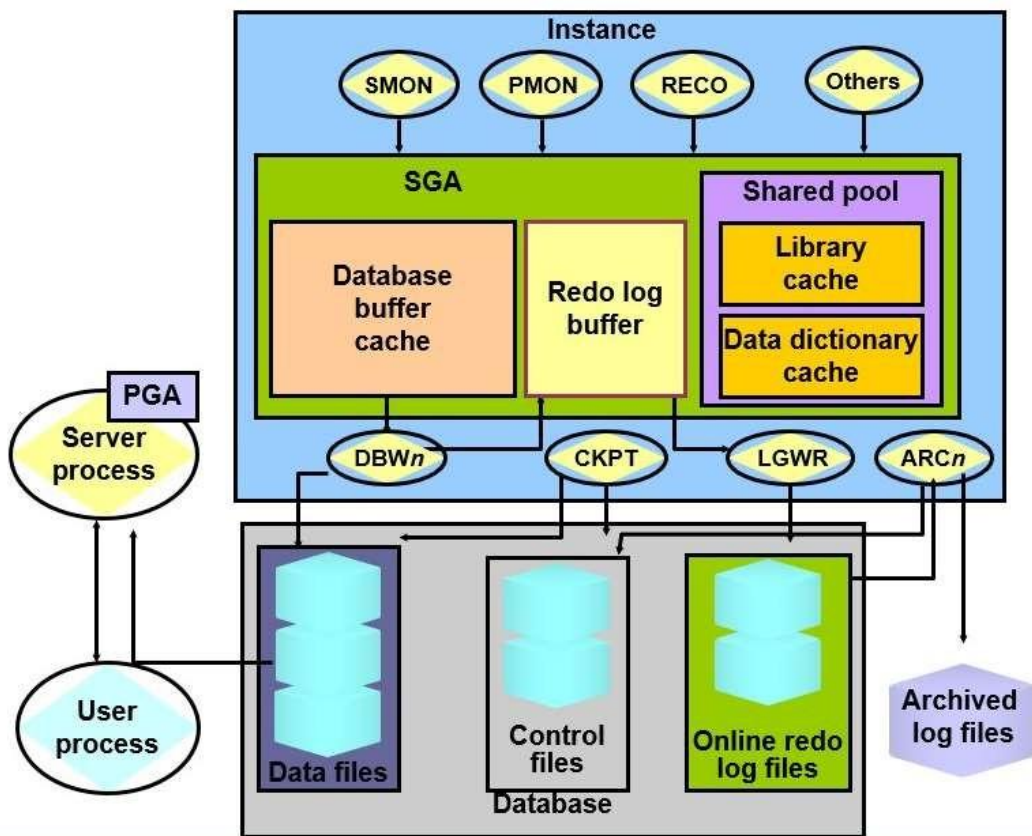


Oracle Database Architecture

Oracle Database Architecture : Overview



Oracle Database Architecture

An Oracle database consists of an instance and its associated databases.

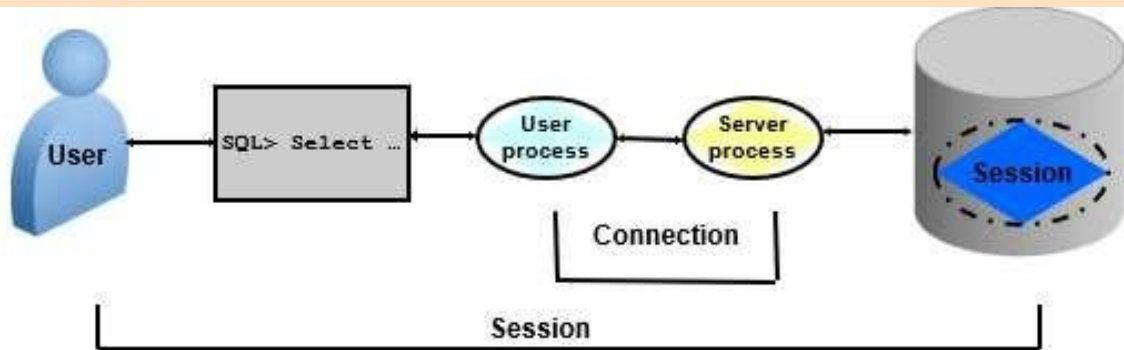
The instance consists of memory structures and background processes.

Every time an instance is started, a shared memory area called the System Global Area (SGA) is allocated and the background processes are started.

The database consists of both physical structures and logical structures.

Because the physical and logical structures are separate, the physical storage of data can be managed without affecting access to logical storage structures.

Connecting to the Database



Connections and sessions are closely related to user processes but are very different in meaning.

A connection is a communication pathway between a user process and an Oracle Database instance.

A communication pathway is established using available interposes communication mechanisms or network software.

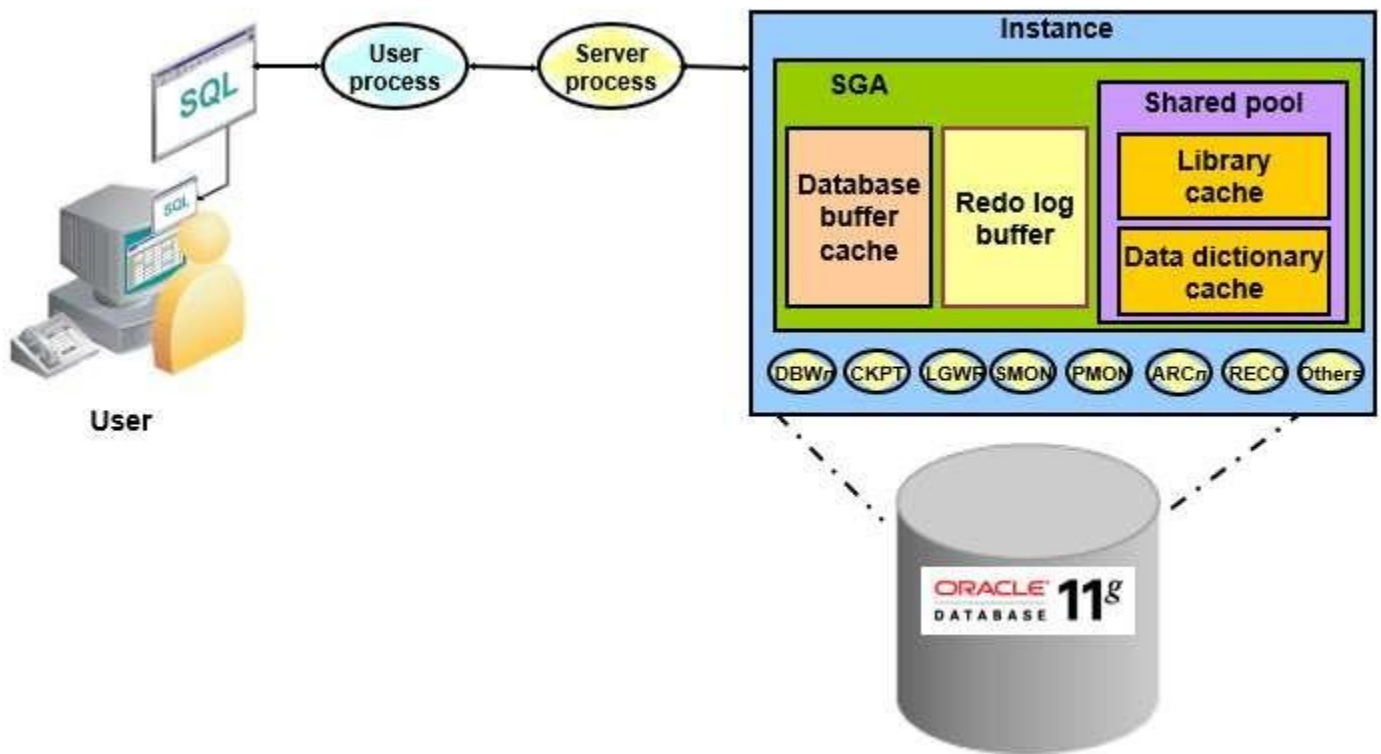
A session represents the state of a current user login to the database instance.

For example, when a user starts SQL*Plus, the user must provide a valid username and password, and then a session is established for that user. A session lasts from the time a user connects until the user disconnects or exits the database application.

In the case of a dedicated connection, the session is serviced by a permanent dedicated process. The session is serviced by an available server process selected from a pool, either by the middle tier or by Oracle shared server architecture.

Multiple sessions can be created and exist concurrently for a single Oracle database user using the same username. For example, a user with the username/password of HR/HR can connect to the same Oracle Database instance several times.

Interacting with an Oracle Database

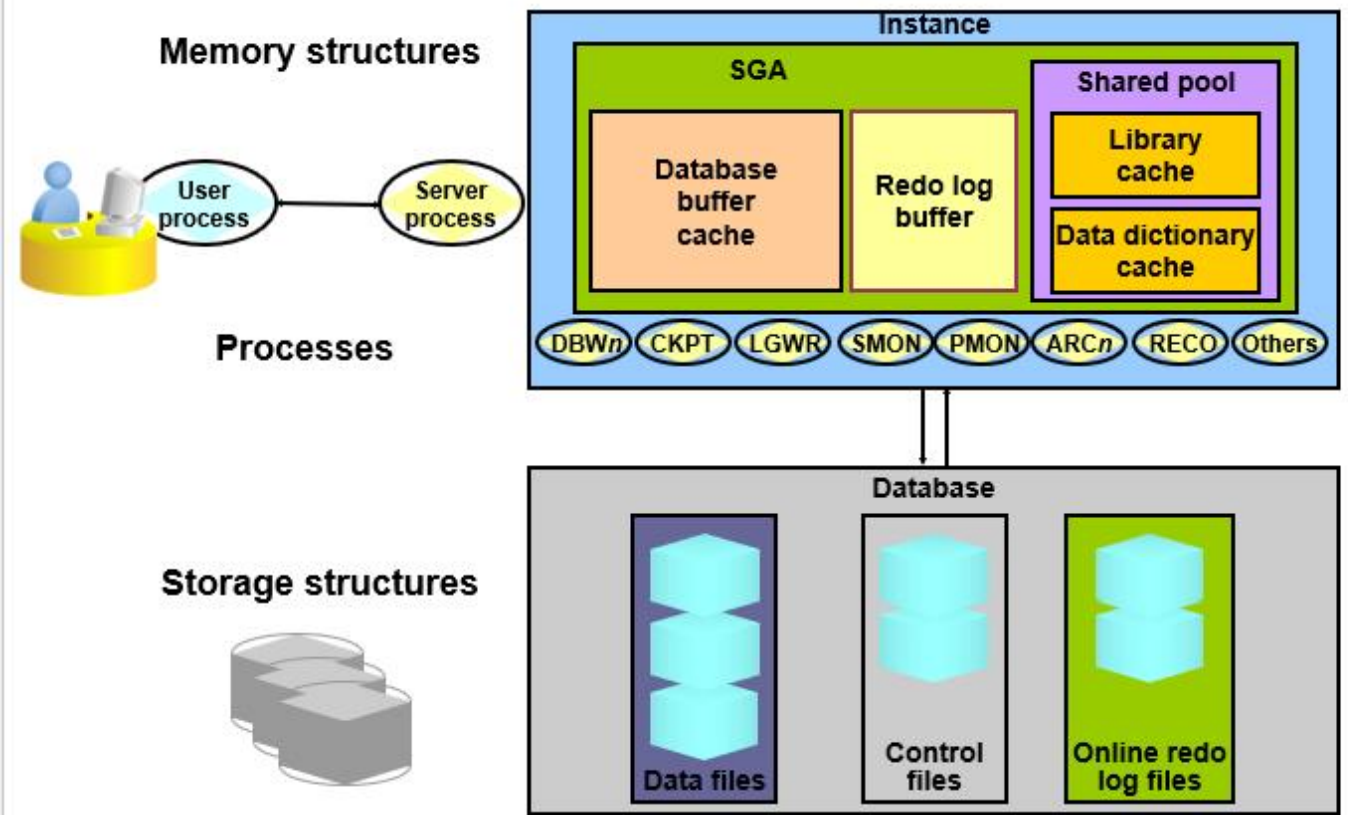


It illustrates an Oracle database configuration in which the user and associated server process are on separate computers, connected through a network.

1. An instance has started on a node where Oracle Database is installed, often called the host or database server.
2. A user starts an application spawning a user process. The application attempts to establish a connection to the server. (The connection may be local, client/server, or a three-tier connection from a middle tier.)
3. The server runs a listener that has the appropriate Oracle Net Services handler. The server detects the connection request from the application and creates a dedicated server process on behalf of the user process.

What is Listener : The Oracle Listener is a server-side process that facilitates communication between client applications and the Oracle Database. It is part of Oracle Net Services and acts as the primary gateway for all incoming connection requests to the database.

Oracle Database Server Structures

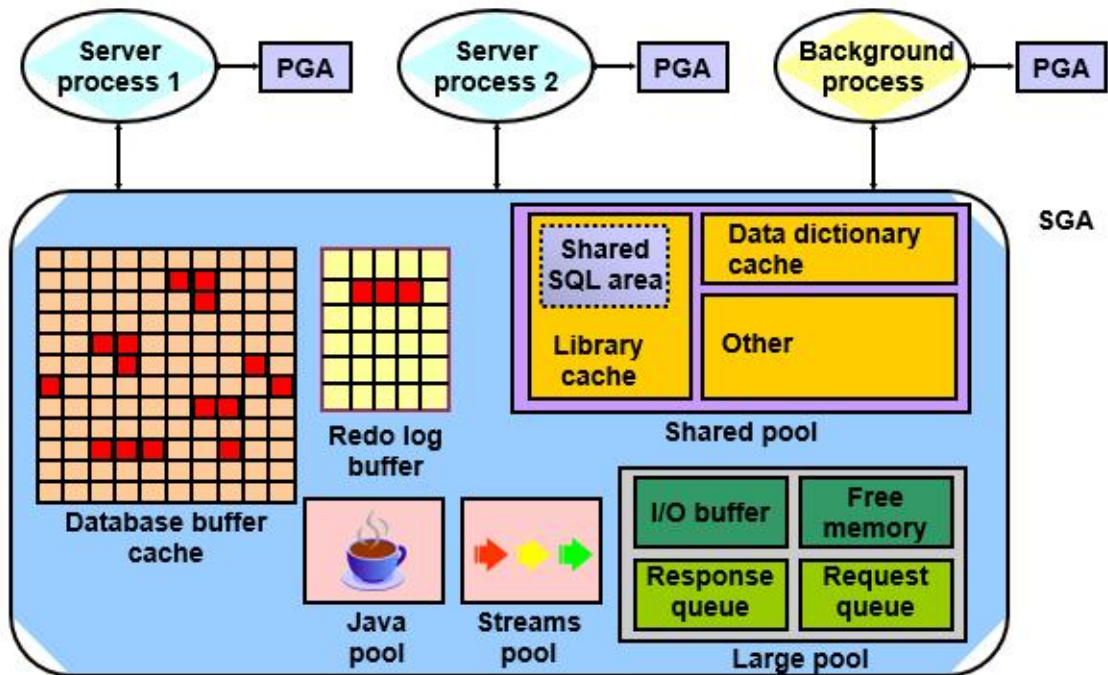


After starting an instance, the Oracle software associates the instance with a specific database. This is called *mounting the database*. The database is then ready to be opened, which makes it accessible to authorized users. Multiple instances can execute concurrently on the same computer, each accessing its own physical database.

You can look at the Oracle Database architecture as various interrelated structural components.

An Oracle instance uses memory structures and processes to manage and access the database. All memory structures exist in the main memory of the computers that constitute the database server. Processes are jobs that work in the memory of these computers. A process is defined as a “thread of control” or a mechanism in an operating system that can run a series of steps.

Oracle Database Memory Structures



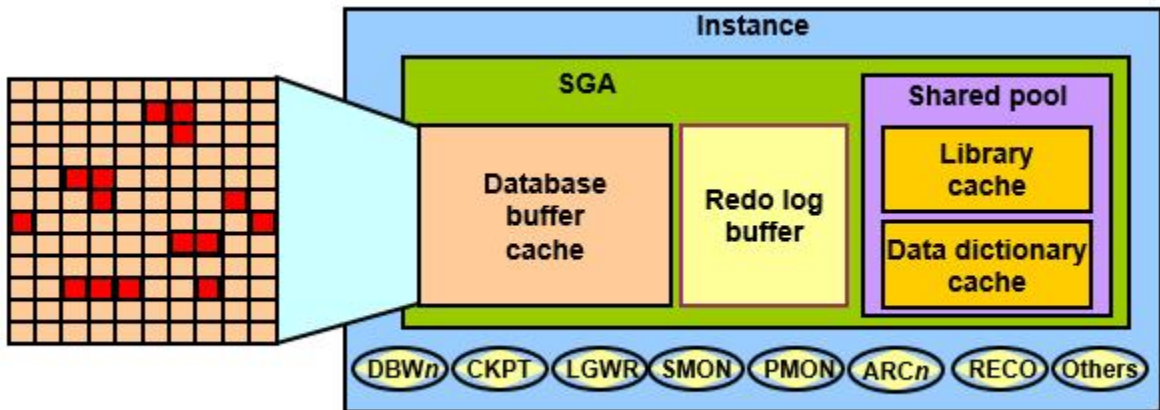
Oracle Database creates and uses memory structures for various purposes. For example, memory stores program code being run, data that is shared among users, and private data areas for each connected user.

Two basic memory structures are associated with an instance:

System Global Area (SGA): Group of shared memory structures, known as SGA components, that contain data and control information for one Oracle Database instance. The SGA is shared by all server and background processes. Examples of data stored in the SGA include cached data blocks and shared SQL areas.

Program Global Areas (PGA): Memory regions that contain data and control information for a server or background process. A PGA is nonshared memory created by Oracle Database when a server or background process is started. Access to the PGA is exclusive to the server process. Each server process and background process has its own PGA.

Database Buffer Cache



The database buffer cache is the portion of the SGA that holds copies of data blocks that are read from data files. All users who are concurrently connected to the instance share access to the database buffer cache.

The first time an Oracle Database user process requires a particular piece of data, it searches for the data in the database buffer cache.

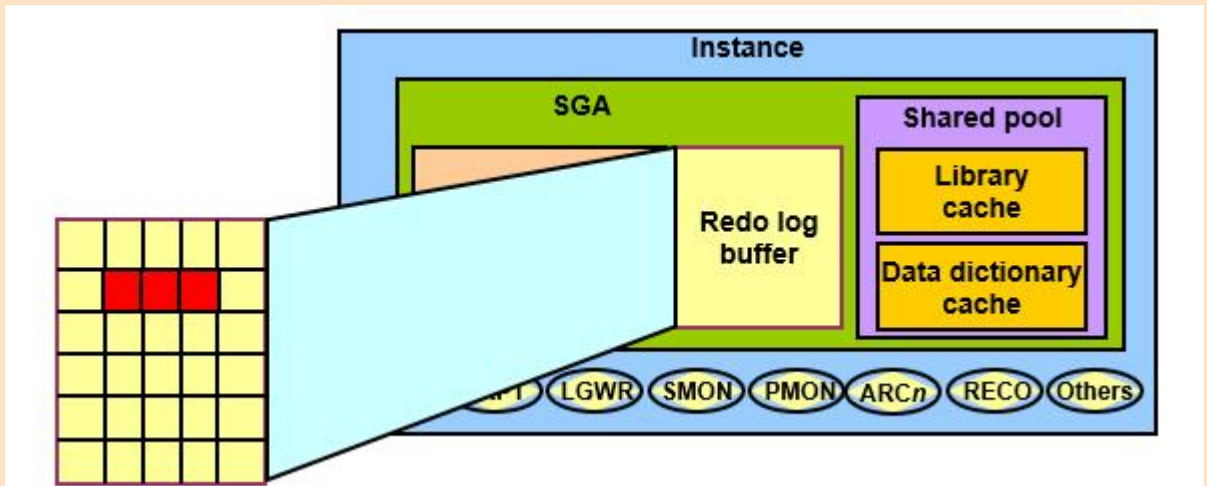
If the process finds the data already in the cache (a cache hit), it can read the data directly from memory.

If the process cannot find the data in the cache (a cache miss), it must copy the data block from a data file on disk into a buffer in the cache before accessing the data.

Accessing data through a cache hit is faster than data access through a cache miss.

The buffers in the cache are managed by a complex algorithm that uses a combination of least recently used (LRU) lists and touch count.

Redo Log Buffer



Redo Log Buffer

The redo log buffer is a circular buffer in the SGA that holds information about changes made to the database. This information is stored in redo entries.

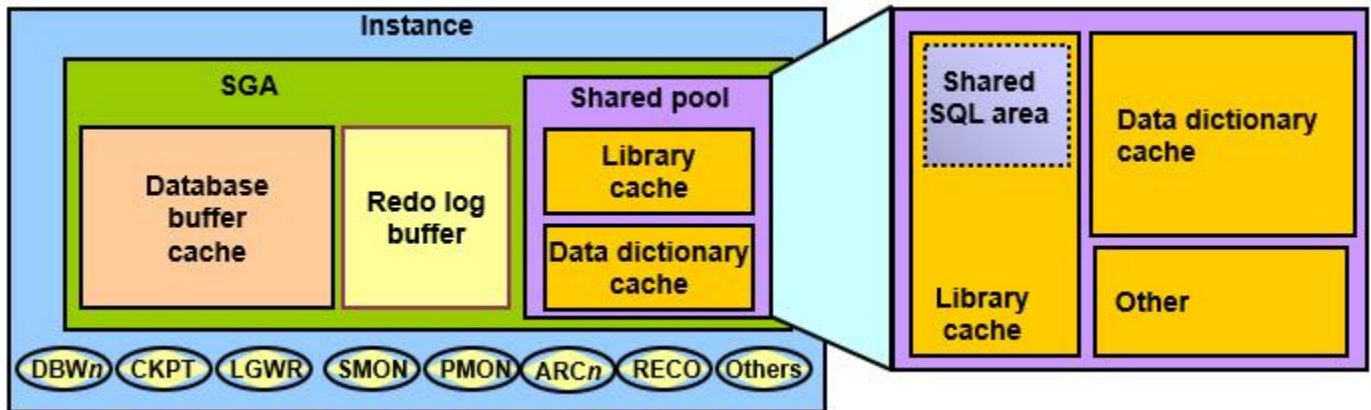
Redo entries contain the information necessary to reconstruct (or redo) changes that are made to the database by DML, DDL, or internal operations.

Redo entries are used for database recovery if necessary.

Redo entries are copied by Oracle Database processes from the user's memory space to the redo log buffer in the SGA.

The redo entries take up continuous, sequential space in the buffer. The LGWR background process writes the redo log buffer to the active redo log file (or group of files) on disk.

Shared Pool



Shared Pool

The shared pool portion of the SGA contains the library cache, the data dictionary cache, the SQL query result cache, the PL/SQL function result cache, buffers for parallel execution messages, and control structures.

The *data dictionary* is a collection of database tables and views containing reference information about the database, its structures, and its users.

Oracle Database accesses the data dictionary frequently during SQL statement parsing. This access is essential to the continuing operation of Oracle Database.

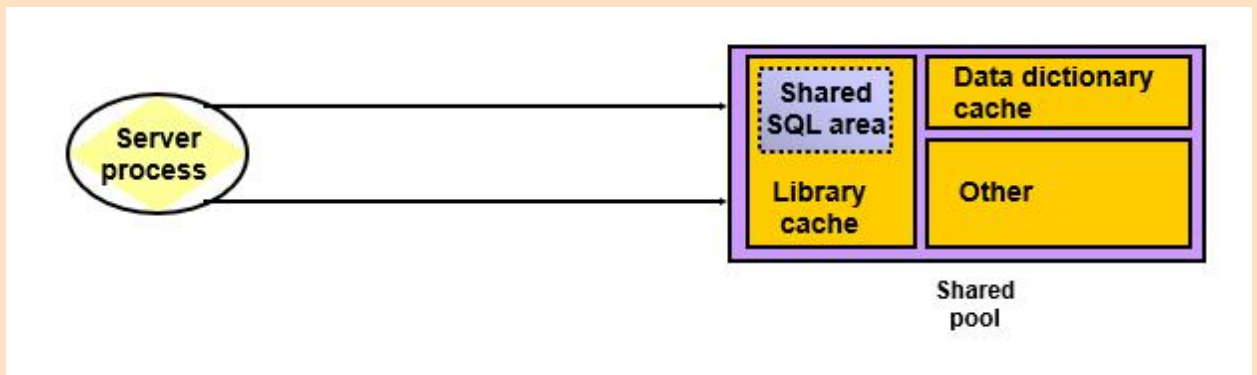
The data dictionary is accessed so often by Oracle Database that two special locations in memory are designated to hold dictionary data.

One area is called the *data dictionary cache*, also known as the row cache because it holds data as rows instead of buffers (which hold entire blocks of data). The other area in memory to hold dictionary data is the *library cache*.

All Oracle Database user processes share these two caches for access to data dictionary information.

Oracle Database represents each SQL statement that it runs with a shared SQL area (as well as a private SQL area kept in the PGA). Oracle Database recognizes when two users are executing the same SQL statement and reuses the shared SQL area for those users.

Allocation and Reuse of Memory in the Shared Pool

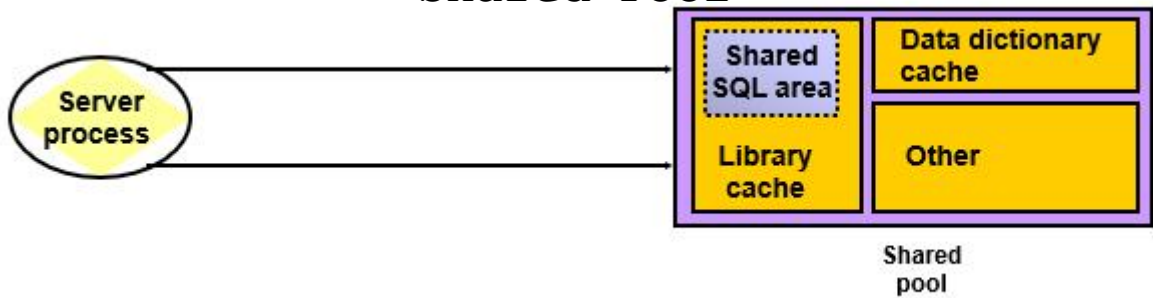


Allocation and Reuse of Memory in the Shared Pool

In general, any item (shared SQL area or dictionary row) in the shared pool remains until it is flushed according to a modified LRU (least recently used) algorithm. The memory for items that are not being used regularly is freed if space is required for new items that must be given some space in the shared pool. A modified LRU algorithm allows shared pool items that are used by many sessions to remain in memory as long as they are useful, even if the process that originally created the item terminates. As a result, the overhead and processing of SQL statements associated with a multiuser Oracle Database system are minimized. When a SQL statement is submitted to Oracle Database for execution, the following memory allocation steps are automatically performed:

1. Oracle Database checks the shared pool to see if a shared SQL area already exists for an identical statement. If so, that shared SQL area is used for the execution of the subsequent new instances of the statement. If there is no shared SQL area for a statement, Oracle Database allocates a new shared SQL area in the shared pool. In either case, the user's private SQL area is associated with the shared SQL area that contains the statement.

Allocation and Reuse of Memory in the Shared Pool



2. Oracle Database allocates a private SQL area on behalf of the session. The location of the private SQL area depends on the type of connection established for the session.

Note: A shared SQL area can be flushed from the shared pool even if the shared SQL area corresponds to an open cursor that has not been used for some time. If the open cursor is subsequently used to run its statement, Oracle Database reparses the statement and a new shared SQL area is allocated in the shared pool.

Oracle Database also flushes a shared SQL area from the shared pool in these circumstances:

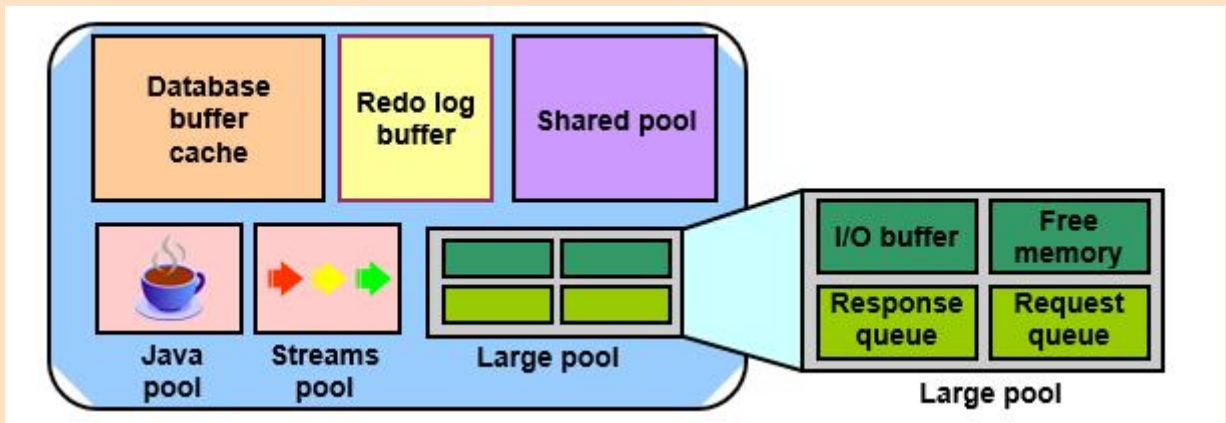
When the DBMS_STATS package is used to update or delete the statistics of a table, cluster, or index, all shared SQL areas that contain statements referencing the analyzed schema object are flushed from the shared pool. The next time a flushed statement is run, the statement is parsed in a new shared SQL area to reflect the new statistics for the schema object.

If a schema object is referenced in a SQL statement and that object is later modified in any way, the shared SQL area is invalidated (marked invalid) and the statement must be reparsed the next time it is run.

If you change a database's global database name, all information is flushed from the shared pool.

The administrator can manually flush all information in the shared pool to assess the performance (with respect to the shared pool, not the data buffer cache) that can be expected after instance startup without shutting down the current instance. The ALTER SYSTEM FLUSH SHARED_POOL statement is used to do this.

Large Pool



Large Pool

The database administrator can configure an optional memory area called the *large pool* to provide large memory allocations for:

Session memory for the shared server and the Oracle XA interface (used where transactions interact with more than one database):

- I/O server processes

- Oracle Database backup and restore operations

By allocating session memory from the large pool for shared server, Oracle XA, or parallel query buffers, Oracle Database can use the shared pool primarily for caching shared SQL and avoid the performance overhead that is caused by shrinking the shared SQL cache.

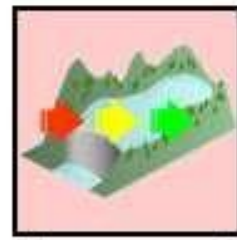
In addition, the memory for Oracle Database backup and restore operations, for I/O server processes, and for parallel buffers is allocated in buffers of a few hundred kilobytes. The large pool is better able to satisfy such large memory requests than the shared pool.

The large pool does not have an LRU list. It is different from reserved space in the shared pool, which uses the same LRU list as other memory allocated from the shared pool.

Java Pool and Streams Pool



Java pool



Streams pool

Java Pool and Streams Pool

Java pool memory is used in server memory for all session-specific Java code and data in the JVM. Java pool memory is used in different ways, depending on the mode in which Oracle Database is running.

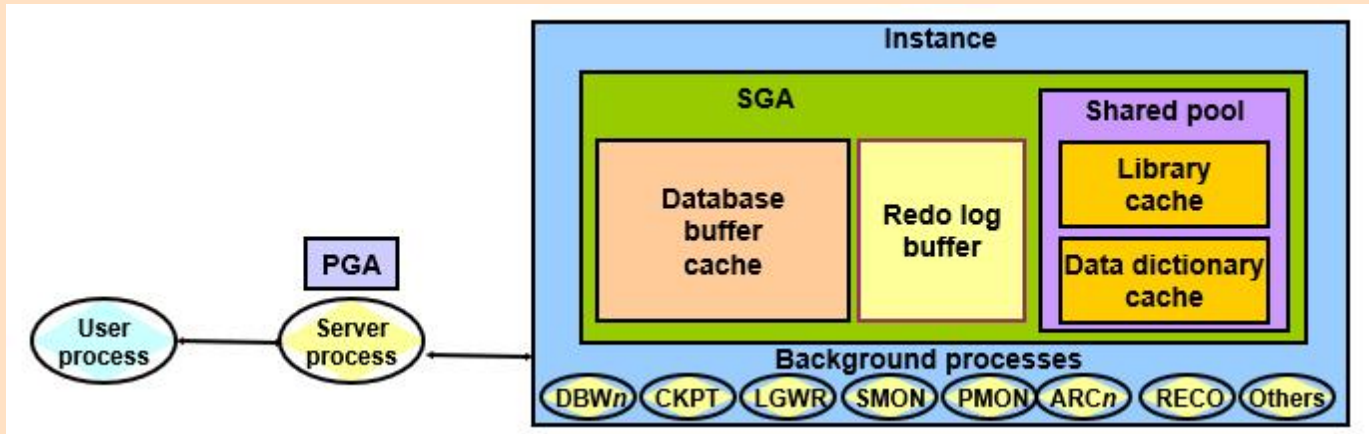
The Java Pool Advisor statistics provide information about library cache memory used for Java and predict how changes in the size of the Java pool can affect the parse rate. The Java Pool Advisor is internally turned on when `statistics_level` is set to `TYPICAL` or higher. These statistics reset when the advisor is turned off.

The Streams pool is used exclusively by Oracle Streams. The Streams pool stores buffered queue messages, and it provides memory for Oracle Streams capture processes and apply processes.

Unless you specifically configure it, the size of the Streams pool starts at zero. The pool size grows dynamically as needed when Oracle Streams is used.

Note: A detailed discussion of Java programming and Oracle Streams is beyond the scope of this class

Process Architecture



Process Architecture

The processes in an Oracle Database system can be divided into two major groups:

- User processes that run the application or Oracle tool code
- Oracle Database processes that run the Oracle database server code (including server processes and background processes)

When a user runs an application program or an Oracle tool such as SQL*Plus, Oracle Database creates a *user process* to run the user's application. Oracle Database also creates a *server process* to execute the commands issued by the user process.

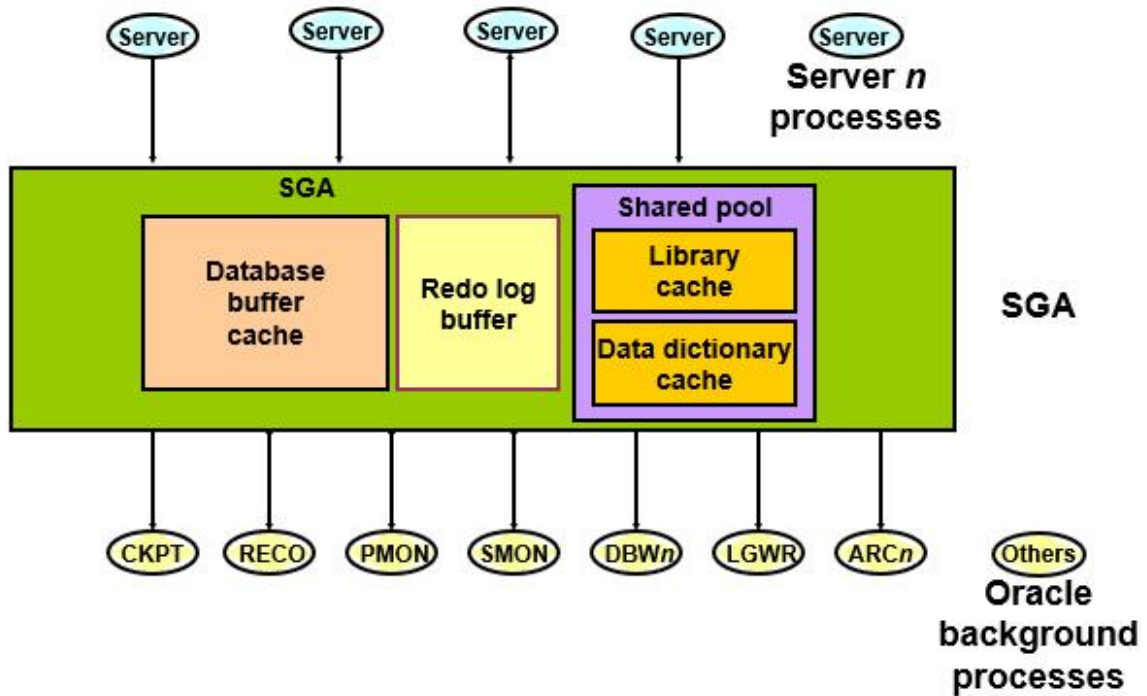
In addition, the Oracle server also has a set of *background processes* for an instance that interact with each other and with the operating system to manage the memory structures, asynchronously perform I/O to write data to disk, and perform other required tasks.

The process structure varies for different Oracle Database configurations, depending on the operating system and the choice of Oracle Database options. The code for connected users can be configured as a dedicated server or a shared server.

Dedicated server: For each user, the database application is run by a user process that is served by a dedicated server process that executes Oracle database server code.

Shared server: Eliminates the need for a dedicated server process for each connection. A dispatcher directs multiple incoming network session requests to a pool of shared server processes. A shared server process serves any client request.

Process Structures



Process Structures

Server Processes

Oracle Database creates server processes to handle the requests of user processes connected to the instance. In some situations, when the application and Oracle Database operate on the same computer, it is possible to combine the user process and corresponding server process into a single process to reduce system overhead. However, when the application and Oracle Database operate on different computers, a user process always communicates with Oracle Database through a separate server process.

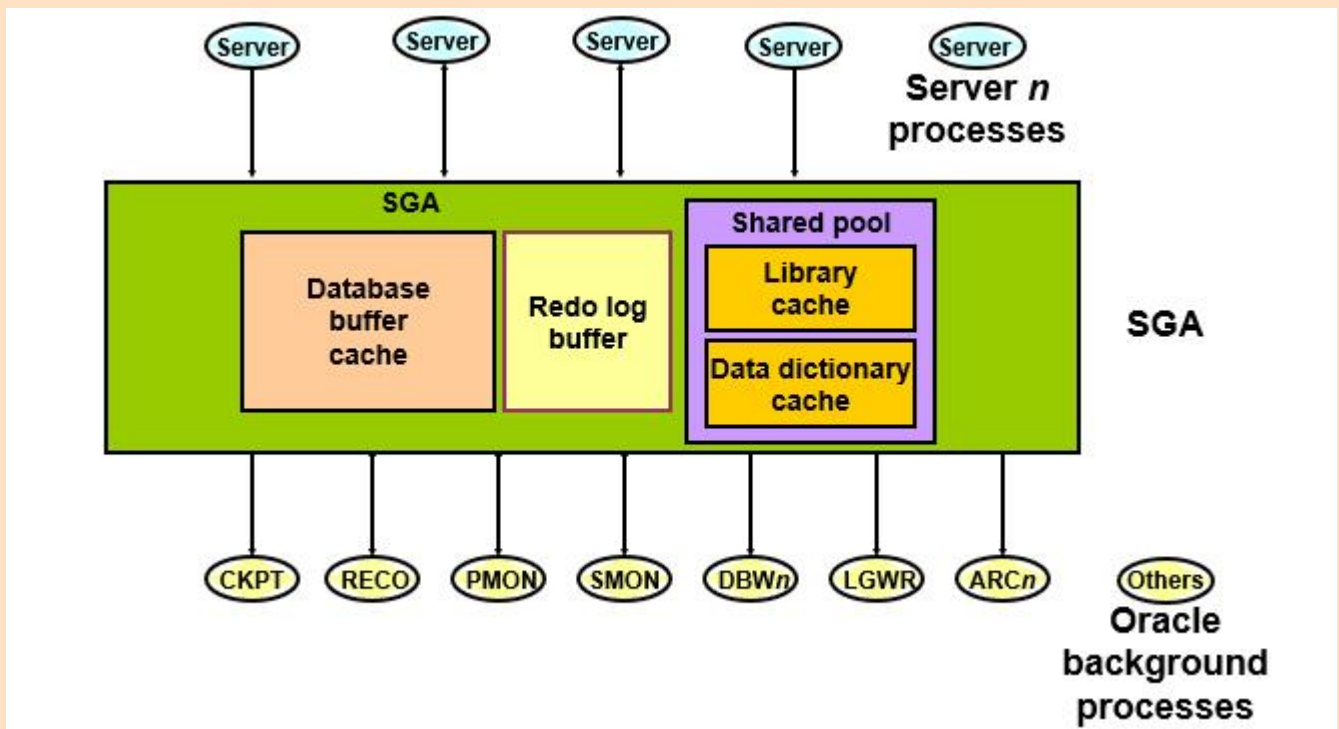
Server processes created on behalf of each user's application can perform one or more of the following:

- Parse and run SQL statements issued through the application
- Read necessary data blocks from data files on disk into the shared database buffers of the SGA (if the blocks are not already present in the SGA)
- Return results in such a way that the application can process the information

Background Processes

To maximize performance and accommodate many users, a multiprocess Oracle Database system uses some additional Oracle Database processes called *background processes*. An Oracle Database instance can have many background processes.

Process Structures



Process Structures (continued)

The background processes commonly seen in non-RAC, non-ASM environments can include the following:

- Database writer process (DBWn)

- Log writer process (LGWR)

- Checkpoint process (CKPT)

- System Monitor process (SMON)

- Process monitor process (PMON)

- Recoverer process (RECO)

- Job queue processes

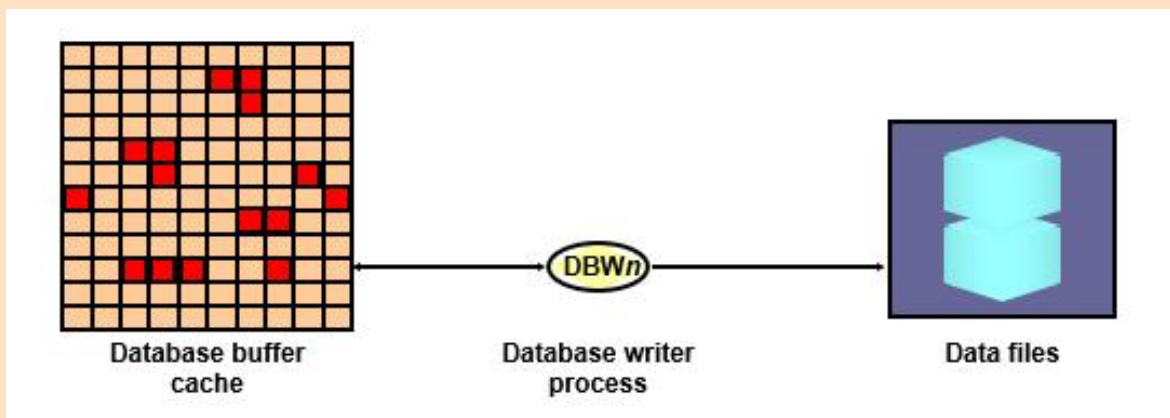
- Archiver processes (ARCn)

- Queue monitor processes (QMn)

Other background processes may be found in more advanced configurations such as RAC. See the V\$BGPROCESS view for more information on the background processes.

On many operating systems, background processes are created automatically when an instance is started.

Database Writer Process (DBWn)

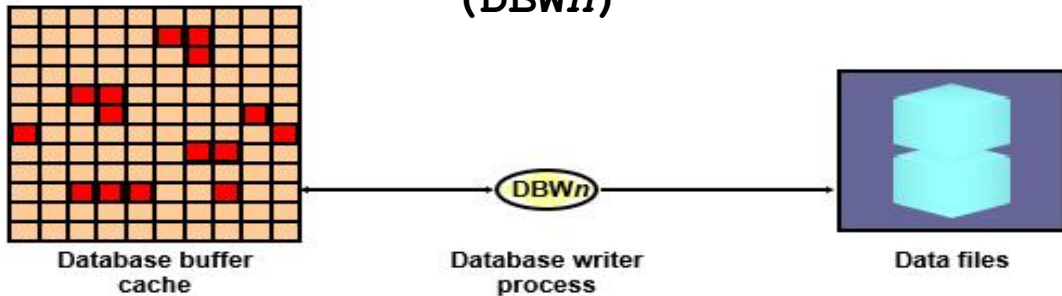


Database Writer Process (DBWn)

The Database Writer process (DBWn) writes the contents of buffers to data files. The DBWn processes are responsible for writing modified (dirty) buffers in the database buffer cache to disk. Although one Database Writer process (DBW0) is adequate for most systems, you can configure additional processes (DBW1 through DBW9 and DBWa through DBWj) to improve write performance if your system modifies data heavily. These additional DBWn processes are not useful on uniprocessor systems. When a buffer in the database buffer cache is modified, it is marked dirty and is added to the LRUW (LRU write) list of dirty buffers that is kept in SCN order. This order therefore matches the order of redo that is written to the redo logs for these changed buffers. When the number of available buffers in the buffer cache falls below an internal threshold (to the extent that server processes find it difficult to obtain available buffers), DBWn writes dirty buffers to the data files in the order that they were modified by following the order of the LRUW list.

Database Writer Process

(DBWn)



Database Writer Process (DBWn) (continued)

The SGA contains a memory structure that has the redo byte address (RBA) of the position in the redo stream where recovery should begin in the case of an instance failure. This structure acts as a pointer into the redo and is written to the control file by the CKPT process once every three seconds. Because the DBWn writes dirty buffers in SCN order, and because the redo is in SCN order, every time DBWn writes dirty buffers from the LRUW list, it also advances the pointer held in the SGA memory structure so that instance recovery (if required) begins reading the redo from approximately the correct location and avoids unnecessary I/O. This is known as *incremental checkpointing*.

Note: There are other cases when DBWn may write (for example, when tablespaces are made read-only or are placed offline). In such cases, no incremental checkpoint occurs because dirty buffers belonging only to the corresponding data files are written to the database unrelated to the SCN order.

The LRU algorithm keeps more frequently accessed blocks in the buffer cache so that, when a buffer is written to disk, it is unlikely to contain data that will soon be useful.

The DB_WRITER_PROCESSES initialization parameter specifies the number of DBWn processes. The maximum number of DBWn processes is 20. If it is not specified by the user during startup, Oracle Database determines how to set DB_WRITER_PROCESSES based on the number of CPUs and processor groups.

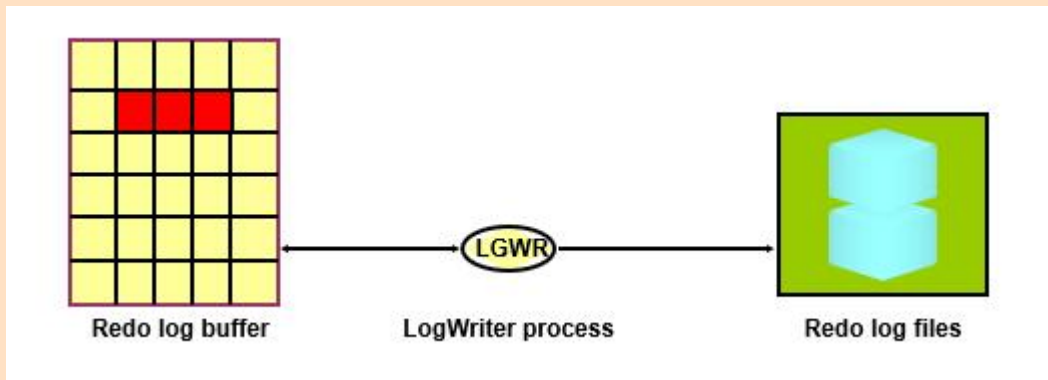
The DBWn process writes dirty buffers to disk under the following conditions:

- When a server process cannot find a clean reusable buffer after scanning a threshold number of buffers, it signals DBWn to write. DBWn writes dirty buffers to disk asynchronously while performing other processing.

- DBWn periodically writes buffers to advance the checkpoint, which is the position in the redo thread (log) from which instance recovery begins. This log position is determined by the oldest dirty buffer in the buffer cache.

In all cases, DBWn performs batched (multiblock) writes to improve efficiency. The number of blocks written in a multiblock write varies by operating system.

LogWriter Process (LGWR)



LogWriter Process (LGWR)

The LogWriter process (LGWR) is responsible for redo log buffer management by writing the redo log buffer entries to a redo log file on disk. LGWR writes all redo entries that have been copied into the buffer since the last time it wrote.

The redo log buffer is a circular buffer. When LGWR writes redo entries from the redo log buffer to a redo log file, server processes can then copy new entries over the entries in the redo log buffer that have been written to disk. LGWR normally writes fast enough to ensure that space is always available in the buffer for new entries, even when access to the redo log is heavy. LGWR writes one contiguous portion of the buffer to disk.

LGWR writes:

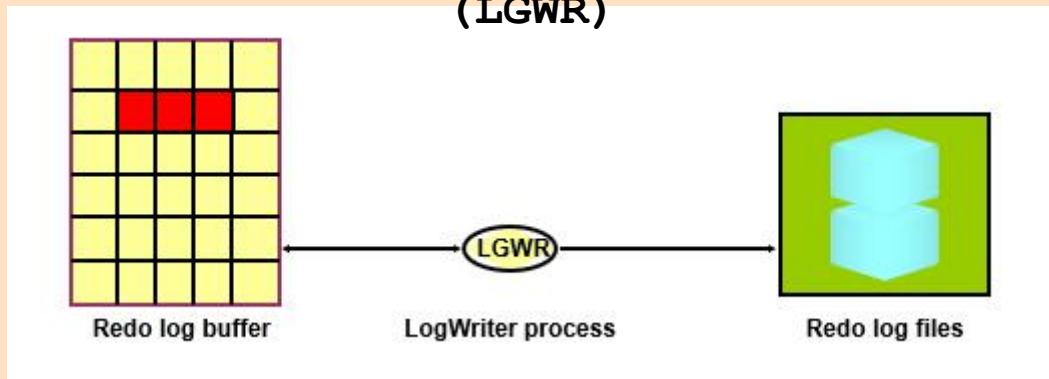
- When a user pLogWriter process commits a transaction

- When the redo log buffer is one-third full

- Before a DBWn process writes modified buffers to disk (if necessary)

LogWriter Process

(LGWR)



Log Writer Process (LGWR) (continued)

Before DBWn can write a modified buffer, all redo records that are associated with the changes to the buffer must be written to disk (the write-ahead protocol). If DBWn finds that some redo records have not been written, it signals LGWR to write the redo records to disk and waits for LGWR to complete writing the redo log buffer before it can write out the data buffers. LGWR writes to the current log group. If one of the files in the group is damaged or unavailable, LGWR continues writing to other files in the group and logs an error in the LGWR trace file and in the system alert log. If all files in a group are damaged, or if the group is unavailable because it has not been archived, LGWR cannot continue to function.

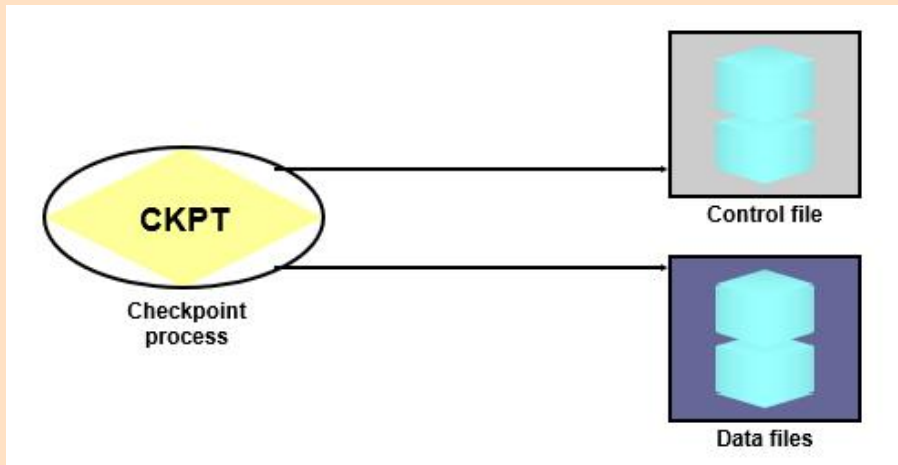
When a user issues a COMMIT statement, LGWR puts a commit record in the redo log buffer and writes it to disk immediately, along with the transaction's redo entries. The corresponding changes to data blocks are deferred until it is more efficient to write them. This is called a *fast commit mechanism*. The atomic write of the redo entry containing the transaction's commit record is the single event that determines whether the transaction has committed. Oracle Database returns a success code to the committing transaction, although the data buffers have not yet been written to disk.

If more buffer space is needed, LGWR sometimes writes redo log entries before a transaction is committed. These entries become permanent only if the transaction is later committed. When a user commits a transaction, the transaction is assigned a system change number (SCN), which Oracle Database records along with the transaction's redo entries in the redo log. SCNs are recorded in the redo log so that recovery operations can be synchronized in Real Application Clusters and distributed databases.

In times of high activity, LGWR can write to the redo log file by using group commits. For example, suppose that a user commits a transaction. LGWR must write the transaction's redo entries to disk. As this happens, other users issue COMMIT statements. However, LGWR cannot write to the redo log file to commit these transactions until it has completed its previous write operation.

After the first transaction's entries are written to the redo log file, the entire list of redo entries of waiting transactions (not yet committed) can be written to disk in one operation, requiring less I/O than do transaction entries handled individually. Therefore, Oracle Database minimizes disk I/O and maximizes performance of LGWR. If requests to commit continue at a high rate, every write (by LGWR) from the redo log buffer can contain multiple commit records.

Checkpoint Process (CKPT)



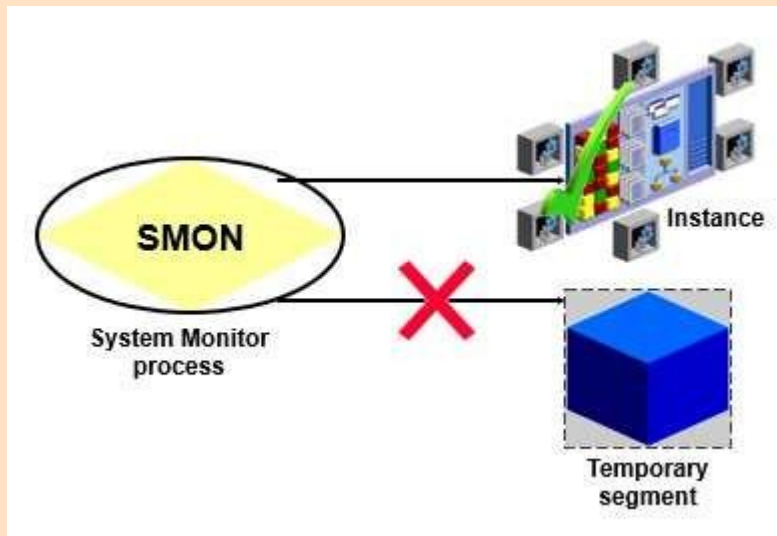
Checkpoint Process (CKPT)

A *checkpoint* is a data structure that defines a system change number (SCN) in the redo thread of a database. Checkpoints are recorded in the control file and in each data file header. They are a crucial element of recovery.

When a checkpoint occurs, Oracle Database must update the headers of all data files to record the details of the checkpoint. This is done by the CKPT process. The CKPT process does not write blocks to disk; DBWR always performs that work. The SCNs recorded in the file headers guarantee that all changes made to database blocks prior to that SCN have been written to disk.

The statistic DBWR checkpoints displayed by the `SYSTEM_STATISTICS` monitor in Oracle Enterprise Manager indicate the number of checkpoint requests that have completed.

System Monitor Process (SMON)

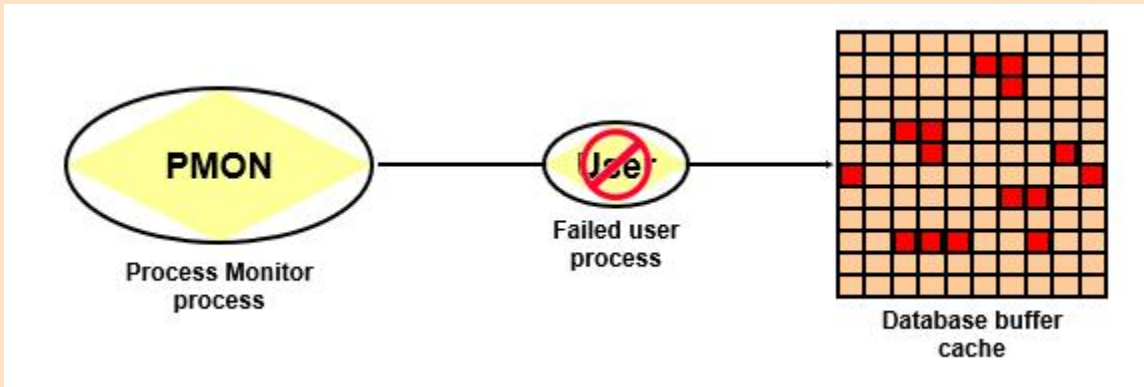


System Monitor Process (SMON)

The System Monitor process (SMON) performs recovery at instance startup if necessary. SMON is also responsible for cleaning up temporary segments that are no longer in use. If any terminated transactions were skipped during instance recovery because of file-read or offline errors, SMON recovers them when the tablespace or file is brought back online.

SMON checks regularly to see whether the process is needed. Other processes can call SMON if they detect a need for it.

Process Monitor Process (PMON)

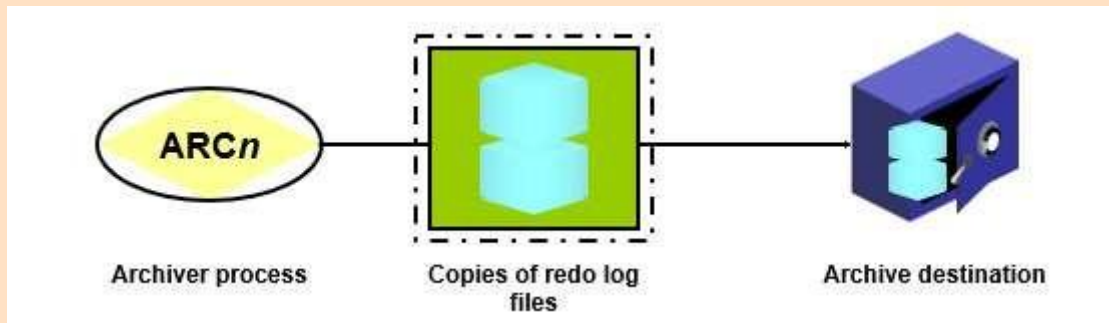


Process Monitor Process (PMON)

The Process Monitor process (PMON) performs process recovery when a user process fails. PMON is responsible for cleaning up the database buffer cache and freeing resources that the user process was using. For example, it resets the status of the active transaction table, releases locks, and removes the process ID from the list of active processes.

PMON periodically checks the status of dispatcher and server processes, and restarts any that have stopped running (but not any that Oracle Database has terminated intentionally). PMON also registers information about the instance and dispatcher processes with the network listener. Like SMON, PMON checks regularly to see whether it is needed; it can be called if another process detects the need for it.

Archiver Processes (ARCn)



Archiver Processes (ARCn)

The archiver processes (ARCn) copy redo log files to a designated storage device after a log switch has occurred. ARCn processes are present only when the database is in ARCHIVELOG mode and automatic archiving is enabled. If you anticipate a heavy workload for archiving (such as during bulk loading of data), you can increase the maximum number of archiver processes with the LOG_ARCHIVE_MAX_PROCESSES initialization parameter. The ALTER SYSTEM statement can change the value of this parameter dynamically to increase or decrease the number of ARCn processes.

OTHER PROCESSES

Other Processes

There are several other background processes that might be running. These can include the following:

The Manageability Monitor process (MMON) performs various manageability-related background tasks, for example:

- Issuing alerts whenever a given metrics violates its threshold value
- Taking snapshots by spawning additional process (MMON slaves)
- Capturing statistics value for SQL objects that have been recently modified

The Lightweight Manageability Monitor process (MMNL) performs frequent tasks related to lightweight manageability, such as session history capture and metrics computation.

The Memory Manager process (MMAN) is used for internal database tasks. It manages automatic memory management processing to help allocate memory where it is needed dynamically in an effort to avoid out-of-memory conditions or poor buffer cache performance.

The Rebalance process (RBAL) coordinates rebalance activity for disk groups in an Automatic Storage Management instance. It performs a global open on Automatic Storage Management disks. ORBn performs the actual rebalance data extent movements in an Automatic Storage Management instance. There can be many of these at a time, named ORB0, ORB1, and so on.

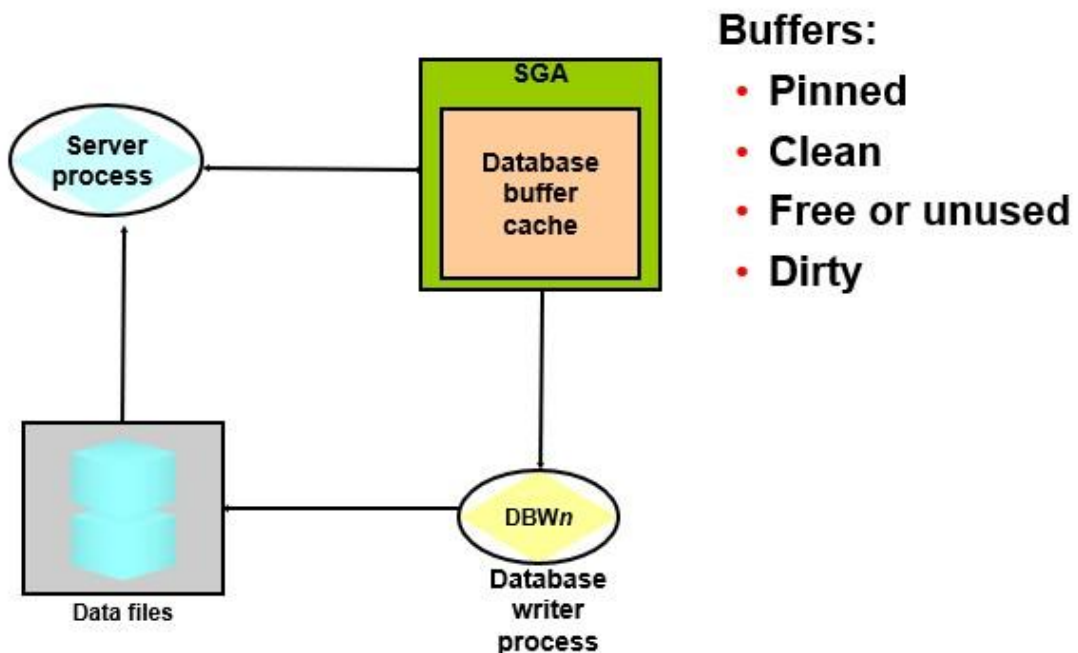
The Automatic Storage Management process (ASMB) is present in a database instance using an Automatic Storage Management disk group. It communicates with the Automatic Storage Management instance.

Job queue processes are used for batch processing. They run user jobs and can be viewed as scheduler services that are used to schedule jobs as PL/SQL statements or procedures on an Oracle Database instance.

The coordinator process, named CJQ0, periodically selects jobs that need to be run from the system JOB\$ table. The CJQ0 process dynamically spawns job queue slave processes (J000 through J999) to run the jobs. The job queue process runs one of the jobs that was selected by the CJQ0 process for execution. The processes run one job at a time.

The Queue Monitor process (QMNx) is an optional background process for Oracle Streams Advanced Queuing, which monitors the message queues. You can configure up to 10 queue monitor processes.

Server Process and Database Buffer Cache



Server Process and Database Buffer Cache

When a query is processed, the Oracle server process looks in the database buffer cache for images of any blocks that it needs. If the block image is not found in the database buffer cache, the server process reads the block from the data file and places a copy in the database buffer cache. Because subsequent requests for the same block may find the block in memory, the requests may not require physical reads. Buffers in the buffer cache can be in one of the following four states:

Pinned: Multiple sessions are kept from writing to the same block at the same time. Other sessions wait to access the block.

Clean: The buffer is now unpinned and is a candidate for immediate aging out, if the current contents (data block) are not referenced again. Either the contents are in sync with the block contents stored on the disk, or the buffer contains a consistent read (CR) snapshot of a block.

Free or unused: The buffer is empty because the instance has just started. This state is very similar to the clean state, except that the buffer has not been used.

Dirty: The buffer is no longer pinned but the contents (data block) have changed and must be flushed to the disk by DBWn before it can be aged out.