

Sistema Gerador de Geradores de Código para Arquiteturas Superescalares

Mariza Andrade da Silva Bigonha¹
José Lucas Mourão Rangel Netto²

Resumo

Para obter um bom desempenho nas arquiteturas RISC, em particular, as arquiteturas superescalares, é necessário haver um bom relacionamento entre as tecnologias dos compiladores e as arquiteturas dos processadores. Baseado neste fato, arquiteturas recentes de computadores motivam pesquisas por técnicas de implementação de compiladores mais eficientes. Isto, contudo, acarreta maior complexidade dos compiladores porque estes novos computadores alcançam a eficiência, delegando aos compiladores a solução de problemas complexos de geração e otimização de código. Mostramos neste trabalho o projeto de um sistema gerador de geradores de código apropriado para arquiteturas superescalares. Mostramos também resultados de um estudo sobre vários problemas relacionados com a geração de código para estes processadores.

Abstract

A good marriage between compiler technology and processor architecture is necessary to achieve high performance on RISC, specifically, superscalar architectures. Based on that, new computer architectures have motivated research for more efficient compiler techniques. These new architectures, however, delegate the solution of the most complicated problems in code generation to the compiler. The focus of this paper is the design of a retargetable code generator system for superscalar architectures. We also show the results of the study of several problems related to code generation to these processors.

¹MSc (UFMG - 1985), DSc (PUC/RJ - 1994). É professora do DCC - UFMG. Áreas de interesse: linguagens de programação, compiladores, arquitetura. E-mail: mariza@dcc.ufmg.br

²MSc Eng. Elétrica (COPPE/UFRJ), Ph.D. em Ciência da Computação (Wisconsin, E.U.A.). É professor do Inst. de Matemática da UFRJ e do DI da PUC/RJ. Áreas de interesse: engenharia de software, compiladores, linguagens de programação, teoria da computação. E-mail: rangel@inf.puc-rio.br

1 Introdução

O alvo do sistema gerador de geradores de código descrito neste trabalho são os processadores superescalares. Estes processadores são uma evolução das arquiteturas RISC (*Reduced Instruction Set*). As principais características destas arquiteturas são: (a) facilidade de executar mais de uma instrução por ciclo e (b) incorporação de várias unidades funcionais que podem operar em paralelo permitindo a execução simultânea de instruções em unidades individuais. A forma como esta operação concorrente é extraída de um conjunto de instruções essencialmente seqüencial varia de arquitetura para arquitetura. Contudo, independentemente do mecanismo utilizado para buscar e executar instruções, o compilador deve administrar de forma eficiente a alocação de registradores e o escalonamento de instruções.

A alocação consiste no mapeamento das variáveis declaradas no programa fonte e das variáveis temporárias, geradas durante o processo de compilação, em registradores da máquina. Ela é considerada ótima se as variáveis permanecerem nos registradores todo seu ciclo de vida. Escalonamento é o processo de decidir a ordem das instruções de forma que elas possam ser executadas em unidades funcionais distintas com o objetivo de minimizar o tempo total de processamento do programa. Dois pontos são fundamentais nesta movimentação de instruções: boa utilização da máquina em análise e preservação da semântica do programa. Ou seja, um escalonamento válido deve preservar a ordem de execução descrita pelas arestas do grafo que descreve as dependências entre instruções.

O maior problema, que ainda permanecia sem uma solução satisfatória até esta data, está relacionado ao nível de comunicação entre a alocação e escalonamento. Questiona-se até que ponto estas duas funções devem interagir para se produzir um bom código e se o escalonamento de instruções ineficiente está associado à falta de comunicação entre eles. Trabalhos nesta linha [2, 3, 1, 19] e [4] sugerem algumas soluções.

O objetivo deste trabalho é apresentar o projeto de uma ferramenta de auxílio à implementação de compiladores de boa qualidade, cuja especificação inclui: (a) o projeto e definição formal da sintaxe e semântica de uma linguagem para descrever o conhecimento embutido nas arquiteturas de computadores (LDA), tendo em vista arquiteturas superescalares [5]; (b) o projeto de um Gerador de Geradores Redirecionáveis de Código Otimizado (GGCO) que, a partir da descrição de uma arquitetura de máquina específica, gere, automaticamente, tabelas que possibilitem dirigir o núcleo básico do gerador. Um gerador de código redirecionável é aquele que pode ser redefinido automaticamente para gerar código para novas máquinas, a partir de descrições de arquiteturas. (c) A identificação do nível de interdependência satisfatória entre o escalonamento e a alocação.

2 Metodologia de Construção de Compiladores

Os geradores e otimizadores de código desenvolvidos na década passada deram mais ênfase à seleção de instruções para as arquiteturas CISC (*Complex Instruction Set Computers*), como pode ser visto no sistema PO [9] e sucessores, no sistema CODEGEN [17], no sistema AutoCode [8]. Como as arquiteturas CISC implementavam as operações mais comuns de várias maneiras, escondendo detalhes de *pipelining* e outras características, geradores de código para estas arquiteturas se baseavam em especificações de máquina que permitiam, na maioria das vezes, que instruções fossem selecionadas utilizando-se de casamento de padrões. Estes geradores e

otimizadores não usavam técnicas de escalonamento e, muitas vezes, não efetuavam alocação global de registradores. Em contrapartida, arquiteturas RISC implementam a maior parte das operações de uma única forma e expõem ao gerador e otimizador de código a estrutura da *pipeline* e os custos das unidades funcionais. Conseqüentemente, a preocupação foi redirecionada para o escalonamento e a alocação, ao invés da seleção de instruções. Redirecionamento dos geradores e otimizadores de código continua importante para estas novas arquiteturas. Primeiro porque o desempenho de novos processadores continua sendo aprimorado, induzindo os usuários a melhorarem seus sistemas. Segundo porque o volume de *software* que deve ser transportado para uma nova máquina é grande o suficiente para proibir o uso intensivo de código em linguagem de montagem. Portanto, para que uma nova arquitetura seja competitiva no mercado, seu compilador deve ser capaz de gerar código de boa qualidade. Terceiro porque podem haver diferenças acentuadas entre diferentes famílias de processadores, tais como Intel i860 e MC88100 [18, 20]; ou mesmo diferenças significativas entre processadores de uma mesma família, como MC68000 e MC88100 [21, 20]. Um sistema redirecionável pode tornar possível a tarefa de construir geradores de código de boa qualidade em um período de tempo relativamente curto.

Devido à mudança de ênfase, problemas relacionados à geração de código para arquiteturas RISC são diferentes daqueles para arquiteturas CISC. Em primeiro lugar, para um compilador RISC redirecionável, a especificação de máquina deve capturar informações de escalonamento, incluindo a latência das operações e conflitos de recursos. Em segundo lugar, como a seleção de instruções nas arquiteturas RISC é relativamente simples, a interação entre a alocação e a seleção de instruções é considerada menos importante. Por outro lado, a interação entre a alocação e o escalonamento é significativa e o escalonador necessita registradores para sobrepor a execução de operações independentes.

Praticamente não existem sistemas de geração de código redirecionável especificamente projetado para arquiteturas RISC. Muito menos ainda para as arquiteturas superescalares. Os sistemas Gnu [24] e Marion [2] são uns dos poucos mais conhecidos. Até esta data Marion [2], é o único sistema que possui uma linguagem de descrição de máquina, contudo as primitivas de sua linguagem descrevem apenas os componentes básicos dos processadores RISC. O sistema Gnu, até pouco tempo, não possuía um escalonador de instruções. Atualmente, existem pelo menos duas de suas versões que já incorporam um método para escalonar instruções. Um deles usa o algoritmo de Gibbon et al. [13], onde a alocação é feita antes de escalonar as instruções e não existe nenhum tipo de comunicação entre eles. Outra versão de Gnu utiliza o algoritmo de Tiemann [26] onde latências dependentes da arquitetura alvo e informações sobre recursos computacionais são encapsulados.

3 O Sistema GGCO

A arquitetura de GGCO ilustrada na Figura 1 compreende as seguintes partes: (a) MD.c, arquivo contendo as tabelas e funções geradas a partir da especificação da arquitetura (veja Seção 5.1); (b) *gen-mdc*, módulo contendo a semântica de uma descrição em LDA (veja Seção 5); (c) MD.h, arquivo predefinido com definições das principais estruturas de dados e tipos usados em MD.c; (d) *front-end*, módulo que corresponde ao programa lcc [12]; (e) *back-end* [2], módulo que engloba as estratégias de escalonamento e alocação (veja Seção 3.2). A entrada para GGCO é a especificação de um processador em uma linguagem de descrição de arquitetura, LDA (veja Seção 4) e gera uma série de tabelas e funções que representam o resultado do processamento das informações extraídas das principais diretivas de LDA. Os arquivos MD.c, MD.h, o *front-end* e o *back-end* mostrados na Figura 1 são processados pelo programa Make e produzem o compilador

mcc para a arquitetura em questão. O compilador mcc recebe como entrada um programa fonte em linguagem C, o transforma para uma linguagem intermediária e então passa o controle para o gerador de código que gera um programa objeto semanticamente equivalente ao programa fonte.

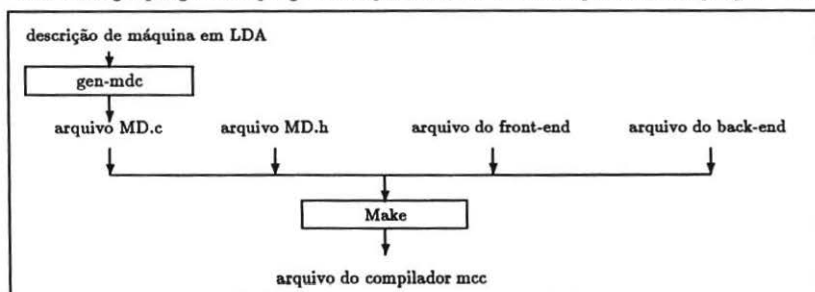


Figura 1: Estrutura do Gerador de Gerador de Código Otimizado-GGCO

3.1 O *Front-end*

Com a interface do programa lcc o módulo *front-end* independente de arquitetura e o módulo do gerador de código dependente de arquitetura estão bem acoplados. A vantagem desta abordagem é que a mesma produz compiladores eficientes e compactos. A desvantagem, como é sabida, complica a manutenção do compilador, porque consertos no *front-end* afetam o *back-end*. Contudo, como a linguagem para a qual pretende-se gerar código é o *ANSI C* e o foco principal deste trabalho é a questão da geração de código certamente existirão poucas mudanças na linguagem, então esta complicação não é tão importante.

O código executável é especificado por DAGs (grafo acíclico dirigido), ou seja, o *front-end* do compilador emite código para uma linguagem DAG. Esta linguagem DAG corresponde a linguagem intermediária gerada a partir do programa fonte e provê a configuração inicial para o DAG de código apresentado na Seção 3.2.2. Os vértices de um DAG representam os possíveis operadores do conjunto de instruções de uma arquitetura.

O *front-end* transforma todos os operadores de fluxo de controle, tais como, "for", "while" e "if" em operações de mais baixo nível, como comparação e desvio. Todos os operadores da linguagem C com efeitos colaterais, por exemplo, operadores de atribuição, operadores condicionais, etc., são transformados em explícitas operações aritméticas, lógicas, de atribuição e desvios. Dois conjuntos de transformações são aplicados na linguagem intermediária antes de efetuar a seleção de código. O primeiro conjunto efetua alocação de dados usando algumas informações dependentes de máquina. Para cada objeto declarado pelo usuário, o alocador de dados aloca: (a) um pseudo-registrador³ se o objeto é uma variável automática que pode ser atribuída a um registrador; (b) uma posição na pilha do sistema se o objeto é uma variável automática, mas não pode ser atribuída a um registrador; (c) uma posição na área de dados global, se o modelo de máquina possui tal área e o objeto é estático; (d) um endereço relocável. No segundo conjunto, os endereços dos operadores são expandidos ou substituídos por referências a pseudo-registradores e os vértices são inseridos para retornar valores em registradores ou posições da pilha, etc.

³registradores criados para conter valores intermediários de expressões e não existe entre eles e o hardware dependência alguma.

3.2 Gerador de Código

Cada gerador de código é construído a partir da descrição de uma determinada arquitetura e efetua a seleção de código, o escalonamento de instruções e a alocação de registradores. GGCO efetua a seleção de instruções antes de passar o controle para a estratégia de geração de código. A estratégia de geração de código dirige a chamada e o grau de comunicação entre a alocação e o escalonamento, além de incluir o algoritmo de escalonamento.

A composição da estratégia de geração de código de GGCO é a mesma do sistema Marion [2]. Ela possui duas partes: (a) uma parte independente da estratégia usada e (b) uma parte dependente. A parte independente da estratégia compreende três componentes lógicos: o construtor do DAG de código (Seção 3.2.2), o alocador de registradores e o suporte para o escalonamento. O construtor do DAG de código é responsável pela construção do DAG das instruções de máquina para cada bloco básico. Bloco básico é uma seqüência de instruções na qual só existe uma entrada e uma saída, isto é, não existem instruções de desvios para dentro do bloco ou para fora, exceto a última instrução que pode ser de desvio. O alocador determina como é usado o conjunto de registradores físicos. O suporte para escalonamento é responsável pela verificação de conflitos de recursos e pela manutenção da lista de instruções aptas a serem escalonadas sem provocar atrasos. A parte dependente da estratégia engloba: algoritmo de escalonamento, tabelas e rotinas geradas a partir da descrição da máquina. Sua estrutura modular permite a substituição ou inserção de novas estratégias de escalonamento e alocação no módulo de geração de código. Baseado nesta modularidade, propusemos a incorporação de uma nova estratégia em GGCO.

3.2.1 Seletor de Código

O seletor de código utiliza o método de seleção descendente recursivo por *pattern matching* para casar vértices das árvores objetos⁴ com vértices das árvores padrões⁵. Para cada instrução da máquina alvo existe um padrão e um símbolo de substituição que a identifica. O algoritmo básico da seleção de código examina os padrões em uma dada ordem, selecionando o primeiro que casa com operadores da árvore objeto, e então tenta casar as sub-árvores, começando pela sub-árvore à esquerda. Se foi bem sucedido, junta-se o padrão ao vértice do DAG da linguagem intermediária. Se não for bem sucedido o seletor recua e examina o próximo padrão da árvore de padrão. Após a seleção, caminha-se no DAG da linguagem intermediária da esquerda para a direita, *bottom-up*, gerando código de acordo com os padrões colocados nos vértices do DAG. Durante a seleção, criam-se *pseudo-registradores* para todas as expressões temporárias. Como a entrada para o seletor de código é um DAG da linguagem intermediária, para acomodá-lo em uma árvore padrão todos os vértices da linguagem intermediária que possuem mais de um pai são forçados em registradores, a menos que sejam constantes que possam ser substituídas por modos de endereçamentos ou operandos imediatos.

3.2.2 Escalonamento de Instruções

A estrutura de dados mais importante no processo de escalonamento é o grafo de escalonamento, o DAG de código. Representa-se nesta estrutura de dados os blocos básicos em que foi dividido

⁴sub-árvores do DAG da linguagem intermediária, que constituem a entrada para o seletor de código.

⁵árvores compostas de operadores da linguagem intermediária, derivadas a partir da descrição da máquina pelo gerador de código.

o programa. Como o escalonamento considerado neste trabalho é feito para blocos básicos, as restrições de precedências consideradas são: restrições baseadas em dependência de dados e dependências de recursos. Dependência de controle existe somente entre blocos básicos e suas correspondentes arestas são derivadas do grafo de fluxo do programa. As dependências de dados entre instruções podem ser verdadeiras ou falsas. As dependências falsas podem ser classificadas em anti-dependência e dependência de saída. Uma dependência de dados de uma instrução u para uma outra v existe se uma das seguintes afirmativas for verdadeira: (a) dependência verdadeira ou de fluxo de dados: registrador definido em u é usado em v ; (b) anti-dependência: registrador usado em u é redefinido em v e destrói o valor usado em u ; (c) dependência de saída: registrador definido em u é redefinido em v , destruindo o valor definido anteriormente em u .

A abordagem mais utilizada para escalonamento de instruções, denominada *lista de escalonamento* [11], [15], [16], [10], [13], [2], [27], [14] funciona da seguinte forma: dado um DAG de código, o escalonador mantém uma lista de instruções que estão aptas para serem escalonadas sem provocar atrasos. A cada iteração utiliza-se uma heurística para selecionar uma instrução pronta para ser escalonada e então atualiza-se a lista. Em geral, esta abordagem possui, no pior caso, um tempo de execução equivalente a $O(e)$, onde e refere-se ao número de arestas no DAG, mas a heurística pode crescer em complexidade.

Em uma lista de escalonamento, normalmente escolhe-se o vértice que possui a prioridade mais elevada em relação aos outros na lista. A heurística mais utilizada para atribuir prioridades denomina-se *distância máxima* e é definida como o comprimento do caminho⁶ mais longo através do DAG de código, partindo do vértice da instrução até o vértice folha. Muitas vezes, a *distância máxima* é também denominada *peso* do vértice dentro do DAG. O raciocínio utilizado é que o vértice mais longe de ser atingido é o mais crítico, vértices menos importantes podem ser escalonados mais tarde. A distância máxima é considerada uma heurística inicial relativamente satisfatória, mas existem situações onde heurísticas alternativas ou complementares são mais vantajosas. Um segundo nível de heurística atribui uma prioridade mais alta para *vértices possuindo mais sucessores*. A idéia por trás desta heurística é que escalonando um vértice com mais sucessores, criam-se mais oportunidades para o escalonador nos ciclos subsequentes e permite-se que mais vértices estejam prontos para serem escalonados mais cedo. Uma outra heurística atribui uma prioridade mais alta ao vértice que possui uma latência maior em relação a um de seus sucessores. O ato de escalonar tal vértice primeiro permite mais oportunidades de preencher o período de latência com o escalonamento de outras instruções. Toda a estratégia de geração de código de GGCO utiliza algoritmos de listas de escalonamento com a distância máxima como primeira heurística.

3.2.3 Alocação de Registradores

Até recentemente, a coloração do grafo de interferência era a abordagem mais usada para resolver o problema da alocação eficiente de valores a registradores físicos em um compilador. Neste grafo, os vértices, representando pseudo-registradores, e as arestas, representando conflitos entre pseudo-registradores em um intervalo de tempo, eram coloridos por um número de cores, derramando para a memória alguns valores quando necessário. Com o advento dos processadores superescalares permitindo o paralelismo entre instruções, algoritmos de coloração de grafos considerados ótimos outrora já não se correlacionam com uma boa utilização dos recursos destas máquinas. Para explorar o paralelismo entre instruções, as mesmas devem ser reordenadas. Quando a alocação é feita antes da reordenação, a seleção de registradores pode

⁶Comprimento de um caminho=soma de todos os valores de rótulos das arestas ao longo do caminho.

limitar as possibilidades de escalonamento devido a introdução de dependências falsas com o reuso dos registradores. Quando o escalonamento é feito antes da alocação o número de pseudo-registradores ativos aumenta dando origem a sérias implicações: (1) seu tempo de vida cresce; (2) são necessários mais registradores físicos; (3) ocorrem mais derramamentos para a memória.

Os métodos existentes usam para as funções de escalonamento e alocação diferentes modelos de grafos para representar o programa fonte, [14], [27], [2] etc. Como o significado dos vértices e arestas nestes grafos são distintas, é difícil obter uma simples combinação dos mesmos.

A estratégia proposta para GGCO é diferente destes métodos citados, ela provê apenas um modelo de grafo, o *grafo de interferência paralelizável*, para representar o programa fonte para estas duas funções [23]. Nesta abordagem, a ênfase é dirigida ao alocador de registradores e o método usado para alocação é baseado no trabalho de Chaitin [7].

O algoritmo [23] funciona da seguinte forma. Para gerar o *grafo de interferência paralelizável* primeiramente inclui todas as dependências e restrições de recursos explicitamente no grafo de escalonamento. Quanto mais arestas estiverem presentes neste grafo melhor será o resultado, porque o algoritmo usa somente as arestas que estão no complemento do grafo construído. Estas arestas representam o paralelismo existente na máquina para um dado programa. Em seguida integra-se o grafo de complemento com o grafo de interferência permitindo que o algoritmo de alocação leve em consideração o paralelismo existente. Com este novo grafo, efetua-se a alocação de registradores. O escalonador propriamente dito é ativado após a alocação.

Como o problema da coloração mínima é NP-completo e em geral o número de registradores é menor que o conjunto mínimo de cores da melhor coloração obtida. Na prática, derramamentos são efetuados, armazenando-se temporariamente os valores de algumas variáveis na memória. Tendo em vista estes fatos, o problema da alocação para arquiteturas superescalares é encontrar um mapeamento de registradores que faça uso de um número mínimo de registradores físicos, que o custo de derramamento seja mínimo e cujo grafo de escalonamento não possua dependências falsas. Para isto aplicam-se sobre o *grafo de interferência paralelizável* as mesmas heurísticas usadas tanto na alocação como no escalonamento. O primeiro tipo de heurística elimina arestas do grafo. A questão de qual aresta deve ser eliminada pode implicar em considerações tanto do escalonador como do alocador. Pode-se considerar a remoção de algumas arestas que evitam dependências falsas, neste caso, está-se fazendo o trabalho do escalonador quando, devido a demanda por registradores, algumas opções de paralelismos são perdidas. Pode ser considerado também preservar algumas arestas que prometem um bom paralelismo e decidir remover arestas de interferência as quais podem ocasionar derramamentos de código⁷.

Os algoritmos de Chaitin [7] e Pinter [23] não discutem o problema de pares de registradores. Pares de registradores são muitas vezes necessários em processadores para representar duas metades de registradores usando precisão dupla para cargas e armazenamentos, para movimentação de registradores, etc. Em seus algoritmos, um pseudo-registrador v pode ser removido do grafo se for garantido existir um registrador físico para ele durante a fase de coloração. Isto significa que o seu grau, ou seja, o número de seus vizinhos, deve ser inferior ao número de registradores disponíveis. A necessidade por mais de um registrador para um dado vértice v muda a definição de seu grau. Um vértice v é considerado sem restrições quando a soma das necessidades de registradores físicos de seus vizinhos mais o número de registradores físicos necessários a ele for menor que o número de registradores alocáveis. Com esta modificação introduzida é possível obter uma coloração do *grafo de interferência paralelizável* mesmo quando a demanda dos vizi-

⁷O derramamento pode ser evitado somente quando se atribuem para dois vértices deste tipo cores diferentes apesar da ausência da aresta de interferência.

nhos de um vértice v seja maior que o número de registradores físicos alocáveis. A reutilização de cores em vizinhos cujas arestas não se interferem neste grafo pode gerar uma coloração de G sem alterar o objetivo do algoritmo proposto por Pinter, que é obter uma ótima alocação de registradores sem introduzir dependências falsas no grafo de escalonamento.

4 A LDA

A linguagem de descrição de arquitetura, LDA, usada em GGCO é composta de três seções: (a) declarações, (b) características da máquina alvo e (c) instruções. Na seção de declarações são especificados registradores, recursos da máquina, constantes e o tamanho da memória disponível entre outras informações. Na seção com as características da máquina alvo são declarados conjuntos de registradores de propósito geral, argumentos, apontadores de pilha e registro de ativação, área global, etc. A seção de instruções introduz a descrição de instruções da máquina, instruções especiais e transformações necessárias para casar padrões da linguagem intermediária com padrões da linguagem da máquina alvo (detalhes sobre a LDA são encontrados em [5]).

A LDA ainda não cobre toda a classe das máquinas superescalares, mas, ela modela a maior parte de suas características, abrangendo também as máquinas RISC. Para ter uma idéia de sua aplicabilidade bem como do sistema GGCO como um todo destacamos alguns aspectos destas arquiteturas ainda não cobertos por outros sistemas:

(A) A LDA oferece ao projetista do compilador facilidades para especificar separadamente os argumentos e parâmetros, o que modela a renomeação de registradores, e, também facilidades para modelar o salvamento dos registradores especificado pelo procedimento chamador e chamado. Com isto cobre janelas de registradores. (B) A descrição de máquina em LDA relaciona as necessidades de recursos para cada instrução da máquina alvo. A partir desta informação, GGCO constroi um *vetor de recursos* para cada instrução. Cada elemento deste vetor contém todos os recursos necessários à instrução em um determinado ciclo de máquina. Para evitar a ocorrência de uma dependência de recurso, o escalonador compara o vetor de recursos usados pela instrução candidata ao escalonamento com um vetor de recursos que é a composição dos recursos necessários a todas as instruções correntemente em execução. Se a interseção for vazia, a instrução candidata é escalonada. Para solucionar dependências mais complicadas, como por exemplo, o esquema de prioridades por *hardware* para contenção de recursos, o GGCO adota o seguinte esquema: (a) permite que uma série de prioridades seja associada com a "declaração de recursos" na descrição da máquina; (b) permite que um elemento do vetor de recursos seja associado com uma "instrução" para indicar sua prioridade; (c) examina prioridades durante a verificação por dependências de recursos e permite que escalonamentos contenham estas dependências, se eles forem causados por recursos prioritizados.

(C) Para evitar as dependências de controle, a LDA especifica o número de *delay slots* dentro da diretiva de instrução. Para evitar o preenchimento dos *slots* com *no-ops* como acontece com o sistema Marion [2] propusemos a implementação da técnica de Hennessy e Gross [16] em GGCO, incluindo-a como um passo separado, intraprocedimental, após o escalonamento. Nesta técnica tenta-se preencher os *slots* com instruções válidas que estão antes da instrução de desvio, com instruções que seguem o objeto do desvio e com instruções que seguem a própria instrução de desvio.

(D) O sistema GGCO foi especificado para dar suporte ao despacho múltiplo de instruções particionando o conjunto de instruções de tal forma que uma instrução de cada partição pode ser despachada a cada ciclo.

(E) A maior parte das arquiteturas superescalares possuem instruções que usam o registrador de código de condição. A solução adotada neste caso é declarar o registrador de código de condição

como sendo um registrador temporal⁸ na seção de declarações da LDA. Declarar na seção de instruções as instruções que o modificam e examinar este registrador durante o escalonamento, garantindo que o mesmo está ligado e corretamente referenciado.

4.1 Mapeamento com a Linguagem Intermediária

Raramente existe um mapeamento direto entre uma linguagem intermediária e um conjunto de instruções da máquina alvo. GGCO incorpora de Marion dois mecanismos que auxiliam o usuário a completar este mapeamento. O primeiro é a transformação "glue" e o segundo são as funções de escape ("*func"). A transformação "glue" é normalmente empregada quando um operador da linguagem intermediária pode ser mapeado em duas ou três instruções de máquina. O exemplo da Figura 2 mostra que para comparar e desviar em uma dada arquitetura é necessário a utilização da instrução "subcc", que coloca o código de condição em um registrador de propósito geral e a instrução de desvio "be", que verifica a condição neste registrador e desvia de acordo com seu valor. Para mapear o operador da linguagem intermediária "==" nesta seqüência de duas instruções, uma transformação "glue" expande a subárvore "(\$1 == \$2)" na subárvore "(((\$1 :: \$2) ==))". O novo operador "::", presente na linguagem intermediária de GGCO, casará a instrução "subcc" e "==" casará a instrução "be". O segundo mecanismo

%instr subcc r, r, r[0] (; clk_cc)	%instr be #rlab
{cc = (\$1 :: \$2); }	{if cc == 0
[IF; ID; IE; IW;]	goto \$1; }
(1, 1, 0)	[IF; ID; IE;]
	(1, 1, 1)
%glue r, #UCONST13	
(((\$1 == \$2) ==> ((\$1 :: \$2) == 0;)	

Figura 2: Instruções em LDA

para auxiliar no mapeamento entre a linguagem intermediária e o conjunto de instruções de máquina alvo, as funções de escape, resolve situações em que são necessários transformações que manipulam metades de registradores e transformações que contém sub-expressões comuns e, em geral, permite ao usuário efetuar mapeamentos complicados sem afetar a geração eficiente de código. Este mecanismo permite ao usuário escrever funções em linguagem C para produzir uma seqüência de instruções individuais escalonáveis. Estas funções podem e devem incluir somente chamadas às primitivas exportadas por GGCO. Estas primitivas criam, manipulam operandos e geram instruções. As funções escritas pelos usuários são invocadas no momento em que o seletor de código casa uma função de escape.

⁸regitradores usados para conter os resultados das sub-operações em arquiteturas que possuem *pipeline de avanço explícito*.

5 Filosofia da Definição Formal de LDA

A definição formal da LDA segue o método denotacional para especificação de semântica. As descrições de arquitetura, que compreendem o domínio sintático⁹ desta descrição formal, são apresentadas na forma de sintaxe abstrata. A Figura 3 mostra o mapeamento da descrição de arquitetura de máquinas para o seu respectivo gerador de código. Este mapeamento é definido pela função `gen-mdc` que define a semântica de uma descrição em LDA. O resultado final deste mapeamento é um trecho de programa em linguagem C, ou seja, o mapeamento de uma descrição de arquitetura para o domínio Mdc, que constitui a parte dependente de máquina do gerador de código para a arquitetura descrita. A função `gen-mdc` recebe como entrada um

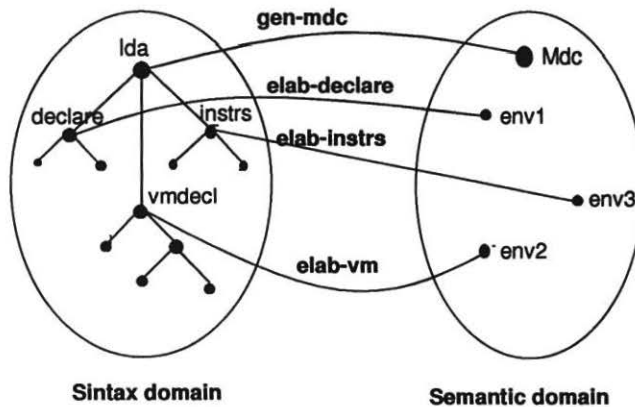


Figura 3: Mapeamento de lda em mdc.

arquivo com a especificação `lda`, na forma de uma árvore de sintaxe abstrata, ativa as funções: `elab-declare`, `elab-vmdecl`, `elab-instr` e `elab-tables` para elaborar as várias partes da descrição de máquina e produz o código dependente de máquina do gerador de código. Esta definição formal pode ser vista como um sistema gerador de código a partir da descrição de arquitetura. O trecho de programa gerado é armazenado no arquivo `MD.c`. A função `elab-declare` elabora as declarações de LDA; a função `elab-vmdecl` elabora as informações sobre a máquina alvo; a função `elab-instr` elabora as instruções de LDA, e a função `elab-tables` elabora as tabelas finais, produzindo o código C.

Na especificação deste mapeamento, destaca-se o domínio semântico "env", que define o ambiente no qual as descrições formais são estabelecidas, ou seja, define uma tupla composta da identificação de todas as entidades auxiliares geradas para a elaboração das tabelas finais. Este domínio compreende um conjunto de tabelas e listas construídas a partir de cada especificação em LDA. Ele é parte da assinatura da maior parte das funções semânticas.

A presente definição formal está escrita em *SCRIPT* [6]. *SCRIPT* é uma linguagem funcional que oferece uma notação simples para descrição da semântica denotacional de uma linguagem

⁹Domínios podem ser vistos como sinônimos de tipo, podendo ser utilizados para modelar propriedades sintáticas ou semânticas em uma linguagem de programação. O domínio de um símbolo terminal é ele mesmo. O domínio de símbolos não-terminais podem ser inferidos, capitalizando sua primeira letra.

de programação de forma modular e legível. A sintaxe abstrata é definida em *SCRIPT* como uma lista de produções de uma gramática livre do contexto. Os não-terminais representam domínios sintáticos e os terminais são domínios de "quotations", contendo apenas o elemento especificado. Os vários domínios semânticos da definição formal são agrupados em módulos, que controlam a visibilidade de suas denotações e provêm os serviços associados a cada domínio. Para cada um dos principais domínios semânticos usados na definição do mapeamento *gen-mdc* existe um módulo *SCRIPT*, que encapsula as suas denotações e provêm os serviços para sua administração.

5.1 Principais Tabelas Geradas

A função de nível mais alto da definição formal de LDA, a função *gen-mdc* mapeia descrição LDA em trechos de programas em C, os quais definem a porção dependente de arquitetura do gerador de código. Este trecho de programa, contido no arquivo MD.c, pressupõe as definições das estruturas de dados e tipos declarados no arquivo pré-definido MD.h (veja Seção 3).

As tabelas mais importantes são a tabela de recursos e a tabela de produções. A tabela de recursos descreve os recursos utilizados pelas instruções e constitui uma das informações mais importantes para o escalonamento. Seu conteúdo é usado essencialmente na rotina básica de escalonamento de instruções, para verificar os conflitos existentes e agrupar instruções.

A tabela de produções contém informações sobre as instruções. Cada linha desta tabela corresponde a uma diretiva de instrução fornecida na descrição da arquitetura e contém o seguinte: (a) uma árvore padrão e um símbolo de substituição derivado da expressão contida na diretiva; (b) uma matriz que indica para cada instrução, os tipos de seus operandos e suas localizações dentro do padrão e árvores objetos; (c) um índice para uma matriz de vetores de recursos e (d) dados referentes a latências, custo, *delay slot*, etc. Informações contidas na tabela de produções são utilizadas em várias funções do gerador de código. Por exemplo, nas rotinas de escalonamento, na construção do DAG de código, na construção do DAG de interferência de registradores, durante a coloração do grafo de interferência, no reconhecimento de padrões, etc. Outras estruturas de dados contêm declarações, informações da máquina virtual, classes, elementos, latências auxiliares, endereços das funções de escape e padrões de transformações "glue".

As rotinas dependentes da arquitetura mais importantes são aquelas que retornam o conjunto de registradores de propósito geral para um dado tipo e aquela que contém informações relativas a sobreposição de registradores. Estas informações são utilizadas durante a criação dos *pseudo-registradores* e a alocação de registradores físicos.

6 Modelos de Pilha de GGCO

Uma informação importante para o projetista do compilador é o modelo de pilha utilizado na implementação de uma dada arquitetura. Quatro modelos de pilha são usados para descrever o modelo de execução da máquina alvo. Estes modelos são definidos de acordo com a direção do par (*stack*, *frame*) apresentados abaixo. O par (*stack*, *frame*) é obtido diretamente das informações especificadas pelas diretivas que manipulam a pilha e o registro de ativação definidos na LDA [5]. (*down,down*): a pilha cresce para baixo, o apontador do registro de ativação é usado para argumentos e variáveis locais; e a área derramada está entre as variáveis locais e o

apontador da pilha.

(down,up): a pilha também cresce para baixo mas o apontador do registro de ativação é usado somente para argumentos e a área derramada está entre as variáveis locais e os argumentos.

(up,down): a pilha cresce para cima, o apontador do registro de ativação é usado somente para os argumentos e a área derramada está entre as variáveis locais e os argumentos. (up,up): a pilha também cresce para cima mas o apontador para o registro de ativação é usado para variáveis locais e argumentos; a área derramada está entre as variáveis locais e o apontador da pilha.

As seguintes convenções são válidas para os quatro modelos apresentados: (a) argumentos de chamada utilizam o apontador da pilha; (b) argumentos de entrada são deslocamentos positivos a partir do apontador do registro de ativação; (c) registradores locais são deslocamentos negativos a partir do registro de ativação; (d) registradores globais são deslocamentos positivos a partir do apontador da área global.

6.1 Modelo de Pilha da Arquitetura SPARC

A SPARC não possui instruções explícitas ou registradores sobre pilhas, portanto o gerenciamento da pilha do usuário é resultado de algumas convenções de software. Por exemplo, (a) na criação de um processo, o sistema operacional reserva uma área para a pilha. Um registrador específico representado por "%sp" aponta para seu topo (posição mais baixa) e um outro registrador específico representado por "%fp" aponta para o registro de ativação corrente; (b) a pilha cresce a partir do endereço de mais alta ordem para o endereço de mais baixa ordem; (c) alguns registradores são preservados entre trocas de janelas; (d) a área reservada para tratar estouros deve iniciar na posição indicada por "%sp" mesmo quando a altura da pilha variar devido a alocação dinâmica de memória; (e) extensão da pilha é de 64 bytes. Os tratadores de estouro pressupõem um alinhamento em múltiplos de oito bytes. Decorre desta convenção que os parâmetros de uma função são encontrados a partir do endereço "fp + 64", as variáveis locais e automáticas presentes na pilha são endereçadas em relação ao "%fp". Temporários e parâmetros de saída são endereçados especificando um deslocamento positivo a partir do "%sp" corrente ("%sp + 64") [25]. No modelo de pilha da SPARC as seguintes convenções são válidas à chamada de uma função em relação a passagem de parâmetros:

(a) Parâmetros Inteiros: até seis parâmetros inteiros usam-se os registradores %o0...%o5. Segundo [22] seis parâmetros é um número mais do que adequado. O número médio de parâmetros, medidos estaticamente e dinamicamente não é maior que 2.1. Parâmetros excedentes, se houver, são passados na pilha em memória.

(b) Parâmetros Agregados: (estruturas, vetores e arranjos) faz cópia na pilha e passa o endereço nos registradores de saída, ou seja, quem chama faz uma cópia dos parâmetros em sua pilha e passa o endereço da cópia.

(c) Parâmetros Endereços: neste caso, o procedimento é semelhante ao caso de agregado, exceto pelo fato de não ser necessário fazer a cópia na pilha em memória. O objetivo deste tipo de passagem de parâmetro é obter o efeito da passagem por referência, logo o endereço passado deve ser do argumento e não de uma cópia sua.

Os registradores locais podem ser usados para parâmetros colocados na pilha do usuário por exceder o valor de seis parâmetros e para automáticos que não tem seu endereço acessado. Se o valor de retorno for um inteiro escalar, ele é colocado em um registrador específico da função chamada. Se o valor de retorno é um real escalar, ele é colocado em registradores de ponto flutuante. Se for agregado, quem chama reserva espaço na pilha de memória para que o valor de retorno seja copiado. O endereço desta área de memória é passado no *hidden*. Note que na

verdade a arquitetura SPARC opera com duas pilhas, uma na memória e outra nos registradores, porque tudo que não é do tipo inteiro é passado na pilha da memória.

A pilha do usuário para um procedimento ativo na versão *SPARC-based Sun-4 workstation* [22] é usada para: (a) colocar as variáveis que devem ser endereçadas, incluindo arranjos e estruturas; (b) alocar espaço de pilha dinamicamente; (c) conter os registradores de ponto flutuante salvos e alguns temporários do compilador; (d) colocar os parâmetros excedentes a seis; (e) conter 6 palavras nas quais o procedimento chamado pode armazenar os registradores com os argumentos; (f) conter apontadores para *buffers* onde procedimentos chamados podem retornar não escalonáveis; (g) salvar janelas que passaram para a memória devido a ocorrência de *overflow*.

Janela de registradores simula esta estrutura de pilha com duas diferenças: (a) a estrutura é implementada nos registradores do chip; (b) o tamanho das áreas para parâmetros de entrada, variáveis locais e parâmetros de saída são fixas. Em sistemas onde a pilha do usuário é alocada na memória, o compilador decide o tamanho e tipo de cada uma destas áreas durante a compilação. As vantagens desta abordagem são: (a) não é necessária referência à memória para se ter acesso aos dados da pilha; (b) como o dado é residente no arquivo de registradores, dois operandos podem ser lidos e um terceiro escrito a cada ciclo de máquina. Contudo, a utilização de janelas de registradores em sistemas multiprocessados deve ser analisada tanto do ponto de vista de um processo quanto do sistema. Em tais sistemas, como é o caso de sistema operacional UNIX, existem vários processos executando simultaneamente (pseudo-parallelismo). Para que cada processo tenha acesso aos recursos físicos do ambiente, por exemplo, CPU, o sistema operacional realiza a troca do contexto de um processo para o de um outro. Caso o processo que estava executando anteriormente, tenha utilizado todas as janelas de registradores, é necessário salvar e invalidar todas janelas, inclusive aquelas que tiveram que ser alocadas na memória para serem utilizadas por outros processos. Desta forma pode ocorrer um aumento no *overhead* de troca de contexto, o que caracteriza uma grande desvantagem do uso de janelas de registradores. Uma opção de utilização das janelas de registradores em tais sistemas seria na base de uma por processo. Neste caso, em um sistema de multiprocessamento poderia ter até oito processos sendo executados, simultaneamente, cada um deles com sua janela de registradores. Se o processador passasse de executar um processo e passasse a executar outro, não seria necessário recarregar todo o seu contexto, o sistema operacional simplesmente trocaria a janela de registradores para a daquele processo e continuaria executando. De certa forma o fato de conceder todo recurso para cada processo ou compartilhá-lo com outros processos é uma escolha do projetista do compilador. Achamos que janela de registradores auxilia do ponto de vista de chamada de função mas implica em uma transferência muito grande de trabalho para este projetista. O ganho é mais real em um sistema monoprocessoado, contudo, ele é questionável em sistemas de multiprocessos.

A pilha do sistema GGCO é uma pilha clássica, ela guarda parâmetros passados, variáveis automáticas e endereço de retorno. Informação suficiente para cobrir grande parte das arquiteturas superescalares. Basicamente, pelo fato das janelas de registradores simularem a estrutura de pilha da SPARC, a pilha do usuário do sistema GGCO deveria fazer uma série de mudanças para tentar acomodá-la. Por exemplo, entre outras informações, ela deveria possuir uma área para salvar os registradores de entrada e locais no caso de estourar a pilha de janelas. Outra informação relevante é o tipo do parâmetro passado, existem algumas regras que têm que ser observadas quanto ao local onde os mesmos são colocados. Se o valor de retorno é um inteiro, este deve ser colocado em um dos seis registradores para inteiro; se o valor não for um inteiro ou real, retorna-se o endereço, o que acarretará por sua vez em uma série de providências; no caso do real, o valor retornado é colocado nos registradores de ponto flutuante. Mas a diferença

principal não está aqui, e sim nos algoritmos que manipulam a pilha. No projeto de GGCO pensou-se inicialmente em utilizar a pilha da SPARC mas, chegou-se a conclusão de que se o fizesse, estaria projetando um sistema totalmente preso a sua arquitetura. Supondo uma implementação nestas bases, os algoritmos projetados de acordo com estas regras manipulariam janelas de registradores com oito entradas, oito saídas e oito locais. Contudo, para executar uma especificação, por exemplo, da arquitetura MC88000 que não possui a facilidade de janela, estes algoritmos não serviriam para nada. Portanto, um sistema com este modelo de pilha não é genérico, pelo contrário, é totalmente ligado a uma máquina específica, inclusive se isto traz um ganho, traz um ganho somente naquele processador porque em máquinas que não possuem esta característica, o ganho seria nulo.

Para implementar um sistema que mantenha estas duas estruturas paralelas, uma forma seria manter os quatro modelos descritos e projetar a pilha da SPARC de tal forma que quando fosse compilar para a SPARC, poderia ser utilizado recurso de pré-processamento para ter acesso ao arquivo onde estaria definido sua pilha. Contudo, esta abordagem resultaria em um sistema ineficiente. Optamos por projetar um novo sistema que seja um superconjunto deste. A parte do sistema GGCO implementado ainda não modela a pilha que engloba janela de registradores, contudo LDA possui as diretivas para tal. Para incorporar janela de registradores, talvez uma opção seja implementar o algoritmo de renomeação, porque na verdade, a troca de janelas salva e renomeia os registradores na ativação de um procedimento: os registradores declarados como saída passam a ser entradas do novo procedimento.

7 Conclusões

Neste trabalho apresentamos uma ferramenta de auxílio à construção de geradores de código otimizados cujas principais contribuições englobam o projeto de um sistema baseado na descrição da arquitetura de máquina que incorpora, de forma que consideramos satisfatória, uma linguagem para descrição de arquiteturas que contém restrições sobre o escalonamento. Esta linguagem está formalmente especificada. Sua definição formal pode ser vista como um sistema gerador de gerador de código a partir da descrição da arquitetura. Ela mostra o mapeamento da descrição de arquitetura de máquina para o seu respectivo gerador de código, cujo resultado final é um trecho de programa em linguagem C, que constitui a parte dependente de máquina do gerador de código para o processador descrito.

O sistema GGCO descrito neste trabalho ainda não está totalmente implementado, contudo a especificação da parte dependente de arquitetura encontra-se definida em [4]. Entendemos que estudos adicionais devem ser realizados, por exemplo, a continuação deste trabalho poderia compreender: (a) implementação e avaliação do algoritmo de alocação de registradores apresentado na Seção 3.2.3; (b) um estudo e implementação de um esquema para modelar múltiplas unidades funcionais idênticas; (c) implementação do algoritmo de Gross e Henessy para preencher os *delay slots* com instruções válidas durante o escalonamento; (d) implementação do mecanismo proposto para permitir a especificação de registradores de código de condição; (e) implementação da solução proposta para atender o esquema de prioridade existente em algumas arquiteturas para controlar o acesso ao barramento de escrita.

Referências

- [1] Benitez, Manuel E. and Davidson, Jack W., *A Portable Global Optimizer and Linker*, ACM Sigplan Notices, 23, 7, July/1988.
- [2] Bradlee, David G. et al., *The Marion System for Retargetable Instruction Scheduling*, ACM Sigplan Conference on Programming Language Design and Implementation, ACM Sigplan Notices 26(7), July/1991.
- [3] Bradlee, David G. et al., *Integrating Register Allocation and Instruction Scheduling for RISCs*, ASPLOS-IV Proceedings - Fourth International Conference on Architectural Support for Programming Languages and Operating Systems, Santa Clara, California, April/1991.
- [4] Bigonha, Mariza A. Silva, *Otimização de Código em Máquinas Superescalares*, Tese de Doutorado, DI-PUC/RJ, Abril/1994.
- [5] Bigonha, Mariza A. S. and Rangel, José Lucas M., *Linguagens para Descrição de Arquiteturas de Computadores*, VI SBAC - Simpósio Brasileiro de Arquitetura de Computadores e Processamento de Auto Desempenho, Caxambu, MG, 137-153, Agosto/1994.
- [6] Bigonha, Roberto S., *SCRIPT An Object Oriented Language for Denotational Semantics (User's Manual and Reference)*, DCC - UFMG, RT 03/94.
- [7] Chaitin, Gregory J., *Register Allocation and spilling via graph coloring*, ACM Sigplan Notices, 17, 6, ACM Sigplan Symposium on Compiler Construction, June/1982.
- [8] Costa, Paulo S. S., *Um Gerador Automático de Geradores de Código*, Tese de Mestrado, PUC-RJ, Fevereiro/1990.
- [9] Davidson, Jack W., *Simplifying Code Generation Through Peephole Optimization*, DCS - The University of Arizona, Tucson Arizona, December/1981.
- [10] Fisher, Joseph A. et al., *Parallel Processing: A smart compiler and a dump machine*, ACM Sigplan Notices, 19, 6, Proceedings of the ACM Sigplan, Symposium on Compiler Construction, June/1984.
- [11] Fisher, Joseph A., *Trace Scheduling: A Technique for Global Microcode Compactation*, IEEE Transactions on Computers, 30, 7, July/1981.
- [12] Fraser, Christopher W. and Hanson, David R., *A Code Generation Interface for ANSI C*, DCS, Princeton University, Research Report, CS-TR-270-90, Last Revised September 1992.
- [13] Gibbons, Phillip B. and Muchnick, Steven S., *Efficient Instruction Scheduling for a Pipelined Architecture*, Proceedings of the ACM Sigplan'86 - Symposium on Compiler Construction, ACM Sigplan Notices 21(7), July/1986.
- [14] Goodman, James R. and Wei-Chung-Hsu, *Code Scheduling and Register Allocation in Large Basic Blocks*, International Conference on Supercomputing - Conference ACM-PRESS Proceedings, St. Malo France, July/1988.
- [15] Hennessy, John and Gross, Thomas, *Code Generation and Reorganization in the Presence of Pipeline Constraints*, Conference Record of the 9th Annual ACM Symposium on Principles of Programming Languages, 128-133, Albuquerque New Mexico, January/1982.
- [16] Hennessy, John and Gross, Thomas, *Postpass Code Optimization of Pipeline Constrains*, ACM Transactions on Programming Languages, 1983, 5(3).

-
- [17] Henry, Robert R., *Code Generation by Table Lookup*, Computer Science Department, University of Washington, Technical Report, # 87-07-07, FR-35 Seattle, WA 98195 USA, July/1987.
- [18] Intel, Inc., *i860 64-bit Microprocessor Programmer's Reference Manual*, Intel, Inc., Santa Clara, California, 1989.
- [19] Kerns, Daniel R., *Balanced Scheduling: Instruction Scheduling when Memory Latency is Uncertain*, Proceedings of the ACM Sigplan'93 Conference on Programming Language Design and Implementation, Albuquerque NM, ACM Sigplan Notices - Vol. 28 number 6 June-1993.
- [20] Motorola, Inc., *MC88100 RISC Microprocessor User's Manual*, Prentice-Hall, Englewood Cliffs, second edition, New Jersey/1990.
- [21] Motorola, Inc., *MC68020 32-Bit Microprocessor User's Manual*, Prentice-Hall, Englewood Cliffs, New Jersey/1984.
- [22] Muchnick, Steven and others, *Optimizing Compilers for the SPARC Architecture An Overview*, IEEE Computer Society International Conference, 284-288, SPRING COMPCON, (33), 1988.
- [23] Pinter, Shlomit S., *Register Allocation with Instruction Scheduling: a New Approach*, Proceedings of the ACM Sigplan'93 Conference on Programming Language Design and Implementation, Albuquerque NM, June/1993, ACM Sigplan Notices - Vol. 28 number 6 June-1993.
- [24] Stallman, Richard M., *Using and Porting GNU C*, Free Software Foundation Incorporation, Cambridge Massachusetts, 1989
- [25] Souza, G. B., *Técnicas de Otimização de Código para Arquiteturas RISC*, Tese de Mestrado, DCC - UNICAMP Campinas, 1992.
- [26] Tiemann, Michael D., *The GNU Instruction Scheduler*, Stanford University, Class Report, CS 343, June/1989.
- [27] Warren Jr, H. S., *Instruction Scheduling for the IBM RISC System/6000 Processor*, IBM Journal of Research and Development, 34, 1, January/1990.