

明示化キャッシュメモリの性能評価

佐藤 正樹 有田 隆也 曽和 将容

名古屋工業大学 電気情報工学科

計算機の最高性能を引き出すためには、キャッシュ・メモリは不可欠なものである。しかし、従来のキャッシュ・システムではキャッシュ・ミスがおこり、計算機の性能向上に対する大きな障害となっている。

本研究は、キャッシュ操作をプログラムの中に埋め込み明示化することによってキャッシュ・ミスを無くそうとするものである。本システムでは起こりうる動作は、コンパイラによって可能な限り静的に解析されており、実行時にはすでにキャッシュ操作に関わるスケジュールがなされてしまっている。それによってキャッシュ・ミスを可能な限り少なくすることができる可能性を持つ。本稿では、明示化キャッシュの概要とソフトウェアシミュレータによる性能評価の結果について述べる。

Performance Evaluation of a Cache under the Control of Program

Masaki SATO Takaya ARITA Masahiro SOWA

Department of Electrical Engineering and Computer Science
Nagoya Institute of Technology

Gokiso, Nagoya 466, Japan

Caches are becoming a more important concern in maximizing overall machine performance. Cache misses are inevitable in traditional cache systems, and they cause a loss of computer system performance.

This study focuses on reducing cache misses using the cache management instructions. The compiler analyzes the behavior of the programs and schedules the cache management instructions at appropriate places to reduce the cache misses. This paper presents the outline of the cache under the control of program and the results of performance evaluation by software simulation.

1. はじめに

近年、プロセッサはクロック周波数の上昇や命令実行時間の短縮などにより性能向上が図られている。プロセッサの最高性能を引き出すためには、プロセッサへの命令やデータの十分な供給が不可欠の問題である。そのためにキャッシュ・メモリが用いられており、その役割は重要性を増してきている。

従来のキャッシュ・メモリ方式では、キャッシュ・ミスを避けることは困難である。それは、実行される各プログラムの性質を考慮せずに、固定的なハードウェアにより動的にキャッシュ操作が行われていることに起因する。また、従来のキャッシュ・メモリ方式では、キャッシュ・ミスが発生した後に、主メモリからキャッシュ・メモリへのデータ転送が開始される。つまり、キャッシュ操作は本質的に要求駆動処理であるといふことができる。このデータ転送期間中プロセッサは待ち状態に入るとともに、パイプライン処理の流れも乱れるため、それだけプログラムの実行が遅くなる。したがって、キャッシュ・ミスをなくすことが計算機のパフォーマンスを向上させることで重要なことである。

本研究はキャッシュ・ミスを、キャッシュ操作をプログラムの中に埋め込み明示化することによって無くそうとするものである。つまり、従来固定的なアルゴリズムでハードウェアによって行っていたキャッシュ操作を、実行される各プログラムに応じてソフトウェアによって行うのである。この方法をキャッシュ操作の明示化と呼び、そのためのプログラムをキャッシュ操作プログラムと呼んでいる。キャッシュ操作プログラムは、高級言語で書かれたプログラムをコンパイル時に静的に解析することによって作られる。作られたキャッシュ操作プログラムは、演算や転送などの通常の処理と並列に実行される。したがって、従来のキャッシュ・メモリ方式のような要求駆動処理ではないので、キャッシュ・ミスの発生しないようにキャッシュ操作を通常の処理より先行させることも可能である。これによって、明示化キャッシュは、キャッシュ・ミスを極限まで少なくできる可能性を持つ。

本論文では、第2章で明示化データ・キャッシュの原理、ハードウェア構成、特徴について述べる。また、第3章では、ソフトウェア・シミュレータによる性能評価の結果について述べる。

2. 明示化データ・キャッシュ

2.1 原理

従来のキャッシュ・システムでは、プロセッサによる命令やデータの参照が起こった後で、キャッシュにヒットしたかどうかの判定が行われている。それに対し、明示化キャッシュ・システムでは通常の処理を行う命令流とキャッシュ操作を行う命令流を持っている。これらはそれぞれ処理命令流、キャッシュ操作命令流と呼ばれており、またこれらのプログラム

は処理プログラム（P P R）、キャッシュ操作プログラム（C H P R）と呼ばれている。

図1に明示化キャッシュ・システムで実行されるプログラムのコントロール・フロー・グラフによる表現例を示す。図1において、左側の流れが通常の処理を行う命令（P P R）であり、右側がキャッシュ操作を行う命令（C H P R）である。左側の丸印はプログラムのブロック（1命令以上からなっている）を表し、右側の丸印は1キャッシュ操作命令を表している。この図で、矢印（アーク）はコントロールの流れ、すなわち実行順序に関する從属関係を表し、複数本のアークが入力されているプログラム・ブロックやキャッシュ操作命令は、その上の命令の実行後に実行可能となることを表す。このことは、コントロール・トークンと呼ばれる制御信号がアーク上を流れることによってなされる。

図1のプログラム

は次のように実行される。まず最初に、
1のキャッシュ操作命令がキャッシュ・メモリにデータを格納する。格納後、実行の終了を示すトークンをアーク上に送出することによって、

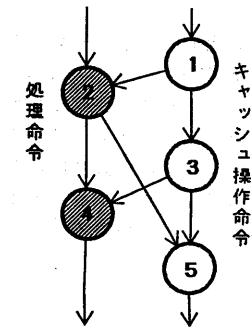


図1 明示化プログラム
のグラフによる表現

2のプログラム・ブ

ロックと3のキャッ

シュ操作命令が実行

可能となり、これら

が並列に実行される。2のプログラム・ブロック中の命令すべてと3のキャッシュ操作命令が終了すると、トークンがそれぞれのアーク上に送出され、4のプログラム・ブロックと5のキャッシュ操作命令の実行が可能となる。以上のように処理プログラムとキャッシュ操作プログラムは、お互いに同期をとりながら並列に実行される。

2.2 ハードウェア構成

データ・キャッシュ操作の明示化を行った計算機（明示化キャッシュ計算機）の構造を図2に示す。P Uはプロセッサ、D C H Mはデータ・キャッシュ・メモリ、Mは主メモリ、C H P Uはキャッシュ操作プロセッサ、C H P U Mはキャッシュ操作プロセッサ用メモリである。主メモリ（M）には処理プログラムやデータが格納されており、キャッシュ操作プロセッサ用メモリ（C H P U M）にはキャッシュ操作プログラム（C H P R）が格納されている。キャッシュ操作プロセッサは、キャッシュ操作プログラム（C H P R）を実行し、P Uは処理プログラム（P P R）を実行する。

データ・キャッシュ・メモリ（D C H M）は、ブロックと呼ばれるmワードからなるメモリn個から構成されている。各ブロックにはD C B 1からD C B nの名前がつけられており、

これをブロックナンバと呼ぶことにする。この1ブロックを構成するmワードがデータ・キャッシュ・メモリと主メモリ間を移動する基本単位であり、この大きさはラインサイズ（あるいはブロックサイズ）と呼ばれている。主メモリ内のデータは、当然のことながらmワード単位で移動させられるので、mワード毎の大きさに分割されている。

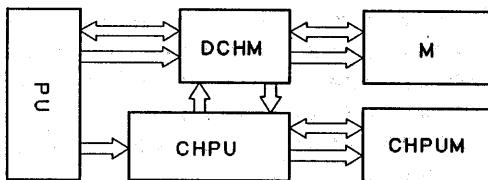


図2 ハードウェアの構成例

2.3 キャッシュ操作命令

次に示す命令が、データ・キャッシュ操作命令である。

(1) 主メモリの内容をデータ・キャッシュ・メモリへロードする命令（ロード命令）

LMC アドレス、キャッシュ・ブロック・ナンバ

(2) データ・キャッシュ・メモリの内容を主メモリへストアする命令（ストア命令）

SCM キャッシュ・ブロック・ナンバ

(3) 分岐命令

JUMP アドレス : 無条件ジャンプ命令

BRANCH アドレス : 条件ジャンプ命令

(4) コール命令

CALL アドレス : コール命令

RETURN : リターン命令

2.4 明示化方法の分類

キャッシュ操作命令のうちロード命令とストア命令とを以下の様に分類することで、明示化方法を3種類検討する。

(1) ロード命令（主メモリ→キャッシュ・メモリ）

ロード命令は、

LMC OP1, OP2

のように表現される。OP1は主メモリのアドレスを示し、OP2はキャッシュ・メモリのブロックナンバを示す。この命令によって、OP1の示す主メモリのブロックからOP2の示すキャッシュ・メモリのブロックへデータの転送を行う。

このオペランドOP1, OP2のとり得る値によって、LMC命令を分類すると、次のようになる。

LMC	OP1	OP2	例
S	S		LMC 100, 1
S	SD		LMC 100
S	D		LMC 100, R1
D	S		LMC R1, 1
D	SD		LMC R1
D	D		LMC R1, R2

この表で、Sはオペランドを静的に決定することを意味し、Dはオペランドを動的に決定することを意味する。また、SDはオペランドを決めるアルゴリズムを静的に決め、そのアルゴリズムによってオペランドを動的に決定することを意味する。表に示した以外の組み合わせも考えられるが、検討する意味はない。また、R1, R2はレジスタを表わし、その値は動的に決定される。

また、本論文でOP2にDを指定する方法についての検討は行っていない。

現在のプログラミング・スタイルでは、S-S型命令 (OP1がSで、OP2がSの命令) のみで明示化プログラムを作成できない。したがって、S型命令 (OP1がSの命令) とD型命令とを組み合わせて、明示化プログラムを作成しなければならない。S型命令とD型命令の組み合せは、

(S-S, D-S), (S-S, D-SD)

(S-SD, D-S), (S-SD, D-SD)

の4通り考えられる。しかし、(S-S, D-SD) と (S-SD, D-S) の組み合せは実現できない。片方の命令のOP2がSであり、もう一方の命令のOP2がSDである場合、Sを決定できないからである (SDは実行時に決定されるので、Sを静的に決定できない)。したがって、実現できる組み合せは、

(S-S, D-S), (S-SD, D-SD)

の2通りである。

(2) ストア命令（キャッシュ・メモリ→主メモリ）

ストア命令は、

SCM OP1, OP2

のように表現される。OP1はキャッシュ・メモリのブロックナンバーを示し、OP2は主メモリのアドレスを示す。この命令によって、OP1の示すキャッシュ・メモリのブロックからOP2の示す主メモリのブロックへデータの転送を行う。

OP2の指定を行うことには意味がない。なぜなら、OP2は静的に決定されるものであって、動的に決定されたり、一定のアルゴリズムによって決定されるものではないからである。オペランドOP1のとり得る値によって、SCM命令を分類すると、次のようになる。

SCM	OP1	例
SCM	S	SCM 100
SCM	SD	
SCM	D	SCM R1

この表中の記号の意味は、LMC命令の場合と同様である。OP1がSDの場合には、SCM命令は省略可能であり、ストアに関してその実行タイミング、実行内容は動的に決定されることになる。

また、本論文でOP1にDを指定する方法についての検討は行っていない。

(3) 明示化方法

(1), (2)よりLMC命令とSCM命令の組み合わせは、

- [1] LMC命令……(S-S, D-S)型
SCM命令……S型
- [2] LMC命令……(S-S, D-S)型
SCM命令……SD型
- [3] LMC命令……(S-SD, D-SD)型
SCM命令……SD型

の3つになる。これらの方法をそれぞれ完全静的決定、半静的決定、動的決定と呼ぶことにする。

(4) キャッシュ操作プログラム例

図3にキャッシュ操作プログラムの例（完全静的決定の場合）を示す。同図(a)は、高級言語による表現であり、同図(b)は処理プログラムであり、同図(c)はキャッシュ操作プログラムである。同図(d)はコントロール・フロー・グラフによる表現である。

2.5 特徴

明示化キャッシュ・システムでは、静的な解析によりキャッシュ操作プログラムを生成しているために、従来のキャッシュ・システムと比べて、以下のような優れた特徴を持っている。

(1) 静的に解析可能な場合には、ヒット率が理論的に最高となる。

コンパイラにより静的に解析可能な場合には、キャッシュ・ミスが発生しないように全てのキャッシュ操作命令の実行スケジューリングを行うことができる。したがって、キャッシュ操作プロセッサは、処理プロセッサの必要とするデータをキャッシュ・メモリに前もって転送することができる。そのため、ヒット率を理論的に最高の値にすることができます。

(2) キャッシュ・メモリの最適な置換が得られる。

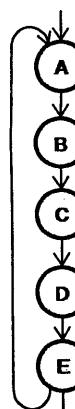
従来のキャッシュ・メモリではハードウェアにより置換を行っているため、固定的なアルゴリズムでしか行うことができない。それに対し、本システムではキャッシュ操作命令により自由に置換を行うことができるため、キャッシュ・メモリを効率よく用いることができる。その例を図4に示す。

同図(a)はループを表すプログラムである。(a)の丸印はプログラム・ブロックを示している。また、ABCDEは、

各プログラム・ブロックで使用されるデータ・ブロックを示している。キャッシュ・ブロックの数は4であるとする。この場合には、(a)のプログラムを実行するとキャッシュ内容の置換が必要となる。同図(b), (c)は、キャッシュ・メモリの状態遷移を表す。キャッシュ・メモリの状態は(1)(2)(3)(4)(5)の順に変化する。また、図中の*は、プロセッサにより実行中のプログラム・ブロックを表す。

(b)はフル・アソシティブ方式で、置換アルゴリズムがLRU法の場合である。この場合をみると、状態(1)から状態(2)への遷移ではAがEに置換され、状態(2)から状態(3)への遷移ではBがAに置換される、というように連続して置換、すなわちキャッシュ・ミスが起こっている。

(c)は明示化キャッシュ・システムの場合である。(c)の場合をみると、状態(2)から状態(3)への遷移でBがEに置換されるだけであり、最小限の置換しか起こっていないことがわかる。データ・ブロックA, C, Dは入れ替えられてはいない。明示化キャッシュ・システムでは、キャッシュ操作が静的にスケジューリングされているために、このようなことが可能である。



(1)	(2)	(3)	(4)	(5)
A	E*	E	E	E
B	B	B*	A*	A
C	C	C	C	B
D*	D	D	D	C*

図4(b) フル・アソシティブ方式、
LRU置換法の場合

(1)	(2)	(3)	(4)	(5)
A*	A	A	A	A
B	B*	E	E	E*
C	C	C*	C	C
D	D	D	D*	D

図4(c) 明示化キャッシュの場合

図4(a) ループを示すコントロール・フロー・グラフ

(3) キャッシュ・ミス時の遅延時間を低減することができる。
従来のキャッシュ・システムはキャッシュ・ヒット、キャッシュ・ミスの判定終了後でなければ、主メモリからのデータ転送が開始されなかった。そのため、キャッシュ・ミス時にはプロセッサが待ち状態に入り、処理が止まってしまった。それに対して、本システムは、キャッシュ・ミスが起こってからデータ転送が開始されるのではなく、時間的に可能な限り早く転送を開始できる（キャッシュ操作を先行して実行することが可能）。そのため、キャッシュ・ミス時の遅延時間を低減できる可能性が大きい。

静的に完全に解析不可能な場合にはキャッシュ・ミスが起こるが、次のようなことを考慮することにより、キャッシュ・ミスを減らすことができる。

(4) 分岐条件の確定を早めるような最適化を行う。

キャッシュ操作プロセッサが、処理プロセッサと同じように条件分岐を行う場合には、分岐条件の確定をはやめることにより、キャッシュ操作プロセッサは先行してキャッシュ操作命令を実行できる。これにより、条件分岐の条件が成立するときも、不成立の時もキャッシュ・ミスをなくすことができる。

(5) 条件分岐の場合に分岐先もキャッシュ・メモリ内に用意する。

条件分岐の場合には実行時でなければ分岐するかどうかは分からない。分岐が成立する場合にキャッシュ・ミスの発生する可能性があるが、分岐先をあらかじめキャッシュ・メモリ内へ転送しておくことにより、キャッシュ・ミスを未然に防ぐことが可能である。また、条件不成立時にキャッシュ・ミスが起らないようにするために、条件分岐命令をプログラム・ブロックの後ろの方におくことを可能な限り避ける必要がある。このようなことをプログラムで指定できることが大きな特徴となっている。

(6) レジスタ間接アドレッシングを行うレジスタの値の確定を早めるような最適化を行う。

図5 (a)に最適化前、(b)に最適化後のプログラムを示す。
最適化の基本的な考え方は、

アドレス生成を行っている命令の実行を早めるように、その命令を前方（プログラムリストの上方）へ移動する。
である。図5の場合、同図(a)中のadd命令（3行目）を前方（同図(b)の2行目）へ移動している。現在のところ局所的な命令順序の入れ替えのみを考慮しているため、大域的にみると最適化にならない場合も存在し得る。

```
begin
  ans1 := 0;
  ans2 := 1;
  for i:=1 to 4 do begin
    ans1 := ans1 + a[i];
    ans2 := ans2 * a[i]
  end
end
```

図3(a) 高級言語によるプログラム

```
01: move x-1,R1
02: move 0,R2
03: move 1,R3
04: add.d R1,1,R1
05: d.load R1,R4
06: add R4,R2,R2
07: mul R4,R3,R3
08: noteq R4,x+3
09: branch 04
0A:d.store R2,y
0B: store.d R3,y+1
```

図3(b) 処理プログラム

```
01:c.LMC.c R1,DCB1
02: branch 01
03: LMC.c y,DCB2
04:c.SCM DCB2
```

図3(c) キャッシュ操作プログラム

ADD.d DataA, R10, R10	ADD.d DataA, R10, R10
d.load R10, R11	ADD.d Temp, R2, R12
ADD.d Temp, R2, R12	d.load R10, R11
d.store R11, R12	d.store R11, R12

図5 (a) 最適化前

ADD.d DataA, R10, R10	ADD.d DataA, R10, R10
d.load R10, R11	ADD.d Temp, R2, R12
ADD.d Temp, R2, R12	d.load R10, R11
d.store R11, R12	d.store R11, R12

図5 (b) 最適化後

3. 性能評価

3.1 ソフトウェア・シミュレータの概要

性能評価を行うために、従来のデータ・キャッシュをもつ逐次計算機用（n c - 1），完全静的決定，半静的決定，動的決定の各場合の明示化キャッシュ計算機用（e c - 1，e c - 2，e c - 3）の4種類のシミュレータを作成した。

これらのシミュレータはC言語で記述され、ワークステーションのX-Wi n d o w上で動作する。

今回作成した4種類のシミュレータでは次のことを仮定している。

- ・プロセッサの命令は、命令キャッシュより供給され、そのヒット率を1とする。
- ・命令キャッシュに関するシミュレーションは行わない。
- ・プロセッサ内はパイプライン方式をとらない。
- ・キャッシュ・メモリには、ライトバッファが存在し、書き込みに要する時間は0であるとする（SCM命令のフェッチ時間は考慮するが、本来の処理（キャッシュ・メモリから主メモリへのデータ転送）を行う時間は0である）。

(1) n c - 1 の仕様

- ・プロセッサは逐次実行を行う。
- ・キャッシュ・メモリの構成方法として、フル・アソシティブ方式、セット・アソシティブ方式、ダイレクト・マップ方式を採用可能。

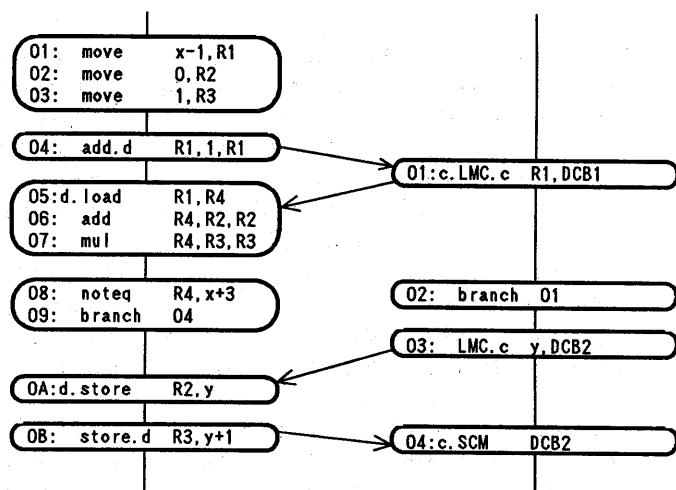


図3 (d) コントロール・フロー・グラフによる表現

- ・キャッシュ・メモリの置換方法として、LRU法、FIFO法、RANDOM法を採用可能。
- ・キャッシュ・メモリの書き込み方式は、コピー・バックのみをインプリメント。
- (2) ec-1, ec-2, ec-3の仕様
- ・プロセッサ、キャッシュ操作プロセッサの並列実行。
- ・命令ストリーム間の同期は機械命令単位で行う。
- ・同期通信のための時間は0とする。
- ・ec-3のみ、nc-1と同様なキャッシュ・メモリを持つ。

3.2 評価プログラム

評価に用いたプログラムは、以下の15種類である。

- ・ソート7種類。
 - (1)バブルソート:bubblesort.
 - (2)分布数えソート:distsort.
 - (3)ヒープソートによって:heapsort.
 - (4)クイックソート(非再帰版):qs.
 - (5)クイックソート(再帰版):quicksort.
 - (6)シェルソート:shellsort.
 - (7)単純挿入法ソート:stringssort.
- ・検索3種類。
 - (8)配列中のデータを二分検索する:binarysearch.
 - (9)行列中のデータを検索する:matrixsearch.
 - (10)配列中のデータを逐次検索する:seqsearch.
- ・その他5種類。
 - (11)ハッシュを用いデータをテーブルへ登録する:hash.
 - (12)行列の乗算(アルゴリズム1):matmult.
 - (13)行列の乗算(アルゴリズム2):matrixmult.
 - (14)データをマージする:merge.
 - (15)素数を求める:sieve.

3.3 評価の前提

明示化データ・キャッシュの評価を行うためにいくつかの条件を設定する。まず最初に、キャッシュ・メモリのブロック数を制限するかどうかによって、

- ・制限なし:キャッシュ操作を明示化する際に、キャッシュ・メモリのブロックを最大限利用する。
- ・制限あり:キャッシュ操作を明示化する際に、キャッシュ・メモリのブロックを2に限定する。

の2つの場合について評価を行う。また、2.4で述べたように、明示化方法を3種類検討する。

- ・完全静的決定:LMC命令の第2オペランド、SCM命令のオペランドを完全に静的に決定する(LMC命令が(S-S, D-S)型であり、SCM命令がS型)。
- ・半静的決定:LMC命令の第2オペランドを静的に、SCM命令のオペランドを動的に決定する(LMC命令が(S-S, D-S)型であり、SCM命令がSD型)。
- ・動的決定:LMC命令の第2オペランド、SCM命令の

オペランドを動的に決定する(LMC命令が(S-SD, D-SD)型であり、SCM命令がSD型)。

である。ラインサイズは2, 4ワードの場合について行った。また、最適化を行った場合と最適化を行っていない場合についても行った。最適化については、2.6参照。キャッシュ・メモリのアクセスタイムと主メモリのアクセスタイムの比が、1:1, 1:2, 1:3, 1:4, 1:5の各場合について評価を行った。

また、逐次計算機では、明示化キャッシュ計算機の場合と同容量のキャッシュ・メモリを使用する。キャッシュ・メモリの方式はフル・アソシエティブ方式であり、置換法はLRU法であり、書き込み方式はコピー・バック方式である。

3.4 評価結果

(1) 実行時間の比較

表1, 2に評価結果の一部を示した。表1はキャッシュ・メモリ・ブロックの数を制限しない場合であり、表2はキャッシュ・メモリ・ブロック数を2に制限した場合である。また、表1, 2ともラインサイズは4であり、最適化はなしであり、キャッシュ・メモリと主メモリのアクセスタイムの比は1:5である。表中の数値は、実行時間比を示している。実行時間比は、従来のキャッシュ・メモリを持つ逐次計算機の実行時間に対する明示化キャッシュ計算機の実行時間の比である。実行時間比が小さいほど、明示化キャッシュ計算機の性能が高いといえる。

[ブロック数の制限をしない場合] 表1参照。

・完全静的決定の場合、明示化キャッシュ計算機の性能の高い例が7ある。その中でもシェルソートは15%程度実行時間が短くなっている。逆にクイックソート(非再帰版)では、実行時間が2倍程度になっている。これは、書き込みをどのタイミングで行うか(実行タイミング)、どのキャッシュブロックを主メモリに書き込むか(実行内容)をSCM命令によって静的に決定していることに起因する。つまり、静的解析では書き込みの実行タイミングと実行内容を特定できない場合が存在するからである。クイックソート(非再帰版)の場合には、多くの書き込みが存在するため、(静的解析で特定できない場合が多く存在するため) SCM命令も多く存在する。そのためSCM命令のフェッチ時間がネックとなっている。また、ラインサイズが2から4へ変化することによって、実行時間比の大きくなったプログラムは7例であり、その他の例では変化ないか、小さくなかった。ラインサイズの変化によって、実行時間比がどのように変化するか一概には決められない。

・半静的決定の場合、明示化キャッシュ計算機の性能の高い例が11ある。完全静的決定の場合と同様に、シェルソートは15%以上実行時間が短くなっている。残りの4例については、10%程度実行時間が長くなっている。これは、キャ

ッシュ・ブロック・ナンバを LMC 命令によって静的に決定していることに起因する。つまり、キャッシュ・ブロック・ナンバを静的に決定することによって、キャッシュ・メモリ全体を有効に使えないからである。また、ラインサイズが 2 から 4 へ変化することによって、実行時間比の大きくなつたプログラムは 1 例のみであり、その他の例では変化ないか、小さくなつた。ラインサイズの変化によって、実行時間比は小さくなる傾向にあるといえる。

・動的決定の場合、明示化キャッシュ計算機の性能の高い例が 9 ある。完全静的決定の場合と同様に、シェルソートは 15 %以上実行時間が短くなつてゐる。残りの 6 例については、従来のキャッシュ・メモリと変わつてない。半静的決定の場合に比べ、実行時間が長くなる例がいくつある。その理由としては、メモリのアクセスパターンによって、従来のキャッシュ・メモリの問題点（無駄な置換がおこる）が発生したことによると思われる。また、ラインサイズが 2 から 4 へ変化することによって、実行時間比の大きくなつたプログラムはなかつた。したがつて、ラインサイズの変化によって、実行時間比は小さくなるといえる。

3 つの場合を比較すると実行時間比は、

動的決定 < 半静的決定 < 完全静的決定

の順に小さくなつてゐる。半静的決定の場合と動的決定の場合との間の大小関係は一概には決まらないが、全体的にみると動的決定の場合の方が優れている。

[ブロック数を 2 に制限した場合] 表 2 参照。

この場合は、ブロック数を制限しない場合とほぼ同様なことがいえ、実行時間比は、

動的決定 < 半静的決定 < 完全静的決定
の順に小さくなつてゐる。

program	完全静的	半静的	動的
binarysearch	1.05	1.05	0.99
bubblesort	0.94	0.88	0.82
distsort	0.94	0.85	0.92
hash	1.00	0.95	0.97
heapsort	1.03	0.93	0.90
matmult	0.99	0.97	1.00
matrixmult	0.99	0.98	1.00
matrixsearch	0.99	0.99	0.99
merge	0.94	0.83	0.88
qs	2.24	1.10	1.00
quicksort	1.23	0.95	0.93
segsearch	1.00	1.00	1.00
shellsort	0.83	0.80	0.80
sieve	1.29	0.99	1.00
strinssort	1.32	1.08	1.00

表 1 実行時間比

(ブロック数の制限なし、ラインサイズ : 4 ワード、最適化なし、アクセスタイムの比 = 1 : 5)

program	完全静的	半静的	動的
binarysearch	1.05	1.05	0.99
bubblesort	0.94	0.88	0.82
distsort	0.99	0.92	0.90
hash	0.97	0.95	0.96
heapsort	0.93	0.92	0.90
matmult	0.95	0.94	1.00
matrixmult	0.95	0.95	1.00
matrixsearch	0.99	0.99	0.99
merge	0.88	0.88	0.88
qs	1.27	1.09	1.00
quicksort	1.00	0.98	0.92
segsearch	1.00	1.00	1.00
shellsort	0.76	0.73	0.80
sieve	1.09	0.99	1.00
strinssort	1.27	1.09	1.00

表 2 実行時間比

(ブロック数の制限あり、ラインサイズ : 4 ワード、最適化なし、アクセスタイムの比 = 1 : 5)

(2) 最適化の効果

表 3, 4 に評価結果の一部を示した。表 3 はキャッシュ・メモリ・ブロックの数を制限しない場合であり、表 4 はキャッシュ・メモリ・ブロック数を 2 に制限した場合である。また、表 3, 4 ともラインサイズは 4 であり、キャッシュ・メモリと主メモリのアクセスタイムの比は 1 : 5 である。表中の数値は、最適化実行時間比を示している。最適化実行時間比は、最適化前の実行時間に対する最適化後の実行時間の比である。最適化実行時間比が小さいほど、最適化の効果が大きいといえる。

[ブロック数の制限をしない場合] 表 3 参照。

・完全静的決定、半静的決定、動的決定の各場合について、最適化実行時間比を平均するとそれぞれ 0.97, 0.97, 0.95 となる。どの場合も、数 % ~ 10 % 程度最適化の効果の出ていることがわかる。しかし、半静的決定の場合、最適化実行時間比が 2 例で 1 をこえている。これは、最適化を局所的に行っていることに起因すると思われる。つまり、最適化する際に大域的な考慮を行っていないので、最適化を行つたことによって関係のない部分に悪影響を与えたものと考えられる。

[ブロック数を 2 に制限した場合] 表 4 参照。

・完全静的決定、半静的決定、動的決定の各場合について、最適化実行時間比を平均するとそれぞれ 0.96, 0.96, 0.95 となる。ブロック数を制限しない場合と同様に、数 % ~ 10 % 程度最適化の効果の出ていることがわかる。

4. おわりに

本論文では、明示化データ・キャッシュの概要を示し、従来のデータ・キャッシュ・メモリを持つ計算機用、明示化データ・キャッシュ・メモリを持つ計算機用のシミュレータを作成し、評価を行つた。

主メモリからキャッシュ・メモリへデータの転送を行う命令(ロード命令), キャッシュ・メモリから主メモリへデータの転送を行う命令(ストア命令)の実現方法を考慮し, 明示化方法を3種類検討した。各方法の異なる点は, ロード命令やストア命令のオペランドの指定方法である。オペランドの値を“静的に決定する”か“動的に決定する”かによって, それらの方法は区別される。1番目の方法(完全静的決定)はロード命令, ストア命令ともオペランドを静的に決定し, 2番目の方法(半静的決定)はロード命令のオペランドを静的に, ストア命令のオペランドを動的に決定し, 3番目の方法(動的決定)はロード命令, ストア命令ともオペランドを動的に決定する。評価の結果, どの方法によって明示化を行っても, 明示化データ・キャッシュはプログラムの実行時間を数%~15%程度短縮できる可能性を持つことがわかった。各方法の優位性について, 一概に言うことはできないが, 完全静的決定の場合より半静的決定, 動的決定の場合の方が優れていることは, 明らかになった。また, プログラムを最適化することによって, 実行時間がさらに数%程度短縮されることも明らかになった。

参考文献

- [1] 佐藤, 有田, 曽和, “キャッシュ操作明示化の提案”, 情報処理学会研究報告, 90-ARC-80 (1990.1).
- [2] 佐藤, 曽和, “並列処理によるキャッシュ操作の明示化”, 並列処理シンポジウムJSPP'90, pp.1-7 (1990.5).
- [3] 佐藤, 曽和, “明示化キャッシュ・システムの性能評価”, 電気関係学会東海支部連合大会 講演論文集, pp.675 (1990.10).
- [4] 佐藤, “キャッシュ・メモリの操作明示化に関する基礎研究”, 昭和63年度名古屋工業大学卒業論文, (1989).
- [5] 佐藤, “キャッシュメモリの操作明示化に関する研究”, 平成元年度名古屋工業大学修士論文, (1991).

program	完全静的	半静的	動的
binarysearch	1.00	1.00	1.00
bubblesort	0.94	0.94	0.93
distsort	0.94	0.94	0.94
hash	1.00	0.98	0.99
heapsort	0.99	1.00	0.99
matmult	0.92	0.92	0.92
matrixmult	0.90	0.90	0.89
matrixsearch	1.00	1.00	1.00
merge	0.98	0.96	0.96
qs	1.00	0.91	0.91
quicksort	0.96	1.00	0.95
seqsearch	1.00	1.00	1.00
shellsort	1.00	1.00	1.00
sieve	0.94	0.93	0.93
strinssort	1.00	0.92	0.91

表3 最適化実行時間比
(ブロック数の制限なし, ラインサイズ: 4ワード,
アクセスタイムの比 = 1:5)

program	完全静的	半静的	動的
binarysearch	1.00	1.00	1.00
bubblesort	0.94	0.94	0.93
distsort	0.95	0.94	0.94
hash	1.00	0.99	0.99
heapsort	0.99	1.00	0.99
matmult	0.92	0.92	0.92
matrixmult	0.90	0.90	0.86
matrixsearch	1.00	1.00	1.00
merge	0.97	0.97	0.97
qs	0.93	0.92	0.91
quicksort	0.96	0.95	0.95
seqsearch	1.00	1.00	1.00
shellsort	1.00	1.00	1.00
sieve	0.94	0.92	0.93
strinssort	0.93	0.92	0.91

表4 最適化実行時間比
(ブロック数の制限あり, ラインサイズ: 4ワード,
アクセスタイムの比 = 1:5)