

分散演算器・データキャッシュを持つクラスタ構成

VLIW計算機

安倍正人 松沢 伸一 岡部 公起 根元義章

東北大学大型計算機センター

〒980-77 仙台市青葉区片平 2-1-1

東北大学大型計算機センター

Tel.022-221-3380 Fax.022-262-3422

分散データキャッシュメモリと分散演算器を数百個あるいは数千個用いたVLIW計算機システムは、例えば各プロセッサエレメントで演算した結果を後処理をすることが必要なアプリケーションのようにプロセッサ間の通信が必要な場合、同期が必要な場合およびバスのアービトレーションが必要な場合でもペナルティが小さいという利点がある。しかし、システムが複数のチップあるいはボードに分かれる場合を想定すると1つのシーケンサで協調させて動作させる必要があることから、MIMD型計算機の $1/5$ ないし $1/10$ 程度のクロック周波数で動作させる必要があり、期待した程の有効性が無い。そこで、外部から与えた比較的遅いクロックに同期させて内部で合成した高い周波数のクロックで常時動作させ、分岐命令は例えば命令の格納アドレスの下位2ビットが11のときのみ実行されるようコンパイラがスケジュールすることにより、チップ間およびボード間のクロックの時間差を吸収し、MIMD型計算機で用いられているものと同じ程度のクロック周波数で動作させることのできるVLIW計算機を提案する。そして、いくつかのリバモアベンチマークにより提案方式の有効性を示す。

VLIW Computer with Distributed Operation Units and Data Cache

Masato Abe

Shinichi Matsuzawa Kouki Okabe Yoshiaki Nemoto

Computer Center, Tohoku University

2-1-1 Katahira, Aoba-ku, Sendai 980-77, Japan

Tel.022-221-3380 Fax.022-262-3422

VLIW computer, which is composed of more than hundreds of operation units, has an advantage that the penalty due to the communication between operation units is small, compared to the MIMD computer. However, since VLIW computer system should be composed of many LSI chips in different PCB boards, a slower clock frequency has to be used in VLIW computer compared in MIMD computer, and the total system performance is not good. Therefore, we propose a method to make the VLIW computer operate with a faster clock as well as a MIMD computer by a compiler technique, by which a branch instruction can only be issued, for example, when the lowest two bits of the instruction address is 11. Then, the difference in the clock, which is fed to each operation unit, can be made negligible.

The effectiveness of the proposed method is confirmed by several benchmark programs.

1 はじめに

計算速度の高速化を図るために、アーキテクチャ的にはスーパースカラ計算機やVLIW計算機が提案され、またディレイドブランチやソフトウェアパイプライン技術といった有用なソフトウェア技術も開発されている。しかし、これらの技術が想定している計算機は複数の演算器は用意しているものの、データキャッシュメモリが1つであるため、演算器へのデータ供給能力が十分でないという問題がある。この問題を解決するために、さまざまな並列システムが研究開発されているが、SIMD計算機は各演算器が全て同一の動作しか行えないので融通性がない。一方、MIMD型計算機においては、システムが複数のチップあるいはボードに分かれた場合を想定すると、各PEは独立に動作できるため、1つのシーケンサで協調させて動作させる必要が無いことから、SIMD計算機の5倍ないし10倍のクロック周波数で動作させることが出来るという利点がある。しかし、例えば各プロセッサエレメントで演算した結果を後処理をすることが必要なアプリケーションのようにプロセッサ間の通信が必要な場合には同期やバスのアービトレーションの点で、有効な結果が得られないという問題がある。そこで、本論文において我々は複数のプロセッサエレメントを用いたMIMD計算機システムに対し、分散データキャッシュメモリと分散演算器を持たせたVLIW計算機システム(DC/PU型VLIW計算機: VLIW Computer with Distributed Caches and Processing Units)を提案する。このとき、1つのシーケンサで複数のチップあるいはボードに分かれた演算器やデータメモリをクロック周波数を下げずに協調させて動作させるために、外部から与えた比較的遅いクロックに同期させて内部で合成した高い周波数のクロックで常時動作させ、分岐命令は例えば命令の格納アドレスの下位2ビットが11のときのみ実行されるようゴンパイラがスケジュールすることにより、チップ間およびボード間のクロックの時間差を吸収する。

本論文では、提案手法の有効性を幾つかのリバモベンチマークテストを用いてMIMD計算機

と比較検討した結果も報告する。

2 DC/PU型VLIW計算機システム構成と命令長の比較

MIMD計算機において、配線量およびバス競合の頻度を考慮して、しばしばクラスタ構成を取るように、本論文で提案するシステムにおいても図1に示すように、クラスタ構成を取ることにする。第4節で述べるシミュレーションにおいては、MIMD計算機では8個のPE(プロセッサエレメント)、提案方式では8個のOPC(Operation Unit with Cache: データキャッシュおよびプロセッシングユニットを合わせたもの)を1つのクラスタとする。(これを1次クラスタと呼ぶ)そして8個の1次クラスタを集めて64個のPEあるいは64個のOPCからなる2次クラスタシステムを構成する。更に、2次クラスタを8個集めることにより、3次のクラスタを構成し、最大512個のPEあるいは512個のOPCからなるシステムを構成する。そして、以上を繰り返すことにより、8のN乗(Nはクラスタの次数)のPEあるいはOPCからなる主システムを構成する。また、最下位のクラスタと同じレベルにDM(分散メモリ)が配置される。更に、主システムと強調して動作するキャッシュおよび分散メモリの管理のためのメモリI/O管理装置を最下位のクラスタと同じレベルそれぞれと最上位のクラスタと同じレベルに1つ配置する。

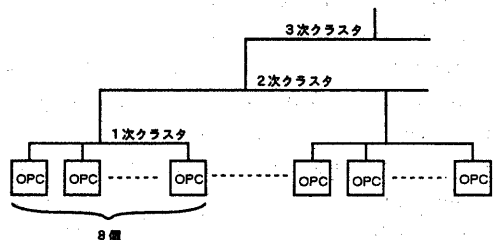


図1: データキャッシュを持つ演算器をクラスタ構成にしたVLIW計算機

OPCは、図2に示すようにデータキャッシュ(DC)、アドレス生成器(AG)、2つのレジスタ

ファイル (RF)、演算器 (PU) で構成され、それぞれが2本の内部データバスに接続されている。ここで、内部バス#1に出力できるのは、データキャッシュ、アドレス生成器付属のレジスタファイル、演算器の3つのうち、一度に1つのみである。一方、内部バス#2に出力できるのは、データキャッシュ、演算器付属のレジスタファイル、演算器の3つのうち、一度に1つのみである。このようにして、各バスへの出力の制御のために必要なVLIW命令のビット数は出力しない場合も含めて、各バスそれぞれに対し2ビットとなる。また、この外に命令キャッシュが別に設けられる。

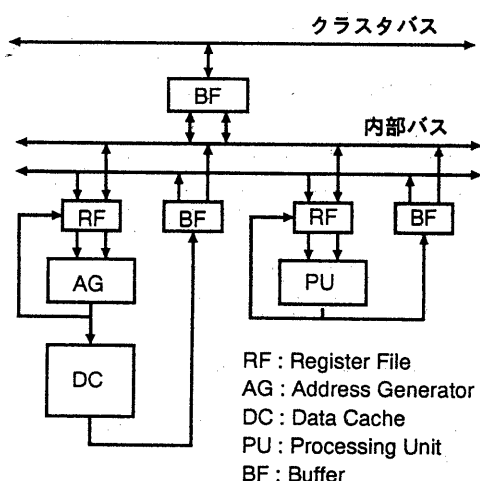


図 2: データキャッシュを持つ演算器OPCの構成

1次のクラスタは、8個のOPCで構成されており、同じクラスタに属する内部バス同士はバスバッファを通して互いに接続されている。これを1次クラスタバスと呼ぶことにする。このクラスタバスに出力するOPCを指定するのに3ビット、バスバッファをイネーブルするかどうかを決定するために1ビット必要とするから、8個のOPCから構成される1次クラスタ1つだけからなるシステムのバス制御のために必要とされるビット数は $(2+2) * 8 + 4 = 36$ ビットとなる。2次のクラスタは8個の1次クラスタから構成され、2

本の1次クラスタバスがバスバッファを経由して互いに接続される。この時、64個のOPCから構成されるシステムのバス制御のために必要とされるビット数は $36 * 8 + 4 = 292$ ビットとなる。最大512個のOPCが接続される3次のクラスタにおいては、同様にして、 $292 * 8 + 4 = 2340$ ビットが必要となる。このことから、1つのOPC当たりに必要なバス制御ビット数は約5 ($= 2340 / 512$) であることが分かる。

アドレス生成器では、付属レジスタファイルの内容そのままの出力、2つのレジスタ間の加減算および1つのレジスタと値1の加減算 (インクリメントとデクリメント) の5通りの演算を行うので、必要な制御ビット数は3となる。2つの付属レジスタファイルの制御には、レジスタ数をそれぞれ64とすると、5ポートレジスタファイルなので、1つのポート当たり6ビットだから、1つのOPC当たりでは $6 * 5 * 2 = 60$ ビット、更に書き込みポートはレジスタ当たり3ポートあるので、その書き込み制御に1レジスタ当たり3ビット必要だから、合計して、2つのレジスタファイルの制御に $60 + 3 * 2 = 66$ ビット必要となる。演算器の制御には、四則演算、論理演算等を考慮して10ビットの制御ビットを割り当てる。最後に、キャッシュメモリの制御には、アドレスおよびデータの一時保持、および保持されたデータの書き込みのために合計3ビット必要とする。以上を総合すると、我々が提案するVLIW型計算機では1つのOPC当たり約90ビットを必要とする。ただし、多くの計算機のように命令の組み合わせを制限し、デコードによって必要な制御信号を作ることにより、必要な命令長を削減することもできる。よって、命令長の立場から見れば、MIMD方式の計算機と大差がないと考えられる。

なお、システム全体として、分散メモリの制御、シーケンサの制御、イミューディエイト生成器の制御のためにさらに幾ビットが必要とするが、例えば64個のOPCからなるシステムの場合、1つのOPC当たり3ビット程度となるのでさほど問題ではない。

3 分散キャッシュメモリの動作

3.1 分散メモリとキャッシュメモリ間のデータ転送

分散メモリは同期式DRAMを用いて実現されるので、連続したアクセスは十分に高速だが、離散的なアクセスは低速である。一方、キャッシュメモリは離散的なアクセスでも十分に高速なSRAMで実現されるが容量が少ない。そこで我々は、演算を制御する主システムとは独立に動作するメモリI/O処理装置(MIO)を用いて、キャッシュメモリを以下に示すように管理することにした。MIOはN個のOPCに分散されたキャッシュメモリを管理するが、各キャッシュメモリは論理アドレス32ビット、1ページ当たり256ワード、全部で64ページのダイレクトキャッシュ方式で管理される。

一方、分散データキャッシュは、Write invalidate方式[1]によりキャッシュのコヒーレンスを保つことにする。このため、図??に示すように分散キャッシュ内のTLBには、アドレスタグAT、有効フラッグVF、他のOPC内の分散キャッシュメモリに同じアドレスページのコピーがあるかどうかを示すコピーフラッグCFおよび、そのページの内容が分散メモリの第一番目のコピーであるかどうかを示すファーストフラッグFFから構成されている。アドレスタグATは該当ページに入っているデータの分散メモリ上での論理アドレスの上位18ビットを保持し、有効フラッグVFは、1ならばキャッシュメモリ中のデータが有効で、0は有効でないことを意味する。一方、8個のOPCあたり1つずつと全体に共通に1つだけ設けられているメモリ入出力装置MIOの構成は図??となっている。

4 OPC間の同期

分岐する場合あるいは自分に所属しないレジスタファイルを参照する場合は、他のOPCと同期を取る、すなわちOPC間のクロックの時間差を補正する必要がある。

4.1 分岐における同期

ここでは、同期用の外部クロックの周期は内部クロックの周期の4倍を仮定し、命令は内部クロックに同期して実行される。図3において、シーケンサ内で外部クロックに従って内部クロックが生成され、外部クロックが立ち上がってから、2つめの内部クロックの立ち下がりで、次に実行すべき4つの命令のアドレスが生成され、全OPCに転送される。各OPCでは、図3に示すように ΔT_1 あるいは ΔT_n の時間差があるが、その時間差が外部クロックの周期の半分程度であれば問題無く動作する。

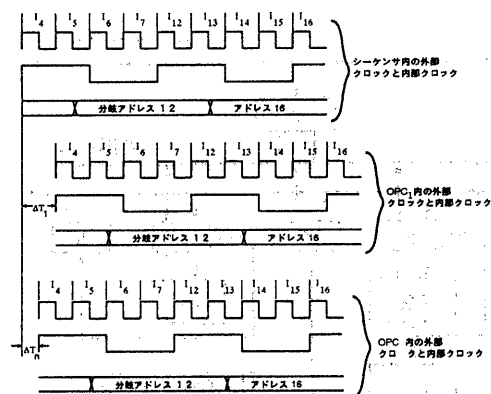


図3: OPC間のクロックの時間差を補正するための分岐シーケンス

この場合、分岐命令が実行できるアドレスが限られ、NOP命令を挿入しなければならないことも生じ、効率が低下する。しかし、全て遅いクロックで動作させるよりは高速化が期待でき、またループの場合はループアンローリングを用いることにより、効率の低下は多少抑制できる。

4.2 自分に所属しないレジスタファイルを参照する場合の同期

例えばOPC₁のレジスタファイルから、OPC_nのレジスタファイルヘデータを転送する場合を考える。この場合もOPC間のクロックの時間差を補正する必要があるので、図4に示すように、例えばOPC₁のレジスタファイルから内部クロック

4つ分の間データを送出し、 OPC_n 内のレジスタファイルへの書き込みはデータが揃ったところを見計らって行われる。この手法では、MIMD計算機におけるメッセージ交換に伴うオーバーヘッドあるいは、バス経由に伴うバスのアービトレーションと同期の問題が無く、高速にOPC間通信を行うことができる。

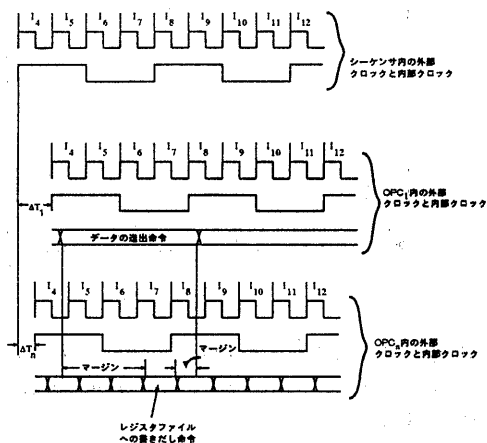


図 4: OPC間のクロックの時間差を補正するためのOPC間通信のシーケンス

5 ベンチマークによる比較

5.1 比較対象としたMIMD計算機の構成

MIMD方式では、さまざまな構成が考えられるが、ハードウェア規模を考慮して我々のと同じクラスタ構成を取ることにした。すなわち、図1においてOPCの代わりにキャッシュメモリ付きのPEがある場合を想定した。また、キャッシュ交換の方式はどちらも同じで、PE間あるいはOPC間の通信方式だけが異なり、提案しているDC/PU方式はバス経由、MIMD計算機はメッセージ交換とした。

5.2 比較条件

4節までの説明ではOPC内のキャッシュメモリの1ページ当たりのデータ数は256ワードとしたが、本シミュレーションでは64、128、256、512ワード(1ワード=4バイト)の4種類について調べた。また、エントリ数は64と

した。更に、キャッシュ入れ替えの際は、1内部クロック当たり8ワード読み込みが出来るものとし、最初のデータを読み込む際のオーバーヘッド、すなわちバスの確保やキャッシュコヒーレンスの確保等に要する時間は内部クロック数で20とし、後続する連続したデータの読み込みに要する時間は内部クロック数で1とした。

一方、MIMD計算機において、PE間でNワードのデータを一度に交換するのに要する時間は、メッセージ交換で行うため、同じクラスタに属する場合は内部クロックで $240 + 10N$ 、異なるクラスタの場合は $490 + 10N$ とした。それに対し、提案するDC/PU計算機では、同期やバスのアービトレーションの問題無く直接バスを経由してデータを交換できるので、内部クロック数で7とした。

6 リバモアベンチマークによるシミュレーション

ここでは、コンパイラが分散メモリにデータを最適に格納しており、各PEあるいはOPC間でデータの競合が出来るだけ起きないようにしている。ただし、ループの回数は100000とした。KERNEL 1、7は、後処理のためのPE間あるいはOPC間通信の必要がないので、スピードアップの結果はDC/PU計算機とMIMD計算機は全く同じである。図5および図6にその結果を示す。

一方、KERNEL 3あるいは24の場合は、各PEあるいはOPCで計算された部分積和あるいは部分最小値を集めて全体の積和あるいは最小値を求める必要がある。このためには、PEあるいはOPC同士が通信し合う必要があり、MIMD計算機ではペナルティが大きい。この様子を図7と図8に示す。

7 むすび

MIMD計算機で問題となる通信によるオーバーヘッドを削減するために、VLW計算機において、1つのシーケンサで多数の分散キャッシュ及び演算器をクラスタ構成にして同時に制御する計

Kernel 1. "HYDRO FRAGMENT"

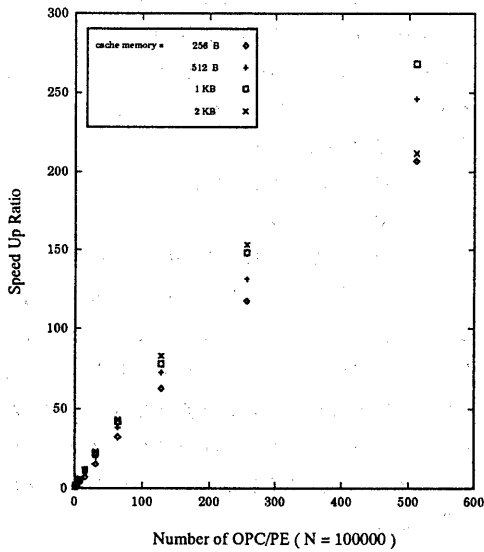


図 5: KERNEL 1 を実行した時の 1 個の CPU に対するスピードアップレシオ (総データ数 100,000)

Kernel 3. "INNER PRODUCT"

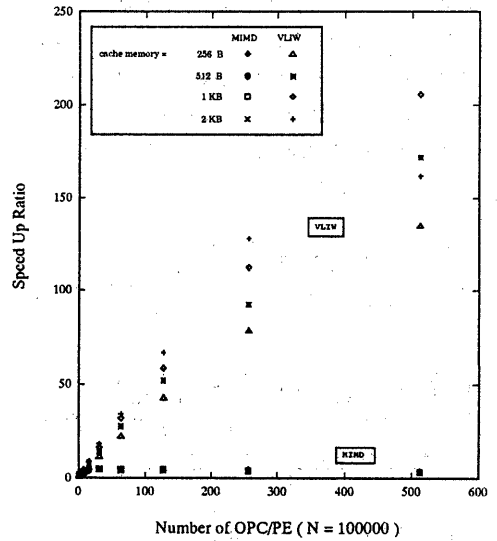


図 7: KERNEL 3 を実行した時の 1 個の CPU に対するスピードアップレシオ (総データ数 100,000)

Kernel 7. "EQUATION OF STATE FRAGMENT"

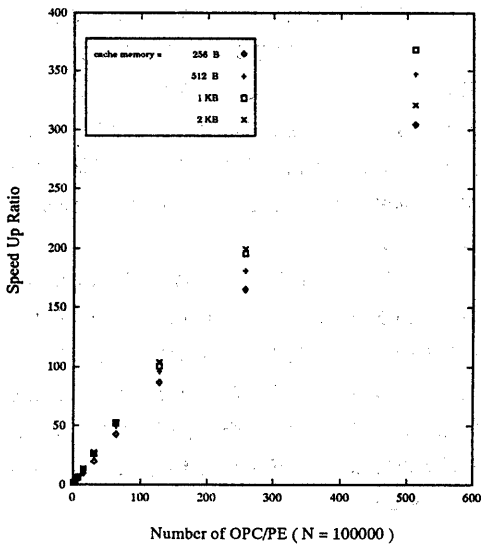


図 6: KERNEL 7 を実行した時の 1 個の CPU に対するスピードアップレシオ (総データ数 100,000)

Kernel 24. "FIND LOCATION OF FIRST MINIMUM IN ARRAY"

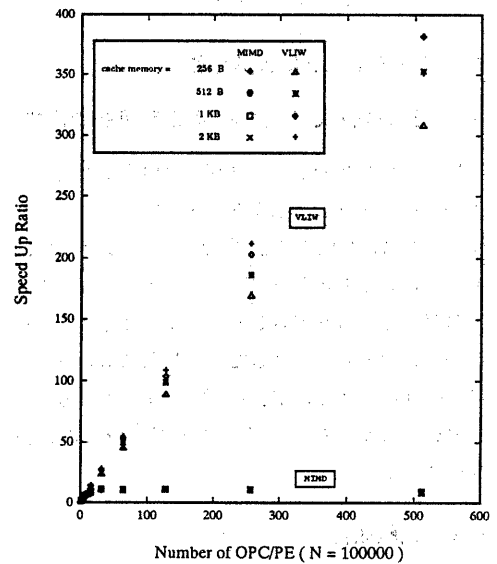


図 8: KERNEL 24 を実行した時の 1 個の CPU に対するスピードアップレシオ (総データ数 100,000)

算機を提案した。この際、複数チップあるいは複

数ボード使用時のクロックの時間差を吸収するた

めに、分岐やPU間の通信は高速な内部クロックではなく比較的低速な外部クロックに同期させることにした。その結果、各PEが独立に動作できることからクロック周波数を上げることで高速化が達成できているMIMD計算機とほとんど変わらない高速化が提案するDC/PU型VLIW計算機で達成できるものと思われる。

リバモアベンチマークのうち並列計算機が有効な4つのプログラムについてシミュレーションを行い、提案するDC/PU型VLIW計算機は、通信が必要なプログラムに対してMIMD計算機より有効であることを示した。

参考文献

- [1] Nitzberg B., Lo. V., "Distributed Shared Memory: A Survey of Issues and Algorithms," IEEE Computer, August 1991, pp.52-60