

Bancos de Dados Distribuídos

Prof. Frank Siqueira
Departamento de Informática e Estatística
Universidade Federal de Santa Catarina



1

Conteúdo

- Introdução aos BDs Distribuídos
- Processamento de Consultas Distribuídas
- Transações
- Controle de Concorrência
- Desenvolvimento de Aplicações

2

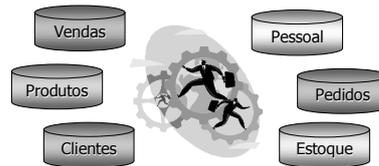
Introdução aos Bancos de Dados Distribuídos

- Motivação
- Conceitos
- Classificação de Bancos de Dados
- Arquitetura de BDs Distribuídos
- Armazenamento dos Dados

3

Motivação

- Bancos de Dados
 - Empregam tecnologia bastante sólida
 - Usados pela imensa maioria das empresas
 - Exercem papel vital na operação da empresa



4

Motivação

- A venda e a utilização de Sistemas de Bancos de Dados crescem constantemente devido à demanda gerada pelos Sistemas de Informação que deles necessitam



5

Motivação

- Sistemas que utilizam Bancos de Dados
 - SAD: Apoio à Decisão
 - SIG: Informações Gerenciais
 - ERP: Planejamento de Recursos
 - CRM: Relacionamento com o Cliente
 - BI: Inteligência de Negócio
 - etc.



6

Motivação

- Benefícios
 - Gerenciar o negócio de modo eficiente
 - Determinar o mercado-alvo de um produto
 - Definir preços, promoções e condições de compra dos produtos
 - Verificar a eficácia de campanhas de marketing
 - Otimizar a quantidade de produtos no estoque
 - Respostas rápidas a mudanças no mercado
- ... ou seja, ganhar eficiência e lucratividade

7

Motivação

- O acesso a Bancos de Dados é possível usando:



Linguagens de Consulta



Aplicações de Rede



Aplicações Gráficas



Páginas e Formulários da Web

8

Motivação

- No entanto...
 - A maioria das empresas utiliza mais de um banco de dados, muitas vezes dispersos em vários locais
 - É preciso ter uma visão integrada dos dados!
 - Operação simultânea de múltiplos BDs pode levar a dados inconsistentes em BDs diferentes
 - É preciso ter uma visão consistente dos dados!



9

Motivação

- Panorama Atual
 - Crescimento do número de usuários
 - Crescimento da quantidade de consultas
 - Automatização de todos os processos dentro de uma empresa ou instituição
 - Maior dependência dos bancos de dados
 - Novos tipos de dados, como som e imagem, exigem maior poder de processamento e armazenamento
- Soluções
 - Usar processamento paralelo e distribuído para processar as consultas em bancos de dados

10

Conceitos

- Processamento Paralelo
 - Consiste em executar simultaneamente várias partes de um mesmo processo ou aplicação
 - Processos são executadas paralelamente:
 - Em um mesmo processador
 - Em uma máquina multiprocessada
 - Em um cluster (máquinas interligadas por uma rede local que se comportam como uma só máquina)
 - Tornou-se possível a partir do desenvolvimento de sistemas operacionais multi-tarefa, multi-thread e paralelos

11

Conceitos

- Processamento Distribuído
 - Consiste em executar processos / aplicações cooperantes em máquinas diferentes
 - Processos são executadas em máquinas diferentes interligadas por uma rede
 - Redes locais
 - Internet
 - Outras redes públicas ou privadas
 - Tornou-se possível a partir da popularização das redes de computadores

12

Conceitos

- Características
 - Acoplamento
 - Sistemas paralelos são fortemente acoplados: compartilham hardware ou se comunicam através de um barramento de alta velocidade
 - Sistemas distribuídos são fracamente acoplados
 - Previsibilidade
 - O comportamento de sistemas paralelos é mais previsível; já os sistemas distribuídos são mais imprevisíveis devido ao uso da rede e a falhas

13

Conceitos

- Características (cont.)
 - Influência do Tempo
 - Sistemas distribuídos são bastante influenciados pelo tempo de comunicação pela rede; em geral não há uma referência de tempo global
 - Em sistemas paralelos o tempo de troca de mensagens pode ser desconsiderado
 - Controle
 - Em geral em sistemas paralelos se tem o controle de todos os recursos computacionais; já os sistemas distribuídos tendem a empregar também recursos de terceiros

14

Conceitos

- Vantagens
 - Usam melhor o poder de processamento
 - Apresentam um melhor desempenho
 - Permitem compartilhar dados e recursos
 - Podem apresentar maior confiabilidade
 - Permitem reutilizar serviços já disponíveis
 - Atendem um maior número de usuários
 - ...

15

Conceitos

- Dificuldades
 - Desenvolver, gerenciar e manter o sistema
 - Controlar o acesso concorrente a dados e a recursos compartilhados
 - Evitar que falhas de máquinas ou da rede comprometam o funcionamento do sistema
 - Garantir a segurança do sistema e o sigilo dos dados trocados entre máquinas
 - Lidar com a heterogeneidade do ambiente
 - ...

16

Classificação de Bancos de Dados

- Bancos de Dados Centralizados
 - Acesso ao banco de dados apenas a partir da máquina na qual o mesmo está localizado
- Bancos de Dados Cliente-Servidor
 - Os dados podem ser acessados por clientes remotos, ligados ao servidor por uma rede de comunicação
- Bancos de Dados Paralelos
 - O servidor utiliza uma máquina paralela (com vários processadores) para processar consultas/transações
- Bancos de Dados Distribuídos
 - Bancos de dados inter-relacionados localizados em diferentes servidores interconectados por uma rede

17

Classificação de Bancos de Dados

- De acordo com sua classificação, o BD requer o uso de um sistema gerenciador apropriado
- Classificação de SGBDs:
 - SGBDs Centralizados: gerenciam BDs acessados apenas por usuários com acesso à máquina local
 - SGBDs Cliente-Servidor: gerenciam o acesso a dados locais a partir de clientes remotos
 - SGBDs Paralelos: gerenciam execução de consultas/transações em paralelo, usando vários processadores
 - SGBDs Distribuídos: gerenciam a execução de consultas em dados distribuídos por vários servidores

18

Classificação de Bancos de Dados

- Bancos de Dados Cliente-Servidor
 - Tornam os dados amplamente acessíveis para os usuários
 - Permitem o acesso remoto a partir de máquinas conectadas ao servidor através da rede
 - Acesso pode ser feito utilizando:
 - Um aplicativo cliente do SGBD (Ex.: console do Oracle, *query analyzer* do SQL Server, etc.)
 - Um programa que acessa o BD usando uma API (Ex.: ODBC, JDBC, ADO, etc.)
 - Uma aplicação Web que acessa os dados no BD (navegador Web → servidor HTTP → BD)

19

Classificação de Bancos de Dados

- Bancos de Dados Paralelos
 - Máquinas paralelas vem sendo usadas para suportar uma carga maior de trabalho dos SDBs
 - Vários processadores executam as operações em paralelo → uso de controle de concorrência
 - O usuário não vê diferença entre um BD paralelo e centralizado → transparência
 - Memória e disco podem ser compartilhados ou não
 - Custo de máquinas multi-processadas está caindo
 - Clusters podem ser montados a um custo reduzido utilizando PCs e software de gerenciamento adequado

20

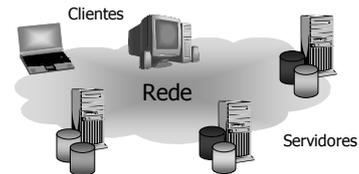
Classificação de Bancos de Dados

- Bancos de Dados Paralelos (cont.)
 - Paralelismo pode ser usado:
 - Na entrada e saída de dados (operações de I/O)
 - No processamento de consultas
 - No processamento de operações individuais
 - Como resultado, pode-se obter um aumento de:
 - Escala: processar mais transações
 - Desempenho: tornar o processamento mais rápido
 - Problema: mais hardware → mais falhas
 - Solução: replicar dados → controlar consistência

21

Classificação de Bancos de Dados

- Bancos de Dados Distribuídos
 - Os dados estão distribuídos por várias localidades
 - Localização dos dados é transparente para o usuário
 - Máquinas se comunicam via rede para acessar os dados e processar as consultas



22

Classificação de Bancos de Dados

- BDs Distribuídos x Centralizados
 - BDs centralizados possuem um ponto único de falha
 - BDs distribuídos podem aumentar a robustez do sistema e a disponibilidade dos dados
 - BDs centralizados são limitados pela capacidade de processamento e armazenamento de uma máquina
 - BDs distribuídos podem crescer em escala adicionando novos servidores ao sistema
 - BDs distribuídos são mais sujeitos a apresentar falhas parciais de funcionamento e de segurança
 - BDs distribuídos são mais difíceis de administrar pois os servidores estão em locais diferentes

23

Classificação de Bancos de Dados

- BDs Distribuídos x Cliente-Servidor
 - BDs distribuídos também podem ser acessados remotamente por clientes
- BDs Distribuídos x Paralelos
 - Ambos podem processar consultas em paralelo usando os vários processadores disponíveis
 - Em BDs paralelos os processadores podem trocar dados usando discos ou memória compartilhada
 - O uso da rede em BDs distribuídos pode prejudicar o desempenho do sistema ao processar consultas
 - BDs paralelos são mais vulneráveis a algumas falhas, e não são tão escaláveis quanto BDs distribuídos

24

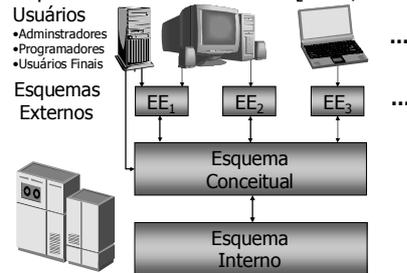
Classificação de Bancos de Dados

- BDs distribuídos podem ser classificados como:
 - Homogêneos: todos os *sites* usam o mesmo *software*
 - Heterogêneos
 - Usam *software* diferente
 - Podem usar esquemas de dados diferentes
 - Linguagem de consulta pode ser diferente
- Dificuldades adicionais em BDs distribuídos
 - Controlar a consistência dos dados armazenados
 - Processar consultas e transações de modo distribuído
 - Controlar a concorrência e resolver conflitos
 - Conciliar as diferenças em sistemas heterogêneos

25

Arquitetura de BDs Distribuídos

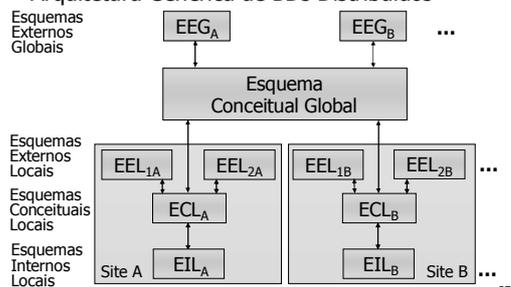
- Arquitetura Genérica de BDs [ANSI/SPARC]



26

Arquitetura de BDs Distribuídos

- Arquitetura Genérica de BDs Distribuídos



27

Arquitetura de BDs Distribuídos

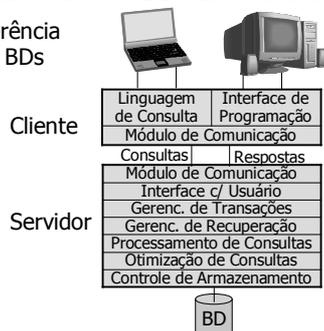
- Usuários de BDs Distribuídos

- Tipos de Usuários
 - Administradores locais ou globais
 - Programadores locais ou globais
 - Usuários finais locais ou globais
- Usuários Globais
 - Possuem visão global do sistema
 - Visualizam um esquema externo global
- Usuários Locais
 - Acessam diretamente o servidor local
 - Visualizam um esquema externo local

28

Arquitetura de BDs Distribuídos

- Gerência de BDs



29

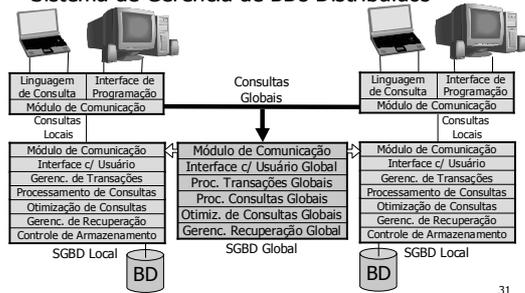
Arquitetura de BDs Distribuídos

- Problemas em BDs distribuídos
 - Otimização e processamento de consultas/transações distribuídas requer algoritmos/protocolos adequados
 - Mecanismos de controle e gerenciamento devem trabalhar de maneira integrada
- Problemas em BDDs heterogêneos
 - Precisamos conciliar as diferenças entre:
 - Modelos lógicos
 - Linguagens de definição e manipulação de dados
 - Formatos de dados: língua, ordenação dos bits, tamanho dos tipos de dados e representação na memória, tabelas de caracteres, etc.

30

Arquitetura de BDs Distribuídos

▪ Sistema de Gerência de BDs Distribuídos



31

Arquitetura de BDs Distribuídos

▪ Construção de BDs com arquitetura distribuída

- Abordagens de projeto
 - *Top-Down*: usada em sistemas construídos do zero, geralmente homogêneos
 - *Bottom-Up*: usada quando os sistemas locais já estão instalados e precisam ser integrados
- Questões a serem respondidas no projeto de BDDs
 - Onde colocar os *sites*?
 - Como fragmentar dados e distribuí-los pelos *sites*?
 - Replicação de dados é necessária?
 - Como integrar diferentes sistemas?

32

Armazenamento dos Dados

- Em BDs distribuídos os dados podem ser:
 - Replicados
 - Fragmentados
 - Replicados e Fragmentados
- Replicação de Dados
 - Uma mesma tabela pode ser armazenada em mais de um servidor
 - Vantagens: aumenta a disponibilidade e o paralelismo
 - Desvantagem: atualizações devem ser feitas em todos os servidores para manter consistência entre réplicas
 - Apresenta bom desempenho nas operações de leitura, mas causa *overhead* nas operações de escrita

33

Armazenamento dos Dados

- Localização das Réplicas
 - Deve levar em conta os locais e usuários que acessam os dados replicados com maior frequência
 - Deve ser transparente para o usuário
- Atualização de Réplicas
 - *Snapshot*: é feita a cópia completa das tabelas replicadas, podendo ocorrer atualizações periódicas
 - Incremental: os dados alterados são transferidos pela rede em um horário programado
 - Transacional: dados são atualizados nas réplicas no instante em que são modificados

34

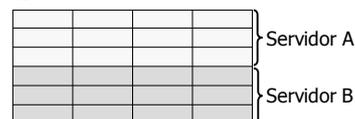
Armazenamento dos Dados

- Fragmentação de Dados
 - Os dados mantidos em uma tabela podem ser divididos em dois ou mais fragmentos
 - Cada fragmento é armazenado em servidor do banco de dados distribuído
 - Fragmentação deve ser transparente para o usuário, que tem visão completa da tabela
 - Métodos usados para fragmentação:
 - Fragmentação horizontal
 - Fragmentação vertical
 - Fragmentação mista

35

Armazenamento dos Dados

- Fragmentação Horizontal
 - Cada fragmento contém um subconjunto das tuplas da relação completa
 - Cada tupla de uma relação precisa ser armazenada em pelo menos um servidor
 - A relação completa pode ser obtida fazendo a união dos fragmentos



36

Armazenamento dos Dados

- Técnicas para fragmentação horizontal:
 - *Round-Robin*
 - Tuplas distribuídas uma a uma entre os servidores
 - Divide os dados igualmente
 - *Hash*
 - Usa uma função *hash* para determinar o servidor
 - Boa para processar seleções de igualdade
 - Por faixa
 - Cada servidor armazena as tuplas que estiverem dentro de uma faixa de valores
 - Boa para processar seleções de igualdade e para procurar por faixas de valores

37

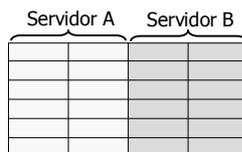
Armazenamento dos Dados

- Tipos de Fragmentação Horizontal
 - Fragmentação Horizontal Primária
 - O(s) atributo(s) usados para fragmentar a tabela fazem parte da mesma
 - Fragmentação Horizontal Derivada
 - O(s) atributo(s) usados para fragmentar a tabela pertencem a outra tabela
 - As tabelas precisam estar relacionadas através de algum atributo chave
 - Pode exigir mais processamento e acesso a disco para execução de consultas

38

Armazenamento dos Dados

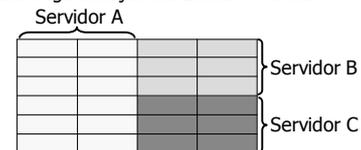
- Fragmentação Vertical
 - Relações são decompostas em conjuntos de atributos mantidos em servidores diferentes
 - Cada fragmento é uma projeção da relação completa
 - A relação completa pode ser obtida fazendo a junção de todos os fragmentos



39

Armazenamento dos Dados

- Fragmentação Mista
 - Combina fragmentação horizontal e vertical



- Fragmentação e Replicação de Dados
 - Dados são fragmentados horizontal ou verticalmente
 - Cada fragmento é mantido em mais de um servidor

40