



**José Eduardo Talavera Herrera**

**Arquitetura para Recomendação de Fontes  
de Dados RDF**

**Dissertação de Mestrado**

Dissertação apresentada ao Programa de Pós-graduação em Informática do Departamento de Informática do Centro Técnico Científico da PUC-Rio como requisito parcial para obtenção do grau de Mestre em Informática.

Orientadora: Prof<sup>ª</sup>. Karin Breitman

Rio de Janeiro  
Outubro de 2012



**José Eduardo Talavera Herrera**

**Arquitetura para Recomendação de Fontes  
de Dados RDF**

Dissertação apresentada ao Programa de Pós-graduação em Informática do Departamento de Informática do Centro Técnico Científico da PUC-Rio como requisito parcial para obtenção do grau de Mestre em Informática. Aprovada pela Comissão Examinadora abaixo assinada.

**Prof<sup>a</sup>. Karin Breitman**

Orientadora

Departamento de Informática — PUC-Rio

**Prof. Marco Antonio Casanova**

Departamento de Informática — PUC-Rio

**Prof. Luiz André P. Paes Leme**

UFF-Rio

**Prof. Antonio Luz Furtado**

Departamento de Informática — PUC-Rio

**Prof. José Eugenio Leal**

Coordenador Setorial do Centro Técnico Científico — PUC-Rio

Rio de Janeiro, 01 de Outubro de 2012

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

### **José Eduardo Talavera Herrera**

Graduou-se em Engenharia de Sistemas Universidad Nacional de San Agustín (UNSA), Arequipa-Perú em 2003. Trabalha como Engenheiro de Aplicações no Banco de Credito do Perú. Ingressou no programa de mestrado do Departamento de Informática em 2010. Atualmente suas áreas de pesquisa são Web Semântica, *Linked Data* e *Cloud Computing*.

#### Ficha Catalográfica

Herrera, José Eduardo

Arquitetura para Recomendação de Fontes de Dados RDF / José Eduardo Talavera Herrera; orientadora: Karin Breitman. — 2012.

112 f. : il. (color); 30 cm

1. Dissertação (mestrado) - Pontifícia Universidade Católica do Rio de Janeiro, Departamento de Informática, 2012.

Inclui bibliografia.

1. Informática – Teses.  
2. Dados Conectados. 3. Computação na nuvem. 4. Recuperação de informação. 5. Descoberta de Links. 6. Similaridade.  
I. Breitman, Karin. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

CDD: 004

## Agradecimentos

Agradeço a Deus por ter me dado força, saúde e por ter estado ao meu lado sempre me guiando. Agradeço aos meus pais, Soledad e Eduardo, pelo amor, dedicação, amizade e apoio e à minha irmã Diana pelo amor, amizade, incentivo, ajuda e apoio em todos momentos. Agradeço a Karin, minha orientadora, por ter me recebido no seu grupo de pesquisa, pelas oportunidades e pela orientação. Agradeço aos meus amigos Renato, Ricardo, Marcelo, Chrystiano e Everton pela amizade, apoio e por tornar os momentos sempre muito divertidos. Agradeço a Michele por ter estado do meu lado e por sempre acreditar em mim, você é e será muito especial para mim. Agradeço aos meus amigos Livia, Thiago e Alessandro pela ajuda na correção deste trabalho, muito obrigado. Agradeço a Amparito, Ximena, Mayra, Edward, Ronald e Adriano pelo incentivo, ajuda e amizade. Enfim, agradeço à todos os meus amigos pois foram peças fundamentais para que eu pudesse chegar até aqui. E Agradeço à CAPES que financiou minha bolsa de mestrado.

## Resumo

Herrera, José Eduardo; Breitman, Karin. **Arquitetura para Recomendação de Fontes de Dados RDF**. Rio de Janeiro, 2012. 112p. Dissertação de Mestrado — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Dentro do processo de publicação de dados na Web recomenda-se interligar os dados entre diferentes fontes, através de recursos similares que descrevam um domínio em comum. No entanto, com o crescimento do número dos conjuntos de dados publicados na Web de Dados, as tarefas de descoberta e seleção de dados tornam-se cada vez mais complexas. Além disso, a natureza distribuída e interconectada dos dados, fazem com que a sua análise e entendimento sejam muito demorados. Neste sentido, este trabalho visa oferecer uma arquitetura Web para a identificação de fontes de dados em RDF, com o objetivo de prover melhorias nos processos de publicação, interconexão, e exploração de dados na *Linked Open Data*. Para tal, nossa abordagem utiliza o modelo de MapReduce sobre o paradigma de computação nas nuvens. Assim, podemos efetuar buscas paralelas por palavras-chave sobre um índice de dados semânticos existente na Web. Estas buscas permitem identificar fontes candidatas para ligar os dados. Por meio desta abordagem, foi possível integrar diferentes ferramentas da web semântica em um processo de busca para descobrir fontes de dados relevantes, e relacionar tópicos de interesse definidos pelo usuário. Para atingir nosso objetivo foi necessária a indexação e análise de texto para aperfeiçoar a busca de recursos na *Linked Open Data*. Para mostrar a eficácia de nossa abordagem foi desenvolvido um estudo de caso, utilizando um subconjunto de dados de uma fonte na *Linked Open Data*, através do seu serviço SPARQL *endpoint*. Os resultados do nosso trabalho revelam que a geração de estatísticas sobre os dados da fonte é, de fato, um grande diferencial no processo de busca. Estas estatísticas ajudam ao usuário no processo de escolha de indivíduos. Um processo especializado de extração de palavras-chave é aplicado para cada indivíduo com o objetivo de gerar diferentes buscas sobre o índice semântico. Mostramos a escalabilidade de nosso processo de recomendação de fontes RDF através de diferentes amostras de indivíduos.

## Palavras-chave

Dados Conectados; Computação na nuvem; Recuperação de informação; Descoberta de Links; Similaridade.

## Abstract

Herrera, José Eduardo; Breitman, Karin (Advisor). **An Architecture for RDF Data Sources Recommendation**. Rio de Janeiro, 2012. 112p. MSc. Dissertation — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

In the Web publishing process of data it is recommended to link the data from different sources using similar resources that describe a domain in common. However, the growing number of published data sets on the Web have made the data discovery and data selection tasks become increasingly complex. Moreover, the distributed and interconnected nature of the data causes the understanding and analysis to become too prolonged. In this context, this work aims to provide a Web architecture for identifying RDF data sources with the goal of improving the publishing, interconnection, and data exploration processes within the Linked Open Data. Our approach utilizes the MapReduce computing model on top of the cloud computing paradigm. In this manner, we are able to make parallel keyword searches over existing semantic data indexes available on the web. This will allow to identify candidate sources to link the data. Through this approach, it was possible to integrate different semantic web tools and relevant data sources in a search process, and also to relate topics of interest defined by the user. In order to achieve our objectives it was necessary to index and analyze text to improve the search of resources in the Linked Open Data. To show the effectiveness of our approach we developed a case study using a subset of data from a source in the Linked Open Data through its SPARQL endpoint service. The results of our work reveal that the generation and usage of data source's statistics do make a great difference within the search process. These statistics help the user within the choosing individuals process. Furthermore, a specialized keyword extraction process is run for each individual in order to create different search processes using the semantic index. We show the scalability of our RDF recommendation process by sampling several individuals.

## Keywords

Linked Data; Cloud computing; Retrieval Information; Link Discovery; Similarity.

# Sumário

1	Introdução	<b>11</b>
1.1	Linked Data	11
1.2	Motivação	12
1.3	Objetivos	13
1.4	Organização da dissertação	14
2	Conceitos e Tecnologias Utilizadas	<b>15</b>
2.1	Tecnologias da Web Semântica	16
2.2	Linked Open Data	18
2.3	Computação na Nuvem	21
2.4	Modelo MapReduce	22
2.5	O Hadoop framework	24
2.6	Web semântica e computação na nuvem	26
3	Trabalhos Relacionados	<b>29</b>
3.1	Nível de Metadados	30
3.2	Nível do Conteúdo	31
3.3	Nível Social	35
4	Processo proposto de Recomendação de Fontes RDF	<b>37</b>
4.1	Etapa de configuração	40
4.2	Etapa de geração de estatísticas	41
4.3	Etapa de escolha de classes	49
4.4	Etapa de escolha de Indivíduos	55
4.5	Etapa de ranking	65
5	Estudo de caso	<b>68</b>
5.1	Etapa de geração de estatísticas	68
5.2	Etapa de escolha classes	70
5.3	Etapa de escolha de indivíduos	74
6	Implementação	<b>81</b>
6.1	Arquitetura	81
6.2	Modelo MapReduce do processamento proposto	90
6.3	Opções de otimização	92
7	Conclusão	<b>95</b>
7.1	Contribuições	95
7.2	Comparação com trabalhos relacionados	96
7.3	Limitações	97
7.4	Trabalhos futuros	98
8	Referências Bibliográficas	<b>99</b>
A	Resultados da etapa da escolha de indivíduos	<b>106</b>

B	Estrutura de documentos no processo de indexação	108
C	Exemplo de resposta do serviço de indexação	111



## Lista de figuras

2.1	Ambiente de Computação na Nuvem	22
2.2	Execução do exemplo de contar palavras	24
2.3	Execução paralela de uma tarefa MapReduce	25
4.1	Processo de recomendação de fontes RDF	38
4.2	Detalhamento da etapa de configuração	40
4.3	Detalhamento da etapa de geração de estatísticas	42
4.4	Fluxo de tarefas	45
4.5	Detalhamento da etapa de escolha de classes	49
4.6	Merge de informação do indivíduo	54
4.7	Detalhamento da etapa de escolha de indivíduos	57
4.8	Tipos de buscas: (a) Combinação independente das palavras-chave, (b) Acúmulo de palavras-chave	64
4.9	Etapa de ranking do processo proposto	65
4.10	Escalabilidade do processo de busca: (a) Busca por combinação independente das palavras-chave, (b) Busca por acúmulo de palavras-chave	66
5.1	Resultados sobre Artigos Acadêmicos: (a) Busca por combinação independente de palavras-chave, (b) Busca por acúmulo de palavras-chave	76
5.2	Resultados sobre o material multimídia: (a) Busca por combinação independente das palavras-chave, (b) Busca por acúmulo de palavras-chave	77
5.3	Ranking de artigo acadêmico: (a) Ranking de classes, (b) Definição de classe bibo:AcademicArticle	79
5.4	Ranking do material multimídia: (a) Ranking de classes, (b) Definição das classes PodCast	80
6.1	Arquitetura proposta	83
6.2	Aplicação Web	83
6.3	Arquitetura conceitual de Solr	87
6.4	Modelo MapReduce da geração de estatísticas	90
6.5	Modelo MapReduce da escolha de classes	91
6.6	Modelo MapReduce da busca de indivíduos	92

## Lista de tabelas

4.1	Relação entre as propriedades do VoiD e consultas SPARQL	44
4.2	Resumo de entrada de dados	47
4.3	Resultado de consultas	48
4.4	Lista de classes e consultas auxiliares	52
4.5	Configuração do processo	53
4.6	Propriedades da classe artigo acadêmico	59
4.7	Configuração do processo	61
4.8	Resultados da busca por combinação (BC) e por acúmulo (BA) de palavras-chave	65
5.1	Lista de classes da fonte data.open.ac.uk	69
5.2	Lista de prefixos usados nas classes	69
5.3	Lista de estatísticas na coleção de documentos	71
5.4	Estatísticas por documento na classe de Artigo Acadêmico	72
5.5	Lista de estatísticas na coleção de documentos	72
5.6	Estatísticas por documento nas classes de material multimídia	73
A.1	Propriedades das classes de material multimídia	107
B.1	Estrutura do documento da classe	108
B.2	Estrutura do documento estatísticas	108
B.3	Estrutura do documento do indivíduo	109
B.4	Estrutura do documento do Síndice	110