



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

## ELEMENTOS DE BASES DE DATOS

Segundo Cuatrimestre 2013

### Clase 19: Sistema de Recuperación: acciones de recuperación - checkpoint

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



## Recuperación del Sistema

- Un sistema está sujeto a fallos.
- Un fallo potencialmente genera pérdida o inconsistencias en la información.
- Ante un fallo, el DBMS debe estar preparado para su 'recuperación' y dejar la base de datos en el estado consistente anterior al fallo.
- Los fallos se clasifican en:
  - Fallo de Transacción
  - Caída de Sistema
  - Rotura de disco

EBD2013\_19 - Mg. Mercedes Vitturini

## Repaso - Tipos de Fallos

- Fallos en la Transacción
    - *Error lógico*: la transacción no puede continuar su ejecución a causa de alguna condición interna.
    - *Error del sistema*: el sistema alcanza un estado no correcto. La transacción puede volver a ejecutarse más tarde.
  - Caída de sistema: por ejemplo errores del hardware o software de base de datos o software del sistema operativo.
  - Fallo de disco: daño físico en el medio de almacenamiento masivo.
- Otras transacciones continúan su ejecución
- Varias transacciones deben ser recuperadas

EBD2013\_19 - Mg. Mercedes Vitturini

## Algoritmos de Recuperación (AR)

- Los algoritmos de recuperación (AR) son técnicas para asegurar la consistencia, atomicidad y durabilidad de las transacciones a pesar de los fallos.
- En un AR se identifican:
  - Acciones a tomar durante el procesamiento normal para asegurar que exista suficiente información para permitir la recuperación de fallos (acciones preventivas).
  - Acciones tomadas a consecuencia de un fallo para asegurar la consistencia, atomicidad y durabilidad de las transacciones (acciones paliativas).

EBD2013\_19 - Mg. Mercedes Vitturini

## AR mediante Bitácora - Repaso

La bitácora es una estructura que se utiliza para guardar información sobre las modificaciones que se realizaron a los datos. La implementación que vimos contiene registros de control y registros con los atributos:

- **Nombre de la Transacción**: el nombre de la transacción que ejecutó la actualización (Write).
- **Nombre del Dato**: el nombre único del dato escrito.
- **Valor Antiguo**: el valor del dato anterior a la escritura.
- **Valor Nuevo**: el valor que tendrá el dato después de la escritura.

EBD2013\_19 - Mg. Mercedes Vitturini

## Repaso – Registros de la Bitácora

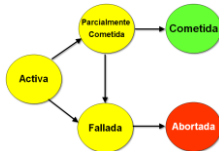
- Cuando la transacción  $T_i$  se inicia, se registra agregando en bitácora:  $\langle T_i, \text{start} \rangle$
- Cada vez que la transacción  $T_i$  ejecuta  $\text{write}(X, x_2)$ , se registra agregando en bitácora:  $\langle T_i, X, v_1, v_2 \rangle$ 
  - donde  $v_1$  representa el valor viejo 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.
- Cuando la transacción  $T_i$  finalizó con su última instrucción agrega en bitácora el registro:  $\langle T_i, \text{commit} \rangle$

EBD2013\_19 - Mg. Mercedes Vitturini

## Esquemas de Modificación

Existen dos esquemas de modificación

- Modificación Inmediata
- Modificación Diferida



EBD2013\_19 - Mg. Mercedes Vitturini

## Modificación diferida

- Implementa atomicidad de la transacción grabando las modificaciones en los registros bitácora y aplazando las operaciones write en la base de datos hasta que la transacción se comete parcialmente.
- La información en la bitácora asociada a la transacción parcialmente cometida se usa durante la ejecución de las escrituras diferidas.
- El esquema de recuperación usa la operación **redo** (rehacer) usando información de la bitácora.
  - **REDO (T)** cuando en la bitácora existen los registros  $\langle T, start \rangle$  y  $\langle T, commit \rangle$ .

EBD2013\_19 - Mg. Mercedes Vitturini

## Modificación inmediata

- Las modificaciones en la base de datos se realizan mientras la transacción está activa, antes de alcanzar el estado de cometida.
- Cuando se produce un fallo, el esquema de recuperación consulta a la bitácora para determinar que transacciones deben deshacerse o rehacerse.
  - **UNDO (T)**: cuando la bitácora contiene el registro  $\langle T, start \rangle$  pero no contiene el registro  $\langle T, commit \rangle$ .
  - **REDO (T)**: cuando la bitácora contiene tanto el registro  $\langle T, start \rangle$  como el registro  $\langle T, commit \rangle$
- Los registros de la bitácora deben ser escritos en memoria estable antes de cualquier modificación de la base de datos.

EBD2013\_19 - Mg. Mercedes Vitturini

## Restricciones impuestas

Hasta ahora asumimos las siguientes *restricciones*:

1. La operación **output de la bitácora es inmediata**. Los registros de la bitácora están "inmediatamente" en memoria estable.
  - La escritura de los registros bitácora en memoria estable siempre es anterior a la escritura de los correspondientes bloques de datos.
  - El **output de los registros de los datos** ocurre en cualquier momento posterior.
2. Durante el proceso de recuperación **se recupera el sistema completo** (desde la primera transacción de la bitácora).
3. **Planificación de las transacciones en serie**.

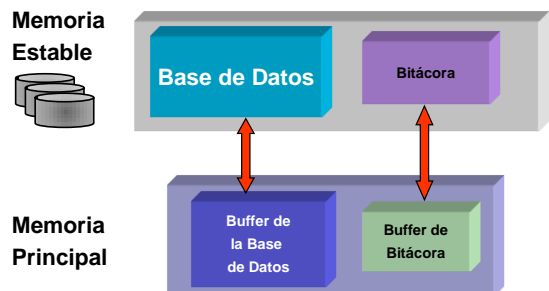
EBD2013\_19 - Mg. Mercedes Vitturini

## Buffering de Registros de Bitácora

- Hasta ahora supusimos que cada registro de bitácora se graba en memoria estable inmediatamente después de cuando se crea.
  - Esto generaría un *overhead extra* Demasiado costoso.
- Cada registro de bitácora es mucho más pequeño que el tamaño del bloque que lo contiene, la grabación de cada registro de bitácora se traduce en una grabación mucho mayor en el nivel físico y por cada registro.

EBD2013\_19 - Mg. Mercedes Vitturini

## Organización de la Memoria



EBD2013\_19 - Mg. Mercedes Vitturini

## Buffering + Registros de bitácora

- Reglas para grabar el buffer de registros de bitácora:
  - El output de registros de bitácora es en el orden en que se crean.
  - La transacción  $T_i$  alcanza el estado *cometida* cuando se grabó en memoria estable  $\langle T_i, commit \rangle$ .
  - Antes de un output de cualquier bloque de datos desde el buffer de memoria principal, todos los registros de bitácora asociados a los datos del bloque deben ser grabados en memoria estable previamente: *write-ahead logging*.

EBD2013\_19 - Mg. Mercedes Vitturini

## Buffering: Base de datos + Bitácora

- La base de datos y la bitácora se almacena en memoria estable.
- Los datos se traen a memoria principal según se requieran mediante un esquema de bloques.
- Si un bloque B1 en memoria principal fue modificado y necesita ser reemplazado por otro bloque B2, entonces:
  - Grabar B1 (Output(B1)).
  - Luego cargar B2 (Input(B2)).
- Este es el concepto estándar de memoria virtual.

EBD2013\_19 - Mg. Mercedes Vitturini

## Buffering: Base de Datos + Bitácora

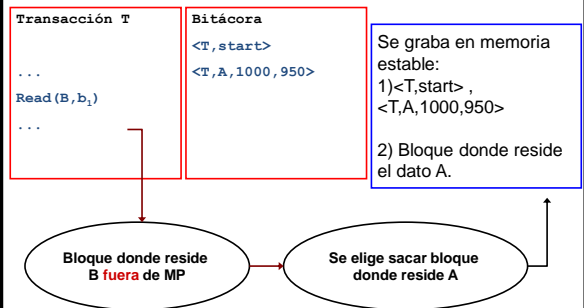
Spongamos que la entrada del bloque B2 hace que el bloque B1 deba salir de memoria principal:

Los pasos a seguir y en ese orden son:

- Grabar en memoria estable hasta todos los registros de bitácora pertenecientes al bloque B1 (ie. grabar en memoria estable los bloques de bitácora hasta cubrir los bloques que contienen registros de B1)
- Grabar el bloque B1 en memoria estable.
- Cargar el bloque B2 desde el disco a memoria principal.

EBD2013\_19 - Mg. Mercedes Vitturini

## Buffering de la base de datos



EBD2013\_19 - Mg. Mercedes Vitturini

## Ejercicio

- Plantear una situación en la que grabar primero los bloques del buffer de datos en memoria estable antes que los registros correspondientes de bitácora provoque una situación de inconsistencia ante un requerimiento de recuperación.

EBD2013\_19 - Mg. Mercedes Vitturini

## El proceso de recuperación

Ante un fallo de sistema las acciones a seguir son:

- Iniciar el proceso de recuperación.
  - No se inicia de la atención de requerimientos de usuarios.
- Analizar la bitácora determinando el conjunto de transacciones a rehacer y deshacer.
  - Deshacer depende de la estrategia de modificación.
- Reiniciar la operación normal del sistema.

EBD2013\_19 - Mg. Mercedes Vitturini

## Recuperación con Registros Bitácora

- Cada vez que ocurre un **fallo de sistema** sería necesario consultar **la bitácora completa** para ver qué transacciones deben rehacerse y cuáles deshacerse.
- **Problemas:**
  - El proceso de recuperación **consume tiempo**.
  - La **mayor parte de las transacciones** que, según el algoritmo, **deben volver a hacerse, seguramente ya escribieron sus actualizaciones** en la base de datos en memoria estable (output).
  - Volver a hacerlas no causa daño, pero **insume tiempo**.

EBD2013\_19 - Mg. Mercedes Vitturini

## Puntos de Verificación (Checkpoints)

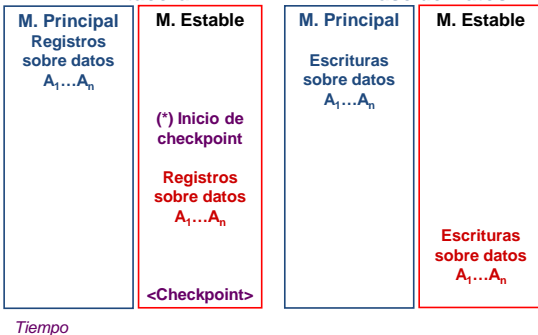
### Solución

- Se agrega al sistema de recuperación un procedimiento de **checkpointing** o puntos de verificación, que sigue los pasos:
  1. Grabar en disco (output B) todos los bloques de **registros de bitácora** que están en memoria principal.
  2. Grabar en disco (output B) los **bloques modificados de los registros intermedios** (buffer de MP).
  3. **Grabar un registro de bitácora <checkpoint>** en memoria estable.
- Durante el proceso de checkpointing se **suspenden todas las otras tareas para realizar sólo tareas de sistema**.

EBD2013\_19 - Mg. Mercedes Vitturini

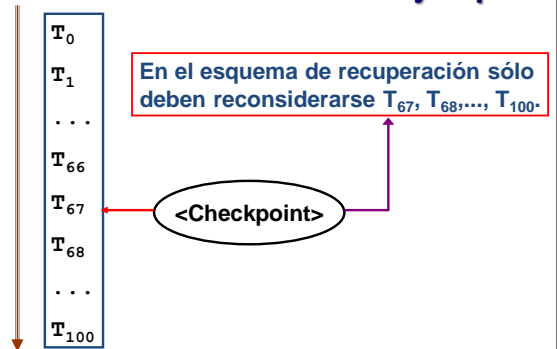
## Checkpoints

### 1. Bitácora → 2. Base de Datos



EBD2013\_19 - Mg. Mercedes Vitturini

## Ejemplo



EBD2013\_19 - Mg. Mercedes Vitturini

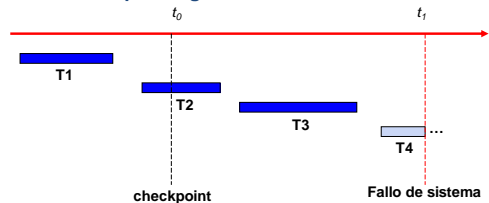
## Recuperación con Checkpoints

- Para cada transacción  $T_i$  tal que aparece en la bitácora el registro  $\langle T_i, \text{commit} \rangle$  **antes** del registro  $\langle \text{checkpoint} \rangle$ , **no se requiere ejecutar un REDO**.
- Después de un fallo de sistema:
  - Examinar la bitácora **hacia atrás** buscando el primer registro  $\langle \text{checkpoint} \rangle$  y cual es el registro  $\langle T_i, \text{start} \rangle$  más cercano (determina la última transacción  $T_i$  que comenzó a ejecutarse antes del *checkpoint*).
  - Luego, se aplica **UNDO o REDO** sobre  $T_i$  y todas las transacciones  $T_k$  que le suceden.

EBD2013\_19 - Mg. Mercedes Vitturini

## Ejemplo

### Checkpointing + Modificación Inmediata



- **T1** se puede **ignorar**.
- **T4** debe ser **deshecha** (undo).
- **T2** y **T3** deben ser **rehеchas** (redo).

EBD2013\_19 - Mg. Mercedes Vitturini

## Checkpoints + Modelo Concurrente

- En un **esquema secuencial** se considera:
  - Transacciones que comenzaron después del checkpoint más reciente (pueden ser varias).
  - La **única transacción activa (si existe)** al momento de dicho checkpoint.
  - La **única transacción activa (si existe)** al momento que se produce un fallo de sistema.
- **Diferencias con un esquema concurrente:**
  - Podrían existir **varias transacciones activas al momento del checkpoint.**
  - Podrían existir **varias transacciones activas al momento de un fallo de sistema.**

EBD2013\_19 - Mg. Mercedes Vitturini

## Recuperación y Concurrencia

### IMPORTANTE

- Un sistema centralizado, aún cuando incorpora facilidades de concurrencia, cuenta con:
  - un **único buffer de registros de disco** y
  - una **única bitácora.**

EBD2013\_19 - Mg. Mercedes Vitturini

## Checkpoint + Modelo Concurrente

El esquema de checkpoint se modifica levemente si existen transacciones concurrentes.

1. Reemplazar el registro **<checkpoint>** por un registro de la forma: **<checkpoint L>**, donde **L** define la **lista de transacciones activas** al momento del checkpoint.
  2. La lista **L** limita la búsqueda hacia atrás del checkpoint en la bitácora.
- Como en el caso anterior, no existen actualizaciones a los bloques de buffer ni a la bitácora durante el proceso de checkpointing.

EBD2013\_19 - Mg. Mercedes Vitturini

## Checkpoints + Modelo Concurrente

### Caída de Sistema (Crash)

- El proceso de recuperación involucra a todas las **transacciones ejecutadas y en ejecución** (posiblemente más de una) a partir del último punto de verificación (**checkpoint**).
- El sistema de recuperación recorre la bitácora para construir dos listas:
  - **undo-list**
  - **redo-list,**

EBD2013\_19 - Mg. Mercedes Vitturini

## Undo-List y Redo-List

1. Recorrer la bitácora **hacia atrás** hasta el primer **<checkpoint, L>**:
  - Por cada registro **<T<sub>i</sub>, commit>**, se agrega T<sub>i</sub> a **redo-list**.
  - Por cada registro **<T<sub>i</sub>, start>**, si T<sub>i</sub> no está en **redo-list** entonces se agrega a la lista **undo-list**.
2. Al llegar **<checkpoint, L>** controlar la lista **L** de transacciones activas:
  - Por cada transacción T<sub>i</sub> de **L** si T<sub>i</sub> no está incluida en **redo-list** se agrega a **undo-list**.

EBD2013\_19 - Mg. Mercedes Vitturini

## Recuperación ante Fallo de Sistema

- 1 Volver a recorrer la **bitácora hacia atrás**, realizando **undo de cada registro** que corresponda a una transacción de la lista **undo-list**. El proceso se detiene al alcanzar en la bitácora el item **<T<sub>i</sub>, start>** de cada transacción T<sub>i</sub> en la lista **undo-list**.
- 2 Avanzar hasta el más reciente **<checkpoint L>** de la bitácora.
- 3 Recorrer la bitácora **hacia adelante** desde el más reciente **<checkpoint L>** y realizar **redo de cada registro** que corresponda a una transacción T<sub>i</sub> en la lista **redo-list**.

EBD2013\_19 - Mg. Mercedes Vitturini

## Recuperación ante Fallos

- Es importante **respetar el orden de ejecución** de las operaciones para la recuperación.
- Las operaciones **UNDO** deben realizarse siempre antes que las operaciones **REDO**.
- Las operaciones de **UNDO** se realizan recorriendo la bitácora desde abajo hacia arriba, esto es, en orden inverso al que se ejecutaron.
- Las operaciones de **REDO** se realizan recorriendo la bitácora desde arriba hacia abajo, esto es, en el mismo orden en que se actualizó.

EBD2013\_19 - Mg. Mercedes Vitturini

## CRASH de Sistema

Dada la siguiente foto de la bitácora antes de un crash de sistema.

```
<T1, start>
<T2, start>
<T1, B, 50, 20>
<T2, A, 20, 30>
<checkpoint, <T1, T2>>
<T1, commit>
<T3, start>
<T3, B, 20, 40>
<T4, start>
<T4, B, 40, 80>
<T4, C, 100, 50>
<T3, commit>
... Crash ...
```

**Undo-List:** T<sub>4</sub>, T<sub>2</sub>

**Redo-List:** T<sub>1</sub>, T<sub>3</sub>

**B: 50,**  
**A: 20,**  
**C: 100,**          **B: 40,**  
**A: 20,**  
**C: 100,**

EBD2013\_19 - Mg. Mercedes Vitturini

## Fallo de Transacciones

- El **fallo de una transacción** también se resuelve usando la bitácora.
- Si una transacción T tiene que ser retrocedida también implica **recorrer la bitácora hacia atrás**: una misma transacción puede realizar más de un cambio sobre el mismo dato.  
 $..., \langle T_5, A, 10, 20 \rangle, \dots, \langle T_5, A, 20, 30 \rangle, \dots$
- Esto es, T<sub>5</sub> primero modifica A de 10 a 20 y luego de 20 a 30. Si no la recorremos hacia atrás, se restauraría A con el valor 20 (cuando debe ser 10).

EBD2013\_19 - Mg. Mercedes Vitturini

## Ejemplo – Fallo de una transacción

T1	T2
read (A)	
write (A)	
	read (A)
	write(B)
	write (A)
read (B)	
...	...

**Estampillas**

R-ts(A)=ts(T<sub>1</sub>)

W-ts(A)=ts(T<sub>1</sub>)

R-ts(A)=ts(T<sub>2</sub>)

W-ts(B)=ts(T<sub>2</sub>)

W-ts(A)=ts(T<sub>2</sub>)

**Bitacora**

<T<sub>1</sub> start>

<T<sub>1</sub>, A, 100, 200>

<T<sub>2</sub> start>

<T<sub>2</sub>, B, 400, 200>

<T<sub>2</sub>, A, 200, 500>

- **Protocolo:** Estampilla de tiempo con  $ts(T_1) < ts(T_2)$
- T<sub>1</sub> debe retroceder por no cumplir con el protocolo
- T<sub>2</sub> debe retroceder en cascada.

EBD2013\_19 - Mg. Mercedes Vitturini

## Checkpoints Difusos (Fuzzy)

- El **proceso para incluir un punto de verificación (checkpointing)** requiere que las actualizaciones a las base de datos se **suspendan** temporalmente.
  - No se atienden requerimientos de usuarios para realizar tareas de sistema
- La técnica de **puntos de verificación difusos (fuzzy checkpointing)** permite que se reinicien las actualizaciones **después de grabar el registro de checkpoint en bitácora en memoria estable pero antes de que los bloques modificados se escriban en disco.**

EBD2013\_19 - Mg. Mercedes Vitturini

## Checkpoints Difusos

1. Se suspenden todas las tareas de usuario.
  - Se graba en memoria estable la bitácora.
  - Se graba **<checkpoint L>** en bitácora en memoria estable.
  - Se identifica la **lista M** de bloques buffer modificados.
2. Se permite a las transacciones proceder.
3. En paralelo se continua con el output de cada uno de los bloques de la lista M
  - No se pueden modificar los bloques de M mientras están en proceso output.
4. En la recuperación usando fuzzy checkpoint, se comienza el recorrido desde el registro **checkpoint** apuntado por **last\_checkpoint**

EBD2013\_19 - Mg. Mercedes Vitturini

## Administración de buffer

El buffer de base de datos se puede manejar de 2 formas:

- El **SMBD se reserva parte de memoria principal para usarla como buffer**.
  - El tamaño del buffer está acotado por los requerimientos de memoria por parte de otras aplicaciones.
  - Aunque no se use el espacio del buffer, no puede usarse para otros fines (menos flexible).
- El **SMBD implementa el buffer en un espacio de memoria virtual provisto por el sistema operativo**.
  - El tamaño del buffer es flexible.
  - El SO decide sacar los bloques de memoria, de modo tal que el SMBD nunca tiene control sobre el número de bloques de buffer.

EBD2013\_19 - Mg. Mercedes Vitturini

## Fallos con pérdida de información

- Aunque el fallo de memoria no volátil es menos frecuente, el sistema debe estar preparado también para este tipo de fallos.
- La técnica es **copiar periódicamente el contenido entero (dump) de la base de datos** a un almacenamiento estable (por ejemplo, cintas magnéticas, dump en servidores remotos).
- Se ejecuta un proceso similar al de checkpointing.
- Ninguna transacción debe estar activa durante el proceso de dump

EBD2013\_19 - Mg. Mercedes Vitturini

## Copias de Seguridad (dump)

### El proceso de Dump de Base de Datos

1. Se **vuelca el buffer de los registros de bitácora** residiendo actualmente en memoria principal al almacenamiento estable.
2. Se **vuelcan todos los bloques de buffer de datos al disco**.
3. Se **copian los contenidos de la base de datos** al medio de almacenamiento estable.
4. Se vuelca un registro de bitácora **<dump>** al almacenamiento estable.

EBD2013\_19 - Mg. Mercedes Vitturini

## Fallo de Disco

Si falla el almacenamiento no volátil:

1. Esto **restaura la base de datos al estado en que se encontraba cuando se realizó la copia de la bases de datos completa a memoria estable (dump)**.
2. Luego, el sistema consulta la bitácora para hacer **redo de todas las transacciones cometidas después del registro <dump>**.

EBD2013\_19 - Mg. Mercedes Vitturini

## Temas de la clase de hoy

- Sistema de Recuperación
  - Puntos de Verificación. Modificación para ambientes concurrentes.
  - Puntos de Verificación Difusos.
  - Recuperación y Buffering.
  - Implementación de memoria estable.
- Bibliografía
  - “Database System Concepts” – A. Silberschatz. Capítulo 17.

EBD2013\_19 - Mg. Mercedes Vitturini