

メモリモジュール上での等間隔アクセス連続化の効果

田邊 昇^{†1} 箱崎 博孝^{†2} 安藤 宏^{†2}
土肥 康孝^{†2} 中條 拓伯^{†3} 宮代 具隆^{†4}
北村 聡^{†4} 天野 英晴^{†4}

パーソナルコンピュータ(PC)のメモリスロットに装着されるプリフェッチ機能を有するメモリモジュールを提案する。このデバイスは、Pentium[®]4などのCOTS(Comercial Off-The-Shelf)型MPUのキャッシュアーキテクチャの弱点を軽減することで、パーソナルスーパーコンピュータに匹敵する実効性能をPC上でも実現可能にすることを目指している。本論文では、プリフェッチ機構付メモリモジュールの基本コンセプトとそのソフト的対応法について述べ、DIMMnet-2プロトタイプにおける実装案を紹介する。主記憶データベースに対するWisconsinベンチマークへの適用を通し、等間隔参照の高速化に関する評価結果を示す。その結果、タプルサイズがキャッシュラインサイズより大きい場合は、最大で25.2倍の加速率を並列処理なしで得られることが判った。

Effects of Conversion of Strided Access into Continuous Access on a Memory Module

NOBORU TANABE,^{†1} HIROTAKA HAKOZAKI,^{†2} HIROSHI ANDO,^{†2} YASUNORI DOHI,^{†2}
HIRONORI NAKAJO,^{†3} TOMOTAKA MIYASHIRO,^{†4} AKIRA KITAMURA^{†4}
and HIDEHARU AMANO^{†4}

Prefetching functions in a memory module plugged into a memory slot of a PC are proposed. The proposed device overcome the defect of cache architecture of a COTS (Commercial Off-The-Shelf) type MPU such as Pentium[®]4 in order to realize personal supercomputer class performance. In this paper, the concept and its software corresponding method of a memory module with prefetching functions and its implementation plan on a DIMMnet-2 prototype are presented. In addition, evaluations for improvement of a stridden access with Wisconsin benchmark programs are shown. As a result, at most 25.2 times performance improvement without parallel computation is gained, when tuple size is bigger than cache line size.

1. はじめに

半導体技術やアーキテクチャの進歩を背景にしたマイクロプロセッサ(MPU)の目覚ましい進歩を遂げている。例えばPentium[®]4などはスーパーコンピュータを凌駕する周波数で、スーパーコンピュータの1CPUに匹敵する演算性能をオフィスや一般家庭に提供している。量産効果に下支えされ、COTS(Comercial Off-The-Shelf)であるMPUおよびパーソナルコンピュータ(PC)の性能向上およびコストパフォーマンスの向上は目覚ましいものがある。

これらのCOTS部品であるMPUはほぼ例外なくキャッシュアーキテクチャに基づいている。キャッシュアーキテクチャは主記憶の脆弱さを隠蔽できることが多いため、低コストなPCにおける主記憶は、ベクトル型スーパーコンピュータのそれとは異なり、キャッシュが効かないアプリケーションに対しては演算能力にバランスしたものにはなっていない。

科学技術計算の分野では、SX-6i¹⁾のようにベクトル型スーパーコンピュータのほぼ1CPU分を切り出したパーソナルスーパーコンピュータ製品上で、PCとの性能差が数十倍に及ぶアプリケーションが多数存在することが報告されている。このようなアプリケーションにおいてはベクトル型スーパーコン

ピュータは依然として存在意義がある。

一方、DB2[®]の父と称されるIBM[®]フェローのLinsay氏は、コンピュータ上の主記憶が増加することに伴ってDBMSの重要な部分が主記憶に常駐するようになり、キャッシュミスの最小化の重要性が高まることを指摘²⁾しており、キャッシュが効かないアプリケーションは科学技術計算にとどまるものではない。この点は将来、MRAMのような不揮発性大容量高速メモリがDRAMに置き換わるような状況になった場合、この分野におけるキャッシュミスの最小化の重要性は加速する。

キャッシュベースのMPUは記憶階層上で主記憶直前にあるキャッシュのラインサイズは伸びる傾向にある。例えばPentium[®]3ではL2キャッシュのラインサイズが32バイトであるのに対し、Pentium[®]4では128バイトである。キャッシュベースのMPUは常に最下位階層キャッシュライン単位で主記憶をアクセスする。

ところが、例えばデータを主記憶に配置した際のリレーショナルデータベース処理のように、ストライドの大きい不連続アクセスが必要なアプリケーションでは、真に必要なデータが128バイトのキャッシュラインの中に1バイトしか存在しないというケースもありうる。

このような状況が、間接指標配列を用いた圧縮型の配列や多次元配列を用いる科学技術計算のみならず、データベースのベンチマークであるWisconsinベンチマークの多くのクエリーで見受けられる。

ベクトル型スーパーコンピュータはデータベース処理には使われてこなかったが、データ全体が乗ってしまうような大容量な主記憶が安価に提供できるのであれば、テーブルスキャンの高速化には向いている。

^{†1} (株) 東芝, 研究開発センター
Corporate Research and Development Center, Toshiba
^{†2} 横浜国立大学
Yokohama National University
^{†3} 東京農工大学
Tokyo University of Agriculture and Technology
^{†4} 慶應義塾大学
Keio University

しかし、ベクトル型スーパーコンピュータは高価である点、製品更新サイクルが長い点、ユーザーが PC クラスタ等に移行している点などから、今後も市場の拡大や低価格化は望めず、提供するベンダーも減少の一途にある。その入手困難性の増加ともあいまって、その低コストな解決の工学的価値は非常に高いと考えられる。

以上のような観点から、キャッシュが効かないアプリケーションに対するアーキテクチャ面からの対応としては、COTS の目覚ましい性能向上を容易に取り込めるものでありつつ、キャッシュアーキテクチャの欠点を補う方式の研究開発が望ましい。

本研究は以上のような背景を鑑み、PC のメモリスロットに装着されるプリフェッチ機能を有するメモリモジュールを提案してきた。このデバイスは、Pentium[®]4 などの COTS 型 MPU のキャッシュアーキテクチャの弱点を軽減することで、パーソナルスーパーコンピュータに匹敵する性能を PC 上でも実現可能にすることを目指している。

さらに、本研究は現在の PC の主流である DDR 型の DIMM スロットにメモリ以外のデバイスを装着するという先例が無い実装形態をとるため、電氣的観点から実現可能であることを実証する必要がある。そこで、本研究の実機検証は DIMM スロットに装着されるネットワークインタフェース (NIC) である DIMMnet-2⁽³⁾ 上で行われており、既に FPGA 版の DIMMnet-2 プロトタイプがメモリとしてホストからアクセス可能であることは確認できている。

本方式を DIMMnet-2 上に適用することで、実質的に高速化された安価なノード PC が高速なネットワークで結合され、全体として巨大かつ高速アクセス可能な共有メモリを有するスーパーコンピュータに見える PC クラスタが構築される。

ただし、本論文における提案技術は必ずしも NIC と共存して並列処理しなければ意味がないものではない。並列化に抵抗感を感じるより多くのユーザーに対しても、PC の単体性能を向上できる高付加価値メモリモジュールという提供形態もありうる。つまり、Myrinet⁽⁵⁾ や QsNET⁽⁶⁾ などの PC クラスタ用ネットワークインタフェースとは対象とするユーザー層や市場規模が異なる。このため、本手法はこれらにはない量産効果向上による低価格化の実現が期待できる。この点でも独創性と経済性を有する。

本報告では、2 章でキャッシュアーキテクチャの問題点を指摘する。3 章ではプリフェッチ機構付メモリモジュールの基本コンセプトとそのソフト的対応法について述べる。4 章ではその DIMMnet-2 における実装例を紹介する。5 章ではキャッシュアーキテクチャが苦手とする配列への等間隔アクセスを多用する応用である Wisconsin ベンチマークを紹介し、その本報告における評価の対象としての位置づけを示す。6 章では提案システム用に改造された Wisconsin ベンチマークプログラムが PC 上で示す性能から、実機上での処理性能を推定する。7 章で関連研究について述べ、8 章でまとめる。

2. キャッシュアーキテクチャの問題点

キャッシュアーキテクチャでは間接参照などの不連続アクセス時にキャッシュラインの利用効率が低下する。

例えば 128 バイトを超える構造体の配列に対して、整数型のあるメンバを全てのインデックスに対してアクセスするようなストライドの大きい不連続アクセスが必要なアプリケーションを 128 バイトのキャッシュラインを有する Pentium[®]4 で実行する場合、真に必要なデータが 128 バイトのキャッシュラインの中に 4 バイトしか存在しない。キャッシュを用いた CPU を使って等間隔参照を行う際の問題点を図 1 に示し、以下に列挙する。

- (1) 有効データが少ないゆえのメモリバンド幅の浪費
- (2) 有効データが少ないゆえのキャッシュエントリの浪費
- (3) 有効データが少ないゆえの TLB エントリの浪費

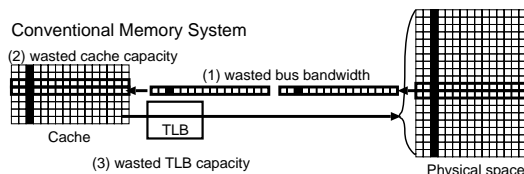


図 1 キャッシュを用いた CPU での等間隔参照の問題点

キャッシュアーキテクチャの上記のような問題点は、例えば主記憶上にデータを保持する主記憶データベースにおける Wisconsin ベンチマークの多くのクエリーにおいても顕著に発生する。半導体メモリの容量拡大と価格低下の着実な進歩に伴い、処理ネックがハードディスクから主記憶アクセスに移動する傾向が強まるため、主記憶上での等間隔アクセスの高速化のニーズが高まってきている。

一方、ベクトル型スーパーコンピュータの等間隔ベクトルアクセス命令によれば、ベクトルレジスタ上に必要なデータの圧縮格納されるため、本来、このような問題に対する適応性が高い。しかし、ベクトル型スーパーコンピュータは大変高価なハードウェアである上、HPC 市場の主流がベクトル型から PC クラスタなどの並列型の計算機に移行し、ソフトサポートもユーザー数の多いそれらのプラットフォーム側に移行してきていることに伴い、今後の安定した入手性を疑問視する意見もある。

上記のような問題が、例えば NAS CG ベンチマークや Wisconsin ベンチマークで代表される重要なアプリケーション群で見受けられ、それらがベクトル型スーパーコンピュータという受け皿を失いつつある現状と合わさって、安価な PC 上でのそれらの問題の解決は工学的な価値が高い。

3. プリフェッチ機能付メモリモジュール

3.1 基本コンセプト

メモリ空間にマップされたメモリモジュール側にあるバッファ(プリフェッチバッファ)へのプリフェッチコマンドを発行し、ホスト CPU から利用確率が高い状態に整えられたデータ群に対してブロックアクセスを行う。

その結果、キャッシュ・TLB・FSB・メモリバスの利用効率が向上する。メモリモジュールは着脱可能なので、MPU やチップセットを改造することなく、高性能かつ低価格な COTS を HPC 向けコンピュータとして有効に活用できる。

なお、プリフェッチコマンドには種々の実装法がありうるが、本研究ではベクトル転送命令について検討する。そのうち本論文で検討するのはベクトル等間隔ロード命令で、配列の等間隔ロードをベクトルレジスタに対して行う命令である。

図 2 は提案メモリモジュールにおける書き込み用 (Write window) と読み出し用 (Read Window) のバッファを有するプリフェッチ機能付き Window メモリであり、これらがベクトルレジスタとして用いられる。これらは通信にも用いることができる。ホストはこれらの window のデータ格納状態をフラグを通じて確認できる。例えばホスト MPU の L2 キャッシュライン分のデータ毎にフラグが 1bit ついていて、ステータスレジスタを読むとこのフラグを確認できる。

これらのハード的な仕組みを使って等間隔配列参照を行った場合の動作と高速化原理を図 3 に示す。前章で列挙した (2) ~ (4) の問題点が全て解決される。

3.2 ソフトウェア的対応

本機構は、キャッシュメモリ等とは異なり、ソフト的に透過なハードではない。利用するためにはプログラムの変更またはコンパイラによる対応が必要である。

本機構を Pentium[®]4 などのキャッシュベースの MPU から利用する際のリード時における指針を以下に示す。

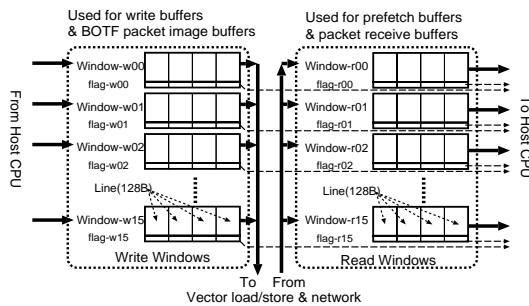


図 2 プリフェッチ機能付き Window メモリ

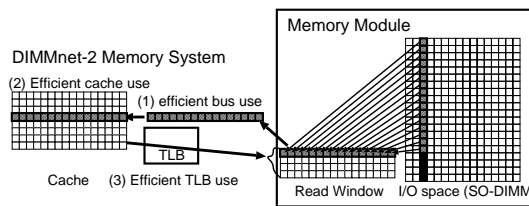


図 3 提案システムにおける等間隔配列参照時のキャッシュ挙動

- (1) 主記憶の WriteBack 属性の領域にスクラッチパッド領域 A を確保する。
- (2) 本機構を利用するためのプリフェッチコマンドを実際にデータを利用するより前に発行する。
- (3) Read window から領域 A にブロックコピーを行うコードをプリフェッチが完了する頃に実行する。Pentium[®]4 の場合は時間的局所性を持たないデータ用のプリフェッチ命令である PREFETCHNTA 命令を適切なタイミングで実行しておくが良い。
- (4) Read window のアドレスに対応するキャッシュラインを再利用前にフラッシュする命令 (Pentium[®]4 の場合は CLFLUSH 命令) を実行する。
- (5) 領域 A (実体はキャッシュ上) にあるデータを用いて、所望の処理を行う。プリフェッチと処理の時間差が少なければキャッシュへのアクセスだけで必要なデータが得られる。

一方、ライト時における高速化の指針は DIMMnet-1 における BOTF における Window メモリへの CPU からの書き込みの高速化指針⁷⁾ とほぼ同様であり、Window メモリにコピーし終えたら、送信を指示するキック操作の代わりにベクトルストアコマンドを起動する。

4. DIMMnet-2 プロトタイプ上での実装

提案されたアーキテクチャの検証用として DIMMnet-2 プロトタイプ上に実装して実機検証を行なうことにしている。本報告では主に等間隔アクセスに関して評価を行なう。

4.1 資源モデル

ローカルな SO-DIMM 上にある配列への等間隔ベクトルロードやベクトル間接ロードを実行するのに必要な、DIMMnet-2 プロトタイプ上に実装されているプログラマから見える資源は以下のものがある。

- (1) SO-DIMM 領域 (4GB×1~16)
ベクトルアクセスされる対象の配列データを格納する。仕様上は最大 4GB のページを 16 枚まで持てるが、このサイズは搭載している SO-DIMM の容量や枚数に依存する。
- (2) RW : Read Window (64bit 幅 64 語 × N 本)
等間隔ベクトルロードコマンドを用いて (1) から必要なデータをロードする。ホストからはユーザー空間にマップされており、キャッシュ可能なメモリ領域として見えている。
- (3) RWSR : RW Status Register (64bit 幅 1 語)

(2) 上へのデータのロードのライン単位での完了状況が反映される。ホストからはユーザー空間にマップされており、キャッシュ不可能なメモリ領域として見えている。

- (4) CMR : Command Register (64bit 幅 16 語)
ベクトルコマンドをホストから書き込むことで等間隔ベクトルロードコマンドなどが起動される。ホストからはユーザー空間にマップされており、キャッシュ不可能なメモリ領域として見えている。

- (5) CMSR : Command Status Register (16bit 幅 1 語)
(4) の完了状況が反映される。ホストからはユーザー空間にマップされており、キャッシュ不可能なメモリ領域として見えている。

- (6) CPR : Current Page Register (16bit 幅 4 語)
(1) のどのページに対するベクトルコマンドなのかを指定する。ホストからはユーザー空間にマップされており、キャッシュ不可能なメモリ領域として見えている。

4.2 ベクトル転送命令

ローカルな SO-DIMM 上にある配列への等間隔ベクトルロードやベクトル間接ロードを実行するのに必要な DIMMnet-2 プロトタイプ上に実装されたコマンドは以下のものがある。コマンドは 64bit 形式であり、ホストから CMR に書き込まれる。

以下、XW は WW と RW の総称、type はデータ型、iter は繰り返し回数、displace は XW 上での開始位置、top は SO-DIMM 領域を集めて構成される大域仮想空間上の offset、stride は要素間間隔 (byte 数)、indexXW は index ベクトルを保持する window である。

- (1) ベクトル等間隔ロード:
VLS(type,iter,displace,RWi,top,stride)
- (2) ベクトル間接ロード:
VLI(type,iter,displace,RWi,top,indexXWj)

4.3 等間隔ベクトルロード時間

現在、DIMMnet-2 プロトタイプは FPGA 内部の論理設計途上であり、一部の機能しか利用できず、使える場合でも性能チューニングが十分な状態ではない。本報告で評価するベクトル等間隔ロードコマンドについては性能チューニングが十分な状態ではないながら、動作が確認できており、Verilog レベルでの現状の所要サイクル数が判っている。現状の回路における等間隔ベクトルロードコマンドのデータ型 type が 4 バイト、要素間間隔 stride が 60 バイトという条件における所要時間を表 1 に示す。ここでは FPGA 版を 100MHz 動作、ASIC 版を 200MHz 動作するものとしている。

繰り返し回数 iter	所要サイクル数	FPGA	ASIC
64	83	830ns	415ns
128	162	1620ns	810ns

5. 評価対象とするアプリケーション

本報告では提案システムの等間隔アクセス主体のアプリケーションの高速化の評価を目的として、主記憶上にデータを格納するデータベースを対象に検討する。半導体メモリの大容量化とビット単価の減少および検索処理への高速化のニーズを背景に、主記憶上にデータを格納する主記憶データベースへの期待感が高まってきており、本報告ではこれを扱う。

各属性データの格納形式はホスト CPU の基本アクセス単位である大きさ (Pentium[®]4 の場合は 32bit CPU なので 32bit) で格納し、それを超える大きさのオブジェクト (文字

列等)についてはポイントとして構成されるタブルの配列で格納されるものとする。4バイトを超えるオブジェクトの内容で検索する必要がある場合のみ、上記のタブルに記載されたポイントをたどり、別途格納されたオブジェクトをアクセスするものとする。すなわち、基本的な検索ではデータのミスアラインの影響という問題はない。

データベースの性能評価には Wisconsin ベンチマークが用いられることがあり、文献⁹⁾でも主記憶データベース処理のSDTによる高速化の評価を目的に Wisconsin ベンチマークが用いられている。本報告での評価でもこれをベースにした評価を行なう。文献⁹⁾の記述では本方式がデータベース処理には向かないということになっているが、それが正しいかどうかを実験的に証明するために、できるだけ文献⁹⁾に記載の評価対象に近いもので評価を行なうことを目指す。

Wisconsin ベンチマークでは15の属性からなるタブルが10K個または1K個から構成されるデータベース(10K個のtenk, 1K個のonek)が検索対象になる。1個のタブルは上記の主記憶に配置される形式で格納すると15個の属性は4バイトのデータまたは4バイトのポイント(合計60バイト)からなる。そのうちの2つのポイントで指し示される2つの文字列データからなる。Wisconsin ベンチマークにおけるクエリー(問合せ)では上記の60バイトの部分が検索に用いられるので、本報告ではこの部分全体がSO-DIMMで構成される領域上に乗るようにする。

このように表形式のデータベースではある属性に対する検索処理を行なう際には等間隔アクセスとなってしまう。上記の Wisconsin ベンチマークでは、ストライド値60バイトで不連続に格納されている4バイトデータに対する等間隔アクセスとなる。

例えば Pentium[®]4のL2 キャッシュラインは128バイトなので1ラインに2個強の4バイトデータしか有効ではない。もし、検索対象の総データサイズがキャッシュ容量よりはるかに大きければ、このキャッシュラインリフィルに対する時間が実行時間の大半を占めることが予想される。この場合、データベースの先頭から順にスキャンするとしても、メモリバス上では4バイトのデータを取得するために平均60バイトを読み出すこととなり、必要なメモリバンド幅は15倍にも膨れ上がる。

文献⁹⁾上での評価ではキャッシュ容量を8KBとして上記のベンチマークそのものをシミュレータ上で実行させている。しかし、近年のMPUは内蔵するキャッシュの容量が増加してきており、例えば Pentium[®]4ではL1およびL2キャッシュの合計が520KBにもなるため、本来の Wisconsin ベンチマークで規定されたタブル数(10K個)では大半がキャッシュに乗ってしまう。それでは主記憶データベースの評価としては不適切であるので、タブル数をキャッシュの増加分である65倍に増やして評価を行なうものとする。

6. 性能評価

6.1 評価環境

表2に性能評価を行ったマシンの仕様を示す。この環境は単にシミュレーションを行った機械という意味合いだけでなく、後述する評価方法を用いているために評価対象そのものを構成する。つまり、キャッシュラインサイズやキャッシュ容量はもちろんのこと、メモリバンド幅や命令のCPU内部での実装の細部に至るまで評価対象そのものを構成し、提案メモリモジュールと組み合わせた際の性能を左右する。これらを仮想的に変更することができない代わりに、通常のシミュレータを用いた評価手法に比べて高速性と精度が高い。

6.2 評価方法

性能評価は Wisconsin ベンチマークにおけるデータの挿入、更新、削除以外の問合せ(クエリー)を利用して行なう。基本

表2 評価環境

機種名	Dell TM Precision TM 360
CPU	Pentium [®] 4
FSB 周波数	800MHz
コア周波数	2.4GHz
L1 キャッシュ容量	8KB
L2 キャッシュ容量	512KB
L1 キャッシュラインサイズ	64B
L2 キャッシュラインサイズ	128B
メモリ種類	PC3200 (DDR SDRAM)
メモリバス本数	2
メモリ容量	2GB
OS	Linux 2.4.20-8
コンパイラ	gcc 3.2.2
最適化オプション	-O3

的には文献⁹⁾と同じクエリーを用いる。ただし Pentium[®]4を用いた場合に On キャッシュデータベースではなく主記憶データベースの評価とするために、タブル数は文献⁹⁾における評価対象のキャッシュサイズと Pentium[®]4のキャッシュサイズの比率である65倍に増やし、520KBものオンチップキャッシュにデータの大半が収まってしまうようにして評価を行なう。

まず Wisconsin ベンチマークで規定された表形式のデータをC言語で記述したプログラムで主記憶上に作成する。以下のクエリーをC言語で記述したものをOriginalとし、この実行時間を1とした際の加速率を測定する。本来の Wisconsin ベンチマークではtenkは10K個のタブルからなるが、本評価では650K個のタブルからなる。クエリーの番号は文献⁹⁾と合わせている。

```
(Q1) select * from tenk1 where (unique2 > 301) and
(unique2 < 402)
(Q2) select * from tenk1 where (unique1 > 647) and
(unique1 < 65648)
(Q3) select * from tenk1 where unique2 = 2001
(Q4) select * from tenk1 t1,tenk1 t2 where (t1.unique2 =
t2.unique2) and (t2.unique2 < 1000)
(Q6) select t1.*,o.* from onek o,tenk1 t1,tenk1 t2 where
(o.unique2 = t1.unique2) and (t1.unique2 = t2.unique2)
and (t1.unique2 > 1000) and (t2.unique2 > 1000)
(Q7) select MIN(unique2) from tenk1
```

このOriginalプログラムをベースに、提案メモリモジュール用の改造を行ったもので行った。ハードウェアで実行する部分を記述したエミュレーション版を作成し、その実行結果の一致による改造の妥当性をチェックする。

その後、ハード部の記述をしている関数の中身をコメントアウトすることで、ハード部の遅延がゼロになった状態の実行時間を再現して、ベクトルコマンド起動やフラグチェックなどのソフトウェアオーバーヘッドを含んだ性能を再現する。その際、本来ハードによって設定されるRead Window上のデータが更新されないため実行結果自体はOriginalとは異なったものになるが、その選択率や実行時間が狂わないように初期値を設定したRead Window上のダミーデータを用いて処理をすることで実行時間を再現する。

その上で、フラグのポーリング部の遅延時間としてホスト側に見えるハード部の遅延時間を Pentium[®]4に内蔵されるパフォーマンスカウンタを用いて挿入し、その変動により全体の実行時間がどのように変化するかを観測する。

さらに、実際の現場で用いられているデータベースのようにもっと大きなタブルサイズを有するテーブルに対する性能を評価するために、140バイトのタブルサイズを有するテーブルを作成し、同様の評価を行なう。

6.3 実験結果

6.3.1 ハードによる SO-DIMM からの読み出しを理想化した場合の性能

エミュレーション版のプログラムのハード部の記述をしている関数の中身をコメントアウトした際に観測される実行時間が、ハードによる SO-DIMM からの読み出しを理想化した場合の性能に対応する。その測定結果を図 4 に示す。各クエリーの Original での実行時間との相対値である加速率を示している。ここで、Read Window のサイズは 512 バイト (4 バイトデータ 128 個分) に固定し、Simple は Read Window を単に 1 枚用いた場合の加速率、Simple unroll は Simple にループアンローリングを施した場合の加速率、Pipelining は Read Window を 2 枚用いてソフトウェアパイプラインを行なった場合の加速率、Pipelining unroll は Pipelining にループアンローリングを施した場合の加速率、NOCLFLUSH はキャッシュラインをフラッシュする CLFLUSH 命令の影響を排除できた場合の加速率である。

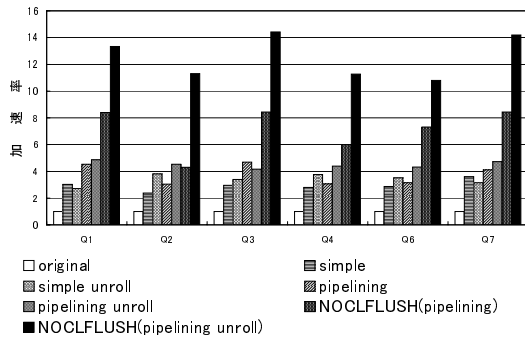


図 4 Wisconsin ベンチマークにおける提案システムの効果

この結果 Pipelining unroll の利用形態において 4.16 倍 ~ 4.87 倍の加速率が得られることがわかった。この結果は逐次性能向上における加速率としては高い加速率ではあるが、NOCLFLUSH(Pipelining unroll) が 11.27 倍 ~ 14.42 倍という高い加速率を示しているように、提案メモリモジュールの潜在能力的にはもっと高い加速率を期待できる。Wisconsin ベンチマークの場合は 4 バイトの有効なデータを取得するのに 60 バイトという 15 倍のメモリバンド幅を浪費することから、15 倍が提案メモリモジュールの加速率の限界であるが、今回の結果は 15 倍に迫る加速率が NOCLFLUSH(Pipelining unroll) で得られることは理に適っている。

このように今回の評価で加速率を制限した主な原因は CLFLUSH 命令の影響と考えられる。文献⁴⁾における評価でも測定に用いた Pentium[®]4 の CLFLUSH 命令の実行時間やその副作用はかなり大きく、実効性能を半減させてしまうことが判っており、今回の実験結果でも同様の傾向が観測された。CLFLUSH 命令がこのような影響を与えることは、Intel[®] 社内以外では詳細不明なブラックボックス状態にある Pentium[®]4 を Intel[®] 外部の人間がシミュレーションによって発見することは事実上困難であるが、本評価手法はこの問題点の発見を可能とした。

なお、SDT⁹⁾ ではほとんど加速できなかったマルチウェイの結合演算 (クエリー Q6) も、提案メモリモジュールでは高速化が達成されている。その理由は間接ベクトルロードが高速化されるためである。

6.3.2 プリフェッチにかかる遅延を変動させた場合の性能

SO-DIMM から Read window までの等間隔ベクトルロード実行にかかる時間を、1 回のポーリング時間を越える範囲で変化させた時の性能の変化を測定した。遅延時間の生成に

は Pentium[®]4 の性能カウンタを用いて正確な遅延時間を生成してフラグポーリングの部分に挿入した。その結果を図 5 に示す。本測定における対象のクエリーは Q7 である。

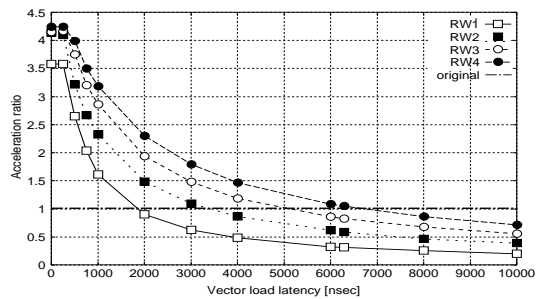


図 5 プリフェッチにかかる遅延を変動させた場合の Q7 の処理時間

この測定は 512B の Read window が 1 本の場合 (Simple unroll) と 2~4 本の場合 (Pipelining unroll) で上限性能が分離する。本数は多いほど等間隔ベクトルロードコマンド実行にかかる時間に対する耐性が強くなる。複数の本数を用いれば概ね 1 μ s 程度の遅延を隠ぺいできることがわかる。現状の DIMMnet-2 プロトタイプ上で SO-DIMM から Read window までの等間隔ベクトルロードコマンド実行にかかる時間は、512 バイトの Read window を一杯にする 128 ワードに対して、前述のように FPGA 版で 1.62 μ s、ASIC 版で 810ns である。つまり半分のクロック周波数で動作する FPGA 版では Pipelining unroll の利用形態をとっても性能低下が発生することがわかる。一方、ASIC 化して 200MHz 動作をさせれば Pipelining unroll の利用形態でほぼ隠ぺいできる。

6.3.3 キャッシュラインより大きなタプルサイズにおける性能

実際の現場で用いられている多くのデータベースのように Wisconsin ベンチマークより大きなタプルサイズを有するテーブルに対する問合せ性能を評価するために、キャッシュラインより大きな 140 バイトのタプルサイズを有するテーブルを作成し、同様の評価を行なった。その結果を図 6 に示す。

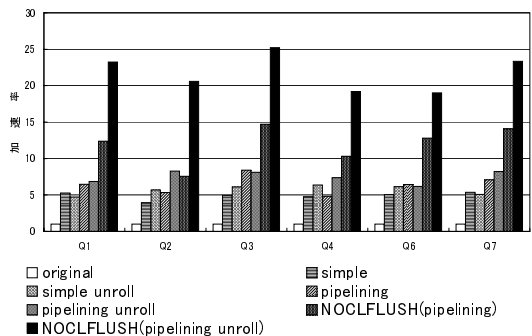


図 6 キャッシュラインより大きなタプルサイズにおける提案システムの効果

この結果から明らかなように、クエリー Q3 で 8.13 倍、CLFLUSH 命令の影響を排除すると 25.2 倍と、極めて高い加速率が得られることもわかった。その理由は、Pentium[®]4 を用いる場合は L2 キャッシュライン 128 バイト中に有効なデータが 4 バイトしか無くなり、Wisconsin ベンチマーク (60 バイト中 4 バイトが有効) よりもメモリバンド幅の浪費の程度が高まるため、Original の性能低下が起こるのに対し、提案メモリモジュールを用いた場合はそのような問題が発生せず、実行時間が Wisconsin ベンチマークの場合と比べてほとんど増加しないためである。

7. 関連研究

通常のキャッシュアーキテクチャベースのMPUのメモリアクセス部分の改良として、SCIMA⁸⁾や、SDT⁹⁾が提案されている。しかしこれらはMPUそのものを改造する必要があるため、Intel[®]やAMDなどのMPUベンダーが採用しない限りPCクラスタに利用できるCOTSとはならない。さらにSDTでは書き込みには対応していない。SCIMAでは等間隔アクセスのみ命令でサポートし、間接参照には対応していない。これに対し本方式は間接参照や書き込みにも対応しており、MPUやマザーボードに手を加える必要がないので、COTSのメリットを最大限享受できる。データベース応用を想定した場合は、SDTだけではマルチウェイの結合演算は加速率が上がらないことがわかっているが、本方式は間接参照もベクトル化されるためにそのようなクエリーの高速化にも対応できる。

メモリコントローラ(ノースブリッジ)の改良としてImpulse¹⁰⁾が提案されている。Impulseはエイリアス上の連続アクセスを等間隔ベクトルアクセスや間接ベクトルアクセスに変換する点で機能が類似している。しかし、ホストCPU上の仮想空間上に設けられる実体と同サイズのエイリアスへのアクセスとなるために、メモリコントローラ側でのアドレス変換が多段階でオーバーヘッドが増加し、かつ32bit CPUではその空間の狭さに縛られる。さらに、ホストCPUのFSB(Front Side Bus)に接続されるためホスト側のメモリコントローラを変更しない限り実装できない。このため、COTSのマザーボードを利用することはできず、PCクラスタ用技術としては適用できない。これに対し本方式はホストの仮想空間とは独立した巨大なページサイズ(DIMMnet-2では4GB)を有する一種のI/O空間に対するアクセスであるため、自ノードのSO-DIMMへのアクセスについては実質的にアドレス変換そのものが排除されているためベクトル型メモリアクセス部と相まって高速化が可能である。さらにノースブリッジやマザーボードに手を加える必要がないので、COTSのメリットを最大限享受できる。

8. おわりに

本報告では、キャッシュアーキテクチャの問題点を指摘し、その解決に向けて考案されたプリフェッチ機構付メモリモジュールの基本コンセプトとそのソフト的対応法について述べた。さらにそのDIMMnet-2における実装例を紹介した。

キャッシュアーキテクチャが苦手とする等間隔アクセスを多用する応用として主記憶データベースをとりあげ、プリフェッチ機能を有するメモリモジュール用に改造されたWisconsinベンチマークがPC上で示す性能から、実機上での処理性能を推定した。

その結果、4.16倍~4.87倍の加速率が得られることがわかった。この結果は逐次性能向上における加速率としては高い加速率ではあるが、CLFLUSH命令の影響を排除することで提案メモリモジュールは例えばクエリーQ3で14.42倍もの高い加速率が期待できることがわかった。さらに、SDT⁹⁾ではほとんど加速できなかったマルチウェイの結合演算(クエリーQ6)も等間隔ベクトルロードと間接ベクトルロードの組み合わせにより高速化が達成された。

Simple Scalar Tool Setなどの命令レベルシミュレータを使う従来手法では内部詳細構造が非公開なPentium[®]4をホストとして用いる場合のCLFLUSH命令のような特殊な命令のアプリケーション性能への影響を測定することは事実上困難であるが、本評価手法はPentium[®]4を対象として精密に結果に反映できている点で画期的な手法である。

また、DIMMnet-2プロトタイプ上で実現されている等間隔ベクトルロードコマンド実行時間を考慮しても、その遅

延時間は2枚のRead Windowを用いることでASIC版のDIMMnet-2では十分に隠ぺいできることもわかった。

さらに、実際の現場で用いられている多くのデータベースのようにWisconsinベンチマークより大きなダブルサイズを有するテーブルに対する問合せに対してはクエリーQ3で8.13倍、CLFLUSH命令の影響を排除すると25.2倍と、極めて高い加速率が得られることもわかった。

このような高い加速率がCOTSベースのPCに提案メモリモジュールを装着するだけで得られるということは、本方式はCLFLUSH命令問題の克服というさらなる改良の余地を残しつつも、文献⁹⁾における本手法はデータベースに向かないという指摘は適切ではなかったと考える。

DIMMnet-2の実機での評価および他のPC向けMPUにおけるCLFLUSH命令問題の確認とその克服、およびコンパイラによる対応は今後の課題である。

謝辞 本研究は総務省戦略的情報通信研究開発制度の一環として行われたものである。SDTに関して教授いただいた北陸先端科学技術大学院大学の田中助教に感謝いたします。DIMMnet-2の開発に関する議論にご参加いただいている慶應義塾大学の西講師、渡辺氏、大塚氏、伊豆氏、伊沢氏、宮部氏、東京農工大学の並木助教、浜田氏、三橋氏、荒木氏、木立氏、森氏、立命館大学の国枝教授、和歌山大学の齋藤講師、日立IT社の上嶋氏、今城氏、岩田氏に感謝いたします。

Trademarks: Pentium[®]はIntel[®] Corporationの登録商標です。本書に記載の商品の名称は、それぞれ各社が商標および登録商標として使用している場合があります。

参考文献

- 1) 萩原, 梅澤: “パーソナルスーパーコンピュータSX-6i”, 情報処理, Vol.44, No.3, pp.277-282 (Mar. 2003)
- 2) B.Lindsay: “DB2の父ブルース・リンゼイ博士が語る, データベース・テクノロジーの現在と未来”, <http://db2.jp/interview/special/lindsay.html>
- 3) 田邊, 濱田, 中條, 天野: “メモリスロット装着型ネットワークインタフェースDIMMnet-2の構想”, 情報処理学会計算機アーキテクチャ研究会, 2003-ARC-152, pp.61-66 (Mar. 2003)
- 4) 田邊, 中武, 箱崎, 土肥, 中條, 天野: “プリフェッチ機能付きメモリモジュールによる不連続アクセスの連続化”, 情報処理学会計算機アーキテクチャ研究会, 2004-ARC-157, pp.139-144 (Mar. 2004)
- 5) Myricom[®], Inc., <http://www.myri.com/>
- 6) Quadrics[®], Ltd., <http://www.quadrics.com/>
- 7) 田邊, 山本, 濱田, 中條, 工藤, 天野: “DIMMスロット搭載型ネットワークインタフェースDIMMnet-1とその高バンド幅通信機構BOTF”, 情報処理学会論文誌, Vol.43, No.4, pp.866-878 (Apr. 2002)
- 8) 近藤, 中村, 朴: “SCIMAにおける性能最適化手法の検討”, 情報処理学会論文誌ハイパフォーマンシングシステム, Vol.42 No.SIG12, pp.37-48 (Nov. 2001)
- 9) 宮崎, 府川, 田中: “ストライドデータアクセスによる主記憶データベースの問合せ処理の評価”, 情報処理学会データベースシステム研究会, 2004-DBS-134(II), pp.361-367 (Jul. 2004)
- 10) Carter, Hsieh, Stoller, Swanson, Zhang, Brunvand, Davis, Kuo, Kuramkote, Parker, Schaelicke and Tateyama: “Impulse: Building a Smarter Memory Controller”, International Symposium on High Performance Computer Architecture (HPCA-5), pp.70-79 (Jan. 1999)