Recuperação após Falha

Banco de Dados II

Cap. 18 (Ramakrishnan) Cap. 23 (Elmasri) Cap. 17 (Silberschatz)

Prof. Denio Duarte



- A recuperação de falha significa que o BD é restaurado ao estado consistente mais recente antes do momento da falha
- SGBD precisar manter informações sobre as mudanças que estão sendo aplicadas aos itens pelas diversas transações
 - Armazenadas no log do sistema

- Tipos de falhas
 - Transação
 - Sistema
 - Meio de armazenamento
 - Problemas físicos e catástrofes

- Transação: uma transação ativa termina de forma anormal
 - Exemplos: violação de Restrição de Integridade, lógica da transação mal definida, deadlock ou cancelamento pelo usuário.
 - Não compromete a memória principal nem a memória secundária (disco, em geral)
 - Falha de maior probabilidade de ocorrência
 - Tempo de recuperação pequeno

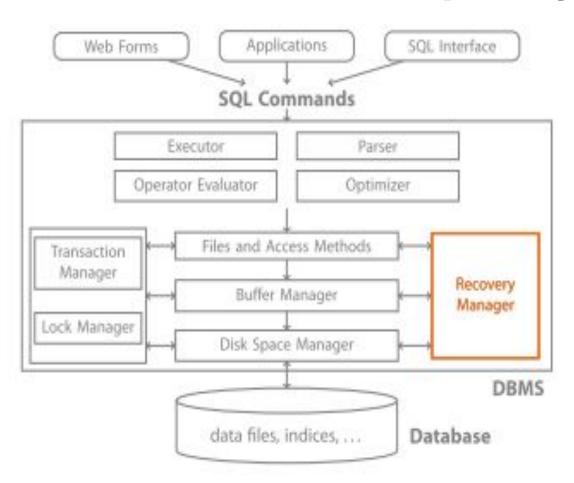
- Sistema: o SGBD encerra a sua execução de forma anormal
 - Interrupção de energia, falha no SO, erro interno no
 - SW do SGBD, falha de HW, por exemplo
 - Compromete a memória principal mas não compromete o disco
 - Falha com probabilidade média de ocorrência
 - Seu tempo de recuperação é médio

- Meio de armazenamento: o BD se torna total ou parcialmente inacessível
 - Setores corrompidos no disco, falha no cabeçote de leitura/gravação, entre outras
 - Não compromete a memória principal mas compromete o disco
 - Falha com menor probabilidade de ocorrência
 - Seu tempo de recuperação é grande

- Problemas físicos e catástrofes: erros do ambiente
 - Falha de energia, ar condicionado, incêndio, alagamento, roubo
 - Falha com menor probabilidade de ocorrência
 - Seu tempo de recuperação é grande



Gerenciador de Recuperação



Recuperação de Falhas

- Garante a atomicidade e a durabilidade das transações
 - Um SGBD deve ser tolerante a falhas
- Tolerância a falhas
 - Capacidade de retornar o BD a um estado passado consistente, após a ocorrência de uma falha que o deixou em um estado inconsistente
 - Baseia-se em redundância de dados
 - Não é um mecanismo 100% seguro
 - Responsabilidade do subsistema de recuperação (recovery) do SGBD

Recuperação

- Estratégia básica
 - Caso dano extensivo (catástrofe, falha disco): restaura backup e reconstrói um estado mais recente, reaplicando ou refazendo as operações das transações no log.
 - Caso dano não extensivo (não houve falha no disco, etc): refazer as operações das transações efetivadas e desfazer as operações das transações abortadas

- Controles
 - Durante a operação do SGBD
 - Manter informações sobre quais itens do BD foram atualizados
 - Realizar cópias periódicas do BD
 - Após a falha
 - Executar ações para retornar o BD a um estado consistente
 - Ações básicas
 - Undo: desfazer uma atualização
 - Redo: refazer uma atualização

- As técnicas de recuperação (dano não extensivo) podem ser baseadas em
 - Reescrita no disco
 - Cópia dos dados (paginação de sombra shadowing)

- Técnica de reescrita no disco
 - Baseada em log
 - Estratégias
 - Atualização adiada (deferred update)
 - Atualização imediata (immediate update)

- Atualização adiada (deferred update)
 - Não atualiza fisicamente o BD até que a transação atinja o commit
 - Atualização são registradas no espaço de trabalho da transação ou nos buffers (através dos logs)
 - Quando a transação entra em Parcialmente
 Efetivada as atualizações do log são efetivadas

- Atualização adiada (deferred update)
 - Em caso de falha antes do fim (system crash ou aborto), as informações no log são ignoradas
 - Não é necessário o UNDO das ações mas pode ser preciso o REDO
 - Algoritmo de recuperação, neste caso, é chamado NO-UNDO/REDO

- Atualização adiada (deferred update)
 - Funcionamento
 - Antes de T iniciar sua execução, a tupla
 T, start> é registrada no log
 - A cada write(item), uma tupla <T, item, novo valor> é registrada também
 - Como UNDO é desnecessário, a tupla contém apenas o novo valor do item
 - Ao fim da transação, registra <T, commit>
 - Antes de escrever o valor do item no banco de dados (output (item)) o registro do log deve ser gravado em um meio estável de armazenamento

Atualização adiada (deferred update)

```
T_1 (C=700)
     (A=100, B=200)
                              read(C)
read(A)
A = A - 50
                          C = C - 100
                              write(C)
write(A)
read(B)
                                                    log
B = B + 50
                                                    (T_0 - > T_1)
write(B)
                                                    <T<sub>0</sub>, start>
                                                    <T_0, A, 50>
                                                    <T_0, B, 250>
                                                    <T_0, commit>
                                                    <T<sub>1</sub>,start>
                                                    <T_1, C, 600>
                                                    <T_1, commit>
```

- Atualização adiada (deferred update)
 - Em caso de falha o procedimento redo (T) é executado e todos os novos itens da transação T são confirmados no banco
 - redo (T) deve ser idempotente, ou seja, pode ser executada várias vezes para a mesma transação mas o resultado deve ser sempre o mesmo
 - O Uma transação T' é refeita se possuir as tuplas <T', start> e <T', commit> no log.

Atualização adiada (deferred update)

$\log (T_0^->T_1^-)$ <T₀,start> <T_o,start> <T_o,start> $<T_0, A, 50>$ $<T_0, A, 50>$ $<T_0, A, 50>$ $<T_0, B, 250>$ $<T_0, B, 250>$ $<T_0, B, 250>$ <T $_{0}$, commit> $<T_0$, commit> crash <T₁,start> <T₁,start> $<T_1, C, 600>$ $<T_1, C, 600>$ <T $_1$, commit>crash crash

Atualização adiada (deferred update)

```
\log (T_0 -> T_1)
                             <T<sub>o</sub>,start>
<T<sub>0</sub>,start>
                                                               <T<sub>o</sub>,start>
<T_0, A, 50>
                             <T_0, A, 50>
                                                               <T_0, A, 50>
<T_0, B, 250>
                             <T_0, B, 250>
                                                               <T_{\circ}, B, 250>
                             <T_{0}, commit>
                                                               <T_{0}, commit>
crash
                                                               <T<sub>1</sub>,start>
                             <T_1, start>
                             <T_1, C, 600>
                                                               <T_1, C, 600>
                                                               <T_{_{1}}, commit>
                             crash
                                                          drash
                             redo(T<sub>0</sub>)
                                                               redo(T<sub>0</sub>)
Nada a fazer.
                             Log atualizado.
                                                          redo(T<sub>1</sub>)
log atualizado.
```

- Atualização imediata (immediate update)
 - Pode atualizar o BD antes da transação atingir o commit (atualizações não comitadas)
 - As operações são registradas no log
 - Se uma transação falhar antes do commit, seus efeitos devem ser desfeitos (UNDO)
 - Se houver uma falha após a commit, os efeitos da transação devem ser refeitos (REDO)
 - Algoritmo de recuperação UNDO/REDO

- Atualização imediata (immediate update)
 - Funcionamento
 - Antes de T iniciar sua execução, <T, start> é registrada no log
 - A cada write(item), uma tupla <T, item, valor antigo, novo valor> é registrada também
 - Como UNDO é necessário, a tupla deve conter o valor antigo do item
 - Ao fim da transação: <T, commit> é registrada
 - Antes de escrever o valor do item no banco de dados (output (item)) o registro do log deve ser gravado em um meio estável de armazenamento

Atualização imediata (immediate update)

```
\begin{array}{l} \text{log (T}_{0}\text{->} \text{T}_{1}\text{)} \\ <\text{T}_{0}, \text{start>} \\ <\text{T}_{0}, \text{A,100,50>} \\ <\text{T}_{0}, \text{B,200,250>} \\ <\text{T}_{0}, \text{commit>} \\ <\text{T}_{1}, \text{start>} \\ <\text{T}_{1}, \text{C,700,600>} \\ <\text{T}_{1}, \text{commit>} \end{array}
```

- Atualização imediata (immediate update)
 - Em caso de falha os procedimentos redo(T) e undo(T) são executados e todos os novos itens da transação T são confirmados no banco ou os antigos retornados, respectivamente.
 - undo (T) deve ser idempotente, ou seja, pode ser executada várias vezes para a mesma transação mas o resultado deve ser sempre o mesmo
 - redo (T) deve ser idempotente, ou seja, pode ser executada várias vezes para a mesma transação mas o resultado deve ser sempre o mesmo

- Atualização imediata (immediate update)
 - O Uma transação T' é refeita se possuir as tuplas <T', start> e <T', commit> no log.
 - O Uma transação T' é desfeita se possuir a tupla <T', start> e não a tupla <T', commit>.

Atualização imediata (immediate update)

```
\log (T_0^->T_1^-)
                             <T<sub>0</sub>,start>
<T<sub>o</sub>,start>
                                                               <T<sub>o</sub>,start>
<T_0, A, 100, 50>
                             <T_0, A, 100, 50>
                                                               <T_0, A, 100, 50>
<T_0, B, 200, 250>
                                                               <T_0, B, 200, 250>
                             <T_0, B, 200, 250>
                             <T_{0}, commit>
                                                               <T_{0}, commit>
crash
                             <T<sub>1</sub>,start>
                                                               <T<sub>1</sub>,start>
                             <T_1, C, 700, 600>
                                                               <T_1, C, 700, 600>
                                                           \langle T_1, commit \rangle
                             crash
                                                          drash
```

Atualização imediata (immediate update)

```
\log (T_0^->T_1^-)
                            <T<sub>0</sub>,start>
<T<sub>o</sub>,start>
                                                             <T<sub>0</sub>,start>
<T_0, A, 100, 50>
                            <T_0, A, 100, 50>
                                                             <T_0, A, 100, 50>
<T_0, B, 200, 250>
                                                             <T_0, B, 200, 250>
                            <T_0, B, 200, 250>
                            <T_{0}, commit>
                                                             <T_{0}, commit>
crash
                            <T<sub>1</sub>,start>
                                                             <T<sub>1</sub>,start>
                            <T_1, C, 700, 600>
                                                             <T_1, C, 700, 600>
                                                             <T<sub>1</sub>,commit>
                            crash
                                                             crash
                            redo(T_0) e undo(T_{1})
                                                             redo(T<sub>0</sub>)
undo (T_0)
log atualizado Log atualizado.
                                                   redo(T<sub>1</sub>)
```

CheckPoint

- Quando ocorre uma falha é necessário percorrer o log para procurar transações para realizar o undo ou redo
 - O processo de procura consome tempo
 - A maioria das transações que precisariam ser refeitas já tiveram seus dados gravados no banco (overhead)
- Checkpoints resolvem esses problemas

CheckPoint

- Checkpoints
 - O sistema, periodicamente, executa checkpoints
 - Coloca em memória estável todos os registros que estão na memória RAM
 - Coloca no disco todos os buffers modificados
 - Coloca a tupla <checkpoint> no log
 - Quando a ação está sendo executada, nenhuma outra ação no banco pode ocorrer
 - O sistema ao encontrar um <checkpoint> "sabe"
 que as ações anteriores já foram confirmadas

CheckPoint

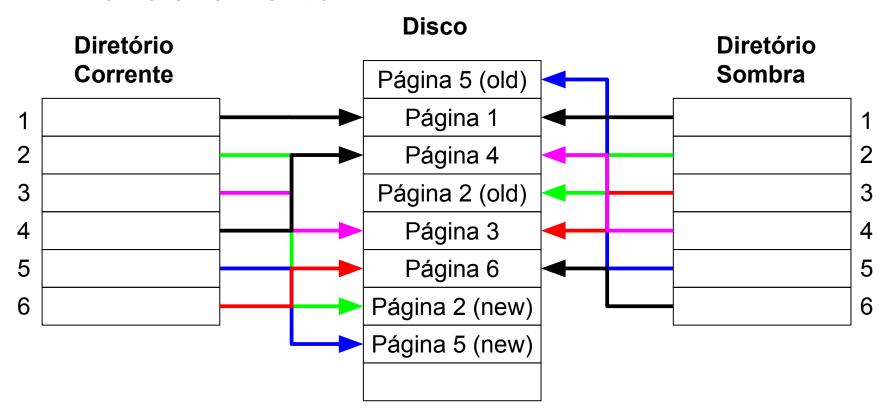
```
<start T3>
<write T3,B,15,12>
<start T2>
<write T2,B,12,18>
<start T1>
<write T1,D,20,25>
<commit T1>
<write T2,D,25,26>
<checkpoint>
<write T3,A,10,19>
<commit T3>
CRASH
```

- Técnica alternativa para recuperação que não necessita de log
- O BD é considerado como uma série de páginas no disco (ou blocos) de tamanho fixo
 - Um diretório com n entradas (n=número de páginas) no qual a i-ésima entrada aponta para a i-ésima página
 - Quando uma transação começa a ser executada, o diretório atual é copiado para um diretório sombra (que ficará em um meio não volátil)

- Funcionamento
 - Durante a execução da transação, o diretório sombra não pode ser atualizado
 - Quando um write(X) é executado, a página onde X está e copiada e o diretório atual apontará para a página com X modificado e o diretório sombra apontará para a página com X original
 - O undo é feito copiando as cópias apontados pelo diretório sombra
 - O redo é feito ignorando as cópias do diretório sombra

- Funcionamento
 - Uma transação T' executa write(X) e X está na iésima página:
 - Se a i-ésima página não está no buffer, o sistema copia para o buffer
 - Se é o primeiro write executado na i-ésima página pela transação, o sistema o modifica o diretório corrente:
 - Procura uma página livre no disco
 - Cópia a i-ésima página para a página livre (torna-se ocupada)
 - Aponta o diretório corrente para a nova página
 - O diretório sombra aponta para a página antiga

Funcionamento



- Vantagens sobre log
 - Overhead na gravação do log é eliminado
 - Recuperação de falha é mais rápida (pois não tem undo nem redo)
- Desvantagens
 - O commit de cada transação exige a atualização das página (commit overhead)

- Desvantagens (cont ...)
 - Fragmentação: as páginas do banco mudam frequentemente de lugar
 - Páginas do diretório fantasma devem ser liberadas
- A técnica que utiliza log é conhecida como Atualização no Local

Vamos complicar? Transações Concorrentes

Recuperação com Escalonamentos Concorrentes

- Um único log para as transações
 - As entradas no log são compartilhadas
- Um único buffer
 - As páginas do buffer são compartilhadas pelas transações

Recuperação com Escalonamentos Concorrentes

- O esquema de recuperação depende do esquema de controle de concorrência
 - T' executa write(X) e é abortada (o valor antigo de X deve ser restaurado)
 - T" executa write(X) sobre o X atualizado por T', como T' abortou T" terá que abortar também
 - 2PL evita esse problema

Rollback de Transações

- Uma transação T é cancelada utilizando o log
- O sistema procura (do fim para o início) as tuplas
 <T,X,Va, Vn>, e restaura X com o valor antigo Va.
 A busca para quando encontrar <T,start>
- A busca do fim para o início evita inconsistências
 - o <T,X,10,20>
 - o <T,X,20,30>
 - O valor correto de X é 10.
- Devido ao 2PL não é necessário realizar o rollback em cascata

Checkpoints

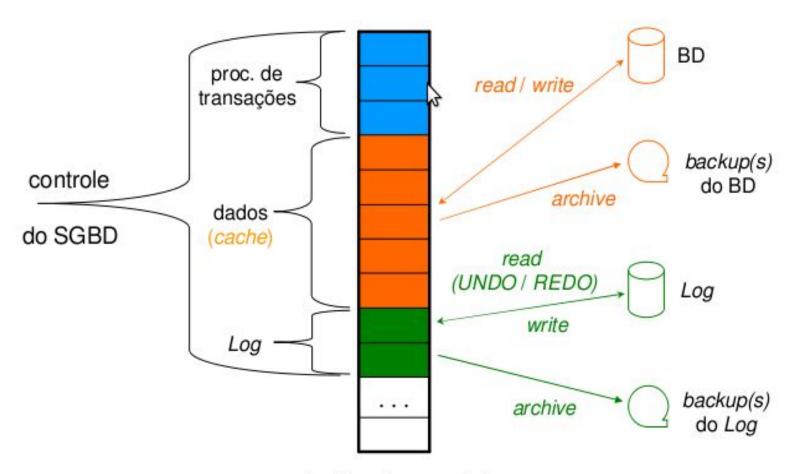
- No caso sem concorrência era considerado
 - Apenas as transações que que haviam começado após o checkpoint mais recente
 - Apenas uma transação poderia estar ativa durante um checkpoint
- Para o caso de concorrência a tupla do checkpoint deve ser modificado
 - <checkpoint, L>, onde L é a lista de transações efetivadas até o checkpoint.

Buffer

Gerenciamento de Buffer

- Buffer
 - Conjunto de blocos (páginas) da memória principal
- O SGBD gerencia os buffers
 - De dados
 - De Logs
 - Para processamento de transações
- Solicita ao SO apenas os serviços de escrita/leitura dos blocos

Gerenciamento de Buffer



buffers de memória

Fonte: Prof. Ronaldo (UFSC)

Gerenciamento de Buffer

- Os buffers de dados e log devem ser sincronizados
 - Um bloco atualizado na cache pode ser gravado apenas no BD após o histórico dos dados atualizados neste bloco ter sido gravado no Log em disco
 - Técnica chamada Write-Ahead Log (WAL)
 - Uma transação T' só pode passar para o estado efetivada (committed) após todas as suas atualizações terem sido gravadas no BD segundo o princípio WAL
- O gerenciamento, aqui apresentado, é utilizado na técnica Atualização no Local

- STEAL (Roubo)
 - O gerenciador de buffer pode roubar a página do banco
 - Pode escrever uma página arbitrária no disco e utilizar tal página para outra coisa do disco
 - O SGBD não controla a política de troca de páginas do buffer
 - A página pode conter escritas feitas por transações que não comitaram
 - Pode ser utilizado por razões de desempenho e não há necessidade de manter blocos bloqueados por transações

- NO-STEAL (sem roubo)
 - Um bloco na cache utilizado por uma transação T' não pode ser gravado antes do commit de T'
 - O bloco possui um bit de status (dirty bit) indicando se foi (1) ou não (0) modificado
 - vantagem: processo de recovery mais simples evita dados de transações inacabadas sendo gravadas no BD

- FORCE (Forçar)
 - O SGBD força todas as atualizações das transações a serem gravadas no disco antes da efetivação das mesmas
 - Os blocos que mantêm dados atualizados por uma transação T' são imediatamente gravados no BD
 - quando T' é efetivada (após o commit)
 - Deve-se saber quais os blocos que T' atualizou dados
 - Vantagem: garante a durabilidade de T' o mais cedo possível - permite o REDO de T' em caso de falha

- NO-FORCE (Sem Forçar)
 - Os blocos que mantêm dados atualizados por T' não são imediatamente gravados no BD quando T' é efetivada (após o commit)
 - Vantagem: blocos atualizados podem permanecer na cache e serem utilizados por outras transações, após a efetivação de T' (reduz custo de acesso a disco)

- Política preferida STEAL/NO-FORCE
 - Esta combinação é mais complicada porém traz os melhores desempenhos
 - STEAL (difícil implementar a durabilidade)
 - Uma transação que gravou em disco abortou
 - O sistema pode falhar antes da transação finalizar
 - Deve lembrar o valor antigo de um item (para suportar o UNDO)

- Política preferida STEAL/NO-FORCE
 - NO-FORCE (difícil implementar a atomicidade)
 - Uma transação que gravou em disco abortou
 - O sistema pode falhar antes da transação finalizar
 - Deve lembrar o valor antigo de um item (para suportar o UNDO)

Gerenciamento (Sumário)

No-Force

Force

	เทบ-อเธลเ	Sical
No-Force		Mais Rápido
Force	Mais Lento	

No Stool

Desempenho

Staal

No-Steal	Steal	
Sem Undo	Undo	
Redo	Redo	
Sem Undo	Undo	
e Redo	Sem Redo	

Implicação na recuperação e no log

Leitura

- Chapter 18 Crash Recovery (Database Management Systems - third edition – Ramakrishnan & Gerhke)
 - Verificar o capítulo equivalente na versão em português.