

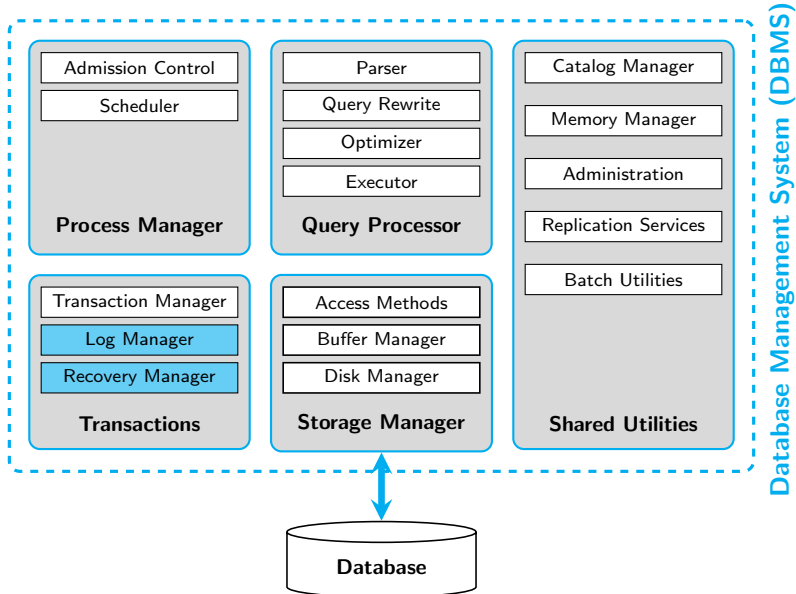
# Undo-redo logging

Clase 21

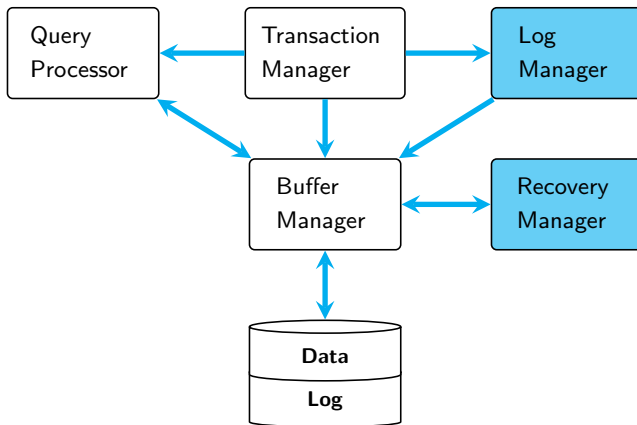
IIC 3413

Prof. Cristian Riveros

# Recuperación de fallas en el sistema



## Zoom a la arquitectura de las transacciones



# Propiedades “ácidas”

- A**tomicity: Se ejecuta todos los pasos de una transacción o no se ejecuta nada.
- C**onsistency: Al terminar una transacción los datos deben estar en un estado consistente.
- I**solation: Cada transacción se ejecuta sin ser interferida por otras transacciones.
- D**urability: Si un transacción hace commit, sus cambios sobrevivirán cualquier tipo de falla.

Queremos asegurar **atomicity** y **durability**.

¿qué puede salir mal para **atomicity** o **durability**?

### Ejemplo

T	$t$	Mem $B_1$	Mem $B_2$	$B_1$	$B_2$
PIN( $B_1$ )		1000		1000	1000
READ( $B_1, t$ )	1000				
$t := t - 100$	900				
WRITE( $B_1, t$ )		900			
UNPIN( $B_1$ )					
FLUSH( $B_1$ )				900	
PIN( $B_2$ )			1000		
READ( $B_2, t$ )	1000				
$t := t + 100$	1100				
WRITE( $B_2, t$ )			1100		
UNPIN( $B_1$ )					
ERROR	×	×	×	900	1000

# Steal y force protocolos (recordatorio)

**Steal** frame:

*¿Puede un elemento  $X$  en memoria ser escrito a disco antes de que la transacción  $T$  termine?*

**Force** page:

*¿Es necesario que una transacción haga **FLUSH** de todos los elementos modificados inmediatamente antes o después de un **COMMIT**?*

- **NO-Steal** + **Force**: fácil de recuperar.
- **Steal** + **NO-Force**: mayor performance.

Queremos un **Steal** + **NO-Force** buffer manager!

# Tres tipos de logging/recovery

1. *undo* - logging.
2. *redo* - logging.
3. *undo/redo* - logging.

# Outline

Undo logging

Redo logging

Undo-redo logging

Concurrencia y recuperación



# Outline

Undo logging

Redo logging

Undo-redo logging

Concurrencia y recuperación

# Undo logging

Log records:

- $\langle \text{START } T \rangle$
- $\langle \text{COMMIT } T \rangle$
- $\langle \text{ABORT } T \rangle$
- $\langle T, X, t \rangle$ : donde  $t$  es el valor **antiguo** en  $X$ .

Undo-log registra el valor antiguo!

# Reglas de undo-logging

## Regla $U_1$ :

*"Si T modifica X, el log record  $\langle T, X, t \rangle$  debe ser escrito antes que X sea escrito en disco."*

## Regla $U_2$ :

*"Si T hace commit, el log record  $\langle \text{COMMIT } T \rangle$  debe ser escrito justo después que todos los datos modificados por T esten almacenados en disco"*

# Reglas de undo-logging

Reglas traducidas:

1. Escribir en el log  $\langle T, X, t \rangle$ .
2. Escribir los datos a disco.
3. Escribir  $\langle \text{COMMIT } T \rangle$ .

¿cómo nos aseguramos que  
el log se escribe antes que la página en disco?

(usamos/definimos el método FLUSH-LOG)

# Reglas de undo-logging

## Ejemplo

<b>T</b>	<i>t</i>	Mem <i>B</i> <sub>1</sub>	Mem <i>B</i> <sub>2</sub>	<i>B</i> <sub>1</sub>	<i>B</i> <sub>2</sub>	LOG
						< START <b>T</b> >
READ( <i>B</i> <sub>1</sub> , <i>t</i> )	1000	1000		1000	1000	
<i>t</i> := <i>t</i> − 100	900					
WRITE( <i>B</i> <sub>1</sub> , <i>t</i> )		900				< <b>T</b> , <i>B</i> <sub>1</sub> , 1000 >
READ( <i>B</i> <sub>2</sub> , <i>t</i> )	1000		1000			
<i>t</i> := <i>t</i> + 100	1100					
WRITE( <i>B</i> <sub>2</sub> , <i>t</i> )			1100			< <b>T</b> , <i>B</i> <sub>2</sub> , 1000 >
FLUSH-LOG						
FLUSH( <i>B</i> <sub>1</sub> )				900		
FLUSH( <i>B</i> <sub>2</sub> )					1100	
						< COMMIT <b>T</b> >
FLUSH-LOG						

# Recuperación de fallas con undo-logging

¿qué pudo andar mal?

Miremos el log:

Log en disco: ... < START **T** > ... < COMMIT **T** > ...



... < START **T** > ... < ABORT **T** > ...



... < START **T** > ...



# Recuperación de fallas con undo-logging

Procesamos el log desde el **final** hasta el **principio**:

1.  $\langle \text{COMMIT } T \rangle$ : marcar  $T$  como realizada.
2.  $\langle \text{ABORT } T \rangle$ : marcar  $T$  como realizada.
3.  $\langle T, X, t \rangle$ : restituir  $X := t$  en disco si  $T$  no fue realizada.
4.  $\langle \text{START } T \rangle$ : ignorar.

# Recuperación de fallas con undo-logging

## Ejemplo

<b>T</b>	<b>t</b>	Mem $B_1$	Mem $B_2$	$B_1$	$B_2$	LOG
						< START <b>T</b> >
READ( $B_1, t$ )	1000	1000		1000	1000	
$t := t - 100$	900					
WRITE( $B_1, t$ )		900				< <b>T</b> , $B_1$ , 1000 >
READ( $B_2, t$ )	1000		1000			
$t := t + 100$	1100					
WRITE( $B_2, t$ )			1100			< <b>T</b> , $B_2$ , 1000 >
FLUSH-LOG						
FLUSH( $B_1$ )				900		
FLUSH( $B_2$ )					1100	
						< COMMIT <b>T</b> >
FLUSH-LOG						

¿qué pasa si el sistema falla en alguna instrucción?



# Recuperación de fallas con undo-logging

Algunas preguntas sin resolver:

1. ¿Qué ocurre si el sistema falla otra vez en el momento del recovery?
2. ¿Hasta donde tenemos que leer el log?
3. ¿Es posible “truncar” el log en algún minuto?

# Checkpoints

Checkpoint de la base de datos (periodicamente):

1. Dejamos de recibir nuevas transacciones.
2. Esperamos hasta que todas las transacciones actuales terminen.
3. FLUSH del log a disco.
4. Escribir < CKPT > al log y hacer FLUSH-LOG.
5. Reanudar las transacciones.

¿cómo usamos los checkpoints?

# Checkpoints

Uso del log record < CKPT >:

- Hacemos undo recovery hasta que veamos < CKPT >.
- No es necesario seguir con undo recovery después de este punto.

¿Algún problema?

- Es casi necesario “**apagar**” el DBMS para hacer checkpoint.

¿cómo hacemos checkpoint sin apagar nuestro DBMS?

# Nonquiescent checkpointing

**Quiescent** = ser tranquilo, quedarse quieto, o descansar.

Non-**quiescent** = lo contrario  
(permitir nuevas transacciones)

# Nonquiescent checkpointing

1. Escribimos un log record:

$\langle \text{START CKPT}(\mathbf{T}_1, \dots, \mathbf{T}_n) \rangle .$

- $\mathbf{T}_1, \dots, \mathbf{T}_n$  es la lista de transacciones activas.

2. Esperamos hasta que  $\mathbf{T}_1, \dots, \mathbf{T}_n$  terminen.

- Sin restringir el comienzo de nuevas transacciones.

3. Cuando  $\mathbf{T}_1, \dots, \mathbf{T}_n$  hayan terminado escribimos el log record:

$\langle \text{END CKPT} \rangle .$

# ¿cómo usamos el checkpoint para hacer el recovery?

## Ejemplo (1)

Considere el siguiente undo-log después de una falla:

< START  $T_1$  >  
<  $T_1, A, 5$  >  
< START  $T_2$  >  
<  $T_2, B, 10$  >  
< START CKPT( $T_1, T_2$ ) >  
<  $T_2, C, 15$  >  
< START  $T_3$  >  
<  $T_1, D, 20$  >  
< COMMIT  $T_1$  >  
<  $T_3, E, 25$  >  
< COMMIT  $T_2$  >  
< END CKPT >

# ¿cómo usamos el checkpoint para hacer el recovery?

## Ejemplo (2)

Considere el siguiente undo-log después de una falla:

$\langle \text{START } T_1 \rangle$   
 $\langle T_1, A, 5 \rangle$   
 $\langle \text{START } T_2 \rangle$   
 $\langle T_2, B, 10 \rangle$   
 $\langle \text{START CKPT}(T_1, T_2) \rangle$   
 $\langle T_2, C, 15 \rangle$   
 $\langle \text{START } T_3 \rangle$   
 $\langle T_1, D, 20 \rangle$   
 $\langle \text{COMMIT } T_1 \rangle$   
 $\langle T_3, E, 25 \rangle$

# Implementación de ROLLBACK

¿cómo hacemos ROLLBACK?

Usamos **undo**-logo para implementar ROLLBACK:

- Optimización: usar referencia entre logs records.

LSN      =      Log Sequence Number



# Outline

Undo logging

**Redo logging**

Undo-redo logging

Concurrencia y recuperación

# Redo logging

¿inconvenientes con **undo**-logging?

- No es posible hacer COMMIT antes de almacenar los datos en disco.
- Estamos obligados a usar una política FORCE pages.

# Redo vs undo

	undo	redo
<b>Trans. incompletas</b>	cancelarlas	ignorarlas
<b>Trans. commiteadas</b>	ignorarlas	repetirlas
<b>Escribir COMMIT</b>	después de alm. en disco	antes de alm. en disco
<b>Update log record</b>	valores antiguos	valores nuevos

# Log records para redo logging

Log records:

- $\langle \text{START } T \rangle$
- $\langle \text{COMMIT } T \rangle$
- $\langle \text{ABORT } T \rangle$
- $\langle T, X, t \rangle$ : donde  $t$  es el valor **nuevo** en  $X$ .

# Regla principal de redo-logging

## Regla $R_1$ :

*"Antes de modificar cualquier elemento  $X$  en disco, es necesario que todos los logs records estén almacenados en disco, incluido el COMMIT log record."*

En otras palabras:

1. Escribir en el log  $\langle \mathbf{T}, X, t \rangle$ .
2. Escribir  $\langle \text{COMMIT } \mathbf{T} \rangle$ .
3. Escribir los datos a disco.

Esto es al revés de undo-logging!

# Reglas de redo-logging

## Ejemplo

<b>T</b>	<b>t</b>	Mem $B_1$	Mem $B_2$	$B_1$	$B_2$	LOG
						< START <b>T</b> >
READ( $B_1, t$ )	1000	1000		1000	1000	
$t := t - 100$	900					
WRITE( $B_1, t$ )		900				< <b>T</b> , $B_1$ , 900 >
READ( $B_2, t$ )	1000		1000			
$t := t + 100$	1100					
WRITE( $B_2, t$ )			1100			< <b>T</b> , $B_2$ , 1100 >
						< COMMIT <b>T</b> >
FLUSH-LOG						
UNPIN( $B_1$ )				900		
UNPIN( $B_2$ )					1100	

# Recuperación de fallas con redo-logging

De nuevo miramos el log:

Log en disco: ... < START **T** > ... < COMMIT **T** > ...



... < START **T** > ... < ABORT **T** > ...



... < START **T** > ...



# Recuperación de fallas con redo-logging

1. Identificamos las transacciones que hicieron COMMIT.
2. Hacemos un scan desde el principio.
3. Para cada log record  $\langle T, X, v \rangle$ :
  - Si  $T$  no hizo COMMIT: no hacer nada.
  - Si  $T$  hizo COMMIT: reescribir con el valor  $v$ .
4. Para cada transacción incompleta  $T$ , escribir en el log  $\langle \text{ABORT } T \rangle$ .



# Recuperación de fallas con redo-logging

## Ejemplo

$T$	$t$	Mem $B_1$	Mem $B_2$	$B_1$	$B_2$	LOG
						< START $T$ >
READ( $B_1, t$ )	1000	1000		1000	1000	
$t := t - 100$	900					
WRITE( $B_1, t$ )		900				< $T, B_1, 900$ >
READ( $B_2, t$ )	1000		1000			
$t := t + 100$	1100					
WRITE( $B_2, t$ )			1100			< $T, B_2, 1100$ >
						< COMMIT $T$ >
FLUSH-LOG						
UNPIN( $B_1$ )				900		
UNPIN( $B_2$ )					1100	

¿qué ocurre si el sistema falla en alguna instrucción?

# Checkpoints en redo-logging

- ¿Qué diferencia hay con undo-logging?
- ¿Qué debiéramos asegurar en un checkpoint?

# Checkpoints en redo-logging

1. Escribir  $\langle \text{START CKPT}(\mathbf{T}_1, \dots, \mathbf{T}_n) \rangle$  y hacer FLUSH-LOG.
  - $\mathbf{T}_1, \dots, \mathbf{T}_n$  es la lista de transacciones activas y sin COMMIT.
2. Escribir a disco **solo** los **objetos** modificados por transacciones en estado COMMIT.
  - Sin restringir el comienzo de nuevas transacciones.
3. Escribir  $\langle \text{END CKPT} \rangle$  y hacer FLUSH-LOG.

¿qué hacemos con las páginas  
que comparten transacciones comiteadas/no-comiteadas?

# ¿cómo usamos el checkpoint para hacer el recovery?

## Ejemplo

Considere el siguiente redo-log después de una falla:

< START  $T_1$  >  
<  $T_1, A, 5$  >  
< START  $T_2$  >  
< COMMIT  $T_1$  >  
<  $T_2, B, 10$  >  
< START CKPT( $T_2$ ) >  
<  $T_2, C, 15$  >  
< START  $T_3$  >  
<  $T_3, E, 25$  >  
< END CKPT >  
< COMMIT  $T_2$  >  
< COMMIT  $T_3$  >

# Comparación undo vs redo

- Undo logging: política **STEAL**/FORCE.
- Redo logging: política **NO-STEAL**/**NO-FORCE**.

¿cómo logramos una política STEAL/NO-FORCE?

# Outline

Undo logging

Redo logging

**Undo-redo logging**

Concurrencia y recuperación

# Undo-redo logging

Log records:

- $\langle \text{START } T \rangle$
- $\langle \text{COMMIT } T \rangle$
- $\langle \text{ABORT } T \rangle$
- $\langle T, X, t_O, t_N \rangle$ : donde  $t_O$  es el valor **antiguo** y  $t_N$  es valor **nuevo**.

# Reglas para loggings anteriores

## Regla $U_1$ :

*"Si  $T$  modifica  $X$ , el log record  $\langle T, X, t \rangle$  debe ser escrito antes que  $X$  sea escrito en disco."*

## Regla $U_2$ :

*"Si  $T$  hace commit, el log record  $\langle \text{COMMIT } T \rangle$  debe ser escrito justo después que todos los datos modificados por  $T$  estén almacenados en disco."*

---

## Regla $R_1$ :

*"Antes de modificar cualquier elemento  $X$  en disco, es necesario que todos los logs records estén almacenados en disco, incluido el COMMIT log record."*



# Reglas de undor-redo logging

## Regla $UR_1$ :

*"Antes de modificar cualquier elemento en disco  $X$ , es necesario que el log record  $\langle T, X, t_O, t_N \rangle$  este registrado en el log."*

$UR_1$  es la "intersección" de los dos conjuntos de reglas.

# Reglas de undor-redo logging

## Ejemplo

<b>T</b>	<b>t</b>	Mem $B_1$	Mem $B_2$	$B_1$	$B_2$	LOG
						< START <b>T</b> >
READ( $B_1, t$ )	1000	1000		1000	1000	
$t := t - 100$	900					
WRITE( $B_1, t$ )		900				< <b>T</b> , $B_1$ , 1000, 900 >
READ( $B_2, t$ )	1000		1000			
$t := t + 100$	1100					
WRITE( $B_2, t$ )			1100			< <b>T</b> , $B_2$ , 1000, 1100 >
FLUSH-LOG						
UNPIN( $B_1$ )				900		
						< COMMIT <b>T</b> >
UNPIN( $B_2$ )					1100	

# Recuperación de fallas con undo-redo logging

Después de la falla de un sistema:

- hacer **redo** de todas las transacciones que hicieron commit.
  - desde el principio hasta el final del log.
- hacer **undo** de todas las transacciones incompletas.
  - desde el final hasta el principio del log.

¿cuál hacemos primero? ¿redo-undo? ¿o undo-redo?

# Recuperación de fallas con undo-redo logging

## Ejemplo

T	t	Mem $B_1$	Mem $B_2$	$B_1$	$B_2$	LOG
						< START T >
READ( $B_1, t$ )	1000	1000		1000	1000	
$t := t - 100$	900					
WRITE( $B_1, t$ )		900				< T, $B_1$ , 1000, 900 >
READ( $B_2, t$ )	1000		1000			
$t := t + 100$	1100					
WRITE( $B_2, t$ )			1100			< T, $B_2$ , 1000, 1100 >
FLUSH-LOG						
UNPIN( $B_1$ )				900		
						< COMMIT T >
UNPIN( $B_2$ )					1100	

¿qué ocurre si el sistema falla en alguna instrucción?

# Checkpoints en undo-redo logging

1. Escribir  $\langle \text{START CKPT}(\mathbf{T}_1, \dots, \mathbf{T}_n) \rangle$  y hacer FLUSH-LOG.
  - $\mathbf{T}_1, \dots, \mathbf{T}_n$  es la lista de transacciones activas y no commiteadas.
2. Escribir a disco **todas** las páginas marcadas en **dirty**.
  - Sin restringir el comienzo de nuevas transacciones.
3. Escribir  $\langle \text{END CKPT} \rangle$  y hacer FLUSH-LOG.

# ¿cómo usamos el checkpoint para hacer el recovery?

## Ejemplo

Considere el siguiente undo-redo log después de una falla:

< START  $T_1$  >  
<  $T_1, A, 4, 5$  >  
< START  $T_2$  >  
< COMMIT  $T_1$  >  
<  $T_2, B, 9, 10$  >  
< START CKPT( $T_2$ ) >  
<  $T_2, C, 14, 15$  >  
< START  $T_3$  >  
<  $T_3, D, 19, 20$  >  
< END CKPT >  
< COMMIT  $T_2$  >  
< COMMIT  $T_3$  >

# Outline

Undo logging

Redo logging

Undo-redo logging

Concurrencia y recuperación

# Concurrencia y recuperación en caso de fallas

Undo-redo **logging**:

- No hace diferencia para schedules no-serializables.

Políticas de control de **concurrencia**:

- No se preocupan de la labor del recovery manager.



# ¿qué podría fallar en nuestro algoritmo de recuperación?

## Ejemplo

$T_1$	$T_2$	$A$	$B$
READ( $A, x$ )		1000	1000
WRITE( $A, x - 100$ )		900	
	READ( $A, y$ )		
	WRITE( $A, y * 1.1$ )	990	
READ( $B, x$ )			
WRITE( $B, x + 100$ )			1100
	READ( $B, y$ )		
	WRITE( $B, y * 1.1$ )		1210
	COMMIT	990	1210
ABORT			



Necesitamos hacer rollback en cascada!

# Schedules recuperables

## Definición

Un schedule  $S$  es **recuperable** si cada transacción hace COMMIT solo después del COMMIT de todas las transacción de las cuales leyó sus datos.

## Ejemplo

- ¿Es este schedule recuperable?

T <sub>1</sub>	T <sub>2</sub>
	$w_2(A)$
$w_1(B)$	
$w_1(A)$	
	$r_2(B)$
COMMIT	
	COMMIT

# Schedules recuperables

## Definición

Un schedule  $S$  es **recuperable** si cada transacción hace COMMIT solo después del COMMIT de todas las transacción de las cuales leyó sus datos.

## Ejemplo

- ¿y este otro schedule? ¿és recuperable?

T <sub>1</sub>	T <sub>2</sub>
$w_1(A)$	
$w_1(B)$	
	$w_2(A)$
	$r_2(B)$
	COMMIT
COMMIT	

# Evitar rollbacks en cascadas

## Definición

Un schedule  $S$  se dice que **evita rollbacks en cascadas** (ACR: *avoid cascading rollback*) si cada transacción lee solo valores escritos por transacciones que ya hicieron COMMIT.

## Ejemplo

- ¿Es este schedule del tipo ACR?

$T_1$	$T_2$
$w_1(B)$	
$w_1(A)$	
	$w_2(A)$
	$r_2(B)$
COMMIT	
	COMMIT

ACR  $\Rightarrow$  recuperable.

# Queremos control de concurrencia que producen schedules ACR

¿Qué políticas de control de concurrencia producen schedules ACR?

- **Strict 2PL** produce schedules ACR.
- CC basado en **timestamps** produce schedules ACR.