

# Além do Hadoop, uma análise sobre os novos paradigmas de Computação em Nuvem

Otávio M. de Carvalho, Philippe O. A. Navaux  
Universidade Federal do Rio Grande do Sul  
Grupo de Processamento Paralelo e Distribuído  
{omcarvalho,navaux}@inf.ufrgs.br

## Introdução

A iniciativa principal do processamento online de grandes volumes de dados, caracterizado pelo conceito atual de *Big Data*, partiu principalmente dos trabalhos iniciais da Google, com a publicação dos seus artigos sobre o *Google File System* e o *MapReduce*.

A partir desses trabalhos, foi desenvolvido o projeto *Apache Hadoop*, que propôs novas versões dessas abordagens, baseadas em software livre. Esse projeto foi inicialmente desenvolvido por engenheiros do Yahoo e posteriormente entregue à fundação Apache. Desde então, esse projeto cresceu e se tornou o modelo predominante para o processamento de grandes volumes de dados.

Grande parte da adoção destas ferramentas se deve também a outro trabalho da Google, chamado *Google BigTable*. Este trabalho tornou possível o desenvolvimento de ferramentas de armazenamento distribuído de dados operando sobre a infraestrutura do Hadoop. A partir desse trabalho, surgiram ferramentas como o *Facebook Cassandra*, que deram origem ao denominado segmento dos bancos de dados NoSQL.

A partir da percepção das limitações existentes no modelo *MapReduce*, surgiram diversas alternativas, focando cada vez mais no tipo de dado a ser processado.

### Iniciativas de extensão do Modelo MapReduce e suas principais características

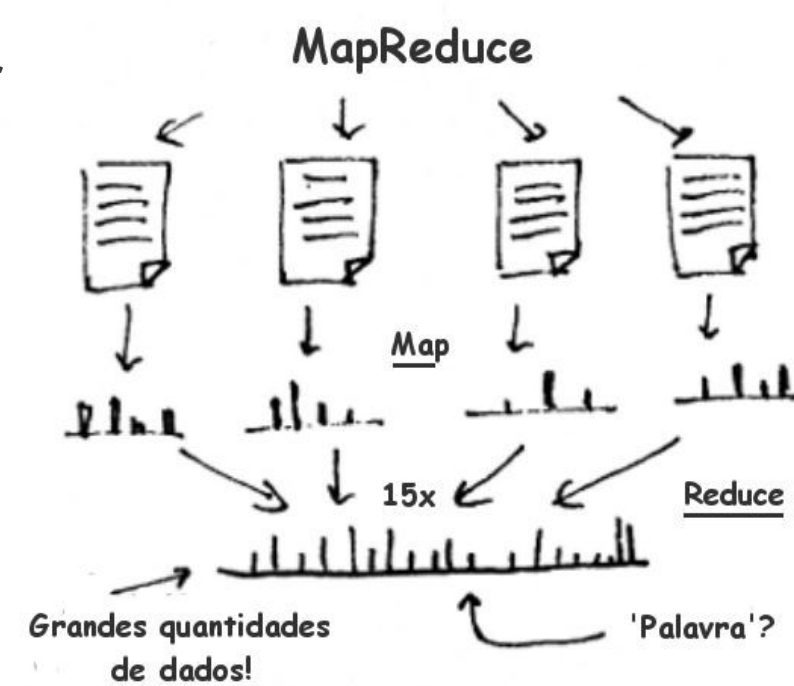
Nome	Ano	Descrição	Batch	Interativo	Tempo Real
Teradata Aster	2013	Banco de dados MPP	●	●	
Pivotal HD	2013	Conjunto de ferramentas de processamento distribuído	●	●	●
Google Photon	2013	Sistema para o processamento distribuído de fluxos contínuos de dados			●
AMPLab BDAS	2012	Conjunto de ferramentas de processamento distribuído em memória	●	●	●
Google Spanner	2012	Primeiro banco de dados distribuído com transações externamente consistentes		●	
Actian ParAccel	2012	Banco de dados MPP	●	●	
Cloudera Impala	2012	Sistema para o processamento de consultas interativas		●	
StreamBase CEP	2012	Ferramenta comercial de processamento complexo de eventos			●
Apache Giraph	2012	Ferramenta para o processamento distribuído de grafos	●		
Apache Drill	2012	Ferramenta para o processamento de consultas interativas		●	
Apache Flume	2012	Ferramenta para o processamento de fluxos contínuos de dados			●
Apache YARN	2011	Evolução do Apache Hadoop	●	●	●
SAP HANA	2011	Banco de dados em memória		●	
Google Megastore	2011	Banco de dados distribuído que precedeu o Google Spanner		●	●
Apache Storm	2011	Ferramenta para o processamento de eventos complexos			●
Apache Kafka	2011	Sistema para o processamento de fluxos contínuos de dados			●
MapR M5	2011	Conjunto de ferramentas de processamento distribuído	●	●	●
Hortonworks HDP	2011	Conjunto de ferramentas de processamento distribuído	●	●	●
Google Pregel	2010	Sistema distribuído para o processamento de grafos	●		
Google Percolator	2010	Sistema distribuído para processamento incremental	●		
Google Dremel	2010	Ferramenta para a análise interativa de dados		●	
AMPLab Spark	2010	Sistema de processamento de dados distribuído que opera em memória	●		
VoltDB	2010	Sistema de banco de dados em memória		●	
Apache S4	2010	Ferramenta para o processamento de fluxos contínuos de dados			●
HP Vertica	2010	Banco de dados MPP	●	●	
Apache Hive	2009	Ferramenta para o processamento de consultas interativas		●	
Cloudera CDH	2009	Conjunto de ferramentas de processamento distribuído	●	●	●
Apache Cassandra	2009	Sistema de armazenamento de dados distribuído		●	
Google BigTable	2006	Sistema de armazenamento de dados distribuído		●	
Apache Hadoop	2005	Sistema de processamento de dados distribuído	●		
Google MapReduce	2004	Sistema de processamento distribuído que deu origem ao Hadoop	●		

## Caracterização

### Processamento Batch

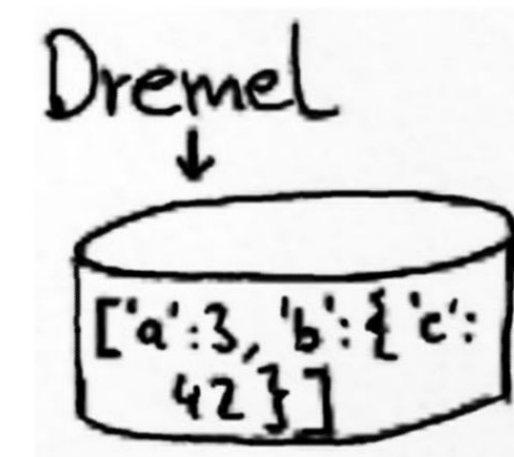
O Processamento Batch consiste no conjunto de aplicações que tem por objetivo o processamento de grandes volumes de dados de forma sequencial, sem que ocorra nenhuma intervenção durante o processamento.

Representado pelo Apache Hadoop e pelas ferramentas comerciais que surgiram a partir dele, o grupo possui uma intersecção com os sistemas de MPP (*Massively Parallel Processing*), que normalmente também executam tarefas de processamento batch por trás de suas arquiteturas.



### Processamento Interativo

O Processamento interativo é caracterizado pela necessidade de processar um tamanho intermediário de dados, uma vez que as plataformas dessa categoria operam sobre dados em sistemas de armazenamento distribuídos.

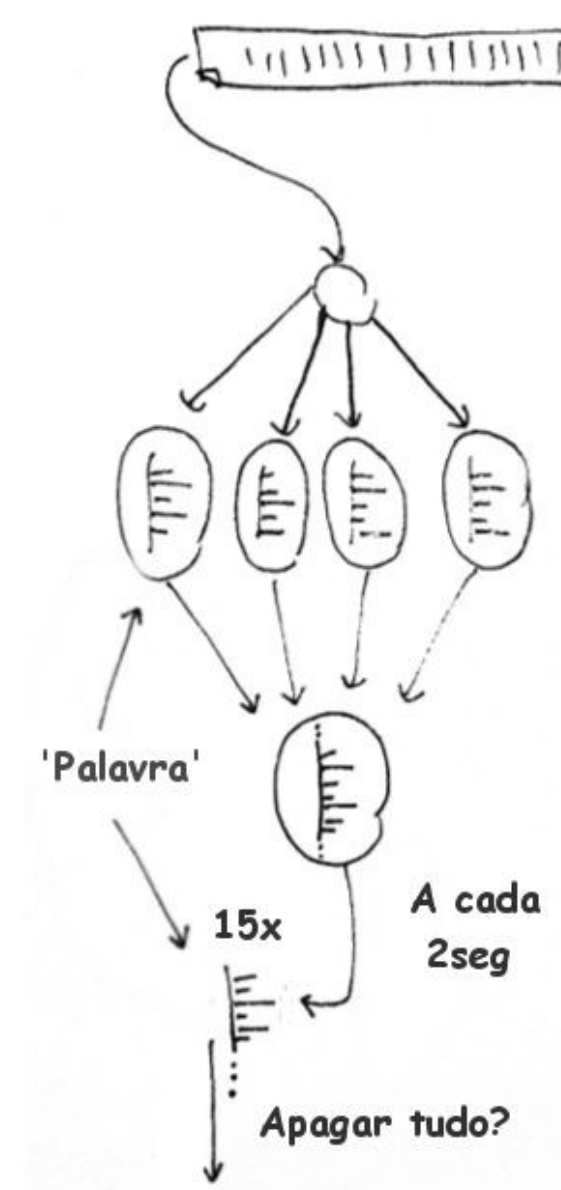


Esta categoria inclui os sistemas de armazenamento de dados distribuídos, incluindo os bancos de dados MPP, os bancos de dados NoSQL e os mais recentes bancos de dados NewSQL.

### Processamento em Tempo Real

O Processamento em Tempo Real visa suprir as novas necessidades de processamento, uma vez que o volume requisições e dados sendo produzidos é tão grande que, é mais interessante analisá-los como um fluxo contínuo de dados, do que depender de sistemas de armazenamento para realização de posterior de análises sobre os dados.

Esses sistemas são representados principalmente pelas ferramentas de processamento de fluxo (*Stream Processing*) e ferramentas de processamento de eventos complexos (*Complex Event Processing*).



## Conclusões

Este trabalho mostra que o ambiente de aplicações distribuídas para o processamento na nuvem não limita-se ao *Hadoop*, e está sendo constantemente estendido.

A proposta de caracterização, nos três grandes grupos sugeridos, facilita o processo de seleção das ferramentas e ajuda a determinar quais apresentam potencial para serem utilizadas por aplicações distribuídas na nuvem.

Entretanto, não é possível afirmar se as implementações atuais convergirão para grandes ferramentas, que ofereçam características de processamento diferentes para cada tipo de conjunto de dados e necessidade de tempo de resposta, ou se evoluirão para um conjunto ainda mais heterogêneo de modelos e ferramentas específicas para cada tipo de problema.