

# Failure Recovery

## ▲ Failure Classification [故障分类]

### Transaction failure:

Logical errors: transaction cannot complete due to some internal error condition

System errors: the database system must terminate an active transaction due to an error condition

### System crash: a power failure or other hardware or software failure causes the system to crash.

Fail-stop assumption: non-volatile storage contents are assumed to not be corrupted by system crash

Database systems have numerous integrity checks to prevent corruption of disk data

### Disk failure: a head crash or similar disk failure destroys all or part of disk storage

Destruction is assumed to be detectable: disk drivers use checksums to detect failures

## ▲ Recovery algorithms have two parts:

Actions taken during normal transaction processing to ensure enough information exists to recover from failures [在正常事务处理期间为确保存在足够的信息以从故障中恢复而采取的操作]

Actions taken after a failure to recover the database contents to a state that ensures atomicity, consistency and durability [在将数据库内容恢复到确保原子性、一致性和持久性的状态失败后所采取的操作]

## ▲ Storage Structure [存储结构]

### Volatile storage [易失性存储器]:

does not survive system crashes

examples: main memory, cache memory

### Non-volatile storage [非易失性存储]:

survives system crashes

examples: disk, tape, flash memory, non-volatile (battery backed up) RAM

but may still fail, losing data

### Stable storage [稳态存储]:

a mythical form of storage that survives all failures

usually approximated by maintaining multiple copies on distinct nonvolatile media

## ▲ Data Access [数据存取]

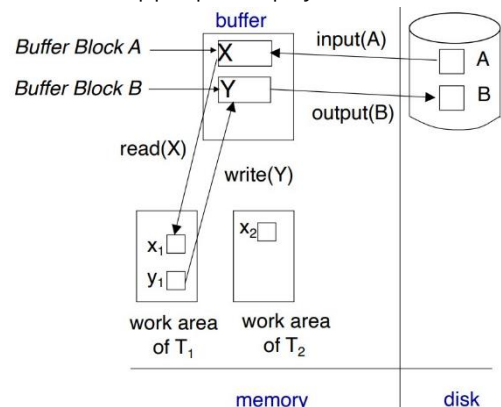
Physical blocks are those blocks residing on the disk.

Buffer blocks are the blocks residing temporarily in main memory.

Block movements between disk and main memory are initiated through the following two operations:

`input(B)` transfers the physical block B to main memory.

`output(B)` transfers the buffer block B to the disk, and replaces the appropriate physical block there.



We assume, for simplicity, that each data item fits in, and is stored inside, a single block.

Each transaction  $T_i$  has its private work-area in which local copies of all data items accessed and updated by it are kept.

$T_i$ 's local copy of a data item  $X$  is called  $x_i$ .

Transferring data items between system buffer blocks and its private work-area done by:

$read(X)$  assigns the value of data item  $\{X\}$  to the local variable  $x_i$ .

$write(X)$  assigns the value of local variable  $x_i$  to data item  $\{X\}$  in the buffer block.

Note:  $output(BX)$  need not immediately follow  $write(X)$ . System can perform the output operation when it deems fit.

Transactions

Must perform  $read(X)$  before accessing  $X$  for the first time (subsequent reads can be from local copy)

$write(X)$  can be executed at any time before the transaction commits

### ▲ Recovery and Atomicity [恢复与单元性]

To ensure atomicity despite of failures, we first output information describing the modifications (e.g. logs) to stable storage without modifying the database itself.

A log is kept on stable storage.

The log is a sequence of log records, and maintains a record of update activities on the database.

When transaction  $T_i$  starts, it registers itself by writing a  $\langle T_i \text{ start} \rangle$  log record

Before  $T_i$  executes  $write(X)$ , a log record  $\langle T_i, X, V1, V2 \rangle$  is written, where  $V1$  is the value of  $X$  before the write (the old value), and  $V2$  is the value to be written to  $X$  (the new value).

When  $T_i$  finishes its last statement, the log record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$  is written.

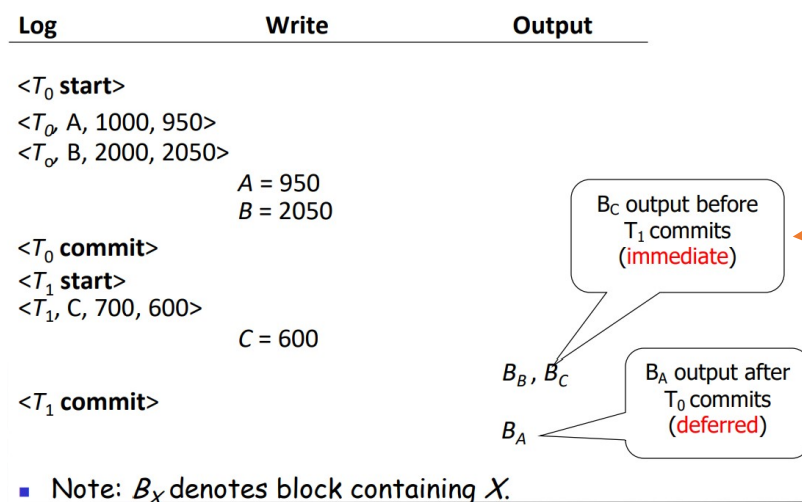
### Two approaches using logs:

Deferred [延迟] database modification

performs updates to buffer/disk only at the time of transaction commit.

Immediate [立即] database modification

allows updates of an uncommitted transaction to be made to the buffer, or the disk itself, before the transaction commits



### ▲ Transaction Commit

A transaction is said to have committed when its commit log record is output to stable storage.

all previous log records of the transaction must have been output already.

Writes performed by a transaction may still be in the buffer when the transaction commits, and may be output later.

### ▲ Concurrency Control and Log-Based Recovery [并发控制&基于日志恢复]

With concurrent transactions, 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 that if a transaction  $T_i$  has modified an item, no other transaction can modify the same item until  $T_i$  has committed or aborted, i.e. using the strict two-phase locking protocol. (不可同时操作同一 item)

Log records of different transactions may be interspersed (散布) in the log.

**Undo [撤销]** of a log record  $\langle T_i, X, V_1, V_2 \rangle$  writes the old value  $V_1$  to  $X$

Undo ( $T_i$ ) restores the values of all data items updated by  $T_i$  to their old values, going backwards from the last log record for  $T_i$ .

(1) Each time a data item  $X$  is restored to its old value  $V$ , a special log record  $\langle T_i, X, V \rangle$  is written out.

(2) When undo of a transaction is complete, a log record  $\langle T_i \text{ abort} \rangle$  is written out.

**Redo [重做]** of a log record  $\langle T_i, X, V_1, V_2 \rangle$  writes the new value  $V_2$  to  $X$  (again)

Redo ( $T_i$ ) sets the value of all data items updated by  $T_i$  to the new values, going forward from the first log record for  $T_i$ .

(1) No additional logging is done in this case

### ▲ Recovering after failure:

Transaction  $T_i$  needs to be **undone** if the log

contains the record  $\langle T_i \text{ start} \rangle$ ,

but does not contain either the record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$ .

Transaction  $T_i$  needs to be **redone** if the log

contains the records  $\langle T_i \text{ start} \rangle$

and contains the record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$

**Tips:** 如果事务  $T_i$  在更早的时候被 undone,  $\langle T_i \text{ commit} \rangle$  记录写到日志中, 然后发生了失败, 那么从失败中恢复  $T_i$  时会选择 redone, 这样的 redone 会重做所有原始的操作, 包括恢复旧值的步骤, 称为重复历史。看起来很浪费, 但极大地简化了恢复算法。

**Example:**

$\langle T_0 \text{ start} \rangle$	$\langle T_0 \text{ start} \rangle$	$\langle T_0 \text{ start} \rangle$
$\langle T_0, A, 1000, 950 \rangle$	$\langle T_0, A, 1000, 950 \rangle$	$\langle T_0, A, 1000, 950 \rangle$
$\langle T_0, B, 2000, 2050 \rangle$	$\langle T_0, B, 2000, 2050 \rangle$	$\langle T_0, B, 2000, 2050 \rangle$
	$\langle T_0 \text{ commit} \rangle$	$\langle T_0 \text{ commit} \rangle$
	$\langle T_1 \text{ start} \rangle$	$\langle T_1 \text{ start} \rangle$
	$\langle T_1, C, 700, 600 \rangle$	$\langle T_1, C, 700, 600 \rangle$
		$\langle T_1 \text{ commit} \rangle$
(a)	(b)	(c)

**Recovery actions** in each case above are:

(a) undo ( $T_0$ ): B is restored to 2000 and A to 1000, and log records  $\langle T_0, B, 2000 \rangle$ ,  $\langle T_0, A, 1000 \rangle$ ,  $\langle T_0, \text{abort} \rangle$  are written out

(b) redo ( $T_0$ ) and undo ( $T_1$ ): A and B are set to 950 and 2050 and C is restored to 700. Log records  $\langle T_1, C, 700 \rangle$ ,  $\langle T_1, \text{abort} \rangle$  are written out.

(c) redo ( $T_0$ ) and redo ( $T_1$ ): A and B are set to 950 and 2050 respectively. Then C is set to 600.

### ▲ Checkpoints [检验点]

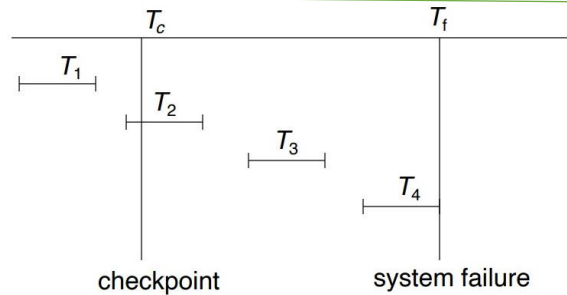
Redoing/undoing all transactions recorded in the log can be very slow (processing the entire log is time-consuming if the system has run for a long time). So, we might unnecessarily redo transactions which have

already output their updates to the database long time ago. (节省时间不必做重复的工作)

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  $\langle \text{checkpoint } L \rangle$  onto stable storage where  $L$  is a list of all transactions which are active at the time of checkpointing.
- (4) All updates are stopped while doing checkpointing.

example



回溯查询 checkpoint

In this example,  $T_1$  can be ignored (updates already output to disk due to checkpoint),  $T_2$ ,  $T_3$  redone and  $T_4$  undone.

### ▲ Recovery Algorithm

**Logging** (during normal operation)

- $\langle T_i \text{ start} \rangle$  at transaction start
- $\langle T_i, X_j, V_1, V_2 \rangle$  for each update, and
- $\langle T_i \text{ commit} \rangle$  at the end of transaction

**Transaction rollback** (during normal operation)

Let  $T_i$  be the transaction to be rolled back

Scan log backwards from the end, and for each log record of  $T_i$  of the form  $\langle T_i, X_j, V_1, V_2 \rangle$

perform the undo by writing  $V_1$  to  $X_j$ ,

write a log record  $\langle T_i, X_j, V_1 \rangle$  (such log records are called compensation (补偿) log records)

Once the record  $\langle T_i \text{ start} \rangle$  is found stop the scan and write the log record  $\langle T_i \text{ abort} \rangle$

### Recovery from failure have two phases:

[1] Redo phase: replay updates of all transactions, whether they committed, aborted, or are incomplete, at and after checkpoint

[2] Undo phase: undo all incomplete transactions

### Redo phase:

Find last  $\langle \text{checkpoint } L \rangle$  record, and set the undo-list to  $L$  (undo-list =  $L$ )

Scan forward from above  $\langle \text{checkpoint } L \rangle$  record:

Whenever a record  $\langle T_i, X_j, V_1, V_2 \rangle$  is found, redo it by writing  $V_2$  to  $X_j$

Whenever a log record  $\langle T_i \text{ start} \rangle$  is found, add  $T_i$  to undo-list

Whenever a log record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$  is found, remove  $T_i$  from undo-list

### Undo phase:

Scan log backwards from the failure point:

Whenever a log record  $\langle T_i, X_j, V_1, V_2 \rangle$  is found where  $T_i$  is in undo-list (same as transaction rollback):  
perform undo by writing  $V_1$  to  $X_j$ .

write a log record  $\langle T_i, X_j, V1 \rangle$

Whenever a log record  $\langle Ti \text{ start} \rangle$  is found where  $T_i$  is in undo-list:

Write a log record  $\langle Ti \text{ abort} \rangle$

Remove  $T_i$  from undo-list

Stop when undo-list is empty! 🗑️

i.e.  $\langle Ti \text{ start} \rangle$  has been found for every transaction in undo-list

After undo phase completes, normal transaction processing can commence (开始) again.

### example

