.



No part of this publication may be reproduced or transmitted in any form or for any purpose without the express permission of SAP AG. The information contained herein may be changed without prior notice.

Some software products marketed by SAP AG and its distributors contain proprietary software components of other software vendors.

SAP, R/3, mySAP, mySAP.com, xApps, xApp, SAP NetWeaver, Duet, Business ByDesign, ByDesign, PartnerEdge and other SAP products and services mentioned herein as well as their respective logos are trademarks or registered trademarks of SAP AG in Germany and in several other countries all over the world. All other product and service names mentioned and associated logos displayed are the trademarks of their respective companies. Data contained in this document serves informational purposes only. National product specifications may vary.

The information in this document is proprietary to SAP. This document is a preliminary version and not subject to your license agreement or any other agreement with SAP. This document contains only intended strategies, developments, and functionalities of the SAP® product and is not intended to be binding upon SAP to any particular course of business, product strategy, and/or development. SAP assumes no responsibility for errors or omissions in this document. SAP does not warrant the accuracy or completeness of the information, text, graphics, links, or other items contained within this material. This document is provided without a warranty of any kind, either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement.

SAP shall have no liability for damages of any kind including without limitation direct, special, indirect, or consequential damages that may result from the use of these materials. This limitation shall not apply in cases of intent or gross negligence.

The statutory liability for personal injury and defective products is not affected. SAP has no control over the information that you may access through the use of hot links contained in these materials and does not endorse your use of third-party Web pages nor provide any warranty whatsoever relating to third-party Web pages

Weitergabe und Vervielfältigung dieser Publikation oder von Teilen daraus sind, zu welchem Zweck und in welcher Form auch immer, ohne die ausdrückliche schriftliche Genehmigung durch SAP AG nicht gestattet. In dieser Publikation enthaltene Informationen können ohne vorherige Ankündigung geändert werden.

Einige von der SAP AG und deren Vertriebspartnern vertriebene Softwareprodukte können Softwarekomponenten umfassen, die Eigentum anderer Softwarehersteller sind.

SAP, R/3, mySAP, mySAP.com, xApps, xApp, SAP NetWeaver, Duet, Business ByDesign, ByDesign, PartnerEdge und andere in diesem Dokument erwähnte SAP-Produkte und Services sowie die dazugehörigen Logos sind Marken oder eingetragene Marken der SAP AG in Deutschland und in mehreren anderen Ländern weltweit. Alle anderen in diesem Dokument erwähnten Namen von Produkten und Services sowie die damit verbundenen Firmenlogos sind Marken der jeweiligen Unternehmen. Die Angaben im Text sind unverbindlich und dienen lediglich zu Informationszwecken. Produkte können länderspezifische Unterschiede aufweisen.

Die in diesem Dokument enthaltenen Informationen sind Eigentum von SAP. Dieses Dokument ist eine Vorabversion und unterliegt nicht Ihrer Lizenzvereinbarung oder einer anderen Vereinbarung mit SAP. Dieses Dokument enthält nur vorgesehene Strategien, Entwicklungen und Funktionen des SAP®-Produkts und ist für SAP nicht bindend, einen bestimmten Geschäftsweg, eine Produktstrategie bzw. -entwicklung einzuschlagen. SAP übernimmt keine Verantwortung für Fehler oder Auslassungen in diesen Materialien. SAP garantiert nicht die Richtigkeit oder Vollständigkeit der Informationen, Texte, Grafiken, Links oder anderer in diesen Materialien enthaltenen Elemente. Diese Publikation wird ohne jegliche Gewähr, weder ausdrücklich noch stillschweigend, bereitgestellt. Dies gilt u. a., aber nicht ausschließlich, hinsichtlich der Gewährleistung der Marktgängigkeit und der Eignung für einen bestimmten Zweck sowie für die Gewährleistung der Nichtverletzung geltenden Rechts.

SAP übernimmt keine Haftung für Schäden jeglicher Art, einschließlich und ohne Einschränkung für direkte, spezielle, indirekte oder Folgeschäden im Zusammenhang mit der Verwendung dieser Unterlagen. Diese Einschränkung gilt nicht bei Vorsatz oder grober Fahrlässigkeit.

Die gesetzliche Haftung bei Personenschäden oder die Produkthaftung bleibt unberührt. Die Informationen, auf die Sie möglicherweise über die in diesem Material enthaltenen Hotlinks zugreifen, unterliegen nicht dem Einfluss von SAP, und SAP unterstützt nicht die Nutzung von Internetseiten Dritter durch Sie und gibt keinerlei Gewährleistungen oder Zusagen über Internetseiten Dritter ab.

Alle Rechte vorbehalten.