

Centro Federal de Educação Tecnológica de Minas Gerais

Departamento de Computação

Algoritmos paralelos de ordenação em Hadoop

Mariane Raquel Silva Gonçalves

Orientadora: Cristina Duarte Murta

Belo Horizonte

16 de março de 2012

Sumário

1	Introdução	2
1.1	Tema do projeto	2
1.2	Relevância	2
1.3	Objetivos	3
1.4	Resultados esperados	4
1.5	Metodologia	4
1.6	Infraestrutura necessária	4
1.7	Cronograma de trabalho	5
	Referências.....	7

1 Introdução

1.1 Tema do projeto

A ordenação é o processo de reordenar uma sequência de entrada e produzir uma saída ordenada de acordo com um atributo. A ordenação paralela realiza esse processo através de múltiplas unidades de processamento, que ordenam em conjunto a sequência de entrada. Na criação de algoritmos de ordenação paralela, é ponto fundamental ordenar coletivamente os dados de cada processo individual, de forma a utilizar todas as unidades de processamento e minimizar os custos de redistribuição de chaves entre os processadores. [Kale e Solomonik 2010]

1.2 Relevância

Com a recente mudança de foco no desenvolvimento de algoritmos sequenciais para paralelos, conhecer os modelos de programação paralela se tornou uma grande necessidade na computação. Os algoritmos paralelos ainda são um ramo pouco explorado, devido a maior complexidade no desenvolvimento e recentes aplicações em sistemas *multicore*.

A arquitetura paralela é um conceito conhecido há várias décadas

Uso crescente de computação paralela para sistemas computacionais gera a necessidade de algoritmos de ordenação inovadores, desenvolvidos para dar suporte a essas aplicações.

Um grande número de aplicações paralelas possui uma fase de computação intensa, na qual uma lista de elementos deve ser ordenada com base em algum de seus

atributos. Um exemplo é o algoritmo de Page Rank [Page et al. 1999] da Google: as páginas de resultado de uma consulta são rankeadas de acordo com sua relevância, e então precisam ser ordenadas de maneira eficiente. [Kale e Solomonik 2010]

Fatores como movimentação de dados, balanço de carga, latência de comunicação e distribuição inicial das chaves são considerados ingredientes chave para o bom desempenho da ordenação paralela, e variam de acordo com o algoritmo escolhido como solução. No exemplo do Page Rank, o número de páginas é enorme, e elas são recolhidas de diversos servidores da Google; é uma questão fundamental escolher algoritmo paralelo com o melhor desempenho dentre as soluções possíveis.

O MapReduce[Dean e Ghemawat 2008] é um modelo de programação paralela criado pela Google para processamento e geração de grandes volumes de dados em *clusters*. Esse modelo propõe simplificar a computação paralela e ser de fácil uso, abstraindo conceitos complexos da paralelização - como tolerância a falhas, distribuição de dados e balanço de carga - e utilizando duas funções principais: map e reduce. A complexidade do algoritmo paralelo não é vista pelo desenvolvedor, que pode se ocupar em desenvolver a solução proposta. O Hadoop [White 2009] é uma das implementações do MapReduce, um *framework* open source que provê o gerenciamento de computação distribuída. Foi desenvolvido por

1.3 Objetivos

Implementar e realizar a análise de desempenho, em termos do número de dados a serem ordenados e de máquinas utilizadas, de algoritmos de ordenação paralela, implementados de acordo com o modelo MapReduce no ambiente Hadoop.

(i) Estudar a programação paralela, seus algoritmos e suas possibilidades de implementação em ambiente paralelo multicore; (ii) Implementar e avaliar o desempenho de um algoritmo de ordenação paralela; (iii) Estudar e implementar a solução no modelo MapReduce, com o software Hadoop; (iv) Comparar os resultados obtidos com os algoritmos de ordenação e analisar seu desempenho com relação à quantidade de dados a serem ordenados, variabilidade dos dados de entrada e número máquinas utilizadas.

1.4 Resultados esperados

Com a realização do trabalho, busca-se ampliar e consolidar conhecimentos adquiridos na área de computação paralela, assim como a capacidade de análise e desenvolvimento de algoritmos paralelos no modelo MapReduce.

Ao final do trabalho, espera-se obter a implementação de algoritmos de ordenação paralela em ambiente Hadoop e uma análise comparativa de desempenho entre tais algoritmos.

1.5 Metodologia

O início do projeto será destinado ao estudo mais detalhado da computação paralela, em especial os algoritmos de ordenação paralela, dos fatores que influenciam o desempenho de tais algoritmos, o modelo MapReduce e a plataforma Hadoop. O passo seguinte é conhecer detalhadamente o algoritmo paralelo a ser implementado e definir as estratégias para sua implantação ambiente Hadoop. O algoritmo implementado deve ser cuidadosamente avaliado para verificar um funcionamento adequado com diferentes entradas e número de máquinas. Em seguida, serão realizados experimentos para testes de desempenho dos algoritmos com relação à quantidade de máquinas, quantidade de dados e conjunto de dados. Os resultados obtidos serão analisados e permitirão comparar a performance dos algoritmos em cada situação.

1.6 Infraestrutura necessária

A infra estrutura necessária ao desenvolvimento do projeto será fornecida pelo Laboratório de Redes e Sistemas (LABORES) do Departamento de Computação (DECOM). Esse laboratório possui um cluster formado por cinco máquinas Dell Optiplex 380, que serão utilizadas na realização dos testes dos algoritmos. Os algoritmos serão desenvolvidos em linguagem Java, de acordo com o modelo MapReduce, no ambiente Hadoop.

O cluster a ser utilizado apresenta as seguintes características:

- 5 nodos
- Processador Intel Core 2 Duo de 3.0 GHz
- Disco rígido SATA de 500 GB 7200 RPM
- Memória RAM de 4 GB
- Placa de rede Gigabit Ethernet
- Sistema operacional Linux Ubuntu 10.04 32 bits (kernel 2.6.XX)
- Sun Java JDK 1.6.0 19.0-b09
- Hadoop 0.20.2

1.7 Cronograma de trabalho

Na Tabela 1.7 está descrito o cronograma esperado para o desenvolvimento do projeto. Cada atividade foi descrita para se adequar da melhor maneira ao tempo disponível do projeto, mas é possível que o cronograma seja refinado posteriormente, com a inclusão de novas atividades ou redistribuição das tarefas existentes.

Citações:

[Kale e Solomonik 2010]

[Manferdelli, Govindaraju e Crall 2008]

[Dean e Ghemawat 2008]

[Asanovic et al. 2009]

Atividade	M	A	M	J	J	A	S	O	N	D
Definição do tema de trabalho.										
Pesquisa bibliográfica sobre o tema.										
Escrita da proposta de projeto.										
Estudo mais detalhado dos algoritmos de ordenação paralela, modelo MapReduce e Hadoop.										
Escrita da introdução e do referencial teórico do projeto.										
Configuração do ambiente Hadoop no laboratório.										
Testes iniciais para conhecer o funcionamento do Hadoop.										
Realização de experimentos com o algoritmo de ordenação paralela encontrado na literatura ou desenvolvido de acordo com a metodologia proposta.										
Descrição dos experimentos e da metodologia no texto do projeto.										
Finalização e entrega do projeto.										
Desenvolvimento algoritmos de ordenação.										
Desenvolvimento e testes dos algoritmos desenvolvidos										
Análise dos resultados obtidos até o momento, em busca de pontos de melhorias no projeto										
Teste final dos algoritmos, análise e escrita dos resultados										
Escrita e revisão do projeto final.										
Entrega e apresentação.										
Revisão nas observações realizadas pela banca.										

Tabela 1: Cronograma proposto para o projeto

Referências

- [Asanovic et al. 2009]ASANOVIC, K. et al. A view of the parallel computing landscape. *Commun. ACM*, ACM, New York, NY, USA, v. 52, n. 10, p. 56–67, out. 2009. ISSN 0001-0782.
- [Dean e Ghemawat 2008]DEAN, J.; GHEMAWAT, S. Mapreduce: simplified data processing on large clusters. *Commun. ACM*, ACM, New York, NY, USA, v. 51, n. 1, p. 107–113, jan. 2008. ISSN 0001-0782.
- [Kale e Solomonik 2010]KALE, V.; SOLOMONIK, E. Parallel sorting pattern. In: *Proceedings of the 2010 Workshop on Parallel Programming Patterns*. New York, NY, USA: ACM, 2010. (ParaPLoP '10), p. 10:1–10:12. ISBN 978-1-4503-0127-5.
- [Manferdelli, Govindaraju e Crall 2008]MANFERDELLI, J. L.; GOVINDARAJU, N. K.; CRALL, C. Challenges and opportunities in Many-Core computing. *Proceedings of the IEEE*, v. 96, n. 5, p. 808–815, maio 2008. ISSN 0018-9219.
- [Page et al. 1999]PAGE, L. et al. *The PageRank Citation Ranking: Bringing Order to the Web*. 1999.
- [White 2009]WHITE, T. *Hadoop: The Definitive Guide*. first edition. O'Reilly, 2009. Disponível em: <<http://oreilly.com/catalog/9780596521981>>.