

# **Armazenamento distribuído de dados e checkpointing de aplicações paralelas em grades oportunistas**

Raphael Yokoingawa de Camargo

TESE APRESENTADA  
AO  
INSTITUTO DE MATEMÁTICA E ESTATÍSTICA  
DA  
UNIVERSIDADE DE SÃO PAULO  
PARA  
OBTENÇÃO DO TÍTULO DE DOUTOR  
EM  
CIÊNCIAS

Área de Concentração: Ciência da Computação  
Orientador: Prof. Dr. Fabio Kon

Durante a elaboração deste trabalho o autor recebeu o auxílio financeiro do CNPq

São Paulo, 30 de Agosto de 2007

# **Armazenamento distribuído de dados e checkpointing de aplicações paralelas em grades oportunistas**

Este exemplar corresponde à redação  
final da tese devidamente corrigida  
e defendida por Raphael Yokoingawa de Camargo  
e aprovada pela Comissão Julgadora.

São Paulo, 30 de agosto de 2007

Banca Examinadora:

- Prof. Dr. Fabio Kon (orientador) - IME-USP.
- Prof. Dr. Marco Dimas Gubitoso - IME-USP
- Prof. Dr. Luiz Eduardo Buzato IC - Unicamp
- Prof. Dr. Markus Endler - PUC-Rio
- Prof. Dr. Francisco Brasileiro - UFCG

Dedico este trabalho à Selma, meu grande amor,  
cuja compreensão, dedicação e carinho trouxeram equilíbrio e grandes alegrias à minha vida.

## Agradecimentos

Agradeço de modo especial ao professor Fabio Kon que, com sua competente orientação, constantemente me motivou a superar os desafios da carreira acadêmica. Ao CNPq, pelo seu apoio financeiro durante todos estes anos.

A todos meus amigos do LCPD, das listas aabalados e psychhunters e da ETFSP. Dos simples almoços no bandejão às longas baladas noturnas, sua companhia foi essencial para tornar minha vida mais alegre. E um agradecimento especial àqueles que ainda contribuíram com meu trabalho de doutorado.

Aos meus pais, por todo o apoio e carinho que possibilitaram o meu amadurecimento e a conclusão deste doutorado, e às minhas irmãs e meus familiares, por sua constante presença a meu lado.

À Selma, por me apoiar nos momentos mais difíceis. Seu afeto e compreensão me deram forças para superar este grande desafio que é o doutorado.

## RESUMO

Grades computacionais oportunistas utilizam recursos ociosos de máquinas compartilhadas para executar aplicações que necessitam de um alto poder computacional e/ou trabalham com grandes quantidades de dados. Mas a execução de aplicações paralelas computacionalmente intensivas em ambientes dinâmicos e heterogêneos, como grades computacionais oportunistas, é uma tarefa difícil. Máquinas podem falhar, ficar inacessíveis ou passar de ociosas para ocupadas inesperadamente, comprometendo a execução de aplicações. Um mecanismo de tolerância a falhas que dê suporte a arquiteturas heterogêneas é um importante requisito para estes sistemas. Neste trabalho, analisamos, implementamos e avaliamos um mecanismo de tolerância a falhas baseado em checkpointing para aplicações paralelas em grades computacionais oportunistas. Este mecanismo permite o monitoramento de execuções e a migração de aplicações entre nós heterogêneos da grade.

Mas além da execução, é preciso gerenciar e armazenar os dados gerados e utilizados por estas aplicações. Desejamos uma infra-estrutura de armazenamento de dados de baixo custo e que utilize o espaço livre em disco de máquinas compartilhadas da grade. Devemos utilizar somente os ciclos ociosos destas máquinas para armazenar e recuperar dados, de modo que um sistema de armazenamento distribuído que as utilize deve ser redundante e tolerante a falhas.

Para resolver o problema do armazenamento de dados em grades oportunistas, projetamos, implementamos e avaliamos o middleware OppStore. Este middleware provê armazenamento distribuído e confiável de dados, que podem ser acessados de qualquer máquina da grade. As máquinas são organizadas em aglomerados, que são conectados por uma rede *peer-to-peer* auto-organizável e tolerante a falhas. Dados são codificados em fragmentos redundantes antes de serem armazenados, de modo que arquivos podem ser reconstruídos utilizando apenas um subconjunto destes fragmentos. Finalmente, para lidar com a heterogeneidade dos recursos, desenvolvemos uma extensão ao protocolo de roteamento em redes *peer-to-peer* Pastry. Esta extensão adiciona balanceamento de carga e suporte à heterogeneidade de máquinas ao protocolo Pastry.

## ABSTRACT

Opportunistic computational grids use idle resources from shared machines to execute applications that need large amounts of computational power and/or deal with large amounts of data. But executing computationally intensive parallel applications in dynamic and heterogeneous environments, such as opportunistic grids, is a daunting task. Machines may fail, become inaccessible, or change from idle to occupied unexpectedly, compromising the application execution. A fault tolerance mechanism that supports heterogeneous architectures is an important requisite for such systems. In this work, we analyze, implement and evaluate a checkpointing-based fault tolerance mechanism for parallel applications running on opportunistic grids. The mechanism monitors application execution and allows the migration of applications between heterogeneous nodes of the grid.

But besides application execution, it is necessary to manage data generated and used by those applications. We want a low cost data storage infrastructure that utilizes the unused disk space of grid shared machines. The system should use the machines to store and recover data only during their idle periods, requiring the system to be redundant and fault-tolerant.

To solve the data storage problem in opportunistic grids, we designed, implemented and evaluated the OppStore middleware. This middleware provides reliable distributed storage for application data, which can be accessed from any machine in the grid. The machines are organized in clusters, connected by a self-organizing and fault-tolerant peer-to-peer network. During storage, data is codified into redundant fragments, allowing the reconstruction of the original file using only a subset of those fragments. Finally, to deal with resource heterogeneity, we developed an extension to the Pastry peer-to-peer routing substrate, enabling heterogeneity-aware load-balancing message routing.

# Sumário

<b>1</b>	<b>Introdução</b>	<b>1</b>
1.1	Motivação . . . . .	4
1.2	Objetivos . . . . .	5
1.3	Organização da Tese . . . . .	5
<b>2</b>	<b>Checkpointing e Tolerância a Falhas</b>	<b>7</b>
2.1	Checkpointing de aplicações paralelas . . . . .	8
2.1.1	Conceitos e terminologia . . . . .	9
2.1.2	Protocolos de <i>checkpointing</i> . . . . .	11
2.2	Abordagens para <i>checkpointing</i> de processos . . . . .	14
2.2.1	Checkpointing no nível do sistema . . . . .	14
2.2.2	Checkpointing no nível da aplicação . . . . .	15
2.2.3	Obtenção do estado de um processo no nível da aplicação . . . . .	16
2.3	Execução de aplicações paralelas em arquiteturas heterogêneas . . . . .	18
2.3.1	O modelo de programação BSP . . . . .	18
2.3.2	Comunicação portátil entre processos . . . . .	21
2.4	Armazenamento de checkpoints . . . . .	21
2.4.1	Estratégias de armazenamento . . . . .	22
2.5	Resumo . . . . .	25
<b>3</b>	<b>Tolerância a Falhas de Aplicações no Sistema InteGrade</b>	<b>27</b>
3.1	Trabalhos relacionados . . . . .	27
3.1.1	Checkpointing no nível da aplicação . . . . .	28
3.1.2	Checkpointing de aplicações paralelas . . . . .	29
3.1.3	Checkpointing em grades computacionais . . . . .	29
3.1.4	Estratégias de armazenamento . . . . .	30
3.2	Visão geral de recuperação por retrocesso no InteGrade . . . . .	31
3.2.1	Arquitetura do InteGrade . . . . .	31

3.2.2	Recuperação por retrocesso baseada em <i>checkpointing</i> . . . . .	33
3.3	Implementação dos principais componentes . . . . .	35
3.3.1	Módulos de Gerenciamento . . . . .	35
3.3.2	Componentes para provedores de recursos . . . . .	39
3.3.3	Reinicialização de processos e aplicações . . . . .	39
3.4	Pré-compilador e biblioteca de <i>checkpointing</i> . . . . .	42
3.4.1	Pré-compilador . . . . .	42
3.4.2	Instrumentação do código de uma aplicação . . . . .	43
3.4.3	Criação de <i>checkpoints</i> e recuperação do estado da aplicação . . . . .	49
3.4.4	Transformação de dados . . . . .	51
3.4.5	Armazenamento e recuperação de checkpoints . . . . .	52
3.5	Aplicações paralelas BSP . . . . .	53
3.5.1	Execução de aplicações BSP em máquinas de diferentes arquiteturas . . . . .	53
3.5.2	<i>Checkpointing</i> de aplicações BSP . . . . .	54
3.6	Resumo . . . . .	54
<b>4</b>	<b>Avaliação do Mecanismo de <i>Checkpointing</i></b> . . . . .	<b>55</b>
4.1	Sobrecarga causada pelo mecanismo de <i>checkpointing</i> . . . . .	56
4.2	Execução de uma aplicação BSP na presença de falhas . . . . .	58
4.3	Execução de uma aplicação BSP em nós heterogêneos . . . . .	59
4.4	Recuperação do estado contido em um checkpoint . . . . .	60
4.5	Estratégias de armazenamento . . . . .	61
4.5.1	Codificação e decodificação de dados . . . . .	61
4.5.2	Sobrecarga sobre o tempo de execução . . . . .	62
4.6	Resumo . . . . .	64
<b>5</b>	<b>Armazenamento Distribuído de Dados</b> . . . . .	<b>67</b>
5.1	Sistemas <i>peer-to-peer</i> . . . . .	68
5.1.1	Redes <i>peer-to-peer</i> não-estruturadas . . . . .	69
5.1.2	Redes <i>peer-to-peer</i> estruturadas . . . . .	70
5.1.3	Balanceamento de carga em DHTs . . . . .	72
5.2	Armazenamento distribuído de dados em redes <i>peer-to-peer</i> . . . . .	74
5.2.1	CFS . . . . .	74
5.2.2	PAST . . . . .	75
5.2.3	OceanStore . . . . .	75
5.2.4	Outros sistemas de armazenamento <i>peer-to-peer</i> . . . . .	76
5.3	Armazenamento de dados em grades computacionais . . . . .	76



---

5.3.1	Grades de dados . . . . .	76
5.3.2	Transporte de dados . . . . .	77
5.3.3	Armazenamento e gerenciamento de réplicas . . . . .	78
5.3.4	Armazenamento com redes <i>peer-to-peer</i> . . . . .	79
5.4	Resumo . . . . .	80
<b>6</b>	<b>Identificadores Virtuais</b>	<b>81</b>
6.1	Trabalhos relacionados . . . . .	82
6.2	Pastry . . . . .	83
6.2.1	Protocolos . . . . .	84
6.3	Identificadores virtuais . . . . .	85
6.3.1	Visão geral dos protocolos e tabelas . . . . .	86
6.4	Protocolos sobre Pastry . . . . .	88
6.4.1	Protocolo de roteamento virtual . . . . .	88
6.4.2	Protocolo de partição do espaço virtual . . . . .	89
6.4.3	Protocolo de ingresso de nós . . . . .	91
6.4.4	Protocolo de saída de nós . . . . .	92
6.4.5	Protocolo de atualização de capacidade . . . . .	92
6.5	Implementação sobre o FreePastry . . . . .	93
6.5.1	FreePastry . . . . .	93
6.5.2	Arquitetura . . . . .	94
6.5.3	Conjunto de folhas e vizinhos virtuais . . . . .	95
6.5.4	Protocolos . . . . .	96
6.6	Resumo . . . . .	99
<b>7</b>	<b>OppStore: Middleware para Armazenamento Distribuído de Dados</b>	<b>101</b>
7.1	Trabalhos Relacionados . . . . .	101
7.2	Projeto do middleware OppStore . . . . .	103
7.2.1	Armazenamento e recuperação de dados . . . . .	105
7.2.2	Identificadores virtuais no OppStore . . . . .	108
7.3	Principais componentes . . . . .	109
7.3.1	Gerenciador de repositórios de dados do aglomerado (CDRM) . . . . .	109
7.3.2	Repositório autônomo de dados (ADR) . . . . .	110
7.3.3	Intermediador de acesso ( <i>Access Broker</i> ) . . . . .	111
7.4	Gerenciamento de dados armazenados . . . . .	112
7.4.1	Disponibilidade de dados . . . . .	112
7.4.2	Gerenciamento de arquivos . . . . .	113

---

7.4.3	Gerenciamento de fragmentos . . . . .	114
7.4.4	Otimização de desempenho . . . . .	115
7.5	Implantação sobre grades computacionais . . . . .	115
7.5.1	Interface com grades computacionais . . . . .	115
7.5.2	InteGrade . . . . .	116
7.5.3	Condor . . . . .	117
7.6	Implementação . . . . .	118
7.6.1	Principais componentes . . . . .	118
7.6.2	Gerenciamento de dados . . . . .	122
7.6.3	Interface com o InteGrade . . . . .	124
7.7	Resumo . . . . .	125
<b>8</b>	<b>Avaliação do middleware OppStore</b>	<b>127</b>
8.1	Simulações . . . . .	127
8.1.1	Ambiente de simulação . . . . .	127
8.1.2	Avaliação do uso de identificadores virtuais . . . . .	129
8.1.3	Recuperação de dados armazenados no OppStore . . . . .	133
8.2	Experimentos . . . . .	135
8.2.1	Armazenamento de dados . . . . .	136
8.2.2	Recuperação de dados . . . . .	137
8.3	Resumo . . . . .	137
<b>9</b>	<b>Conclusões</b>	<b>139</b>
9.1	Trabalhos Futuros . . . . .	140
9.2	Contribuições do Trabalho . . . . .	141
9.2.1	Contribuições tecnológicas . . . . .	141
9.2.2	Contribuições científicas . . . . .	142
9.3	Publicações durante o doutorado . . . . .	143

## Gracias por visitar este Libro Electrónico

Puedes leer la versión completa de este libro electrónico en diferentes formatos:

- HTML(Gratis / Disponible a todos los usuarios)
- PDF / TXT(Disponible a miembros V.I.P. Los miembros con una membresía básica pueden acceder hasta 5 libros electrónicos en formato PDF/TXT durante el mes.)
- Epub y Mobipocket (Exclusivos para miembros V.I.P.)

Para descargar este libro completo, tan solo seleccione el formato deseado, abajo:

