

Rafael Silva Pereira

**A Cloud based Real-Time Collaborative
Filtering Architecture for Short-Lived Video
Recommendations**

TESE DE DOUTORADO

DEPARTAMENTO DE INFORMÁTICA

Programa de Pós-Graduação em Informática

Rio de Janeiro
December 2015



Rafael Silva Pereira

**A Cloud based Real-Time Collaborative Filtering
Architecture for Short-Lived Video Recommendations**

TESE DE DOUTORADO

Thesis presented to the Programa de Pós-Graduação em
Informática of the Departamento de Informática, PUC-Rio
as partial fulfillment of the requirements for the degree of
Doutor em Ciências - Informática

Advisor: Prof. Hélio Côrtes Vieira Lopes

Rio de Janeiro
December 2015



Rafael Silva Pereira

A Cloud based Real-Time Collaborative Filtering Architecture for Short-Lived Video Recommendations

Thesis presented to the Programa de Pós-Graduação em Informática, of the Departamento de Informática do Centro Técnico Científico da PUC-Rio, as partial fulfillment of the requirements for the degree of Doutor.

Prof. Hélio Côrtes Vieira Lopes

Advisor

Departamento de Informática – PUC-Rio

Prof. Marco Antonio Casanova

Departamento de Informática – PUC-Rio

Profa. Karin Koogan Breitman

EMC

Profa. Giseli Rabello Lopes

UFRJ

Prof. José Viterbo Filho

UFF

Prof. Marcus Vinicius Soledade Poggi de Aragão

Departamento de Informática – PUC-Rio

Prof. Luiz André Portes Paes Leme

UFF

Prof. Márcio da Silveira Carvalho

Coordinator of the Centro Técnico Científico - PUC-Rio

Rio de Janeiro, December 11th, 2015

All rights reserved

Rafael Silva Pereira

Graduated in Electronics and Computer Engineering at Universidade Federal do Rio de Janeiro – UFRJ in 2006. Masters in Computer Science at Pontifícia Universidade Católica do Rio de Janeiro - PUC-Rio in 2011, and Engineering Manager from Globo.com since 2006.

Bibliographic data

Pereira, Rafael Silva

A Cloud based Real-Time Collaborative Filtering Architecture for Short-Lived Video Recommendations / Rafael Silva Pereira ; advisor: Hélio Côrtes Vieira Lopes. – 2015.

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

Tese (doutorado)–Pontifícia Universidade Católica do Rio de Janeiro, Departamento de Informática, 2015.

Inclui bibliografia

1. Informática – Teses. 2. Computação na nuvem. 3. Filtragem colaborativa. 4. Recomendação. 5. Sistemas distribuídos. 6. Arquiteturas orientadas a serviços. I. Lopes, Hélio Côrtes Vieira. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

CDD: 004

To my parents and my wife,
for your unconditional love and support

Acknowledgements

First of all, to my family, specially my parents that offered me all the support and encouragement to always continue independently of the obstacles, and to afford all my primary education. To my wife Sofia Gross, that always supported me in this journey since the beginning, even in the moments where all my attention was directed to this work.

To my advisor Prof. Hélio Lopes for its kindness and for the constant support and technical knowledge, always available when necessary. To all professors of the Department of Informatics from PUC-Rio for their knowledge and help, especially to Karin Breitman who was my advisor during my Masters at PUC-Rio, and is who encourage me to start this work. I also thank Professor Simone Barbosa who helped me to solve the issues with my registration during the course. To all my colleagues from PUC-Rio.

To the professors responsible for the evaluation of this work, for your availability and dedication reading and contributing to this research.

To Globo.com, my employer, that paid all expenses of this Doctoral course and gave me all resources needed for this research.

Finally, to Brazilian Government, for its public, high quality, and free federally funded higher education system, that allowed me to complete my graduation degree.

Abstract

Pereira, Rafael Silva; Lopes, Hélio Côrtes Vieira (Advisor). **A Cloud based Real-Time Collaborative Filtering Architecture for Short-Lived Video Recommendations**. Rio de Janeiro, 2015. 86p. DSc. Dissertation – Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

This dissertation argues that the combination of collaborative filtering techniques, particularly for item-item recommendations, with emergent cloud computing technology can drastically improve algorithm efficiency, particularly in situations where the number of items and users scales up to several million objects. It introduces a real-time item-item recommendation architecture, which rationalizes the use of resources by exploring on-demand computing. The proposed architecture provides a real-time solution for computing online item similarity, without having to resort to either model simplification or the use of input data sampling. This dissertation also presents a new adaptive model for implicit user feedback for short videos, and describes how this architecture was used in a large scale implementation of a video recommendation system in use by the largest media group in Latin America, presenting results from a real life case study to show that it is possible to greatly reduce recommendation times (and overall financial costs) by using dynamic resource provisioning in the Cloud. It discusses the implementation in detail, in particular the design of cloud based features. Finally, it also presents potential research opportunities that arise from this paradigm shift.

Keywords

Cloud Computing; Recommendations; Collaborative Filtering; Distributed Systems; Service Oriented Architectures.

Resumo

Pereira, Rafael Silva; Lopes, Hélio Côrtes Vieira (Orientador). **Uma Arquitetura de Filtragem Colaborativa em Tempo Real Baseada em Nuvem para Recomendação de Vídeos Efêmeros**. Rio de Janeiro, 2015. 86p. Tese de Doutorado – Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Esta tese propõe que a combinação de técnicas de filtragem colaborativa, em particular para recomendações item-item, com novas tecnologias de computação em nuvem, pode melhorar drasticamente a eficiência dos sistemas de recomendação, particularmente em situações em que o número de itens e usuários supera milhões de objetos. Nela apresentamos uma arquitetura de recomendação item-item em tempo real, que racionaliza o uso dos recursos computacionais através da computação sob demanda. A arquitetura proposta oferece uma solução para o cálculo de similaridade entre itens em tempo real, sem ter que recorrer à simplificação do modelo de recomendação ou o uso de amostragem de dados de entrada. Esta tese também apresenta um novo modelo de feedback implícito para vídeos de curta duração, que se adapta ao comportamento dos usuários, e descreve como essa arquitetura foi usada na implementação de um sistema de recomendação de vídeo em uso pelo maior grupo de mídia da América Latina, apresentando resultados de um estudo de caso real para mostrar que é possível reduzir drasticamente o tempo de cálculo das recomendações (e os custos financeiros globais) usando o provisionamento dinâmico de recursos na nuvem. Ela discute ainda a implementação em detalhes, em particular o projeto da arquitetura baseada em nuvem. Finalmente, ela também apresenta oportunidades de pesquisa em potencial que surgem a partir desta mudança de paradigma.

Palavras-chave

Computação na Nuvem; Filtragem Colaborativa; Recomendação; Sistemas Distribuídos; Arquiteturas Orientadas à Serviço.

Table of Contents

| | |
|--|----|
| 1 Introduction | 13 |
| 1.1. Goals | 15 |
| 1.2. Main Contributions | 16 |
| 1.3. Related Work | 16 |
| 1.4. Outline | 18 |
| 2 Background | 19 |
| 2.1. Recommender Systems | 19 |
| 2.1.1. The Recommendation Problem | 20 |
| 2.1.2. Recommendation Strategies | 23 |
| 2.1.3. Limitations of Collaborative Filtering | 29 |
| 2.1.4. Comparison of Strategies | 30 |
| 2.2. Cloud Computing | 31 |
| 2.2.1. Cloud Computing Paradigms | 32 |
| 2.2.2. Amazon Web Services Platform | 35 |
| 3 Large Scale Video Recommendations | 40 |
| 3.1. Introduction | 40 |
| 3.2. Item-Item Video Recommendations | 41 |
| 4 A Real-Time Large Scale Collaborative Filtering Architecture | 48 |
| 4.1. A multi-layered architecture | 50 |
| 4.1.1. The first layer – Hits Layer | 51 |
| 4.1.2. The second layer – Persistence Layer | 52 |
| 4.1.3. The third layer – Similarity Layer | 55 |
| 5 An Adaptive Implicit Feedback Model | 57 |
| 6 Implementation and Deploy | 65 |
| 6.1. Implementation Issues | 65 |
| 6.2. Deploy using Amazon Web Services | 67 |

| | |
|--|----|
| 7 Application in Globo.com and Results | 69 |
| 8 Conclusions | 78 |
| 9 References | 81 |

List of Figures

| | |
|--|----|
| Figure 1 - Recommendation Process | 20 |
| Figure 2 - Collaborative Filtering process | 25 |
| Figure 3 - Collaborative Filtering results in Amazon.com | 26 |
| Figure 4 - Neighborhood training process | 26 |
| Figure 5 - Item-based collaborative filtering | 27 |
| Figure 6 - Netflix Star Ratings for explicit feedback | 41 |
| Figure 7 - YouTube "Thumbs up" / "Thumbs Down" explicit feedback | 42 |
| Figure 8 - Similarity recalculation process for a single piece of feedback | 44 |
| Figure 9 - Video views from 5 news videos from the <i>Jornal Nacional</i> | 45 |
| Figure 10 - Multi-layer collaborative filtering process | 51 |
| Figure 11 - Player sending the HTTP call informing that an user started to watch a video | 52 |
| Figure 12 - Layer II storing binary feedback in sets | 53 |
| Figure 13 - Layer II storing items already evaluated by the user in a supporting set structure | 54 |
| Figure 14 - Similarity calculation between items I_1 and I_2 through intersection of sets | 55 |
| Figure 15 - Completion ratio by category of short videos and clips | 58 |
| Figure 16 - Consumption pattern of short videos of the National Journal compared the the average of its category | 59 |
| Figure 17 - Video A and video B have different consumption curves and hence different feedback scales | 61 |
| Figure 18 - Adapted feedback scales for videos A and B | 62 |
| Figure 19 - Redis pools connecting architecture layers | 66 |
| Figure 20 - The recommendation platform deployed on AWS | 68 |
| Figure 21 - Recommendations integrated in the player and presented at the end of the playback | 70 |
| Figure 22 - Recommendations integrated in the web page | 70 |
| Figure 23 - User feedback requests for a 10 day period starting on Monday | 71 |

| | |
|---|----|
| Figure 24 - Variation in number of Layer II nodes according to Queue I size | 72 |
| Figure 25 - Processing Time in Layer II vs. Total Unique Users | 72 |
| Figure 26 - Variation in number of Layer III nodes according to Queue II size | 73 |
| Figure 27 - Persistence Redis pool configuration | 75 |
| Figure 28 - Conversion rate comparison | 76 |

List of Tables

| | |
|---|----|
| Table 1 - Rating Matrix for movies | 22 |
| Table 2 - Total monthly cost for Globo.com recommendations deployment | 77 |