

# RAMCloud Project

Bruno Padilha

19 de Junho de 2015

# O que é o projeto RAMCloud ?

---

Sistema de armazenamento de dados de propósito geral que mantém os dados em memória RAM (DRAM) o tempo todo.

Objetivos: escalabilidade e baixa latência.

# Motivação

---

- Os discos magnéticos foram (e ainda são) o principal modo de armazenamento de dados de sistemas on-line nos últimos 40 anos;

# Motivação

---

- Os discos magnéticos foram (e ainda são) o principal modo de armazenamento de dados de sistemas on-line nos últimos 40 anos;
- Sistemas de arquivos e SGBDs evoluíram para funcionar com HDs;

# Motivação

---

- Os discos magnéticos foram (e ainda são) o principal modo de armazenamento de dados de sistemas on-line nos últimos 40 anos;
- Sistemas de arquivos e SGBDs evoluíram para funcionar com HDs;
- A evolução do desempenho dos HDs não acompanhou a sua evolução de capacidade de armazenamento;

# Motivação

---

- Os discos magnéticos foram (e ainda são) o principal modo de armazenamento de dados de sistemas on-line nos últimos 40 anos;
- Sistemas de arquivos e SGBDs evoluíram para funcionar com HDs;
- A evolução do desempenho dos HDs não acompanhou a sua evolução de capacidade de armazenamento;
- A dificuldade de escalabilidade da taxa de acesso aleatório aos dados acaba limitando a vazão de dados de um banco de dados relacional;

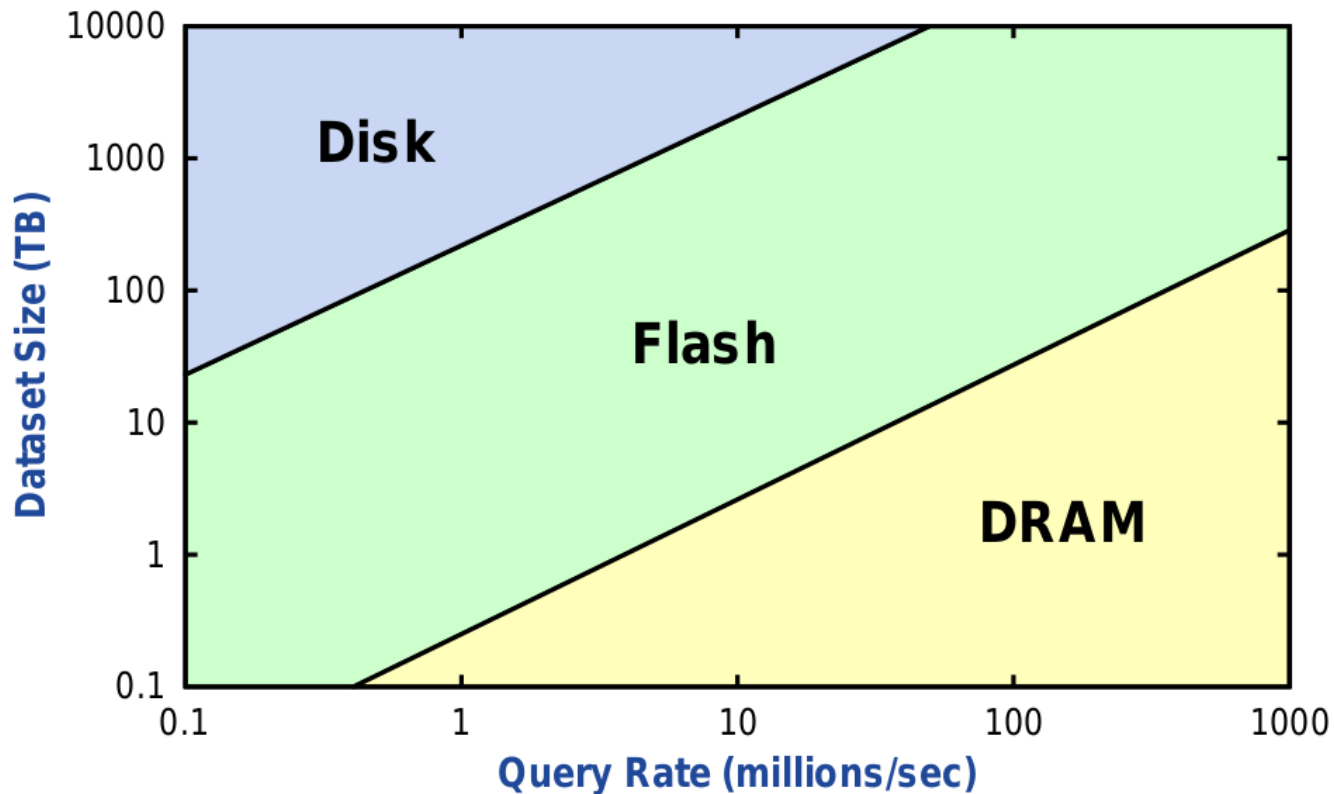
# Evolução do disco rígido (HD)

---

Comparação entre os discos rígidos comuns disponíveis em 2009 e em 1985

	Mid-1980's	2009	Change
<b>Disk capacity</b>	<b>30 MB</b>	<b>500 GB</b>	<b>16667x</b>
<b>Max. transfer rate</b>	<b>2 MB/s</b>	<b>100 MB/s</b>	<b>50x</b>
<b>Latency (seek &amp; rotate)</b>	<b>20 ms</b>	<b>10 ms</b>	<b>2x</b>
<b>Capacity/bandwidth (large blocks)</b>	<b>15 s</b>	<b>5000 s</b>	<b>333x</b>
<b>Capacity/bandwidth (1KB blocks)</b>	<b>600 s</b>	<b>58 days</b>	<b>8333x</b>
<b>Jim Gray's rule</b>	<b>5 min</b>	<b>30 hours</b>	<b>360x</b>

# Regra de Jim Gray

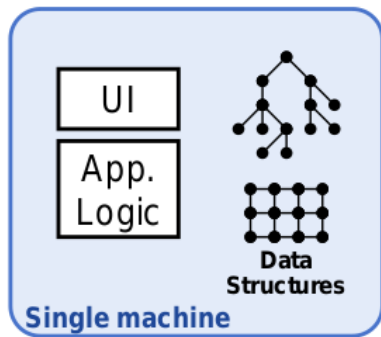




# Por que a latência importa ?

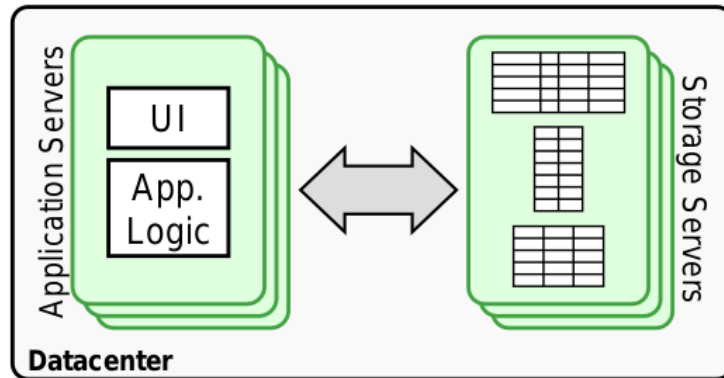
Uma aplicação Web com dezenas de centenas de milhares de usuários não é capaz de ser executada em uma única máquina.

**Traditional Application**



**$\ll 1\mu\text{s}$  latency**

**Web Application**



**0.5-10ms latency**

# Possíveis soluções

---

- Novos sistemas de banco de dados (NoSQL);
- Caching;
- **RAMCloud.**

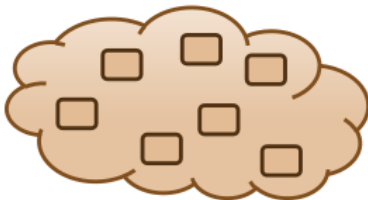
# RAMCloud

---

- Armazenamento em DRAM para DataCenters (chave-valor);
- 1000 a 10000 (possivelmente 100000) servidores commodity;
- 64 a 256 GB de DRAM por servidor;
- Durabilidade e disponibilidade;
- Alta vazão (1M ops/sec/server);
- Baixa latência (5us pra leitura e 15us para escrita);
- Rápida recuperação de falhas (1-2s);
- Custo de cache volátil.

# RAMCloud: Modelo de dados

## Tables



Object

Key ( $\leq 64\text{KB}$ )

Version (64b)

Blob ( $\leq 1\text{MB}$ )

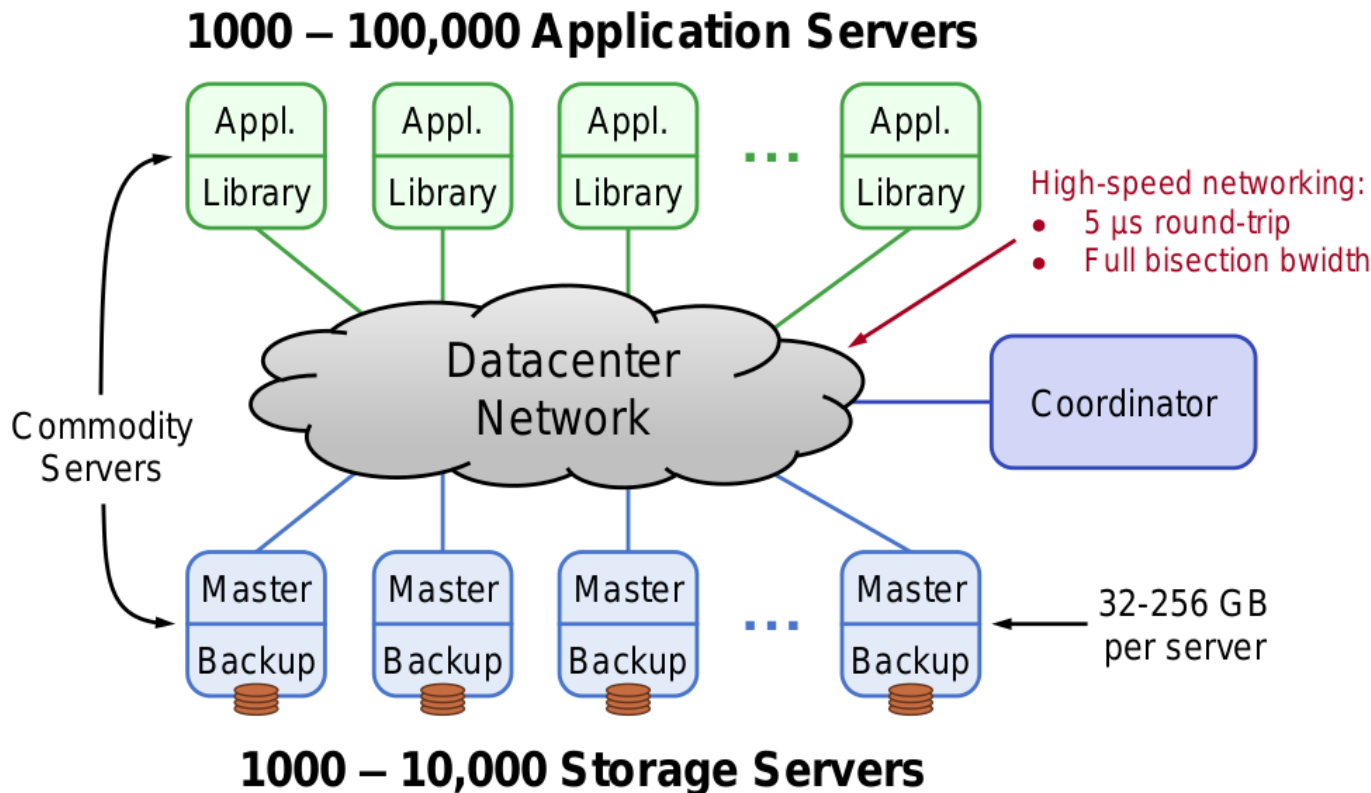
`read(tableId, key) => value, version`

`write(tableId, key, value) => version`

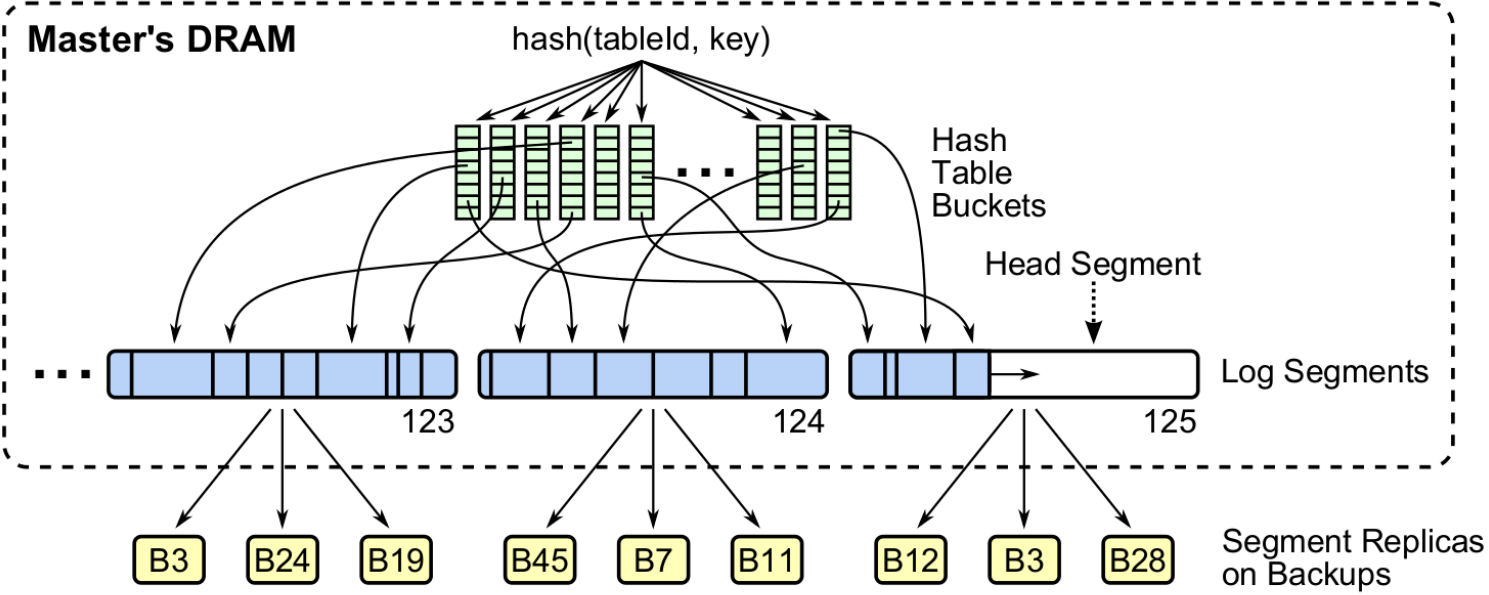
`conditionalWrite(tableId, key, value, condition) => version`

`delete(tableId, key)`

# RAMCloud: Arquitetura

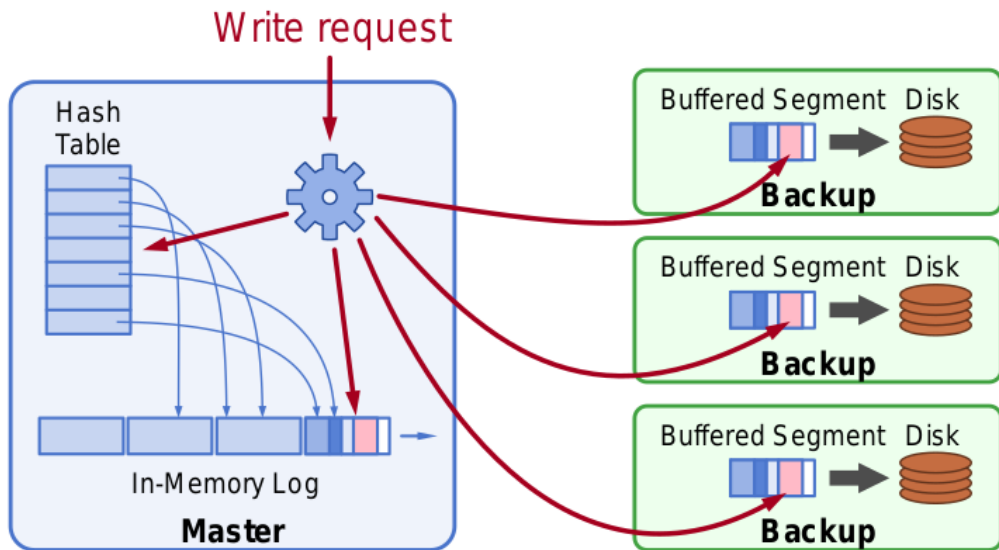


# RAMCloud: Organização

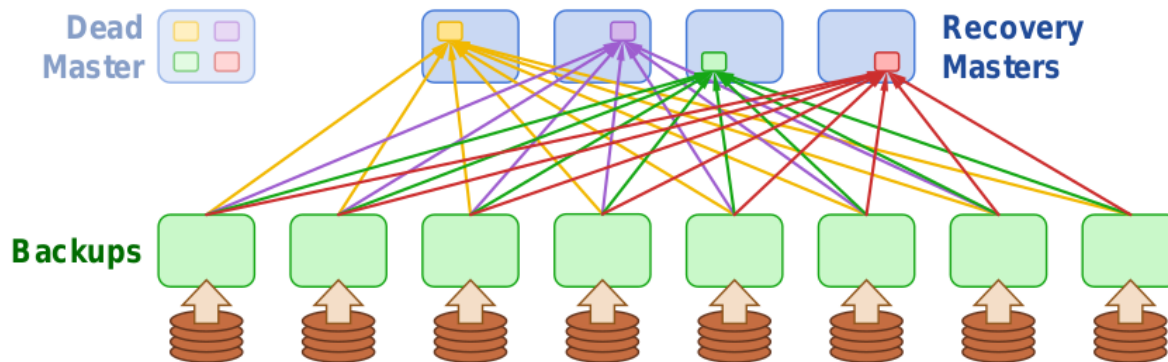


# RAMCloud: Buffered logging

Backup escrito de modo assíncrono (sem espera por I/O).



# RAMCloud: Crash recovery



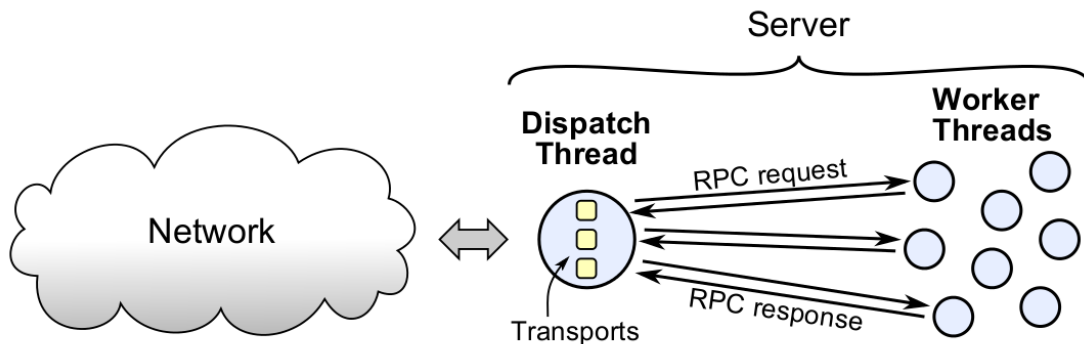
Os dados de um servidor mestre estão divididos em partições.

Em caso de falha, cada partição perdida é recuperada em um servidor distinto



# RAMCloud: Latência

Component	Traversals	2009	Possible 2014	Limit
Network switches	10	100-300 $\mu\text{s}$	3-5 $\mu\text{s}$	0.2 $\mu\text{s}$
Operating system	4	40-60 $\mu\text{s}$	0 $\mu\text{s}$	0 $\mu\text{s}$
Network interface controller (NIC)	4	8-120 $\mu\text{s}$	2-4 $\mu\text{s}$	0.2 $\mu\text{s}$
Application/server software	3	1-2 $\mu\text{s}$	1-2 $\mu\text{s}$	1 $\mu\text{s}$
Propagation delay	2	1 $\mu\text{s}$	1 $\mu\text{s}$	1 $\mu\text{s}$
<b>Total round-trip latency</b>		<b>150-400 <math>\mu\text{s}</math></b>	<b>7-12 <math>\mu\text{s}</math></b>	<b>2.4 <math>\mu\text{s}</math></b>



# RAMCloud: Utilizar

---

- Aplicações que utilizam muitos pedaços pequenos de dados independentes de uma dada requisição e precisam responder em tempo real;
- Aplicações que fazem muitas requisições dependentes que não podem ser feitas em paralelo (algoritmos em grafos). Exemplo: Facebook, Amazon, etc ...

# RAMCloud: Não utilizar

---

- Aplicações que não fazem uso intensivo de dados;
- Aplicações que processam dados em lotes;
- Aplicações cujos dados dependem de propriedades de bancos de dados relacionais .

# Referências

---

As seguintes referências encontram-se disponíveis na página do projeto RAMCloud Project:

<https://ramcloud.atlassian.net/wiki/display/RAM/RAMCloud>

- The RAMCloud Storage System;
- The Case for RAMCLoud;
- RAMCloud presentations (todas as figuras);
- Introductory talk on RAMCloud.

