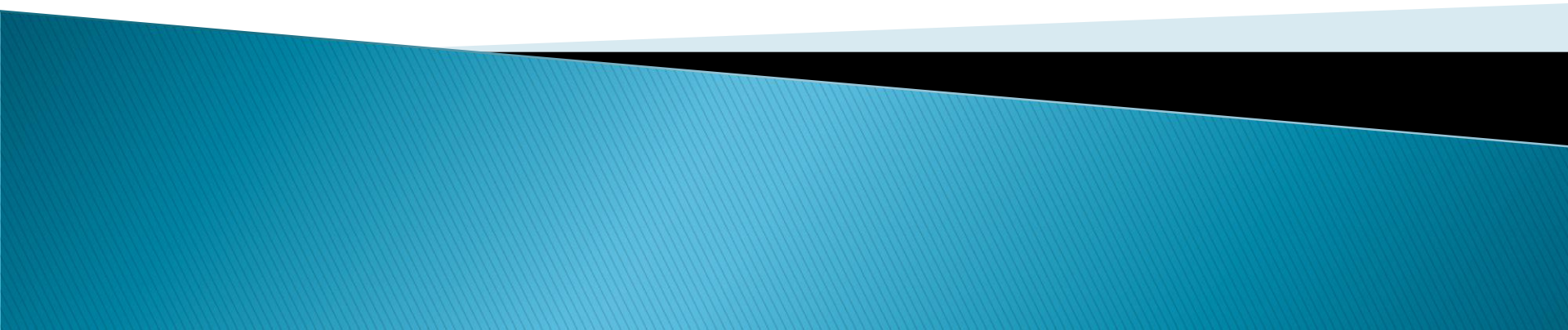


Recovery Mangement– Part 2

Checkpoints



Checkpoints

▶ Problems in recovery procedure as discussed earlier :

1. searching the entire log is time-consuming
2. we might unnecessarily redo transactions which have already
3. output their updates to the database.

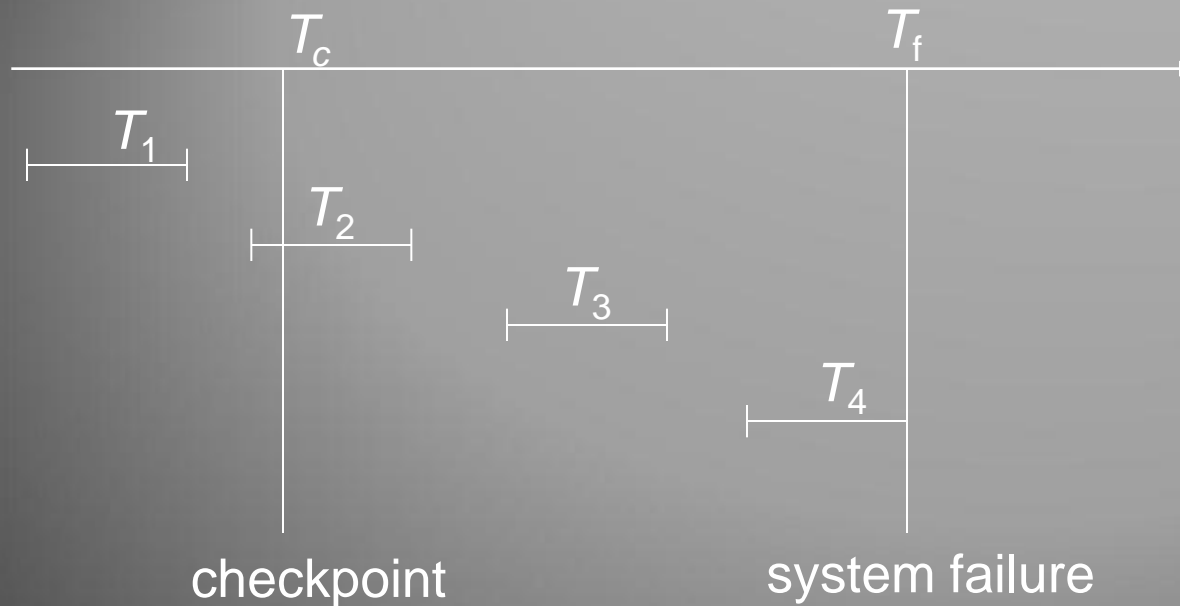
▶ Streamline recovery procedure by periodically performing **checkpointing**

1. Output all log records currently residing in main memory onto stable storage.
2. Output all modified buffer blocks to the disk.
3. Write a log record < **checkpoint** > onto stable storage.

Checkpoints (Cont.)

- ▶ During recovery we need to consider only the most recent transaction T_i that started before the checkpoint, and transactions that started after T_i .
 1. Scan backwards from end of log to find the most recent **<checkpoint>** record
 2. Continue scanning backwards till a record **< T_i start>** is found.
 3. Need only consider the part of log following above start record. Earlier part of log can be ignored during recovery, and can be erased whenever desired.
 4. For all transactions (starting from T_i or later) with no **< T_i commit>**, execute **undo(T_i)**. (Done only in case of immediate modification.)
 5. Scanning forward in the log, for all transactions starting from T_i or later with a **< T_i commit>**, execute **redo(T_i)**.

Example of Checkpoints

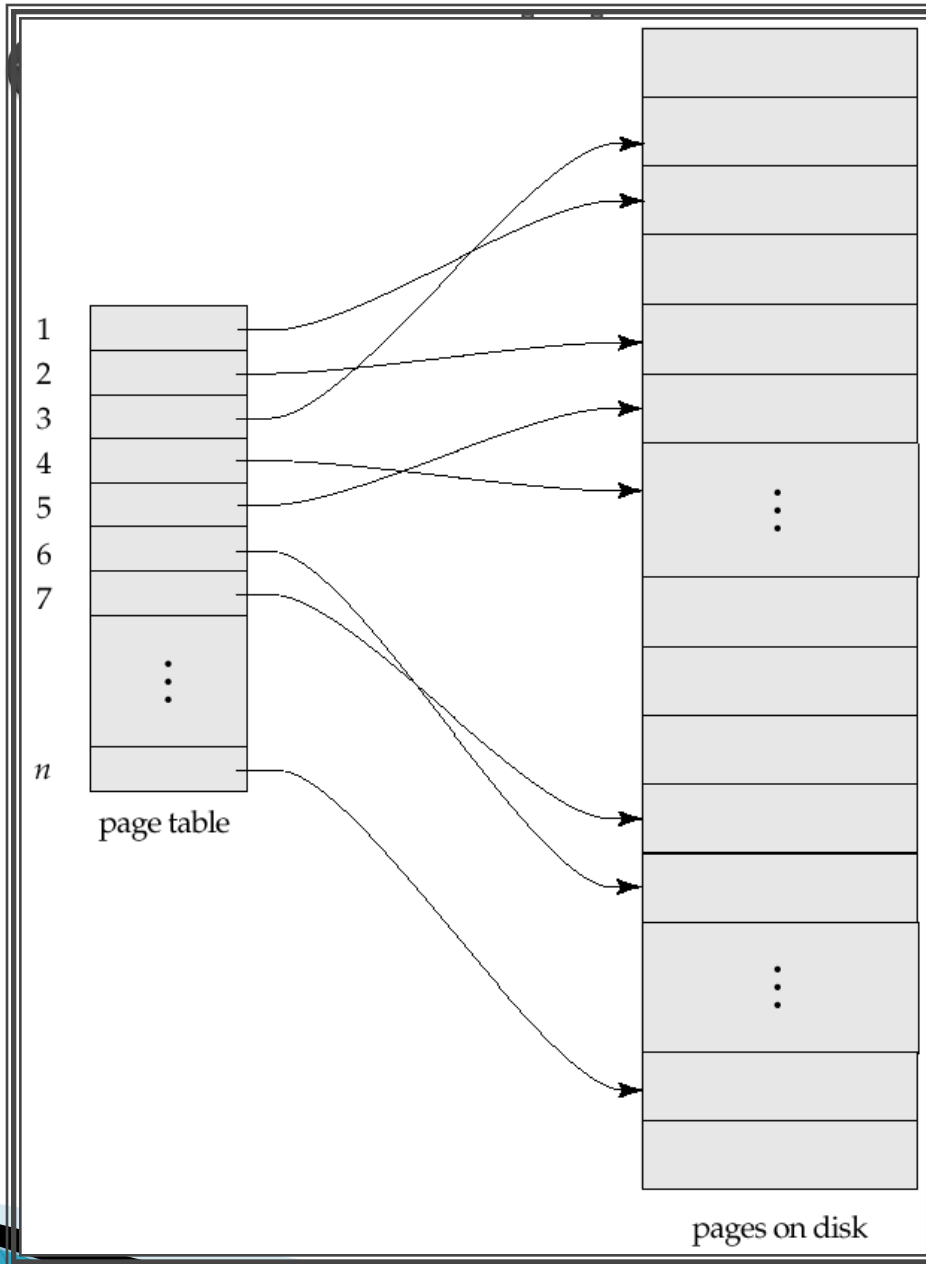


- ▶ T_1 can be ignored (updates already output to disk due to checkpoint)
- ▶ T_2 and T_3 redone.
- ▶ T_4 undone

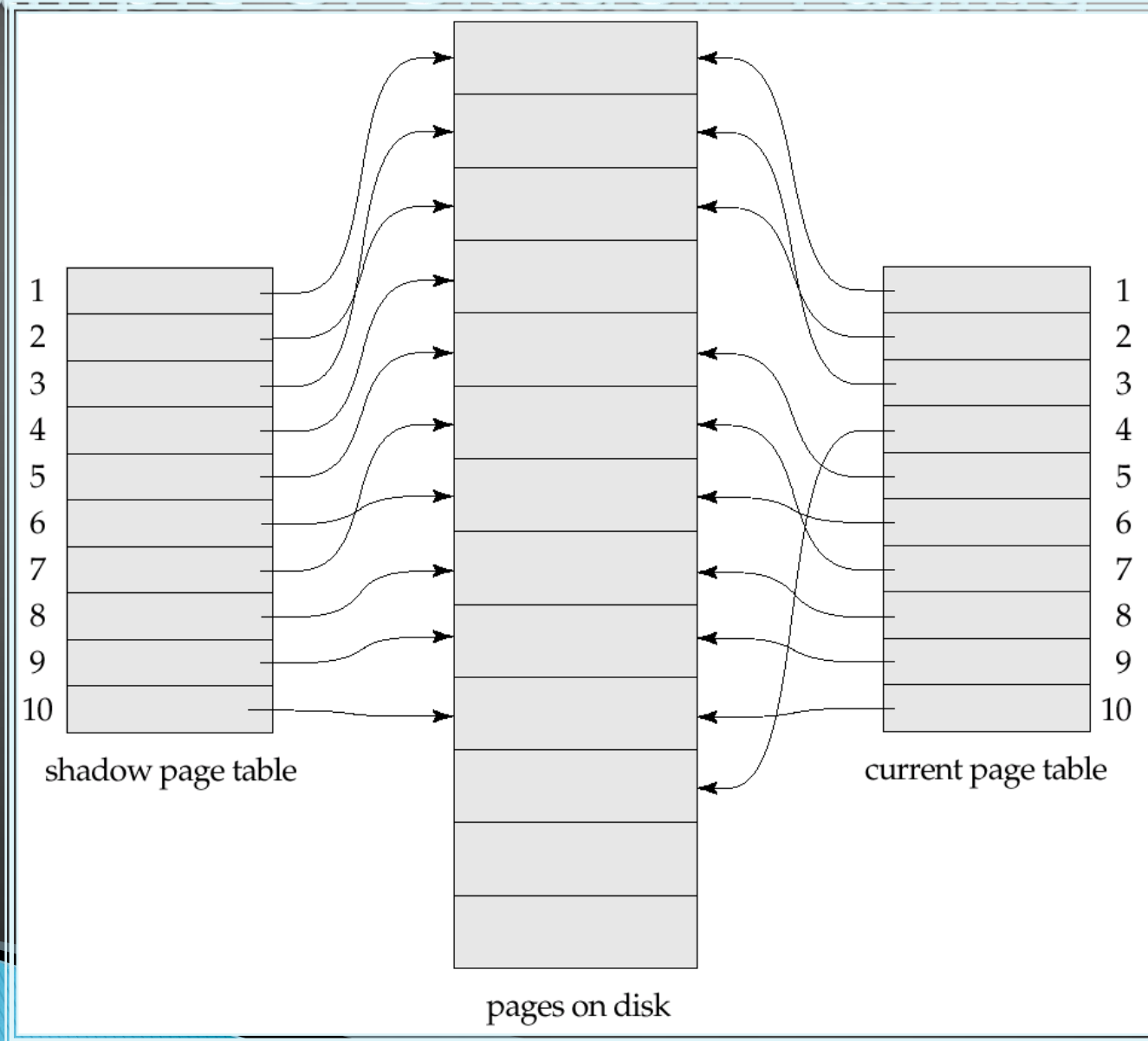
Shadow Paging

- ▶ **Shadow paging** is an alternative to log-based recovery; this scheme is useful if transactions execute serially
- ▶ Idea: maintain *two* page tables during the lifetime of a transaction –the **current page table**, and the **shadow page table**
- ▶ Store the shadow page table in nonvolatile storage, such that state of the database prior to transaction execution may be recovered.
 - Shadow page table is never modified during execution
- ▶ To start with, both the page tables are identical. Only current page table is used for data item accesses during execution of the transaction.
- ▶ Whenever any page is about to be written for the first time
 - A copy of this page is made onto an unused page.
 - The current page table is then made to point to the copy
 - The update is performed on the copy

Sample



Example of Shadow Paging



Shadow Paging (Cont.)

- ▶ To commit a transaction:
 1. Flush all modified pages in main memory to disk
 2. Output current page table to disk
 3. Make the current page table the new shadow page table, as follows:
 - keep a pointer to the shadow page table at a fixed (known) location on disk.
 - to make the current page table the new shadow page table, simply update the pointer to point to current page table on disk
- ▶ Once pointer to shadow page table has been written, transaction is committed.
- ▶ No recovery is needed after a crash — new transactions can start right away, using the shadow page table.
- ▶ Pages not pointed to from current/shadow page table should be freed (garbage collected).

Show Paging (Cont.)

- ▶ Advantages of shadow-paging over log-based schemes
 - no overhead of writing log records
 - recovery is trivial
- ▶ Disadvantages :
 - Copying the entire page table is very expensive
 - Can be reduced by using a page table structured like a B⁺-tree
 - No need to copy entire tree, only need to copy paths in the tree that lead to updated leaf nodes
 - Commit overhead is high even with above extension
 - Need to flush every updated page, and page table
 - Data gets fragmented (related pages get separated on disk)
 - After every transaction completion, the database pages containing old versions of modified data need to be garbage collected
 - Hard to extend algorithm to allow transactions to run concurrently
 - Easier to extend log based schemes

Recovery With Concurrent Transactions

- ▶ We modify the log-based recovery schemes to allow multiple transactions to execute concurrently.
 - All transactions share a single disk buffer and a single log
 - A buffer block can have data items updated by one or more transactions
- ▶ We assume concurrency control using strict two-phase locking;
 - i.e. the updates of uncommitted transactions should not be visible to other transactions
 - Otherwise how to perform undo if T1 updates A, then T2 updates A and commits, and finally T1 has to abort?
- ▶ Logging is done as described earlier.
 - Log records of different transactions may be interspersed in the log.
- ▶ The checkpointing technique and actions taken on recovery have to be changed
 - since several transactions may be active when a checkpoint is performed.

Recovery With Concurrent Transactions (Cont.)

- ▶ Checkpoints are performed as before, except that the checkpoint log record is now of the form $\langle \text{checkpoint } L \rangle$ where L is the list of transactions active at the time of the checkpoint
 - We assume no updates are in progress while the checkpoint is carried out (will relax this later)
- ▶ When the system recovers from a crash, it first does the following:
 1. Initialize *undo-list* and *redo-list* to empty
 2. Scan the log backwards from the end, stopping when the first $\langle \text{checkpoint } L \rangle$ record is found. For each record found during the backward scan:
 - if the record is $\langle T_i \text{ commit} \rangle$, add T_i to *redo-list*
 - if the record is $\langle T_i \text{ start} \rangle$, then if T_i is not in *redo-list*, add T_i to *undo-list*
 3. For every T_i in L , if T_i is not in *redo-list*, add T_i to *undo-list*

Recovery With Concurrent Transactions (Cont.)

- ▶ At this point *undo-list* consists of incomplete transactions which must be undone, and *redo-list* consists of finished transactions that must be redone.
- ▶ Recovery now continues as follows:
 1. Scan log backwards from most recent record, stopping when $\langle T_i \text{ start} \rangle$ records have been encountered for every T_i in *undo-list*.
 - During the scan, perform **undo** for each log record that belongs to a transaction in *undo-list*.
 2. Locate the most recent $\langle \text{checkpoint } L \rangle$ record.
 3. Scan log forwards from the $\langle \text{checkpoint } L \rangle$ record till the end of the log.
 - During the scan, perform **redo** for each log record that belongs to a transaction on *redo-list*

Example of Recovery

- ▶ Go over the steps of the recovery algorithm on the following log:

< T_0 start >

< T_0 , A, 0, 10 >

< T_0 commit >

< T_1 start >

< T_1 , B, 0, 10 >

< T_2 start >

/* Scan in Step 4 stops here */

< T_2 , C, 0, 10 >

< T_2 , C, 10, 20 >

< checkpoint { T_1 , T_2 } >

< T_3 start >

< T_3 , A, 10, 20 >

< T_3 , D, 0, 10 >

< T_3 commit >

Log Record Buffering

- ▶ **Log record buffering:** log records are buffered in main memory, instead of being output directly to stable storage.
 - Log records are output to stable storage when a block of log records in the buffer is full, or a **log force** operation is executed.
- ▶ Log force is performed to commit a transaction by forcing all its log records (including the commit record) to stable storage.
- ▶ Several log records can thus be output using a single output operation, reducing the I/O cost.

Log Record Buffering (Cont.)

- ▶ The rules below must be followed if log records are buffered:
 - Log records are output to stable storage in the order in which they are created.
 - Transaction T_i enters the commit state only when the log record $\langle T_i, \text{commit} \rangle$ has been output to stable storage.
 - Before a block of data in main memory is output to the database, all log records pertaining to data in that block must have been output to stable storage.
 - This rule is called the **write-ahead logging** or **WAL** rule
 - Strictly speaking WAL only requires undo information to be output

Database Buffering

- ▶ Database maintains an in-memory buffer of data blocks
 - When a new block is needed, if buffer is full an existing block needs to be removed from buffer
 - If the block chosen for removal has been updated, it must be output to disk
- ▶ As a result of the write-ahead logging rule, if a block with uncommitted updates is output to disk, log records with undo information for the updates are output to the log on stable storage first.
- ▶ No updates should be in progress on a block when it is output to disk. Can be ensured as follows.
 - Before writing a data item, transaction acquires exclusive lock on block containing the data item
 - Lock can be released once the write is completed.
 - Such locks held for short duration are called **latches**.
 - Before a block is output to disk, the system acquires an exclusive latch on the block
 - Ensures no update can be in progress on the block

Buffer Management (Cont.)

- ▶ Database buffer can be implemented either
 - in an area of real main-memory reserved for the database, or
 - in virtual memory
- ▶ Implementing buffer in reserved main-memory has drawbacks:
 - Memory is partitioned before-hand between database buffer and applications, limiting flexibility.
 - Needs may change, and although operating system knows best how memory should be divided up at any time, it cannot change the partitioning of memory.

Buffer Management (Cont.)

- ▶ Database buffers are generally implemented in virtual memory in spite of some drawbacks:
 - When operating system needs to evict a page that has been modified, to make space for another page, the page is written to swap space on disk.
 - When database decides to write buffer page to disk, buffer page may be in swap space, and may have to be read from swap space on disk and output to the database on disk, resulting in extra I/O!
 - Known as **dual paging** problem.
 - Ideally when swapping out a database buffer page, operating system should pass control to database, which in turn outputs page to database instead of to swap space (making sure to output log records first)
 - Dual paging can thus be avoided, but common operating systems do not support such functionality.

Failure with Loss of Nonvolatile Storage

- ▶ So far we assumed no loss of non-volatile storage
- ▶ Technique similar to checkpointing used to deal with loss of non-volatile storage
 - Periodically **dump** the entire content of the database to stable storage
 - No transaction may be active during the dump procedure; a procedure similar to checkpointing must take place
 - Output all log records currently residing in main memory onto stable storage.
 - Output all buffer blocks onto the disk.
 - Copy the contents of the database to stable storage.
 - Output a record **<dump>** to log on stable storage.
 - To recover from disk failure
 - restore database from most recent dump.
 - Consult the log and redo all transactions that committed after the dump
- ▶ Can be extended to allow transactions to be active during dump; known as **fuzzy dump** or **online dump**
 - Will study fuzzy checkpointing later