

PONTIFÍCIA UNIVERSIDADE CATÓLICA  
DO RIO DE JANEIRO



**Geórgia Regina Rodrigues Gomes**

**Integração de Repositórios de Sistemas de Bibliotecas  
Digitais e de Sistemas de Aprendizagem**

**Tese de Doutorado**

Tese apresentada como requisito parcial para  
obtenção do título de Doutor pelo Programa de  
Pós-Graduação em Informática da PUC-Rio.

Orientadores: Prof. Rubens Nascimento Melo  
Prof. Sean Wolfgang Matsui Siqueira  
Prof<sup>a</sup>. Maria Helena Lima Baptista Braz

Rio de Janeiro  
Setembro de 2006

**Geórgia Regina Rodrigues Gomes**

**Integração de Repositórios de Sistemas de Bibliotecas Digitais  
e Sistemas de Aprendizagem**

Tese apresentada como requisito parcial para obtenção do título de Doutor pelo Programa de Pós-Graduação em Informática da PUC-Rio. Aprovada pela Comissão Examinadora abaixo assinada

**Prof. Rubens Nascimento Melo**

Orientador  
PUC-Rio

**Prof<sup>a</sup>. Maria Helena Lima Baptista Braz**

Co-Orientador  
Instituto Superior Técnico - Lisboa

**Prof. Sean Wolfgang Matsui Siqueira**

Co-Orientador  
UNIRIO

**Prof. Hugo Fuks**

PUC-Rio

**Prof. Emmanuel Piseces Lopes Passos**

IME-RJ

**Prof. Maria Carmen Romcy de Carvalho**

Universidade Católica de Brasília - UCB

**Prof. Luiz Antonio de Moraes Pereira**

Banco Central do Brasil

**Prof. José Eugenio Leal**

Coordenador Setorial do Centro Técnico Científico - PUC-Rio

Rio de Janeiro, 29 de setembro de 2006.

Todos os direitos reservados. É proibida a reprodução total ou parcial do trabalho sem autorização da universidade, do autor e dos orientadores.

### **Geórgia Regina Rodrigues Gomes**

Graduou-se em Matemática pela FAFITA em 1989. Obteve o grau de Mestre em Informática, pela Pontifícia Universidade Católica do Rio de Janeiro (PUC-Rio) em 1999. Trabalhou de 1992 a 2003 como coordenadora da seção de automação da DBD/PUC-Rio. Desenvolveu junto com a PUCPR o sistema Pergamum – Sistema Integrado de Bibliotecas utilizado hoje com mais de 150 instituições no Brasil. Foi coordenadora acadêmica de dois cursos a distância na PUC-Rio. Atualmente é Professora Adjunta e Pesquisadora da Universidade Cândido Mendes em Campos dos Goytacazes, atuando principalmente nas seguintes áreas: Banco de Dados, Bibliotecas Digitais, Ensino a Distância, Integração de Dados, Padrões de Metadados e Recuperação da Informação.

#### Ficha Catalográfica

Gomes, Geórgia Regina Rodrigues

Integração de repositórios de sistemas de bibliotecas digitais e sistemas de aprendizagem / Geórgia Regina Rodrigues Gomes ; orientador: Rubens Nascimento ; co-orientadores: Sean Wolfgang Matsui Siqueira, Maria Helena Lima Baptista Braz. – 2006.

143 f. ; 30 cm

Tese (Doutorado em Informática)–Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2006.

Incluí referências bibliográficas

1. Informática – Teses. 2. Bibliotecas digitais. 3. Educação baseada na web. 4. Objetos de aprendizagem. 5. Integração de dados. 6. Mineração de texto. 7. Ontologia. 8. Banco de dados. I. Melo, Rubens Nascimento. II. Siqueira, Sean Wolfgang Matsui. III. Braz, Maria Helena Lima Baptista. IV. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. V. Título.

CDD: 004

A Deus, pela graça alcançada.  
Ao meu marido Joney e meu filho Diogo.  
Aos meus pais Jorge e Aparecida.  
À minha sogra e meu sogro.

## **Agradecimentos**

A Deus, que ilumina meus caminhos e me deu forças a cada minuto para que chegasse até o fim, sem Ele seria impossível.

Ao meu marido Joney Junior, pelo apoio, ajuda e principalmente, pela compreensão e paciência ao longo da realização deste trabalho.

Ao meu filho Diogo, luz que ilumina minha vida, pela aceitação dos momentos que não pude estar ao seu lado, por causa deste trabalho.

Aos meus pais, Aparecida e Jorge, pelo amor, incentivo e carinho que sempre me dedicaram em todos os momentos da minha vida.

À minha sogra D. Lurdinha e meu sogro Sr. Joney, pela ajuda nos momentos mais importantes deste trabalho, cuidando do meu filho com carinho e amor na minha ausência e pela honra de tê-los como sogros.

As minhas irmãs, Mara, Alzira e Rita, por serem minhas irmãs.

Ao meu orientador, Rubens Nascimento Melo, pelos ensinamentos transmitidos durante todos estes anos, apoio e incentivo nos momentos de desânimo, e por ser este ser humano maravilhoso.

Ao Sean Wolfgang Matsui Siqueira, meu co-orientador, pela grande determinação e ajuda durante todo o processo de desenvolvimento desta tese, sem ele seria quase impossível.

A Maria Helena Lima Baptista Braz, minha co-orientadora, pela força e atenção dispensados a mim, principalmente no mês da entrega deste trabalho, sem ela seria muito difícil.

À minha amiga e irmã de coração Diva de Souza e Silva Rodrigues, que em

todos os momentos difíceis desde o mestrado, estava ao meu lado me apoiando e incentivando.

Aos amigos Roberto Rodrigues (marido de minha amiga Diva) e Elizabeth Vitória (minha grande amiga e irmã), pelo tempo dispensado e a grande boa vontade de revisar o texto final da tese, que Deus os mantenha sempre assim.

Aos amigos do TecBD, Simone Leal de Moura, Carlos Eduardo Portela, Álvaro César Pereira Barbosa, Carolina de Lima Aguiar, Luiz Antônio de Moraes Pereira, Fábio André Machado Porto, Fernanda Lima, Julita Glória Machado Cravo, Paulo Sérgio Simões de Araujo, Cássia Blondet Baruque, Sandra Dias de Souza e Fábio Coutinho.

Aos Professores Hugo Fuks, Emmanuel Piseces Lopes Passos, Maria Carmen Romcy de Carvalho, Luiz Antonio de Moraes Pereira, por aceitarem a participar desta banca.

A todos os Professores do Departamento de Informática pelos valiosos ensinamentos ministrados.

Aos amigos, colegas, professores e funcionários da PUC-Rio, que, a seu modo, mesmo que às vezes sem saber, ajudaram direta ou indiretamente na realização deste trabalho.

À PUC-Rio pela bolsa de isenção concedida até o penúltimo semestre deste curso.

À UCAM-Campos, que me abriu as portas num momento que eu precisava muito, isto também teve influência no término deste trabalho.

Ao meu aluno Igor, pela ajuda na implementação da aplicação de extração de informação.

A todas as pessoas que contribuíram direta ou indiretamente para a realização deste trabalho.

## Resumo

Gomes, Geórgia R. R.. **Integração de Repositórios de Sistemas de Bibliotecas Digitais e Sistemas de Aprendizagem**. PUC-Rio, 2006.143p. Tese de Doutorado – Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Com o uso generalizado das tecnologias de informação no apoio ao ensino, é comum disponibilizar conteúdos digitais, seja através de Sistemas de Bibliotecas Digitais (DLMS) ou de Sistemas de Gerência de Aprendizagem (LMS). No entanto, estes sistemas funcionam de forma independente, têm características diferentes e manipulam tipos diferentes de materiais, sendo seus repositórios com dados e metadados heterogêneos e distribuídos. Os conteúdos destes repositórios seriam melhor aproveitados se estivessem integrados a um ambiente comum, ou fossem acessados de modo integrado a partir dos ambientes de DLMS e LMS. Nesta tese é apresentada uma visão homogênea dos conteúdos de DLMS e LMS. Para esta homogeneização utilizou-se uma extensão da arquitetura de mediadores e tradutores que trata a integração de metadados, assim como ontologias para tratamento semântico. Foram consideradas ontologias locais para descrever os metadados de cada repositório e uma ontologia global para a integração. No entanto, os documentos dos repositórios dos DLMS tendem a ser monolíticos e não têm um enfoque na reutilização( reuso). Assim, foram definidas regras para extração dos conteúdos mais importantes destes documentos, o que possibilita a reutilização. Esta extração envolve técnicas de mineração de texto e utiliza regras para descobrir as definições contidas nos documentos. Foi desenvolvido um protótipo que demonstra a viabilidade do processo. Para facilitar o entendimento do trabalho, é apresentado um estudo de caso que utiliza a técnica proposta e o protótipo desenvolvido. O trabalho facilita e enriquece o desenvolvimento de materiais de aprendizagem, uma vez que torna os conteúdos de documentos das bibliotecas digitais reutilizáveis e integrados aos Objetos de Aprendizagem (LO) existentes.

## Palavras-chave

Bibliotecas Digitais; Educação Baseada na *Web*; Objetos de Aprendizagem; Integração de Dados; Mineração de Texto; Ontologia; Banco de Dados

## **Abstract**

Gomes, Geórgia R. R.. **Integration of Repositories of Digital Library Systems and Learning Management Systems**. PUC-Rio, 2006.143p.  
PhD. Thesis – Computer Science Department, Pontifical Catholic University of Rio de Janeiro, Brazil

With the widespread use of Information Technology for teaching support, it is usual to made digital content available through Digital Library Systems (DLMS) or Learning Management Systems (LMS). These systems, however, work independently, have different characteristics and manipulate different types of materials, and their data and metadata repositories are heterogeneous and distributed. The content of repositories would be better used if it was integrated in the same environment or accessed in an integrated way from DLMS and LMS. This thesis presents a homogeneous view of DLMS and LMS content. In order to provide such homogenization, it is proposed an extension of the mediator and wrapper architecture for dealing with metadata integration and ontologies for treating semantics. Local ontologies are used for describing each metadata repository, and a global ontology for the integration. As documents of DLMS repositories tend to be monolithic and not to follow a reuse approach, rules for extracting the most important content from these documents were developed in order to make them reusable. This extraction includes text mining techniques as well as rules for discovering definitions embedded in the documents. A prototype was developed which implements the extraction and proves the feasibility of this approach. In order to make the work easier to understand, it is presented a case study that uses the proposed technique and the prototype. The work described in this thesis facilitates and enriches the development of learning material by making the content of digital library documents reusable and integrated to existing learning objects.

## **Keywords**

Digital Library; Web-Based Education; Learning Objects; Data Integration; Text Mining; Ontology; Database



## Sumário

1	Introdução	16
1.1.	Motivação	16
1.2.	Objetivos da Tese	19
1.3.	Organização da Tese	19
2	Fundamentação	21
2.1.	Ambientes de Aprendizagem e Bibliotecas Digitais	21
2.1.1.	Ambientes de Aprendizagem	21
2.1.2.	Bibliotecas Digitais	23
2.2.	Integração de Dados Heterogêneos	25
2.2.1.	Mediadores	26
2.2.2.	Heterogeneidade Semântica	28
2.2.2.1.	Metadados	29
2.2.2.2.	Ontologia	31
2.3.	Mineração de Texto	33
2.3.1.	Preparação dos dados textuais	34
2.3.1.1.	Recuperação da Informação	35
2.3.1.2.	Análise dos dados	36
2.3.2.	Processamento dos textos	38
2.3.2.1.	Extração da Informação	39
2.3.3.	Pós-processamento da Mineração	40
3	Preparação de DL para integração	42
3.1.	Extensão da arquitetura do ambiente de DL	42
3.2.	Extração de informação dos DDs	43
3.2.1.	Proposta da CISCO	44
3.2.2.	Extração de Definição	45
4	Integração DLMS e LMS	52

4.1. Arquitetura Proposta	52
4.2. Componentes da arquitetura	54
4.2.1. Camada de Aplicação	54
4.2.2. Camada de Mediação	55
4.2.2.1. Modelo de dados	56
4.2.3. Camada de Tradutores	59
4.3. Caso de Uso do Sistema Integrador	62
5 Estudo de Caso	68
5.1. O Cenário	68
5.2. Extração de RIOs	69
5.3. Consulta Integrada	70
5.3.1. Mediador	72
5.3.2. Tradutores	73
6 Trabalhos Relacionados	77
6.1. DILLEO	77
6.2. ILUMINA	78
6.3. LEBONED	79
6.4. Síntese comparativa	80
7 Conclusão	82
7.1. O Trabalho Apresentado nesta Tese	82
7.2. Contribuições	83
7.3. Trabalhos Futuros	84
Referências Bibliográficas	86
Apêndice A - LOM	95
Apêndice B - MARC	98
Apêndice C – Dublin Core	129



## Lista de figuras

Figura 1 – Arquitetura do Ambiente de Aprendizagem.....	22
Figura 2 - Arquitetura do Ambiente de DL.....	25
Figura 6 – Camadas da Arquitetura de Mediadores .....	27
Figura 7 - Etapas do Processo de Mineração de Texto.....	34
Figura 8 - Arquitetura Sistema IR .....	35
Figura 9- Arquitetura Modificada do Ambiente de DL.....	43
Figura 10 - Arquitetura Proposta .....	53
Figura 11 - Camadas com os componentes da arquitetura.....	55
Figura 12 – Modelo de Dados de Integração do Mediador.....	57
Figura 13 – Mapeamento de assunto do esquema global para os correspondentes termos dos esquemas locais .....	59
Figura 14 - Ontologia Dublin Core.....	60
Figura 15 - Ontologia MARC.....	61
Figura 16- Ontologia LOM .....	62
Figura 17 – Diagrama de Casos de Uso da Arquitetura Proposta .....	63
Figura 18 – Ambiente do estudo de caso .....	71
Figura 19 – Exemplo de Interface da aplicação de consulta.....	73

## Lista de tabelas

Tabela 1- Resultados de extração de definições do Processo 1	50
Tabela 2 - Resultados de extração de definições do Processo 2	50
Tabela 3 – Descrição do Caso de Uso Validar Usuário	65
Tabela 4 - Descrição do Caso de Uso Consultar Objetos	65
Tabela 5 - Descrição do Caso de Uso Incluir Repositório de Dados	66
Tabela 6 - Descrição do Caso de Uso Excluir Repositório de Dados	67
Tabela 7 – Tabela comparativa deste trabalho com o projeto LEBONED	81
Tabela 8 -Tabela com representação das principais <i>tags</i> do MARC	102
Tabela 9 - Tabela com representação das principais tags e subcampos do MARC	105

## Abreviaturas e Siglas

ADL	Iniciativa da Secretaria de Defesa dos EUA no sentido de estabelecer um ambiente distribuído de aprendizagem - <i>Advanced Distributed Learning</i>
ARIADNE	<i>Alliance of Remote Instructional Authoring and Distribution Networks for Europe</i>
DC	<i>Dublin Core</i>
DD	Documento Digital
DL	<i>Digital Library</i> (Biblioteca Digital)
DLF	<i>Digital Library Federation</i>
DLMS	<i>Digital Library Management Systems</i>
DLO	<i>Digital Library Object</i> (Objeto de Bibliotecas Digitais)
DLOMS	<i>Digital Library Object Management Systems</i>
EI	Extração da Informação
F	Medida padrão combinando as métricas P e R
IEEE	<i>Institute of Electrical and Electronics Engineers</i>
IMS	<i>Information Management Systems</i> (Sistemas de Gerência de Informação)
IR	<i>Information Retrieval</i>
KDT	<i>Text Mining ou Knowledge Discovery from Texts</i>
LB	<i>Library of Congress</i> (Biblioteca do Congresso Americano)
LCMS	<i>Learning Content Management Systems</i> (Sistemas de Gerência de Conteúdo de Aprendizagem)
LMS	<i>Learning Management Systems</i> (Sistemas de Gerência da Aprendizagem)
LO	<i>Learning Object</i> (Objeto de Aprendizagem)
LOM	<i>Learning Objects Metadata</i> (Padrão de metadados proposto pelo IEEE para descrição de LOs)
LTSC	<i>Learning Technology Standards Committee</i> (Comitê do IEEE responsável pela padronização da tecnologia de aprendizagem)
MARC	<i>Machine Readable Cataloging</i>
METS	<i>Metadata Encoding and Transmission Standard</i>
OWL	<i>Web Ontology Language</i> (Linguagem para definição/especificação de ontologias para a Web)
P	<i>Precision</i> (Precisão)
PLN	<i>Processamento de Linguagem Natural</i>
PGL	<i>Partnership in Global Learning</i>
R	<i>Recall</i> (Abrangência)
RDA	<i>Remote Data Access</i>
RIO	<i>Reusable Information Objects</i> (Objetos de Informação Reutilizáveis)
RLO	<i>Reusable Learning Object</i>
SCORM	<i>Sharable Content Object Reference Model</i>
SBD	<i>Sistema de Banco de Dados</i>
SBDH	<i>Sistema de Banco de Dados Heterogêneos</i>
SGBD	<i>Sistema de Gerência de Banco de Dados</i>
SGBDH	<i>Sistema de Gerência de Banco de Dados Heterogêneos</i>
SQL	<i>Structured Query Language</i>
TecBD	Laboratório de Tecnologia em Banco de Dados do Departamento de Informática da PUC-Rio
VIO	<i>Very Important Object</i>
XML	<i>Extensible Markup Language</i>
W3C	<i>World Wide Web Consortium</i>
Web	<i>World Wide Web</i>

*“Quando pensares em desistir, lembre-se que Deus está ao seu lado, para te ajudar a prosseguir, e com Ele, o sonho se torna realidade”*

*Geórgia Gomes*