

Integração de Banco de Dados

IN1128/IF694 – Bancos de Dados Distribuídos e Móveis
 Ana Carolina Salgado – acs@cin.ufpe.br
 Bernadette Farias Lôscio – bfl@cin.ufpe.br



Cin.ufpe.br

Formas Básicas de BD Distribuído

- **Abordagem top-down**
 - Adequada para bancos de dados homogêneos
- **Abordagem bottom-up**
 - Adequada para sistemas de múltiplos bancos de dados



Abordagem Top-down

- **Projeto de um novo banco de dados**
 - Busca da **redução dos custos** e da melhoria do **desempenho** das aplicações
 - **Limites da capacidade** de um banco de dados centralizado atingidos
 - Armazenamento
 - Processamento
 - Tempo de acesso



Abordagem Top-down (Distribuição de Dados)

- **Particionar o BD em unidades lógicas (fragmentos)**
 - **Fragmentação**
 - É um processo para o **particionamento de BD em banco de dados menores sem replicação**
 - O resultado é um conjunto de fragmentos de relações, os quais devem ser **alocados aos diferentes bancos de dados locais**
 - **Replicação**
 - Útil para melhorar a disponibilidade de dados
 - Replicação Total x Replicação Parcial



Abordagem Bottom-up

- **Processo de integração de bancos de dados já existentes e independentes**
 - Ambientes onde proliferaram os **bancos de dados departamentais e individuais**
 - Integração gera **bancos de dados globais virtuais**
 - Preserva os investimentos já feitos em aplicações, possibilitando ainda uma **visão integrada** dos dados fisicamente dispersos



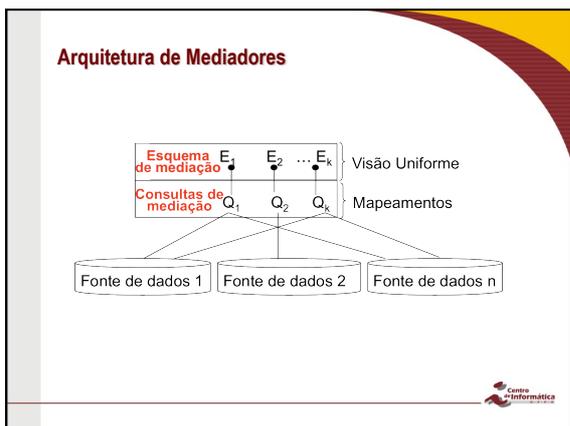
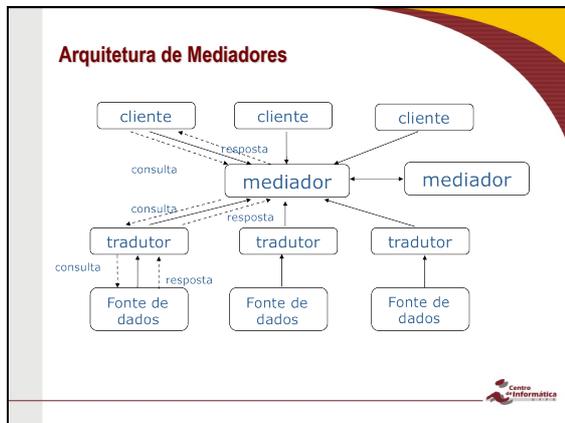
Abordagem Bottom-up

- **Razões**
 - **Sistemas Legados**
 - Organizações possuindo sistemas desenvolvidos há algum tempo (até 1990, aproximadamente)
 - Tecnologias aplicadas bastante diferentes das atuais
 - **Migração entre Plataformas**
 - Evolução de hardware e software disponibilizando novas facilidades
 - **Mudanças Organizacionais**
 - Fusão de duas ou mais empresas
 - Visão integrada da nova organização
 - **Evolução da Tecnologia de Redes de Comunicação**
 - Integração dos ambientes computacionais anteriormente isolados



Abordagem Bottom-up

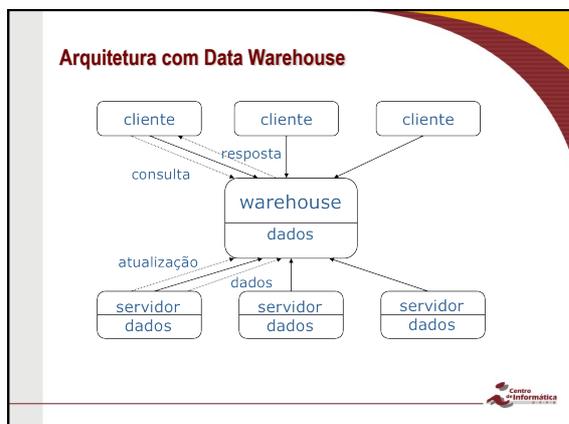
- Integrar os bancos de dados e esquemas locais existentes para um único banco de dados global com seu respectivo esquema global
- Principais arquiteturas
 - Mediadores
 - Data Warehouse

- ### Arquitetura de Mediadores
- O mediador integra dados de diferentes fontes dedados
 - O tradutor é um software que converte os dados das fontes de dados para o modelo de dados comum e converte consultas de aplicações em consultas específicas da fonte de dados correspondente
 - Um tradutor provê uma interface única para uma ou mais fontes de dados que tenham uma interface comum de acesso
 - Fontes de dados baseadas em um mesmo modelo de dados e em uma mesma linguagem de acesso são associados a um mesmo tradutor
- 

- ### Abordagem Virtual
- Os dados são mantidos nas fontes e as consultas são decompostas em tempo real e submetidas às diversas fontes
 - Vantagem: os dados não são replicados e tem-se a garantia de estarem atualizados no momento da consulta
 - Desvantagem: como as fontes de dados são autônomas, são necessários métodos para otimização de consultas para garantir uma performance adequada
- 

- ### Abordagem Virtual
- O Enfoque Virtual é mais apropriado para as seguintes situações:
 - o número de fontes é grande
 - os dados são atualizados frequentemente
 - existe pouco controle sobre as fontes de dados
- 



Abordagem Materializada

- Os dados gerados das diversas fontes são carregados em um repositório (warehouse) e as consultas são aplicadas a estes dados.
 - Vantagem: performance garantida no momento da consulta
 - Desvantagem: atualização do repositório sempre que houver mudança nos dados

Metodologia Bottom-up

- O esquema global pode ser definido a partir dos requisitos dos usuários ou pode ser obtido a partir da **integração dos esquemas locais**
- **Integração de esquemas**
 - É o processo de gerar um esquema integrado (global) a partir de um conjunto de esquemas de entrada (esquemas locais) resolvendo as diversidades estruturais e semânticas existentes entre eles

Integração de Esquemas - Heterogeneidade

- Uma organização pode ter múltiplos SGBDs e diferentes departamentos
 - Dentro de uma mesma organização podemos ter diferentes requisitos para um mesmo conjunto de dados
- Diferentes perspectivas podem levar a diferentes representações (estrutura e restrições) e interpretações de um mesmo dado
- Diferentes modelos de dados oferecem primitivas diferentes para a modelagem dos dados
 - Exemplo: um endereço pode ser modelado como uma entidade em um esquema e como um atributo composto em outro esquema

Integração de Esquemas - Heterogeneidade

- **Problema crítico!**
 - Heterogeneidade terminológica
 - Termos diferentes usados para representar os mesmo conceitos
 - Termos iguais usados para representar conceitos diferente
 - Heterogeneidade estrutural
 - Conceitos similares representados através de construtores divergentes

Integração de Esquemas - Heterogeneidade

- **Heterogeneidade terminológica - Exemplo**
 - Entidades semanticamente relacionadas podem ser representadas com nomes diferentes em bancos de dados distintos

Sinonímia

| Esquema | Definição |
|---------|---|
| S1 | medico1 (num_conselho, nome, especialidade) |
| S2 | doutor2 (CRM, nome) |

Integração de Esquemas - Heterogeneidade

■ Heterogeneidade terminológica - Exemplo

- Dois atributos são semanticamente semelhantes e apresentam nomes diferentes (sinônimos) ou quando atributos com mesmo nome podem não ter qualquer relação semântica (homônimos)

| | Esquema | Definição |
|-----------|---------|---|
| Sinonímia | S1 | medico1 (num_conselho, nome, especialidade) |
| | S2 | doutor2 (CRM, nome) |
| Homonímia | S1 | medico1 (num_conselho, nome, especialidade) |
| | S2 | enfermeiro2 (num_conselho, nome) |



Integração de Esquemas - Heterogeneidade

■ Heterogeneidade estrutural - Exemplo

- Duas entidades são representadas em dois bancos de dados com níveis diferentes de abstração (generalização)

| Esquema | Definição |
|---------|---|
| S1 | enfermeiro1 (num_conselho, nome) |
| S1 | medico1 (num_conselho, nome, especialidade) |
| S2 | agente-saude2 (num_conselho, nome, tipo_agente) |



Integração de Esquemas

- Ferramentas de integração de esquemas são semi-automáticas
- Sugerem conceitos candidatos a serem integrados
 - O especialista aceita ou rejeita



Integração de Esquemas - Etapas

Esquemas locais

Pré-Integração

Esquemas canônicos

Comparação

Esquemas canônicos + equivalência entre conceitos de esquemas

Unificação

Esquemas integrado + mapeamentos para esquemas locais

Reestruturação

Esquemas integrado + mapeamentos para esquemas locais



Pré-Integração

- Nesta fase os esquemas locais são traduzidos para o modelo de dados comum (ex: modelo ER)
- Análise dos esquemas locais a fim de definir a política de integração, por exemplo, a ordem a ser seguida na integração de esquemas, se a integração é parcial ou total



Tradução de Esquemas

- Um sistema de bancos de dados múltiplos deve oferecer suporte para a tradução entre modelos de dados locais e globais
- Quando bancos de dados heterogêneos são integrados os esquemas locais são traduzidos para um modelo de dados comum

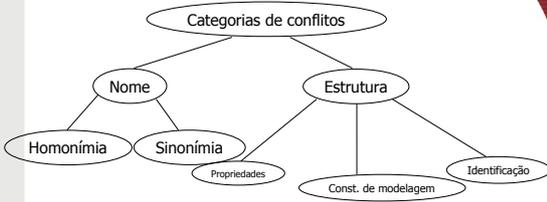


Comparação

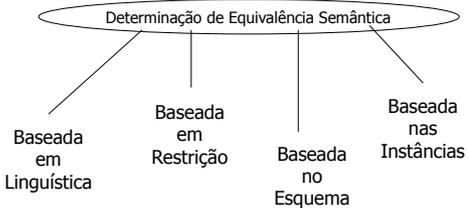
- Identificação de correspondências e conflitos entre esquemas diferentes através da análise de suas características como atributos, relacionamentos e restrições de integridade
- O resultado final desta etapa é um conjunto de correspondências entre os conceitos dos esquemas



Comparação




Comparação




Comparação

- Diferentes técnicas podem ser aplicadas na determinação de equivalências semânticas entre esquemas
- Baseada em linguística: considera afinidade de nomes ou descrições textuais entre conceitos
 - Utiliza informações adicionais fornecidas por um dicionário



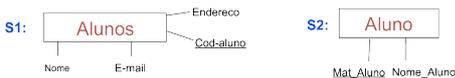
Comparação

- Baseada em restrição: tem como base características estruturais e de domínio entre conceitos e suas propriedades
- Baseada no esquema: apenas informações sobre os esquemas são consideradas
- Baseada nas instâncias: usa informações baseadas no conteúdo dos dados



Comparação - Identificação das correspondências

- Dois elementos de dois bancos de dados são equivalentes se eles descrevem o mesmo objeto do mundo real




Comparação - Identificação das correspondências

- De tipos
- De atributos
- De caminho



Comparação - Identificação das correspondências

- Entre elementos correspondentes com estruturas semelhantes

S1: Alunos (Atributos: Nome, E-mail, Endereço, Cod-aluno)

S2: Aluno (Atributos: Mat_Aluno, Nome_Aluno)



Comparação - Identificação das correspondências

- Correspondência de Entidades

S1: Cursos (Atributos: Nome, Cod-curso, Grade)

S2: Curso (Atributos: desc_curso, turno, cod_curso, Carga_horaria)

S1.Cursos \cong S2.Curso



Comparação - Identificação das correspondências

- Correspondência de Relacionamento

S1: Professores \xrightarrow{n} lotar $\xrightarrow{1}$ Departamentos

S2: Departamento $\xrightarrow{1}$ dep_prof \xrightarrow{n} Professor

S1.lotar \cong S2.dep_prof



Comparação - Identificação das correspondências

- Correspondência de Atributos
- S1.Cursos.Cod-curso \cong S2.Curso.cod_curso
- S1.Cursos.Grade \cong S2.Curso.desc_curso

S1: Cursos (Atributos: Nome, Cod-curso, Grade)

S2: Curso (Atributos: desc_curso, turno, cod_curso, Carga_horaria)



Comparação - Identificação das correspondências

- Correspondência de Caminhos
- S1.Departamento-Dep_centro-Centro-Cod-centro \cong S2.Departamento-Cod-centro

S1: Departamento (Atributos: Nome, Cod-depto, fone) \xrightarrow{n} Dep_centro $\xrightarrow{1}$ Centro (Atributos: Nome, Cod-centro)

S2: Departamento (Atributos: Nome, Cod-depto, Cod-centro)



Comparação - Identificação das correspondências

- Entre elementos correspondentes com estruturas diferentes
 - Entidade e Atributo
 - S1.Autor \equiv S2.Livro.nome_autor

Comparação - Identificação das correspondências

- Entidade e Relacionamento
- S1.RCasamento \equiv S2.Casamento

Comparação - Identificação das correspondências

Passo 1: Identificação de correspondências entre elementos com estruturas correspondentes

- Identificação das entidades correspondentes
 - Identificação dos atributos correspondentes
 - Especificação dos identificadores
 - Identificação da correspondência entre as extensões

Comparação - Identificação das correspondências

(cont.)

- Identificação dos relacionamentos correspondentes
 - Identificação dos atributos correspondentes
 - Identificação das entidades participantes correspondentes
 - Identificação da correspondência entre as extensões

Passo 2: Identificação de correspondências entre elementos com estruturas diferentes

- Ex: Entidade/Relacionamento e Entidade/Atributo

Unificação

- Resolução de conflitos entre os conceitos dos esquemas locais e produção do esquema integrado ou global
- Informações de mapeamento entre os conceitos dos esquemas locais e os conceitos do esquema global são identificadas
 - Estes mapeamentos são importantes para a transformação de consultas

Unificação

Unificação

- **Binária: a unificação ocorre entre pares de esquemas locais**
 - Redução da complexidade de comparação e unificação em cada passo
 - Iterativa: unifica-se seqüencialmente pares de esquemas
 - Balanceada: unifica-se em paralelo pares de esquema



Unificação

- **N-ária**
 - A unificação ocorre entre mais de dois esquemas, podendo ser realizada em um ou vários passos
 - Redução do número de passos de unificação



Unificação

- **Os seguintes critérios devem ser atendidos:**
 - Completude e corretude: o esquema global deve conter todos os conceitos dos esquemas locais de forma semanticamente correta
 - Minimalidade: um conceito representado em mais de um esquema local deve ser representado uma única vez no esquema global
 - Compreensão: o esquema global deve ser fácil de ser compreendido (a representação mais clara deve ser escolhida)



Reestruturação

- **Nesta etapa o esquema global é analisado e reestruturado para remover redundâncias**
- **Ao final do processo, são gerados:**
 - Esquema global
 - Mapeamentos entre o esquema global e os esquemas locais



Abordagens para definição dos mapeamentos entre o esquema global e os esquemas locais

- **Visão Global (Global as View - GAV)**
- **Visão Local (Local as View - LAV)**
- **Abordagens que combinam características de GAV e LAV**
 - GLAV
 - BAV (Both as view)



Abordagens para definição dos mapeamentos

- **Visão Global**
 - Cada elemento no esquema de mediação é definido como uma visão sobre os esquemas locais
 - A reformulação de consultas torna-se mais simples
 - Não é adequada para sistemas em evolução



Abordagens para definição dos mapeamentos

- **Visão Local**
 - Cada elemento em uma fonte local é definido como uma visão sobre o esquema de mediação
 - Torna mais fácil a manutenção das fontes de dados
 - O processo de decomposição de consultas é mais complexo

Abordagens para definição dos mapeamentos – Exemplo

- **Fonte de dados1:**
 - Estudante₁(mat₁, nome₁, curso₁, nota₁)
- **Fonte de dados2:**
 - Estudante₂(mat₂, nome₂, cod_aval₂)
 - Avaliação₂(cod_aval₂, curso₂, aval_escrita₂)
- **Esquema de mediação:**
 - Estudante_m(mat_m, nome_m, curso_m, nota_m, aval_escrita_m)

Exemplo – Visão Global

Exemplo – Visão Global

- **Semântica dos mapeamentos - Fonte 1**

Para cada e em **Estudante_m**
 Existe e_1 em **Estudante₁**
 Com $e_1.mat_1 = e.mat_m$ e
 $e_1.nome_1 = e.nome_m$ e
 $e_1.curso_1 = e.curso_m$ e
 $e_1.nota_1 = e.nota_m$

Exemplo – Visão Global

- **Semântica dos mapeamentos - Fonte 2**

Para cada e em **Estudante_m**
 Existe e_2 em **Estudante₂** e a_2 em **Avaliação₂**
 Onde $e_2.cod_aval = a_2.cod_aval$
 Com $e_2.mat_2 = e.mat_m$ e
 $e_2.nome_2 = e.nome_m$ e
 $a_2.curso_2 = e.curso_m$ e
 $a_2.aval_escrita_2 = e.aval_escrita_m$

Exemplo – Visão Local

Exemplo – Visão Local

■ Semântica dos mapeamentos – Fonte 1

Para cada e_1 em **Estudante₁**
 Existe e em **Estudante_m**
 Com $e.mat_m = e_1.mat_1$ e
 $e.nome_m = e_1.nome_1$ e
 $e.curso_m = e_1.curso_1$ e
 $e.nota_m = e_1.nota_1$



Exemplo – Visão Local

■ Semântica dos mapeamentos - Fonte 2

Para cada e_2 em **Estudante₂** e a_2 em **Avaliação₂**
 Onde $e_2.cod_aval = a_2.cod_aval$
 Existe e em **Estudante_m**
 Com $e.mat_m = e_2.mat_2$ e
 $e.nome_m = e_2.nome_2$ e
 $e.curso_m = a_2.curso_2$ e
 $e.aval_escrita_m = a_2.aval_escrita_2$



O uso de mapeamentos no processo de reescrita de consultas

- Os mapeamentos especificam como obter os elementos do esquema de mediação
- A partir dos mapeamentos podem ser geradas expressões de consulta que facilitam o processo de reescrita de consultas



O uso de mapeamentos no processo de reescrita de consultas

- Passos do processo de execução de uma consulta
 - Identificar quais relações estão sendo consultadas
 - Descobrir as definições das relações (visões)
 - Reformular e submeter a consulta para as fontes locais
 - Integrar os resultados



Reescrita de Consultas - Exemplo

Consulta 1 - Esquema de mediação

For e in **Estudante_m**
 Return ($e.mat_m, e.nome_m, e.curso_m, e.nota_m, e.aval_escrita_m$)

Consulta 1 - Fonte de Dados1

For e_1 in **Estudante₁**
 Return ($e_1.mat_1, e_1.nome_1, e_1.curso_1, e_1.nota_1$)

$Estudante_m.mat_m \equiv Estudante_1.mat_1$
 $Estudante_m.nome_m \equiv Estudante_1.nome_1$
 $Estudante_m.curso_m \equiv Estudante_1.curso_1$
 $Estudante_m.nota_m \equiv Estudante_1.nota_1$



Reescrita de Consultas - Exemplo

Consulta 1 - Esquema de mediação

For e in **Estudante_m**
 Return ($e.mat_m, e.nome_m, e.curso_m, e.nota_m, e.aval_escrita_m$)

Consulta 1 - Fonte de Dados2

For e_2 in **Estudante₂**, a_2 in **Avaliação₂**
 Where $e_2.cod_aval = a_2.cod_aval$
 Return ($e_2.mat_2, e_2.nome_2, e_2.curso_2, e_2.cod_aval_2$)

$Estudante_m.mat_m \equiv Estudante_2.mat_2$
 $Estudante_m.nome_m \equiv Estudante_2.nome_2$
 $Estudante_m.curso_m \equiv Estudante_2.cod_aval_2.Avaliação_2.curso_2$
 $Estudante_m.aval_escrita_m \equiv Estudante_2.cod_aval_2.Avaliação_2.aval_escrita_2$



Integração de dados

- Os resultados das consultas locais devem ser integrados para produzir a visão integrada solicitada pelo usuário
- Conflitos de dados podem surgir durante o processo de integração e devem ser tratados corretamente



Integração de dados

- Tipos de conflito:
 - Valores de dados
 - Dois atributos semanticamente equivalentes com valores distintos
 - Tipos de dados
 - O atributo CPF é representado como string em uma fonte e como inteiro em uma outra fonte
 - Escala
 - O atributo área é definido em m² em uma base de dados e em cm² em uma outra base



Integração de dados

- Tipos de conflito:
 - Precisão dos dados
 - Os atributos nota e conceito (um conceito pode corresponder a um conjunto de notas)
 - Valor default
 - O atributo maioridade pode ter como valor default 18 ou 21 anos
 - Restrições
 - Em uma fonte de dados o intervalo de valores aceitáveis para a frequência cardíaca em um ser humano é 60-100 bpm (batimentos por minuto), em outra fonte o intervalo é 50 a 110bpm



Integração de dados

- Funções de transformação podem ser usadas para resolver os conflitos de dados
- Exemplos:
 - `converteDataParaFormatoPadrao(string formato, string data)`
 - `converteParaMetro(string unidade, int valor)`
 - `converteNotaParaConceito(double nota, string conceito)`

