

スラック予測を用いた 省電力アーキテクチャ向け命令スケジューリング

福山 智久[†] 福田 匡則[†] 三輪 忍[†]
小西 将人^{††} 五島 正裕^{††} 中島 康彦[†]
森 眞一郎[†] 富田 眞治[†]

我々は、命令のスラック (slack) に基づくクリティカルITY予測を提案している。ある命令の実行を s サイクル遅らせてもプログラムの実行時間が増大しないとき、 s の最大値をその命令のスラックという。したがって、いわゆるクリティカルな命令のスラックは 0 サイクルである。前回の実行時のスラックを予測表に登録しておくことによって、それを今回の予測値とすることができる。本稿では、スラックの値が 1 以上である命令を低速 / 低消費電力の演算器で実行するという方法で演算器の省電力化を図った。評価したところ、4.5% の IPC の低下で 19% の EDP を削減できることが分かった。

Instruction Scheduling for Low-Power Architecture with Slack Prediction

TOMOHIISA FUKUYAMA,[†] MASANORI FUKUDA,[†] MIWA SHINOBU,[†]
MASAHITO KONISHI,^{††} MASAHIRO GOSHIMA,^{††} YASUHIKO NAKASHIMA,[†]
SHIN-ICHIRO MORI[†] and SHINJI TOMITA[†]

We proposed an instruction criticality prediction technique based on prediction of instruction slacks. When the execution time of a program doesn't become longer even if an instruction of the program is delayed by s cycles, the maximum of s is referred as the slack of the instruction. Thus the slack of a critical instruction is zero cycles. The slack value is stored to the prediction table to be a predicted value for the next time. This paper describes instruction scheduling with slack. Evaluation result shows EDP is reduced 19% at the cost of 4.5% IPC degradation.

1. はじめに

命令のクリティカルITY (criticality)、すなわち、命令がどれほどクリティカルかを知ることは、スーパースカラ・プロセッサの高性能化と省電力化の両方に効果がある。例えば、命令をスケジューリングするときには、よりクリティカルな命令を優先的に発行した方がよい。小林らは、クラスタ化された演算器を持つプロセッサ¹⁾ に対して、クリティカル・パスの情報をを用いたスケジューリング手法を提案している²⁾。クリティカルITYの情報は、省電力アーキテクチャにおいても有用である。例えば、クリティカルでない命令のみを低速 / 低消費電力の演算器で実行することで、性能を大きく低下させることなく

省電力化を図ることができる³⁾⁻⁷⁾。

さて、従来このような研究の多くは、プログラムのクリティカル・パスに基づいて行われてきた。しかしクリティカル・パスに基づく方法には、以下のような問題点がある：

- (1) 論理的 実行しているプロセッサの物理的な制約が反映されていない。
- (2) 二値的 最もクリティカルな命令を教えるのみで、それ以外の命令がどの程度クリティカルでないのが判定できない。
- (3) クリティカル・パスの判定が困難 実行中のプログラムのクリティカル・パスを判定することはそれほど容易ではない。

一方我々は、クリティカル・パスではなく、命令のスラック (slack)⁸⁾ によって、命令のクリティカルITYを測ることを提案した^{9),10)}。ある命令の実行を s サイクル遅らせてもプログラムの実行時間が増大しないような s の最大値をその命令のスラックという。したがって、クリティカルな命令のスラックは 0 サイクルである。

[†] 京都大学
Kyoto University

^{††} 東京大学
University of Tokyo

^{†††} 大阪工業大学
Osaka Institute of Technology

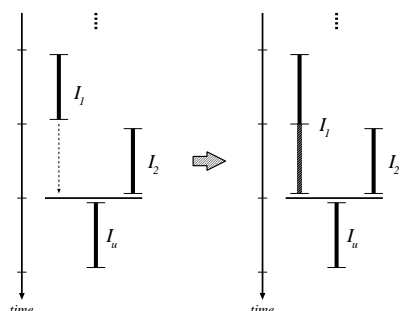


図1 I_1 を1サイクル遅らせた場合

スラック予測器を用いることで、前述した問題点は以下のように解決されると期待できる:

- (1) 実効的 履歴に基づいてスラックを予測するので、物理的、実効的なクリティカルティが反映される。
- (2) 多值的 クリティカルティの大 / 小は、スラックの小 / 大によって多值的に表現される。
- (3) スラックの判定は容易 クリティカル・パスとは異なり、スラックは容易に求めることができる。

スラックの計算

図1に、命令が実行される様子を表わすタイム・チャートを示す。図中、「T」が命令の実行を表し、「T」の長さはその命令の実行レイテンシを表す。上下の「T」の間にある横線は、フロー依存関係を表す。

同図では、命令 I_1 が定義した結果を最初に使用する命令は I_u となっている。すると I_1 のスラック s は、原則的には、定義命令 I_d による定義時刻 t_d と、使用命令 I_u による使用時刻 t_u の差、つまり:

$$s = t_u - t_d - 1 \quad (1)$$

によって得られる。この式に従えば、図1の I_1 のスラックは1サイクル、 I_2 のスラックは0サイクルとなる。なお、以下ではスラックの単位を省略し、「命令 I_1 のスラックは1」のように言うことにする。

本稿では、スラック予測器による予測結果を、実際に命令スケジューリングに応用し、演算器における消費電力を削減する方法について述べる。

スラックの利用

スラックが1以上の命令を低速 / 低消費電力の演算器で実行することによりプロセッサの性能を大きく低下させることなく省電力化を図る。

スラックが1以上の命令は、プログラムの実行時間を増大させることなく、その実行を遅らせることができる。例えば、図1左においてスラックが1と計算された命令 I_1 を、低速な演算器で実行することにより演算実行結果が使用可能になる時刻を1サイクル遅らせ

た様子を図1右に示す。

図1右では、 I_1 と I_2 が前回と同様のタイミングで実行された場合を示している。 I_1 の実行結果が得られる時刻が1サイクル遅れても、 I_u の実行開始は遅れておらず、プログラムの実行時間は増大しないことに注意されたい。

上述したように、本来スラック予測とはクリティカルティを多值的に表現しているものである。しかし、本稿の命令スケジューリングでは、スラック0であるか1以上であるかが重要となる。すなわち、ここで言うスラック予測とはスラック0予測と言えよう。

以下、2章でスラック予測器について、3章でスラック予測を用いて省電力アーキテクチャ向け命令スケジューリングを実現する際のスラックの計算方法について、4章でグローバル分岐履歴と2ビット・カウンタの導入について述べ、5章で本稿で述べた命令スケジューリングの評価結果を示す。

2. スラック予測器

本章では、スラック予測器の基本的な実装について述べる。以下、まず2.1節でスラック予測器のデータ構造についてまとめた後、2.2節で予測器に対する登録、参照といった操作について説明する。

2.1 スラック予測器の構成

スラック予測器は、主に以下の2種の表からなる:

- (1) スラック表 命令のスラックを記録する表。
- (2) 定義表 各データに対し、以下を記録する:
 - (a) 定義時刻 そのデータが定義された時刻
 - (b) 定義命令 そのデータを定義した命令
 - (c) フラグ 命令を遅らせて実行したかどうか

フラグについては、3.4節で詳しく述べる。スラック表は、命令の過去のスラックを記録する予測表本体であり、値予測におけるVHT (Value History Table) に相当する。一方、定義表は、スラック表に記録するスラック自体を計算するために用いられる。

定義表

定義表は、論理的には、レジスタ・ファイルやメモリ上の各データに対して、定義時刻、定義命令、フラグを記録するフィールドを付加したものと考えてよい。データが使用される時、データと同時に定義表に記録された定義時刻とフラグの値を読み出せば、定義命令のスラックを計算することができる。

ただし、システム中のすべてのデータに対して定義時刻、定義命令、フラグを記録することは非現実的である。予測精度とハードウェア・コストのバランスをとるためには、アクセス頻度の高いロケーションに対してエ

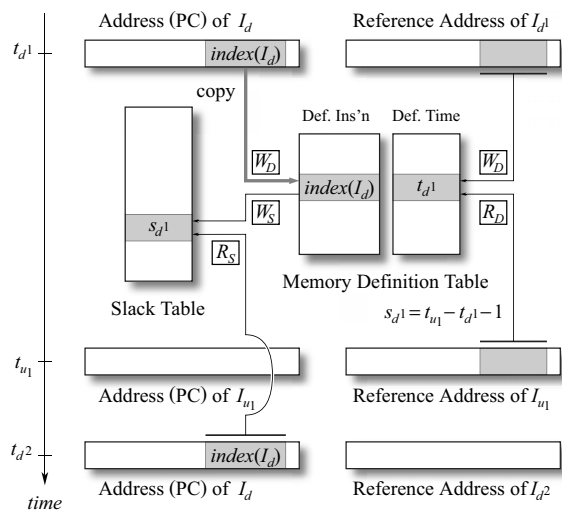


図2 予測器に対する操作(ストア命令の場合)

ントリを提供することが肝要である。

アクセス頻度を考慮して、定義表は、データの格納場所がレジスタかメモリかによって2つに分け、それぞれ以下のように実装する:

- (1) レジスタ定義表 物理レジスタ番号をインデックスとするRAMによって構成する。すなわち、そのエントリ数は物理レジスタ・ファイルに等しく、まさに物理レジスタ・ファイルに定義時刻、定義命令、フラグを格納するフィールドを付加したものと考えるよい。
- (2) メモリ定義表 ロード/ストア命令の参照アドレスをインデックスとするキャッシュとして構成する。

したがって、メモリ定義表ではミスが発生することになる。メモリ定義表のミスや、エントリ数と連想度については後で詳しく述べることにして、次節では、これらの表を用いたスラック予測器の動作について述べる。

2.2 スラック予測器の動作

以下では、命令 I_d の n 回目 ($n \in \mathbb{N}$) の実行を、肩付き数字を用いて、 I_d^n のように表すことにする。また、 I_d^n が定義したデータを最初に使用する命令の実行を I_{u_n} と表す。なお、 I_{d_i} と I_{d_j} ($i, j \in \mathbb{N}, i \neq j$) は同じ命令であるが、 I_{u_i} と I_{u_j} は一般に異なる。

図2では、ストア命令 I_{d^1} とロード命令 I_{u^1} が、それぞれ、時刻 t_{d^1} および t_{u^1} に実行されている。スラック表への登録、および、同スラック表への参照、すなわち、予測は、以下の様に行われる:

- (1) 登録 登録は、以下のように、(a) I_{d^1} がデータを定義するとき、(b) I_{u^1} がそのデータを使用する

ときの2つのフェーズからなる:

- (a) 定義 I_{d^1} がデータを定義するとき、以下の操作が行われる:

W_D 現時刻 t_{d^1} と I_{d^1} 自身(のアドレス)を定義表に書き込む。

- (b) 使用 I_{u^1} がそのデータを使用するとき、以下の2つの操作が連続して行われる:

R_D 定義表を読み出して、定義時刻 t_{d^1} と定義命令 I_{d^1} (のアドレス)を得る。

W_S 定義時刻 t_{d^1} と現時刻 t_{u^1} から、 I_{d^1} のスラック s_{d^1} が、 $s_{d^1} = t_{u^1} - t_{d^1} - 1$ と求まる。スラック表の定義命令 I_{d^1} のエントリに、求めた s_{d^1} を書き込む。

- (2) 予測 命令 I_d が再びフェッチされると、以下の操作が行われる:

R_S I_d のアドレスをインデックスとしてスラック表を直接読み出すことで、前回のスラック s_{d^1} が得られる。

ストア命令以外の場合には、基本的には、メモリ定義表をレジスタ定義表に、参照アドレスを物理レジスタ番号に、それぞれ読み替えればよい。

なお、スラックの計算を行うのは、そのデータが最初に使用されるときのみである。したがって、 R_D において、読み出された定義表のエントリを無効化し、以降同じエントリが繰り返し読み出されないようにする。このことは、メモリ定義表のエントリの有効利用にも効果がある。

2.3 スラック予測器へのアクセス・タイミング

図3と図4に、前述したストア命令の場合と、それ以外の命令の場合の、予測器へのアクセスのタイミングをそれぞれ示す。図中の $W_D \sim R_S$ は、前述した操作と対応している。スラック予測器へのアクセス・タイミングは、以下のとおりである:

- (1) 登録 スストア命令の場合には、定義側、使用側共に、アドレス計算(図3中、A)の次のステージからアクセスを開始できる。

一方、ストア以外の命令の場合には、表へのインデックスとして用いる物理レジスタ番号が既に分かっているため、1サイクル早く、発行(図4中、I)の次のステージからアクセスを開始できる。

- (2) 予測 分岐予測器、値予測器などと同様、命令のアドレスをもってアクセスすればよいので、命令フェッチと同時にアクセスを開始することができる; ただし、得られたスラックは専ら命令スケジューリングに使用されるため、ディスパッチ(図3, 4中、D)に間に合うように読み出せばよい。

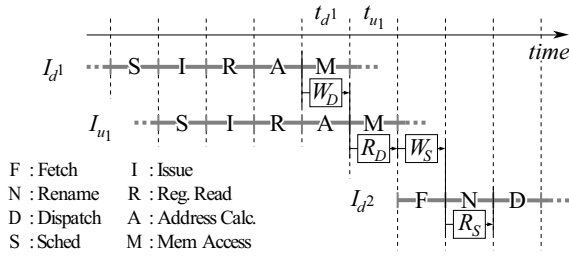


図3 予測器に対する操作のタイミング(ストア命令の場合)

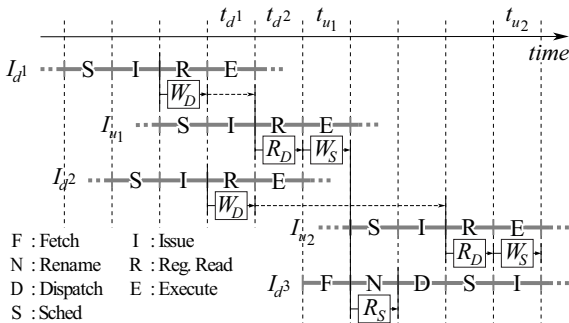


図4 予測器に対する操作のタイミング(ストア命令以外の場合)

2.4 条件分岐命令

分岐命令は、その実行結果を参照する命令が存在しないため、上述の方法でスラックを計算することはできない。そこで、以下に述べる理由により、分岐予測ミスを起こした分岐命令のスラックは0、ヒットした分岐命令のスラックは0または1とする。

分岐予測ミスを起こす分岐命令は、できるだけ早く実行した方がよい。分岐予測ミスを起こす分岐命令より下流にある命令の実行はすべて無駄である。その分岐命令を早く実行すると、この無駄が省かれるとともに、状態回復がそれだけ早く開始されるからである。

一方、分岐予測ヒットする分岐命令は、論理的にはクリティカルにはならない。しかし、以下に述べる理由から、できるだけ早く実行した方がよい；フェッチ後、未完了のまま (pending) にできる分岐命令の数には、分岐予測に関する資源の量に起因して、他の命令より強い制約がある。例えば、MIPS R10000 プロセッサ¹¹⁾ の場合、たかだか4個の条件分岐命令しか未完了にできない。資源が不足した場合には、フロントエンドがストールすることになる。

このように分岐命令は、予測ヒット/ミスに関わらず他の命令より早く実行した方がよく、その中ではミスする命令の方がよりクリティカルである。例えば、分岐命令の実行ユニットが1つしか空いていない場合には、ヒットするものよりミスするものを先に実行した方がよい。したがって、このことを考慮したい場合には、ミ

スすると予測される命令のスラックは0、ヒットすると予測される命令のスラックは1とするとよいと考えられる。

2.5 スラック表、メモリ定義表のミス

メモリ定義表、スラック表は、キャッシュであり、ミスが起こる。メモリ定義表、スラック表共に、 W_D 、 W_S では、ミスが起こっても割り当てられたエントリに上書きするだけであるから、その時点でのミスは性能上影響がない。したがって、 W_D 、 W_S で割り当てられたエントリが、 R_D 、 R_S までにリプレースされてしまった場合のことを考えればよい。

R_D 、 R_S におけるミスは、以下のようにする：

メモリ定義表 定義表のエントリがリプレースされたのだから、定義側の命令を特定することができず、スラック表への登録を行うことはできない。

スラック表 スラック表ミス時には、スラックは0と予測するのが安全である。また、5章で述べるように、スラックが0である命令は全体の半分以上に上り、0としておいても相応の予測ヒット率が期待できるからである。

3. スラックの計算方法

本章では、スラック予測を命令スケジューリングに応用する際に用いる3種類のスラックの計算方法を提案する。以下、3.1節で基本的な計算方法、3.2節で見かけのスラック、3.3節で見かけ上クリティカルになっている命令について述べ、最後に、3.4節で3つの計算がスラック予測器の定義表にフラグを用意することによって可能であることを示した。

3.1 スラックの基本計算

1章で述べたように、スラック s の基本となる計算式は、定義命令 I_d がデータを定義した時刻 t_d とそのデータを最初に参照する命令 I_u がデータを使用する使用時刻 t_u の差、つまり式 (1)：

$$s = t_u - t_d - 1 \quad (1)$$

によって与えられる。

ただし、図5左下において、1サイクル遅らせた命令 I_1 については、そのスラックを再計算する際注意が必要である。式 (1) にしたがって計算すると、遅らせた命令 I_1 のスラックは0となる。すると次回 I_1 をスケジューリングする際には、 I_1 はクリティカルであると判断され、実行時には1サイクル遅らせることなく実行されてしまう。その結果、次回には再び図5左上のように実行され、以降、上下の状態を繰り返すことになる。

しかし、この繰り返しが発生するのは、遅らせた命

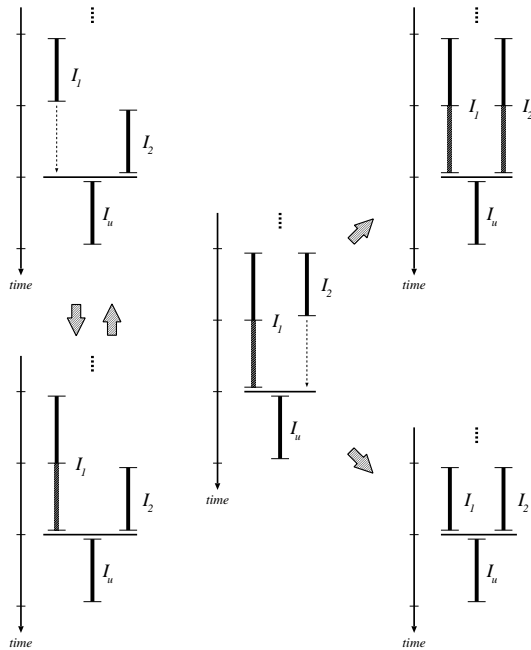


図5 状態の遷移

令の Slack が式 (1) によって 0 と計算される時のみ、つまり、遅らせなかった場合の Slack が 1 の命令のみである。よって、遅らせずに実行した時に Slack が常に 1 となる命令は、毎回遅らせることができるにもかかわらず、2 回に 1 回しか遅らせることができなくなってしまう。

3.2 見かけの Slack

遅らせた命令 I_1 の 0 という Slack は、故意に 1 サイクル遅らせたことによって生ずる「見かけの」の Slack であるといえる。 I_1 の真の Slack は、図 5 左上の場合と同様、遅らせなかった場合の定義時刻を Slack を元に計算すればよい。この遅らせなかった場合の定義時刻を実効定義時刻と呼ぶことにする。つまり、遅らせなかった場合の定義時刻 t_e は：

$$t_e = t_u - 1 \quad (2)$$

によって与えられる。この場合の Slack は、使用時刻と実効定義時刻の差、つまり：

$$\begin{aligned} s &= t_u - t_e - 1 \\ &= t_u - t_d \end{aligned} \quad (3)$$

によって得られる。

このようにすれば、発振が生じ 2 つの状態を繰り返すことはなくなる。

3.3 見かけの Slack (その 2)

前節までの説明では、命令 I_1 と I_2 は、毎回同じタイミングで実行が開始される場合を考えていた。しかし、必ず前回と同じタイミングで実行開始されるとは

限らない。図 5 中に示すように、 I_1 と I_2 が同時に実行開始されることもあり得る。さらには、 I_2 の方が早く実行される可能性もある。このような場合、前回 Slack が 1 以上であったため 1 サイクル遅らせた命令 I_1 が、今回クリティカルになってしまう。このような命令を見かけ上クリティカルになっている命令と呼ぶことにする。

図 5 中のような場合、 I_1 の Slack は実効定義時刻に基づいて、 I_2 の Slack は基本どおりに、それぞれ 1 と計算され、次回実行時には図 5 右上のようになる。しかし実際には、図 5 右下のように I_1 、 I_2 の Slack を共に 0 とし実行した方が I_u を 1 サイクル早く実行することができる。

Slack を共に 1 と計算した図 5 右上の場合は、すべての命令を無条件に 1 サイクル遅らせたのと変わりが無い。また、一旦この状態になると Slack が 1 のまま 0 に戻らないという可能性もある。

結局、図 5 中の I_1 、 I_2 のように見かけ上クリティカルになっている命令がある場合には、次回は図 5 右下のようにするために、 I_1 、 I_2 の Slack を共に 0 とする。

図 5 左下のように、遅らせていない命令（図では I_2 ）がクリティカルになっている状況では、 I_1 は見かけ上クリティカルになっている命令とは言わない。逆に言えば、図 5 中のように、クリティカルになっている命令のすべて（図では、 I_1 と I_2 ）が遅らせた命令である場合に、これらの命令を見かけ上クリティカルになっている命令と言う。

見かけ上クリティカルになっている命令がある場合、命令を遅らせたサイクル数（図では 1 サイクル）だけ無駄に実行が遅れている。したがって各命令の実効 Slack は、式 (1) (3) により得られた Slack の値から遅らせたサイクル数を減ずればよい。

つまり、見かけ上クリティカルになっている命令が生じている場合、遅らせた命令の Slack は式 (3) から 1 減じた式、つまり式 (1) と同じ式により得られ、遅らせなかった命令は式 (1) から 1 減じた式：

$$s = t_u - t_d - 2 \quad (4)$$

によって得られる。

3.4 問題の発見

2 つの見かけの Slack の問題を発見するには、その命令が遅らせたものかどうかを表すフラグが必要になる。そこで、Slack 予測器の定義表にフラグを用意し、命令を遅らせていない場合には 0 を、遅らせた場合には 1 を命令の実行時に記憶する。

まず、式 (1) により暫定的な Slack を求める。この暫定的な Slack の値が 0 となる命令のフラグが全

て1であるとき、見かけ上クリティカルになっている命令がある状況が生じていることとなる。

このように見かけ上クリティカルになっている命令がある状況はフラグを容易することで発見でき、この場合の実効スラックはフラグが0である命令に関しては暫定的なスラックから1減じた値、つまり式(4)で、フラグが1である命令に関しては暫定的なスラックの値、つまり式(1)で得ることができる。

見かけ上クリティカルになっている命令がない場合には、フラグが0の命令は暫定的なスラックの値、つまり式(1)で得られる値が、フラグ1の命令は暫定的なスラックに1足した、つまり式(3)で得られる値がそれぞれ実効スラックとなる。

このようにして得られた実効スラックの値をスラック表に記憶する。

4. 予測精度の向上

本章では、スラック予測の精度の向上を試みるために、スラック予測器にグローバル分岐履歴と飽和型2ビット・カウンタを導入する方法について述べる。

4.1 グローバル分岐履歴の導入

スラック予測器にグローバル分岐履歴を導入することにより、予測精度が向上することが認められている¹⁰⁾。グローバル分岐履歴の導入にはスラック表のインデクスとして、gshare分岐予測器などと同様に命令のアドレスとグローバル分岐履歴の排他的論理和を用いた。ただし以下で述べるように、スラック表はタグを持っていることに注意する必要がある。

分岐予測で用いられるPHT (Pattern History Table) には、競合が発生する; すなわち、通常PHTはタグを保持しておらず、異なる2つの命令が同一のPHTエントリを使用することを許している。

一方、値予測で用いられるVHT (Value History Table) では、タグ比較を行い、競合を許さないことが普通である。

スラック表も、VHTと同様、現在ではタグを保持している。グローバル分岐履歴を導入するにあたっては、インデクスに命令のアドレスとグローバル分岐履歴との排他的論理和を用いるとともに、タグにもグローバル分岐履歴を加え、競合を完全に排除している。そのため、破壊的競合ではなく、ヒット率の低下によって予測精度が低下するおそれがある。

4.2 飽和型2ビット・カウンタの導入

これまで述べてきたスラック予測の方法では、スラッ

クの最新の値 (last value) を記憶している。一方、1章でも述べたようにスラックが1以上の命令を1サイクル遅らせる場合、スラックの厳密な値は必要なく、スラックが0かどうかを記憶しておけばよい。

結局、スラックが0かどうかを予測する上では、これまで述べてきた方法は、いわば1ビット・カウンタによる予測であると言える。

そこで、飽和型2ビット・カウンタを用いることにより、命令のローカルな履歴を反映させることが考えられる。スラックの計算による値が0であればカウンタを1減じ、1以上であれば1足す。そしてスラック予測を行う際に、スラック表から読み出した値が0, 1であればスラックは0と予測し、2, 3であればスラックは1以上と予測する。

5. 評価

シミュレーションにより、3章で提案した3種類のスラックの計算方法を実際に命令スケジューリングに用いた場合について評価する。

なお、シミュレーションではスラック予測器で用いたグローバル分岐履歴の履歴長を全て2ビットとした。スラック予測器にグローバル分岐履歴を用いることで予測精度が向上することが認められている¹⁰⁾が、その度合いはわずかで今回の評価にはほとんど影響が現れなかったためである。

5.1 評価方法

SimpleScalar ツールセット (ver. 3.0) の sim-outorder シミュレータに対して、スラック予測器を実装し、SPECベンチマークを用いて本稿で述べた命令スケジューリングによる省電力アーキテクチャの効果を測定した。測定には表1に示すSPEC CINT95の8つのプログラムを実行した。

コンパイラは、gcc (ver. 2.7.2.3) を用いた。最適化オプションは、-O6 -funroll-loops である。

プロセッサのモデル

評価では、本稿で述べた命令スケジューリング方式を整数演算器に適用する。用いたプロセッサの整数系の

表1 SPEC CINT95 ベンチマーク・プログラム

プログラム	入力セット	実行命令数
099.go	99	132M
124.m88ksim	dcrand.big	120M
126.gcc	genrecog.i	32M
129.compress	10000 q 2131	35M
130.li	train.lsp	183M
132.jpeg	vigo.ppm-GO	26M
134.perl	primes.in	10M
147.vortex	persons.250	157M

タグは、一部省略できる可能性がある¹²⁾。

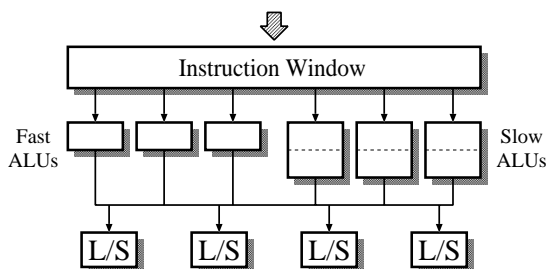


図 6 プロセッサのモデル

ブロック図を図 6 に示す。

同図に示すように、プロセッサは高速な ALU と低速な ALU をそれぞれ 3 個ずつ、合計 6 個持つ。高速な ALU と低速な ALU のレイテンシは、それぞれ、1 サイクル、および、2 サイクルである。低速な ALU は、高速な ALU をパイプライン動作させるものに相当し、レイテンシは 2 に増加するが、スループットは 1 のままである。これらの ALU の電源電圧は、それぞれ、1.1 V、および、0.7 V とした¹³⁾。スラックが 0 と予測された命令を高速な ALU、1 以上と予測された命令を低速な ALU にスケジューリングする。

実際にキャッシュへのアクセスを行うロード/ストア・ユニットは、4 個あり、同図のように接続されている。すなわち、ロード/ストア命令のアドレス計算は、これらの ALU のいずれかで実行される。アドレス計算を実行した ALU とロード/ストア・ユニットの間は、完全結合とした。すなわち、アドレス計算を行った ALU とロード/ストア・ユニットの組み合わせには制限はない。

また、評価の対象としたモデルは、分離 (separate) ロード/ストア方式を採用している。すなわちロード/ストア命令は、ディスパッチ時に、アドレス計算を行う命令と、実際にメモリ(キャッシュ)アクセスを行う命令に分離され、個別にスケジューリングされる。評価では、分離されたそれぞれに対してスラックを予測し、スケジューリングを行っている。したがって、最初からアドレス計算命令とメモリ・アクセス命令の 2 つの命令があったものと考えてよい。

プロセッサのそれ以外のパラメータを表 2 に示す。

テーブルのパラメータ

表 3 に、テーブルのパラメータをまとめる。スラック表、および、メモリ定義表の容量は、それぞれ、1 次命令、および、1 次データ・キャッシュと同じ範囲をカバーできるようにした；すなわち、それぞれ 8K エントリである。ただし連想度は、1 次命令、および、1 次データ・キャッシュがそれぞれ 2 であるのに対して 4 とした。このように、やや大きな容量 / 連想度としたのは、ミ

スによる影響を評価結果から除外するためである。メモリのレイテンシ(18 サイクル)は、メモリ・インタフェースを集積する AMD Athlon プロセッサのものを参考にした。

スラックの計算方法

それぞれの命令をどちらの ALU で実行するかを決定する際に用いるスラックの計算方法は 3 章で提案した次の 3 つとする。

BASE 3.1 節で述べた基本となる計算方法。発振が生じ遅らせることができる命令にもかかわらず遅らせずに実行されてしまう可能性がある。

EDT 3.2 節で述べた実効定義時刻 (EDT: effective definition time) を用いる計算方法。BASE とは逆に命令を遅らせ過ぎる状況が起こる。

ACC 3.3 節で述べた見かけ上クリティカルになっている命令な命令を発見し遅らせ過ぎる状況をなくすことにより、より正確な (ACC: accurate) スラックを求める計算方法。

上記の 3 つのそれぞれの場合について、前回の履歴にのみ依存している場合と飽和型 2 ビット・カウンタを用いた場合の 2 通りをシミュレーションし、それぞれ -1b、-2b として区別する。

それぞれの場合について、性能、ALU における消費電力、低速の ALU で実行した命令の割合を測定し、

表 2 プロセッサの各パラメータ

パラメータ	サイズ
FETCH 幅	8
命令ウインドウ	16
ISSUE 幅	8
COMMIT 幅	8
Functional Units	Int:6, FP:4, MemPort:4
レイテンシ (Latency/Throughput)	iALU 1/1, iMULT 3/1, iDIV 20/19, fADD 2/1, fCMP 2/1, fCVT 2/1, fMULT 4/1, fDIV 12/12, fSQRT 24/24
分岐予測	4K-entry 8-histry-length gshare predictor 6cycle extra misprediction latency 8-entry return address stack

表 3 各表、キャッシュ、メモリのパラメータ

	容量	ライン サイズ	連想度	レイテンシ (cycles)
スラック表	8K命令	—	4	1
レジスタ定義表	64命令	—	64	1
メモリ定義表	8K命令	—	4	1
1 次 命令	8K命令*	8命令*	2	1
1 次 データ	8Kワード	8ワード	2	1
2 次	1MB	64B	2	6
メモリ	—	—	—	18 [†]

*: SimpleScalar ツールセットでは 8B/命令。

†: 最初のワード。後続ワードには 2 サイクル / ワードが必要。

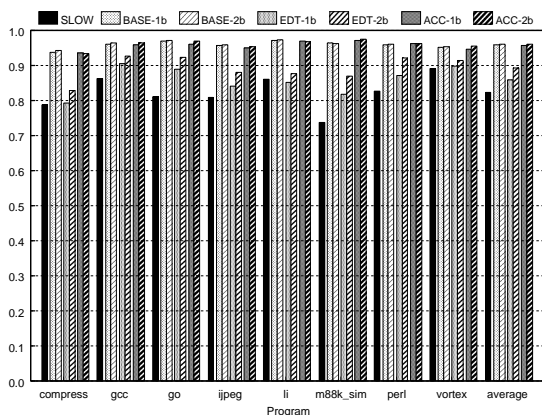


図7 性能の比

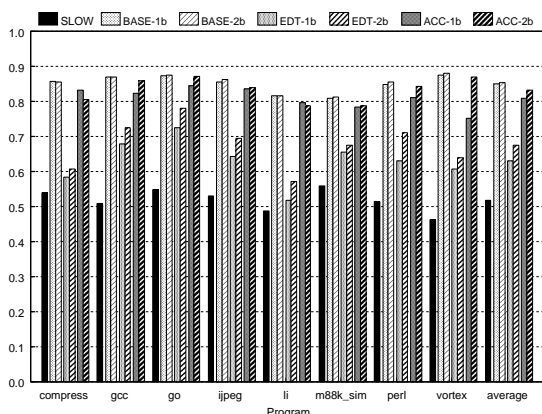


図8 EDPの比

全ての ALU を高速にした場合 (FAST とする) や低速にした場合 (SLOW とする) と比較することによって評価を行う。

なお、ここで言う消費電力とは、ALU のみにおける消費電力でありプロセッサ全体の消費電力ではない。よって、スラック表や各定義表における消費電力は考慮していない。

5.2 評価結果

図7, 図8, 図9に結果を示す。図7は性能を, 図8はEDPを, 図9は低速なALUで実行した命令の割合を示している。図7, 図8には, 9組のバーがあり, 左の8組はSPEC CINT95の8つのプログラムに対応しており, 最右1組はそれらの平均である。図7は, FASTに対するIPCの比を, 図8はEDPの比を表わしている。各組には, バーが7本ずつあり, それぞれ左から前節で述べたSLOW, BASE-1b, BASE-2b, EDT-1b, EDT-2b, ACC-1b, ACC-2bの場合のものを指している。

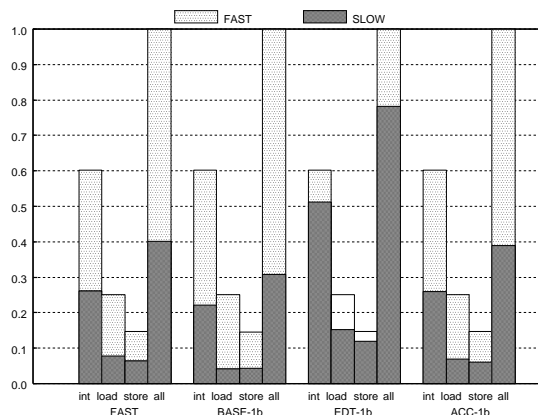


図9 低速なALUで実行した命令の割合

一方, 図9には4組のバーがあり, 左からFAST, BASE-1b, EDT-1b, ACC-1bの場合の8つのプログラムの平均である。それぞれの組の1番左の薄いバーは整数演算器で実行した全命令に対する通常の整数演算命令の比で, 2番目のバーはロードの, 3番目のバーはストアのアドレス計算の比で4番目のバーはそれらの合計である。1番左のFASTの組の濃いバーはそれぞれの命令における全命令に対するスラックが1以上の命令の比を, 薄いバーはスラック0の命令の比を表わしている。また, その他の組の濃いバーは低速なALUで実行した命令の比を, 薄いバーは高速なALUで実行した命令の比を表わしている。

図7, 図8を見ると, SLOWでのIPC低下は18%, EDP削減は48%であるのに対し,

- BASE-1bでは, FASTに比べIPC, EDPがそれぞれ約4.5%, 約14.5%低下した。
- EDT-1bでは, FASTに比べIPCが約14%低下しているものの, 約37%ものEDPを削減できている。
- ACC-1bは, IPCではBASEとほぼ同じ約4.5%であったが, EDPでは約19%とBASEよりも4.5%多く削減できている。
- -2bでは, それぞれの場合の飽和型2ビット・カウンタを用いなかった場合よりもIPC, EDPの比が共に若干減少した程度であった。

一方, 図9を見ると, FASTでは, 整数演算命令の44%, ロード命令のアドレス計算部分の31%, ストア命令のアドレス計算部分の42%がスラック1以上である。3つを合わせると整数演算器を用いる全命令の40%がスラック1以上であることがわかる。

これに対して,

- BASE-1bでは低速なALUで実行した命令はFAST

においてスラック1以上である命令よりも少ない。これは、発振が生じ遅らせることができる命令を遅らせられなかったためである。

- 逆に、EDT-1bは低速なALUで実行した命令が大幅に多い。これは、3.3節で述べた見かけ上クリティカルになっている命令が存在し、命令を遅らせ過ぎているためである。
- ACC-1bはFASTとほぼ一致しており、より正確にスラックの計算がされていると考えられる。

ES-1bは厳密にスラックを計算できているにもかかわらず、全てのALUを高速な演算器にしたときに比べ4.5%の性能が低下している。この原因として次の2つが挙げられる：

- 本当はクリティカルである命令を低速なALUで実行した、つまりスラック予測の予測ミスによる性能低下。
- 高速なALUを6つから3つに減らしたためにクリティカルな命令を1サイクルにつき最大3命令までしか実行できない。それによって、使用するALUに偏りが生じた場合にはALUが不足してしまい、性能低下の原因となる。

そこで、ACC-1bにおいてスラックが1以上の命令を実行する低速なALUのレイテンシを高速なALUと同様1サイクルとしてシミュレーションを行った。こうすることにより、スラックの予測ミスが原因となる性能低下をなくし、使用ALUの偏りが原因となる性能低下のみを図ることができると考えられるからである。このシミュレーションの結果、FASTに対する性能は平均して約98.5%であった。よって、ACC-1bの4.5%の性能低下のうち、約1.5%が使用ALUの偏りによるもので、スラックの予測ミスによる性能低下は約3%程度であると推測できる。

5.3 関連研究

千代延らは文献7)において、小林らが提案しているデータフローグラフの最長パスを特定するためのパス情報テーブル²⁾(PIT: Path info Table)をクリティカルパス予測器として、省電力アーキテクチャの設計に用いている。クリティカルパス予測器によりクリティカルパス上にないと予測された命令を低速/低消費電力のALUで実行し性能と消費電力の計測を行い、その結果を彼ら自身が提案しているクリティカルパス予測器と比較することにより性能と省電力化の評価を行っている。PITを用いた省電力アーキテクチャでは約5%の性能低下で20%のEDPを削減できている。

使用したベンチマーク・プログラムやプロセッサの

モデルが若干異なるものの本稿でのスラック予測を用いた省電力アーキテクチャでは約4.5%の性能低下で19%のEDPを削減できている。PITと同等の性能と省電力化を実現できている。

また、PITは複雑なハードウェアを必要とし、実際に実装することはとても困難である。一方、これまで述べた通りスラックは容易に求めることができることにも注意が必要である。

命令のクリティカル性を判定する際に、クリティカルパスに基づいたPITを用いた場合と、我々が提案する履歴に基づいたスラック予測を用いた場合とで、同様の結果が得られていることも興味深いところである。

6. おわりに

本稿ではスラック予測を用いた省電力アーキテクチャ向け命令スケジューリングについて述べた。評価結果は、スラック予測が省電力アーキテクチャに十分応用できることを示している。

今後は、より厳密なスラックを求めるために以下のことを検討する予定である。

スラックの伝搬

本稿で用いたスラックは文献9)にもあるように近似的なものであり、「命令の実行を s サイクル遅らせてもプログラムの実行時間が増大しないような s の最大値」というスラックの厳密な定義に従っていない。たとえば、図10左において、本稿で用いたスラックの求め方に基づけば、同図の命令 I_1, I_2 のスラックはそれぞれ0, 2となり、本稿の省電力アーキテクチャを適用すると図10中のように I_2 のみが低速なALUで実行され、 I_1 は高速なALUで実行されることとなる。しかし、厳密には I_1, I_2 のスラックは共に2であり、図10右のように両方の命令が共に低速なALUで実行された方がよい。

そこで、我々は I_2 のスラックを I_1 に伝搬させるなどしてより厳密なスラックを求めることを検討している。

謝辞

本研究の一部は、日本学術振興会 科学研究費補助金 基盤研究S(課題番号16100001)、および21世紀COEプログラム(課題番号14213201)による。

参考文献

- 1) Keller, J.: The 21264: A Superscalar Alpha Processor with Out-of-Order Execution, *9th Annual Microprocessor Forum* (1996).
- 2) 小林良太郎, 安藤秀樹, 島田俊夫: データフローグラフの最長パスに着目したクラスタ化スーパー

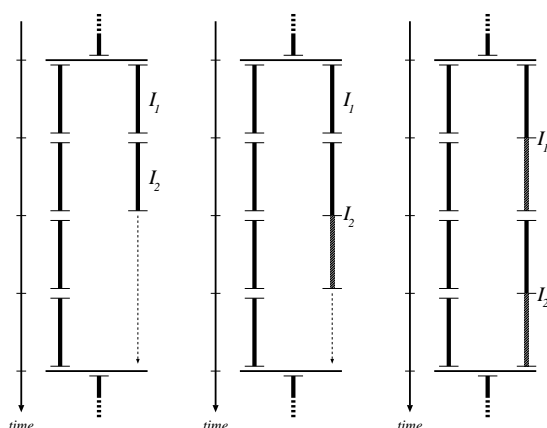


図 10 スラックの伝搬

データ値予測機構のハードウェア量削減, 信学技報 CPSY 2000-3 (2000).

- 13) Levy, M.: SAMSUNG Twists ARM Past 1GHz, *Information Quarterly*, No. 1 (2002).

スカラ・プロセッサにおける命令発行機構, 並列処理シンポジウム JSPP 2001, pp. 31-38 (2001).

- 3) Fields, B. and Blodik, S. R. R.: Focusing Processor Policies via Critical-Path Prediction, *28th Int'l Symp. on Computer Architecture (ISCA-28)*, pp. - (2001).
- 4) Fields, B., Bodik, R. and Hill, M. D.: Slack: Maximizing Performance under Technological Constraints, *29th. Int'l Symp. on Computer Architecture (ISCA-29)* (2002).
- 5) Tune, E., Liang, D., Tullsen, D. M. and Calder, B.: Dynamic Prediction of Critical Path Instructions, *7th Int'l Symp. on High Performance Computer Architecture (HPCA7)* (2001).
- 6) Grunwald, D.: Micro-architecture Design and Control Speculation for Energy Reduction, *Power Aware Computing*, Kluwer, ISBN 0-306-46786-0, chapter 4 (2002).
- 7) 千代延昭宏, 佐藤寿倫: プログラム実行時における命令の重要度決定に関する検討, 情報処理学会研究報告 2003-ARC-154 (SWoPP 2003), pp. 1-6 (2003).
- 8) Casmira, J. and Grunwald, D.: Dynamic Instruction Scheduling Slack, *Kool Chips Workshop (in conjunction with MICRO-33)* (2000).
- 9) 劉小路, 小西将人, 五島正裕, 中島康彦, 森眞一郎, 富田眞治: クリティカリティ予測のためのスラック予測, 先進的計算基盤システムシンポジウム SAC-SIS 2004, pp. 187-196 (2004).
- 10) 福田匡則, 小西将人, 五島正裕, 中島康彦, 森眞一郎, 富田眞治: グローバル分岐履歴を用いたスラック予測器, 情報処理学会研究報告 2004-ARC-159 (SWoPP 2004), pp. 25-30 (2004).
- 11) Yeager, K. C.: The MIPS R10000 Superscalar Microprocessor, *IEEE Micro*, Vol. 16, No. 2, pp. 28-40 (1996).
- 12) 佐藤寿倫, 有田五次郎: タグビット幅を考慮した