

Avaliação de sistemas NoSQL para o gerenciamento de dados em *workflows* científicos

Gabriel Ferreira Guilhoto

Orientadora: Profa. Dra. Kelly Rosa Braghetto

Trabalho de Formatura Supervisionado (MAC0499)

Bacharelado em Ciência da Computação

Instituto de Matemática e Estatística

Universidade de São Paulo

16 de novembro de 2015

Sumário

- ▶ Introdução
- ▶ Sistemas NoSQL
 - ▶ MongoDB
 - ▶ Cassandra
- ▶ Workflows científicos
- ▶ Experimentos
- ▶ Resultados

Introdução

- ▶ *Workflows* científicos são aplicações intensivas em dados
- ▶ Soluções para o gerenciamento de seus dados não satisfazem todas as necessidades
 - ▶ Sistemas de arquivos distribuídos
 - ▶ Sistemas de armazenamento baseado em objetos
 - ▶ Eliminação dos resultados intermediários
- ▶ Sistemas NoSQL podem ser mais adequados

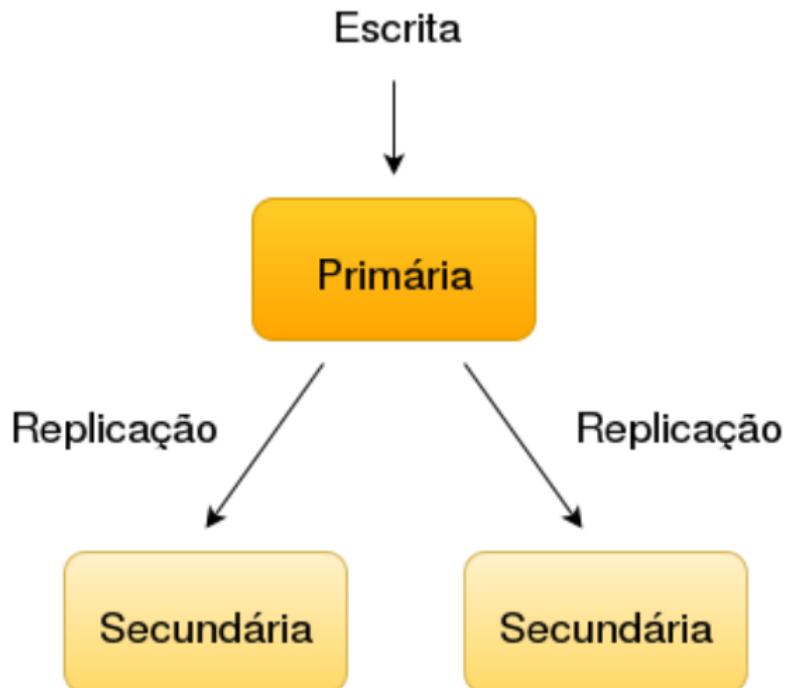
Sistemas NoSQL

- ▶ Projetados para escalarem horizontalmente
- ▶ Não usam o modelo de dados relacional
- ▶ Não têm interfaces padronizadas como a SQL
- ▶ Lidam com bancos de dados sem esquema definido
- ▶ São projetos de código aberto

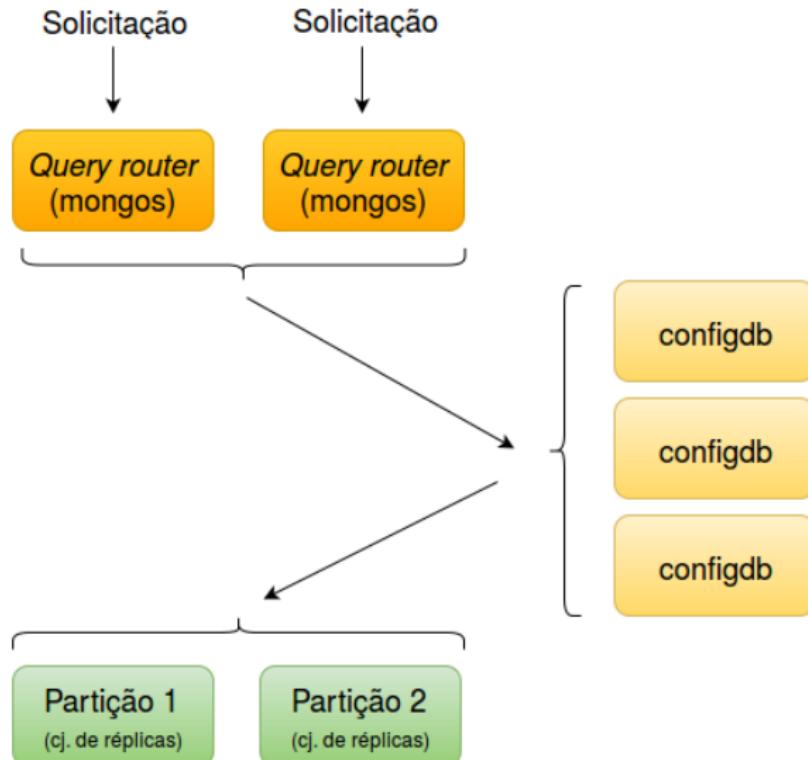
MongoDB

- ▶ Orientado a documentos; armazena seus dados em formato JSON sem estrutura pré-definida
- ▶ Replicação mestre-escravo
- ▶ Leitura opcional nos escravos
- ▶ Partição (*shard*): conjunto de réplicas é responsável por um subconjunto dos dados

Replicação no MongoDB



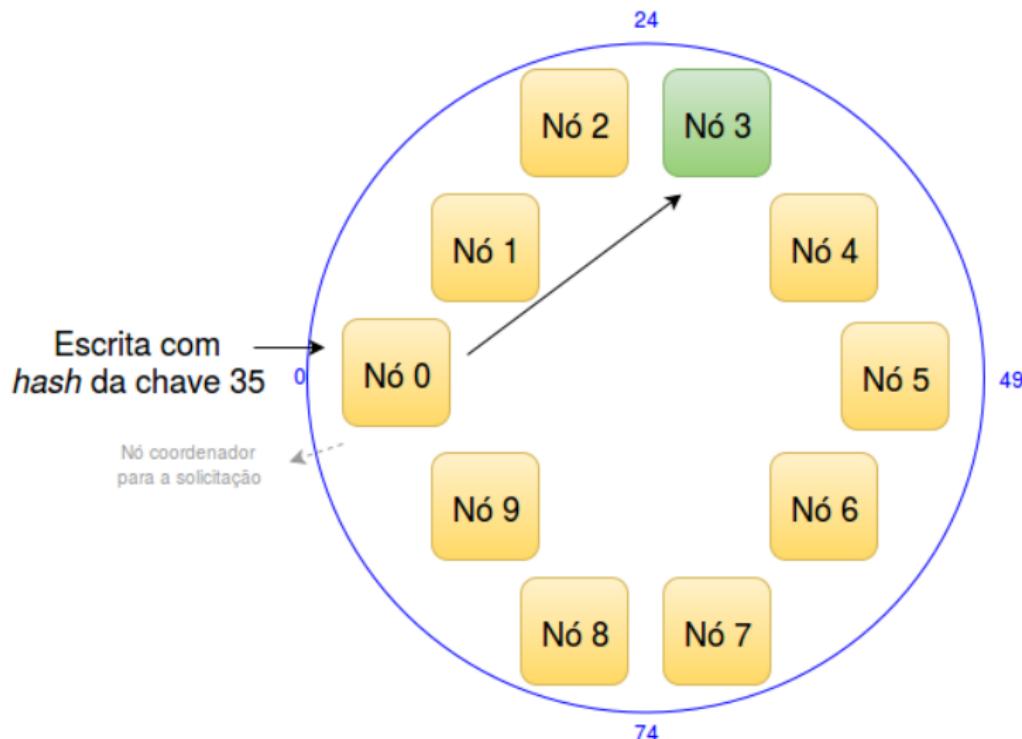
Particionamento no MongoDB



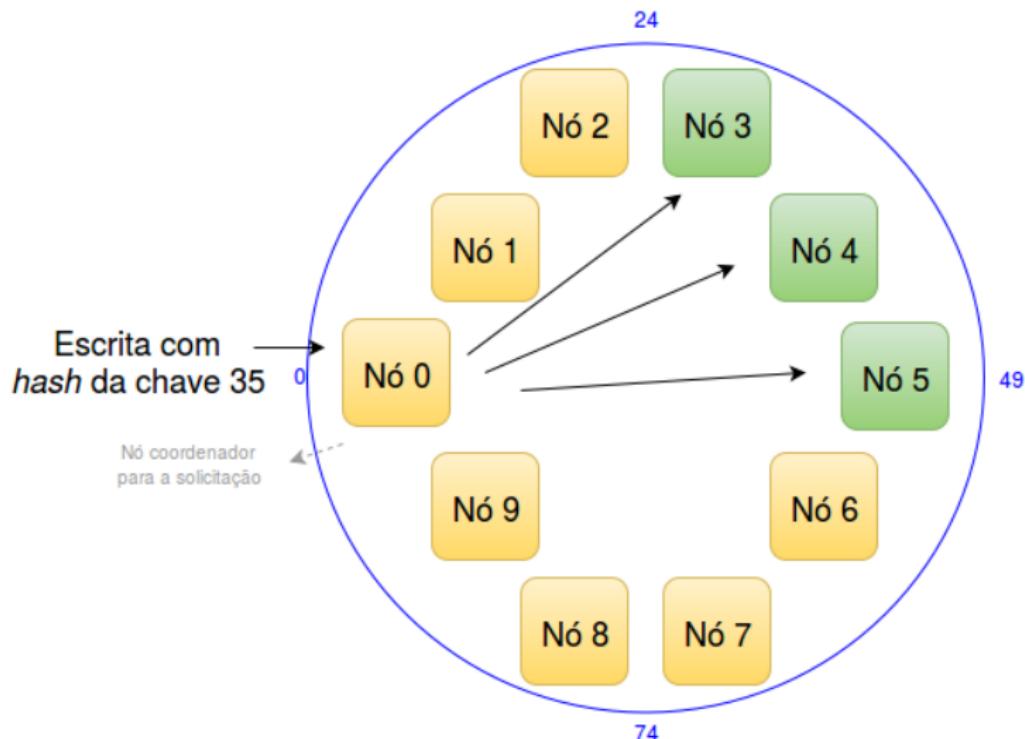
Cassandra

- ▶ Orientado a famílias de colunas
- ▶ Replicação ponto a ponto; não há mestre e todos atendem leituras e escritas
- ▶ Inserção em um nó determinado pela chave, replicado nos próximos nós em sentido horário

Particionamento no Cassandra



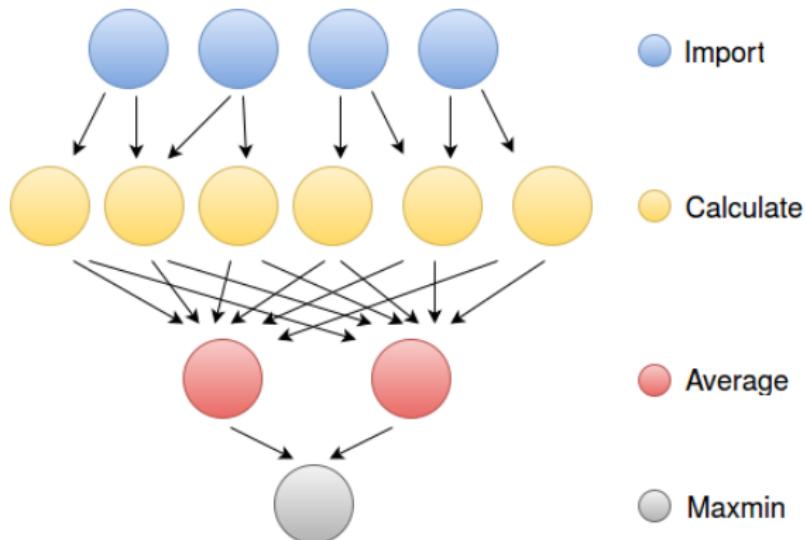
Replicação no Cassandra



Workflows científicos

- ▶ Descrição de um processo computacional de análise de dados
- ▶ Fluxo de dados: atividades produzem dados consumidos por outras

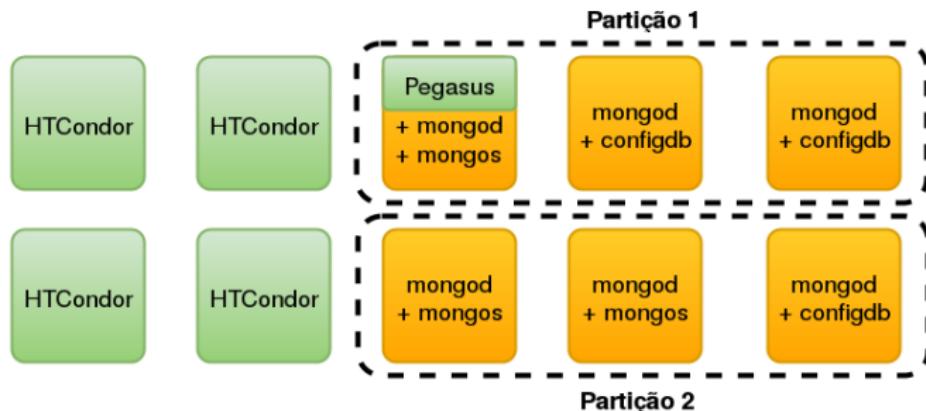
Workflow usado nos experimentos



Máquinas na Nuvem USP

Número da máquina	1	2	3	4	5	6	7	8	9	10
Memória RAM	32 GB									
Clock da CPU	2300 MHz									
Número de cores	8									
Sistema operacional	Ubuntu Server 12.04						Ubuntu 14.04			
Disco rígido	20 GB			120 GB			220 GB			

Máquinas no experimento com o MongoDB

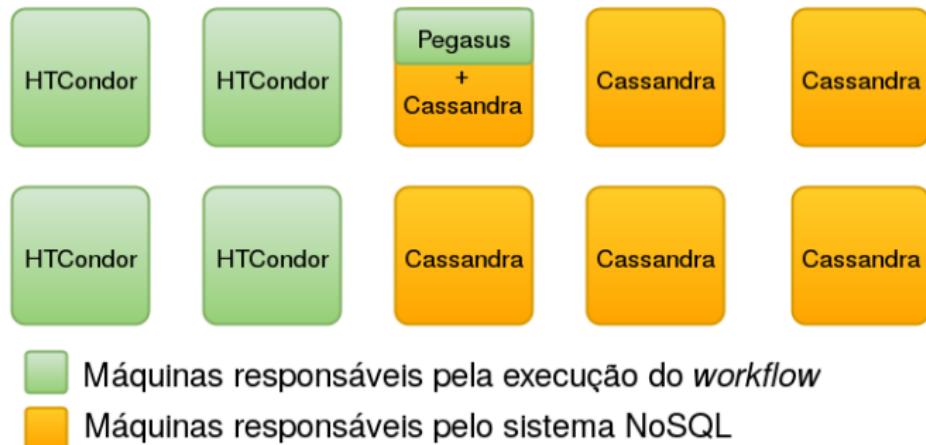


Máquinas responsáveis pela execução do *workflow*

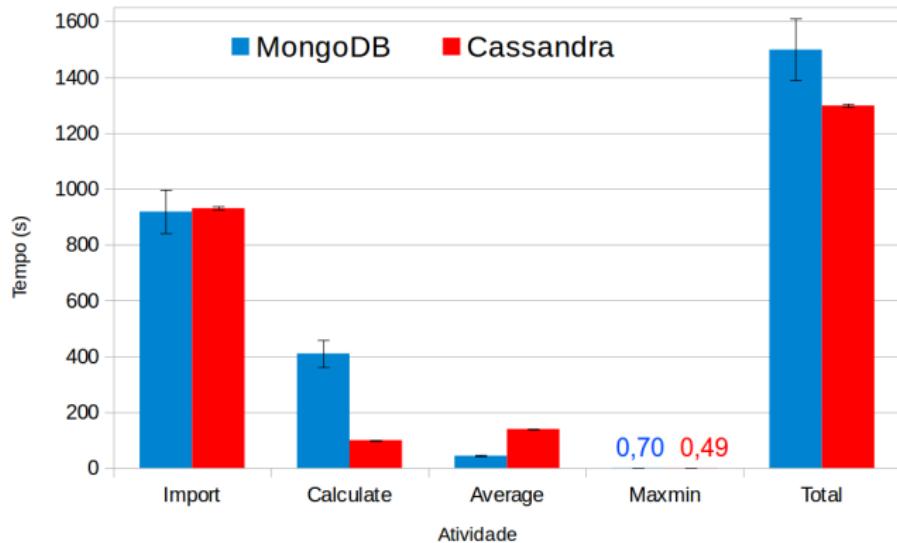


Máquinas responsáveis pelo sistema NoSQL

Máquinas no experimento com o Cassandra



Resultados



Resultados

Tipo de atividade	Tipo de operação	Sistema com melhor desempenho
Import	Inserção	Empate
Calculate	Seleção por faixa de valores e inserção	Cassandra
Average	Seleção por igualdade e inserção	MongoDB

Referências

-  **The Apache Software Foundation.**
The Apache Cassandra Project.
<http://cassandra.apache.org/>.
-  **Kelly Rosa Braghetto e Daniel Cordeiro.**
Introdução à modelagem e execução de workflows científicos.
Atualizações em Informática. 1ª ed. Porto Alegre: SBC, páginas 1–40, 2014.
-  **Ewa Deelman, Gurmeet Singh, Mei-Hui Su, James Blythe, Yolanda Gil, Carl Kesselman, Gaurang Mehta, Karan Vahi, G Bruce Berriman, John Good, et al.**
Pegasus: A framework for mapping complex scientific workflows onto distributed systems.
Scientific Programming, 13(3):219–237, 2005.
-  **Martin Fowler e Pramod J. Sadalage.**
NoSQL essencial: um guia conciso para o mundo emergente da persistência poliglota.
Novatec, 1ª edição.
-  **MongoDB, Inc.**
MongoDB.
<https://www.mongodb.org/>.
-  **Charles Reiss, John Wilkes, e Joseph L. Hellerstein.**
Google cluster-usage traces: format + schema.
Technical report, Google Inc., Mountain View, CA, USA, novembro 2011.
<http://code.google.com/p/googleclusterdata/wiki/TraceVersion2>.