**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

**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

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

# 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**

## List of Figures

## List of Tables

PUC-Rio - Certificação Digital Nº 1122198/CA