



Dpto. Ciencias e Ingeniería de la Computación
 Universidad Nacional del Sur

ELEMENTOS DE BASES DE DATOS

Segundo Cuatrimestre 2013

Clase 18: Sistema de Recuperación: Bitácora, Modificación Inmediata y Diferida

Mg. María Mercedes Vitturini
 [mvitturi@cs.uns.edu.ar]



DBMS y Fallos

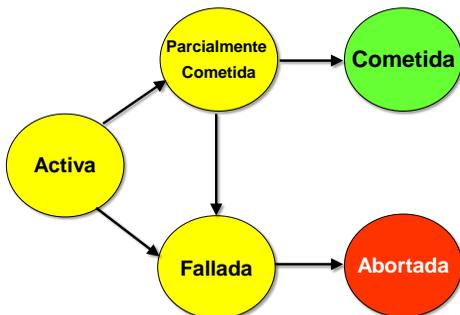
Cualquier sistema está sujeto a fallos: un fallo eventualmente puede ocasionar pérdidas o inconsistencia de información.

“Un fallo se puede presentar en cualquier etapa de la vida de las transacciones que gestiona el SMDB.”

- Ante el fallo, el SMDB se tiene que **‘recuperar’** y dejar sus bases de datos en el estado consistente anterior al fallo.
- Ejemplo:
 - Sea T una transacción que realiza transferencias de fondo de la cuenta A a la cuenta B. El fallo podría ocurrir después de descontar 50\$ de A. Los valores iniciales: A=1000, B=2000. El fallo ocurre después del débito, pero antes de acreditarlos en B. (A=950, B=2000)
 - El DBMS debe recuperarse al estado consistente anterior, antes de iniciar la transferencia (A=1000, B=2000).

EBD2013_18 - Mg. Mercedes Vitturini

Repaso – Estados de una Transacción



EBD2013_13 - Mg. Mercedes Vitturini

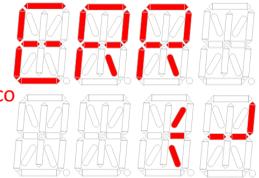
Tipos de Fallos

- Para un SMDB centralizado concurrente se identifican los siguientes tipos de fallos:

– Fallo de Transacción.

– Caída de Sistema.

– Fallo con pérdida de disco



EBD2013_18 - Mg. Mercedes Vitturini

Fallo de Transacción

Fallo de Transacción – cualquier error que cause que una transacción deba retroceder:

- **Error lógico:** la transacción no puede continuar su ejecución a causa de alguna condición interna.
- **Error bajo el entorno del sistema:** el sistema debe terminar con la ejecución de una transacción activa debido a alguna condición de error (deadlock, problemas de estampillas, etc.)
 - La transacción vuelve a ejecutarse más tarde.

Alcance del fallo de transacción – afecta transacción activa y potencialmente a otras transacciones activas por el efecto ‘fallo en cascada o cadena’. Otras transacciones activas no afectadas prosiguen su ejecución normal.

EBD2013_18 - Mg. Mercedes Vitturini

Caída del Sistema

Caída del sistema (crash) – fallos de hardware o software que causan la falla general y caída del sistema.

- Características:
 - Causan **pérdida de información en memoria volátil** (memoria principal y cache).
 - El contenido de la información en la memoria no volátil (disco) permanece intacto.
 - El DBMS posee numerosos chequeos de consistencia para prevenir la corrupción de los datos.

Alcance del fallo de sistema – el fallo afecta a **todas las transacciones activas al momento del error y potencialmente a transacciones recientemente cometidas.**

EBD2013_18 - Mg. Mercedes Vitturini

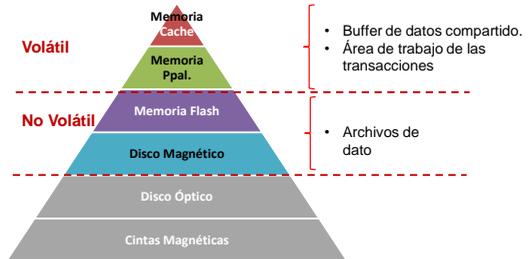
Fallo de Disco

Fallo de disco – i.e. daño físico en el medio de almacenamiento masivo.

- Se asume que el error es detectable (checksum).
- Se requiere de copias de datos de backup en otros discos u otros medios de almacenamiento para recuperar el fallo.
- La recuperación combina los pasos de recuperar un backup con recuperación del sistema.

EBD2013_18 - Mg. Mercedes Vitturini

Estructura de Almacenamiento



EBD2013_11 - Mg. Mercedes Vitturini

Estructura de Almacenamiento

Almacenamiento volátil

- Incluye memoria principal (MP) y memoria cache (MC).
- No sobrevive a las caídas de sistemas.

Almacenamiento no volátil

- Ejemplo: discos.
- Sobrevive a las caídas del sistemas.

Almacenamiento estable

- Organización del almacenamiento que asegura *'nunca'* pierde información.
- Sobrevive a todos los fallos.
- Se basa en replicas en distintos medios de almacenamiento con modos de fallo independiente y actualización de la información controlada.

EBD2013_18 - Mg. Mercedes Vitturini

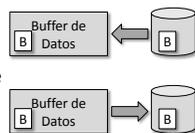
Memoria Estable

- La información que reside en memoria estable **nunca** se pierde.
- Mecanismo para alcanzarlo: **repetir la información (de manera controlada) en varios medios de almacenamiento.**
 - En el caso de **discos espejados**, ambos bloques están físicamente en la misma locación.
 - En el caso de **copias remotas**, uno de los bloques es local y el otro está en un sitio remoto.

EBD2013_18 - Mg. Mercedes Vitturini

Acceso a los datos

- La **base de datos** reside en archivos en memoria estable y se organiza en unidades denominadas **bloques físicos**.
- Los datos se trasladan **desde/hacia** los **bloques físicos** a **bloques de registro intermedio** o **buffer de MP**.
- Las transferencias disco/MP se inician a través de las operaciones:
 - **Input(B)**: transfiere el bloque físico B a memoria principal.
 - **Output(B)**: transfiere el bloque de registro intermedio B (buffer) y sustituye el bloque físico apropiado.



EBD2013_18 - Mg. Mercedes Vitturini

Acceso a los datos

- Las operaciones de entrada y salida desde el disco a la MP se hacen en **unidad de bloque**.
- Cada transacción posee su **área de trabajo (work-area) en MP** donde mantiene las variables locales con los valores de los ítems de datos accedidos y actualizados por ella.
- Las transacciones transfieren ítems de datos entre el **buffer de datos compartido** y su propia **área de trabajo** usando las operaciones: **read (X) y write(X)**.

EBD2013_18 - Mg. Mercedes Vitturini

Transferencias de Datos

READ(X, x_i) asigna el valor del elemento de dato X a la variable local x_i .

Esta operación se ejecuta como sigue:

- Si el bloque B_x en el que reside X **no está en memoria principal**, entonces se ejecuta **INPUT(B_x)**.
- Se asigna a x_i el valor de X del bloque en el registro intermedio.

EBD2013_18 - Mg. Mercedes Vitturini

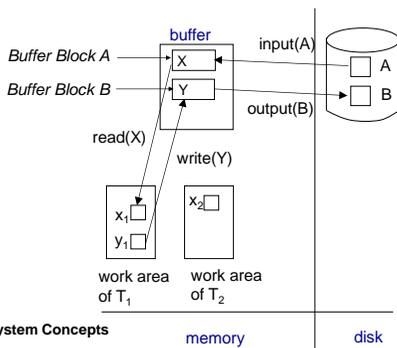
Transferencias de Datos

WRITE(X, x_i): asigna el valor de la variable local x_i al elemento de información X .

Esta operación se ejecuta como sigue:

- Si el bloque B_x en el que reside X **no está en memoria principal**, entonces se ejecuta **INPUT(B_x)**.
- Se asigna el valor de x_i a X en el bloque que contiene a X del **buffer compartido**.
- En algún momento se realiza el **OUTPUT(B_x)**.

EBD2013_18 - Mg. Mercedes Vitturini



Fuente: Database System Concepts
 A. Silberschatz

EBD2013_18 - Mg. Mercedes Vitturini

Write (X, x) – Output(X)

- El Output de bloques modificados a almacenamiento estable tiene lugar en cualquier.
- El orden en el que se realiza el output puede ser diferente del orden en el que se modificaron.
- La operación Output(X) no se siempre ejecuta inmediatamente después de write (X, xi).
 - Por ejemplo, el bloque en el que reside X contiene otros elementos de información a los que se accede todavía.

- Ante un **fallo de sistema después de ejecutar write(X, xi) y antes de ejecutar Output(X)**, existe riesgo de perder el nuevo valor para X .

EBD2013_18 - Mg. Mercedes Vitturini

Sistema de Recuperación



Los sistemas están sujetos a fallas que pueden ocasionar estados inconsistentes

Algoritmo de Recuperación (AR)

Los **algoritmos de recuperación (AR)** son procesos para asegurar **atomicidad y durabilidad** de las transacciones **a pesar de los fallos**, así como la **consistencia** de la BD.

- En un AR se identifican:
 - **Acciones preventivas** a tomar durante el procesamiento normal de las transacciones para asegurar que exista suficiente información para permitir la recuperación de fallos (**acciones preventivas**).
 - **Acciones de recuperación** como consecuencia de un fallo para asegurar la consistencia, atomicidad y durabilidad de las transacciones (**acciones preventivas**).

EBD2013_18 - Mg. Mercedes Vitturini

Acciones Preventivas

- Modificar directamente la base de datos antes de saber si la transacción cometerá puede conducir a estados inconsistentes. Se necesita guardar en memoria estable información sobre las modificaciones realizadas.

La **bitácora** (o registros *log*) se mantiene en **memoria estable** y mantiene una secuencia de **registros con información sobre las actualizaciones realizadas por cada transacción**.

Framework: la bitácora y el buffer de datos **son compartidos por todas las transacciones del sistema**.

Framework inicial: cada uno de los registros de la bitácora se graban en **memoria estable** y **las transacciones se ejecutan en serie**.

EBD2013_18 - Mg. Mercedes Vitturini

Registros de la Bitácora

- Cuando la **transacción T_i se inicia**, se agrega en la bitácora el registro:

$\langle T_i, \text{start} \rangle$

- Cada vez que la transacción T_i ejecuta **write(X)**, se agrega en bitácora un registro:

$\langle T_i, X, v_1, v_2 \rangle$

- donde v_1 representa el valor viejo de X antes de la escritura,
- v_2 el nuevo valor después de la misma,
- T_i y X identifican la transacción y el dato modificado.

EBD2013_18 - Mg. Mercedes Vitturini

Registros de la Bitácora

- Cuando la transacción T_i finalizó con su última instrucción agrega en bitácora el registro:

$\langle T_i, \text{commit} \rangle$

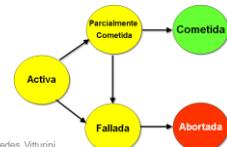
- Si la transacción T_i **aborta** agrega en bitácora el registro:

$\langle T_i, \text{abort} \rangle$

EBD2013_18 - Mg. Mercedes Vitturini

Transacción cometida

- Una transacción se dice que está **cometida** si el registro $\langle T, \text{commit} \rangle$ de bitácora está en **memoria estable**.
 - Todos los registros log anteriores de la transacción también deben haber sido almacenados en memoria estable previamente.
- Sin embargo, las **escrituras realizadas por la transacción podrían haber sido realizadas sólo en el buffer de datos** y el output estar pendiente.



EBD2013_18 - Mg. Mercedes Vitturini

Operaciones Redo y Undo

- **Undo de un registro:** $\langle T_i, X, V_1, V_2 \rangle$ escribe a X con el valor **anterior** V_1 .
- **Redo de un registro:** $\langle T_i, X, V_1, V_2 \rangle$ escribe a X con el valor **nuevo** V_2 .

Undo y Redo de Transacciones

- **undo(T_i)** restaura el valor todos los ítems de datos modificados por T_i a sus valores anteriores, recorriendo la bitácora en sentido hacia atrás desde el último registro para T_i
- **redo(T_i)** escribe el valor de todos los ítems de datos actualizados por T_i con el nuevo valor, recorriendo la bitácora hacia adelante desde el primer registro de T_i

EBD2013_18 - Mg. Mercedes Vitturini

Recuperación con Bitácora

IMPORTANTE:

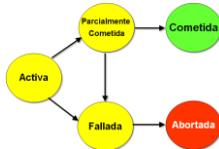
- Los **registros de la bitácora deben siempre ser escritos antes** que la modificación en la base de datos.
 - Asumimos además que la operación output para la bitácora es inmediata.
 - El output de los registros de los datos puede ocurrir en cualquier momento posterior.

EBD2013_18 - Mg. Mercedes Vitturini

Esquemas de Modificación

Existen dos esquemas de modificación

- Modificación Inmediata
- Modificación Diferida



EBD2013_18 - Mg. Mercedes Vitturini

Esquema modificación inmediata

- Bajo el esquema de **modificación inmediata** las modificaciones a la base de datos (*output(X)*) se realizan **en cualquier momento aún mientras la transacción está activa** (modificaciones no cometidas).
- En **caso de que ocurra un fallo**, es necesario **deshacer los cambios** hechos por una transacción no cometida.

EBD2013_18 - Mg. Mercedes Vitturini

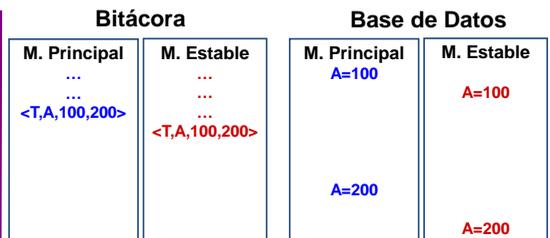
Modificación inmediata: acciones preventivas

- Cuando la transacción T se inicia se agrega en bitácora el registro **<T, start>**
- Cada operación **Write(X)** de T, agrega el registro de bitácora con los **valores anterior y actual** para X,
 $\langle T, X, v_1, v_2 \rangle$
- Luego se actualiza la base de datos.
- Cuando T está cometida se agrega en bitácora el registro **<T, commit>**

Nunca se actualiza la base de datos **antes de que se escriba el registro de bitácora en memoria estable.**

EBD2013_18 - Mg. Mercedes Vitturini

Actualización efectiva de los datos

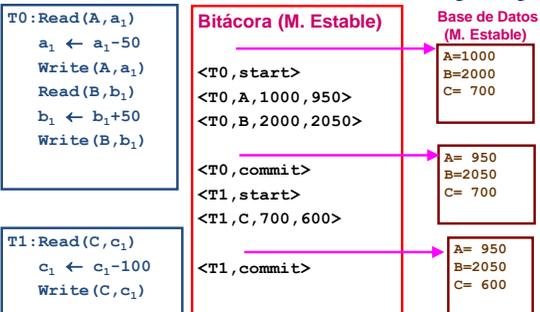


La actualización en memoria principal del valor de A=200 **debe ser posterior** al grabado en memoria estable del registro de bitácora <T, A,100,200>.

↓ Tiempo

EBD2013_18 - Mg. Mercedes Vitturini

Ejemplo



EBD2013_18 - Mg. Mercedes Vitturini

Modificación inmediata: recuperación ante un fallo

- El esquema de recuperación con modificación inmediata necesita de las operaciones **UNDO** y **REDO** usando información de la bitácora.
 - **Undo(T)**: restaura todos los ítems de dato que T actualiza a los valores que tenía anteriormente.
 - **Redo(T)**: asigna los nuevos valores a todos los ítems de dato que actualiza la transacción T.
- Las operaciones undo y redo son **idempotentes**:
 $\text{undo}(T)=\text{undo}(\text{undo}(T))$ y $\text{redo}(T)=\text{redo}(\text{redo}(T))$

EBD2013_18 - Mg. Mercedes Vitturini

Recuperación ante fallos

- Consideremos que se produjo un **fallo del sistema**. El esquema de recuperación inmediata consulta la bitácora para determinar cuáles transacciones deben deshacerse o rehacerse.
- **¿Cuándo aplicar undo(T)?**
 - Cuando la bitácora contiene el registro <T,start> pero no contiene el registro <T,commit>.
- **¿Cuándo aplicar redo(T)?**
 - Cuando la bitácora contiene tanto el registro <T,start> como el registro <T,commit>.

EBD2013_18 - Mg. Mercedes Vitturini

Recuperación de fallos

<T0,start> <T0,A,1000,950> <T0,B,2000,2050>	<T0,start> <T0,A,1000,950> <T0,B,2000,2050> <T0,commit> <T1,start> <T1,C,700,600>	<T0,start> <T0,A,1000,950> <T0,B,2000,2050> <T0,commit> <T1,start> <T1,C,7000,600> <T1,commit>
---	--	--

Undo(T0)

Undo(T1)
Redo(T0)

Redo(T0)
Redo(T1)

Es necesario realizar las operaciones de undo antes que las operaciones redo.

EBD2013_18 - Mg. Mercedes Vitturini

Esquema Modificación Diferida

- En el esquema de **modificación diferida** las modificaciones que realice una transacción (**write**)
 - Durante la ejecución de la transacción sólo se graban en los registros de la bitácora,
 - Las escrituras en la base de datos se demoran hasta que la transacción este cometida.
- La información en la bitácora se utiliza cuando se ejecutan las escrituras diferidas.
- Si ocurre un **fallo de sistema** o **fallo de transacción** antes de que la transacción este cometida se ignora la información de la bitácora.

EBD2013_18 - Mg. Mercedes Vitturini

Modificación diferida: acciones preventivas

- Antes de que comience la ejecución de la transacción T_i , se agrega en bitácora el registro **<T_i, start>**.
- Si T_i realiza una operación **write(X,v)**, se agrega en bitácora el registro **<T_i, X, v>**.
 - Modificación diferida requiere solamente del nuevo valor del dato, omitiendo el valor antiguo.
- Cuando T_i está cometida, se escribe el registro de bitácora **<T_i, commit>**.
- Los registros *log* se usan para las escrituras diferidas.

EBD2013_18 - Mg. Mercedes Vitturini

Modificación Diferida



T1
 Read (A, a1)
 a1 = a1 - 50
 Write (A, a1)
 Read (B, b1)
 b1 = b1 + 50
 Write (B, b1)

Bitacora
 <start T1>
 <T1, A, 950>
 <T1, B, 2050>
 <commit T1>

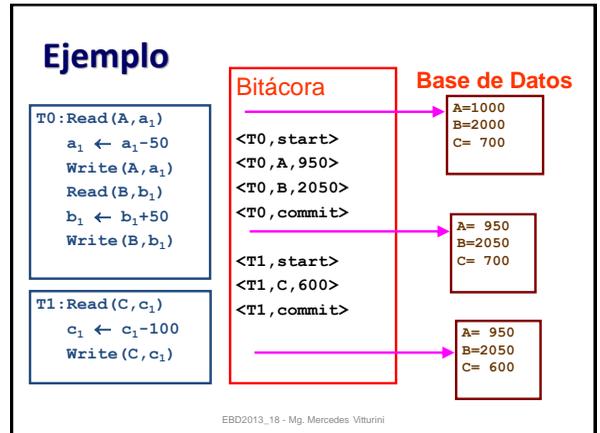
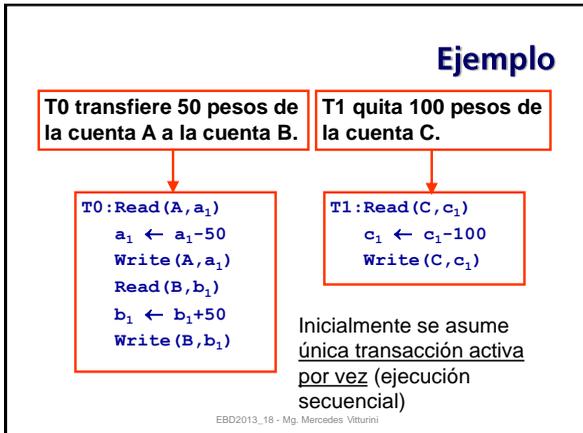


EBD2013_18 - Mg. Mercedes Vitturini

Recuperación y modificación diferida

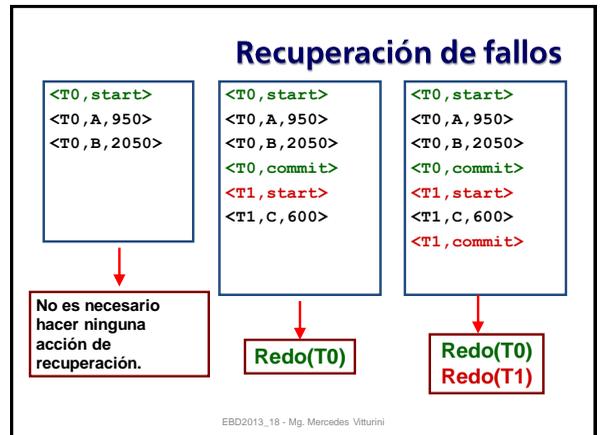
- Supongamos una secuencia de transacciones T_1, T_2, \dots, T_n ejecutándose en serie.
- En cualquier momento puede ocurrir **una caída del sistema**.
- Para recuperar el estado consistente, **necesitarán rehacerse (REDO T_i)** de todas aquellas transacciones que tienen en bitácora **<T_i, start>** y **<T_i, commit>**.
- **¿Y si ocurriese una caída del sistema mientras se realizan tareas de recuperación?**

EBD2013_18 - Mg. Mercedes Vitturini



- ### Recuperación ante Fallos
- El esquema de recuperación ante fallos con modificación diferida usa la operación **REDO** y la **información de la bitácora**.
 - Redo(T)**: asigna los nuevos valores a todos los datos que actualiza la transacción T.
 - La operación redo es **idempotente**, esto es:

$$\text{redo}(T) = \text{redo}(\text{redo}(T))$$
 - ¿Cuándo aplicar redo(T) después de un fallo?
 - Siempre que en bitácora existan los registros <T, start> y <T, commit>.
- EBD2013_18 - Mg. Mercedes Vitturini



- ### Temas de la clase de hoy
- Sistema de Recuperación.
 - Tipos de Fallos.
 - Estructuras de Almacenamiento.
 - Mecanismos de recuperación usando bitácora para ambientes no concurrentes.
 - Modificación Diferida.
 - Modificación Inmediata.
- Bibliografía**
- “Database System Concepts” – A. Silberschatz. Capítulos 17.
- EBD2013_18 - Mg. Mercedes Vitturini