

# 超並列計算機 CP-PACS における NPB Kernel CG の評価

板倉 憲一<sup>†</sup> 松原 正純<sup>†</sup> 朴 泰祐<sup>†</sup>  
中村 宏<sup>††</sup> 中澤 喜三郎<sup>†††</sup>

本論文では NAS 並列ベンチマークの kernel CG を用いて、超並列計算機 CP-PACS の性能評価および解析を行い、CP-PACS における並列プログラムの最適化について述べる。Kernel CG では、広範囲で複雑なデータ転送や短ベクトル処理が要求されるために並列化効率を上げることが難しいが、CP-PACS では柔軟なハイパクロスバ網と擬似ベクトル処理機構によって最適化を行うことが可能である。評価に用いた並列プログラムは基本的な通信性能と CPU 性能に基づいて FORTRAN およびアセンブラ・ソースでのチューニングを行い、CPU clock counter による精密な実行時間測定を基に性能解析を行った。解析結果から通信時間と CPU 処理時間の最適化にはトレードオフが生じるために、そのチューニングが必要であることが分かった。また、このような最適化を行うことで、他の並列計算機に対して CP-PACS が高い台数効果と絶対性能を持つことを示すことができた。

## Performance Evaluation of NPB Kernel CG on CP-PACS

KEN'ICHI ITAKURA,<sup>†</sup> MASAZUMI MATSUBARA,<sup>†</sup> TAISUKE BOKU,<sup>†</sup>  
HIROSHI NAKAMURA<sup>††</sup> and KISABURO NAKAZAWA<sup>†††</sup>

We evaluate the performance massively parallel processor CP-PACS on NAS Parallel Benchmarks Kernel CG. Since Kernel CG requires global data transfer and short vector processing, it is hard to achieve performance in proportion to the number of PU's. However, the pseudo vector processing mechanism implemented in each node processor can effectively handle short vector calculations. Moreover, Hyper-Crossbar Network, the interconnection network of CP-PACS, can provide high-performance global data communication. Owing to these characteristics, CP-PACS can achieve high absolute performance and speed-up ratio when increasing the number of PU's.

### 1. はじめに

超並列計算機 CP-PACS<sup>1),2)</sup>は、主に計算物理学の諸問題を対象としているが、そのアーキテクチャは様々な数値処理を考慮し設計されている。単体ノードプロセッサは擬似ベクトル処理機構<sup>3),4)</sup>を採用し、短ベクトルから大規模ベクトル処理まで十分な実効性能が引き出せる。さらに、3次元ハイパ・クロスバ<sup>5),6)</sup>によるノード間結合網は、数千台規模での実装容易性とともに高い柔軟性を持っており、ユーザの要求する様々なデータ転送パターンに対して高いスループットを提

供することできる。

本論文では、NAS Parallel Benchmarks (以下 NPB) version 1<sup>7)</sup>のうち Kernel CG ベンチマークによって、超並列計算機 CP-PACS の性能評価を行う。NPB は非常に多くの並列計算機において評価がなされており、性能評価の指標として有効である。NPB は、version 2<sup>8)</sup>から MPI によって記述されたベンチマークプログラムでの評価が行われているが、プログラムの実装をユーザに任せられたベンチマークの評価も version 1 として引き続き行われている。

Kernel CG では大規模疎行列とベクトルの積の計算が処理時間の大部分を占め、この時間をいかに短縮するかが高速化のポイントとなる。一般に超並列計算機のプログラム・チューニングでは、いかにデータ転送時間を削減するかが重要である。この点において、隣接転送だけでなく広範囲な転送に対しても高いスループットを提供する CP-PACS は、ユーザに自由な問題マッピングと並列化手法を提供し、データ転送時間を

<sup>†</sup> 筑波大学電子・情報工学系

Institute of Information Sciences and Electronics, University of Tsukuba

<sup>††</sup> 東京大学先端科学技術研究センター

Research Center for Advanced Science and Technology, University of Tokyo

<sup>†††</sup> 明星大学情報学部

Faculty of Computer Science, Meisei University

削減するために様々なアルゴリズムを試みることを可能にする。しかし、このような並列化が内部処理時間に影響を与えることも考慮に入れる必要がある。なぜなら、擬似ベクトル処理機構を持つ CP-PACS のノードプロセッサでは、内部処理時間はベクトル長に依存するからである。これからのハイパフォーマンス・コンピューティングの分野においては、ベクトル処理機構を並列計算機に組み込むことは必須といえ、ノードプロセッサでの処理効率を落とさないように、アルゴリズムを検討する必要がある。

本論文では単にベンチマークプログラムの時間計測を行うだけではなく、処理時間をクロックカウンタを用いて細かく分解し、その解析を行う。特に主なデータ転送によって起きる並列処理のオーバヘッドの時間は CP-PACS の諸元から予想し、実測との比較を行うことでプログラムの想定した最適化が実測に反映していることを示す。また、擬似ベクトル処理機構を持つプロセッサの短ベクトル処理に対する特性を明らかにし、その有効性についても考察する。

## 2. 超並列計算機 CP-PACS

超並列計算機 CP-PACS<sup>1),2)</sup> (図1) は広い範囲の数値計算処理を効率良く実行できるように、その設計・実装段階において様々なアーキテクチャ上の特徴を採り入れている。現在の CP-PACS は計算用ノード (PU) が 2048 台、ディスクなど I/O 装置を接続するためのノード (IOU) が 128 台あり、それらは  $8 \times 17 \times 16$  の 3 次元ハイパ・クロスバ結合網<sup>5),6)</sup> (以下 HXB) によって接続されている。各 PU には 64 MB のメインメモリ、16 KB(I)/16 KB(D) の 1 次キャッシュメモリ、512 KB(I)/512 KB(D) の 2 次キャッシュメモリがある。

3 次元 HXB と PU は NIA (Network Interface Adapter) によって接続されており、そのスループットは 300 MByte/sec である。3 次元 HXB は図 1 に示すように 3 次元直行空間上に PU を並べ、その間をクロスバスイッチ (XB) で結合した間接網である。各次元方向の PU は単一の XB によって完全結合され、それ以外の PU 間のデータ転送はエクスチェンジャ (EX) と呼ばれる小型のクロスバスイッチによって複数次元方向の XB を経由して行う。CP-PACS ではメッセージ転送に wormhole ルーティングを用い、経路の選択には固定ルーティングを用いている。

3 次元 HXB は mesh/torus ネットワークはもちろん HyperCube ネットワークも部分的にエミュレートすることが可能であり、パタフライ転送のような広域

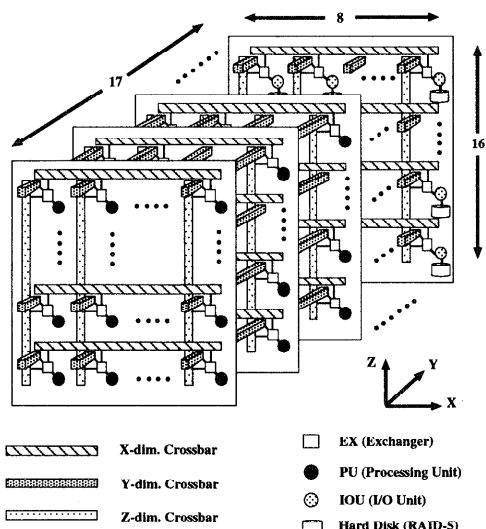


図1 CP-PACS の構成図

Fig.1 Overview of CP-PACS architecture.

な転送も無衝突で行える<sup>9)</sup>。このため、これらの相互結合網を前提にチューニングされた、従来の並列アルゴリズムを HXB に適用することも可能である。

さらに、CP-PACS にはリモート DMA と呼ばれる高速な PU 間データ転送方式がある。これは、PU の NIA がユーザメモリ空間のデータを直接 DMA アクセスして送受信を行うものである。特に受信側では受信命令を発行する必要はなく、メッセージ到着によって自動的に NIA が動作し、ゼロ・コピーでユーザメモリ空間にデータが書き込まれる。つまり OS 領域のバッファを使用しないので、ハードウェアの高速なスループットをそのまま活かすことができる。また、転送の立ち上げオーバヘッドを軽減するために、NIA の起動を OS を介さずにユーザが直接行えるようになっている。

次に PU の特徴について述べる。PU には、既存の RISC プロセッサ (HP PA-RISC 1.1) をベースに擬似ベクトル機構 (PVP-SW: Pseudo Vector Processor based on Slide Window)<sup>3),4)</sup> を付加したノードプロセッサを使用している。超並列計算機は非常に多数のノードプロセッサを用いるため、その量的な制約から 1 チップで構成される RISC マイクロプロセッサを使用することが妥当である。しかし、既存の RISC プロセッサの性能はキャッシュメモリのヒット率に強く依存している。そのためワーキングセットがキャッシュ容量より大きい問題では、長い主記憶アクセス・レイテンシのために性能は著しく低下する。たとえば、多重ループによって使用されるデータ領域をブロック化

し、アクセスパターンを変更することによって、局所的なワーキングセットを限定し、キャッシュヒット率を向上させる手法も提案されている。しかし、このブロック化の最適なサイズはキャッシュメモリのサイズや構成に大きく影響されるため<sup>10)</sup>、注意深い最適化が必要であり、プログラムの可搬性は悪くなる。

これに対し PVP-SW は、キャッシュメモリを介さずに、浮動小数点データをメインメモリから直接レジスタに読み込む preload 命令と、レジスタの内容を直接メインメモリに書き戻す poststore 命令を備える。そして、preload 命令と poststore 命令は non-blocking operation とし、後続命令が load/store 命令であってもパイプライン的に処理される。これによって主記憶アクセス・レイテンシを隠蔽する。さらに、この non-blocking の load/store 命令を有効に動作させるために、メインメモリは 16 バンク構成の擬似パイプライン・メモリを用いる。

具体的には、演算に必要なデータに対する preload 命令を十分前もって先行発行させることで主記憶アクセス・レイテンシを隠蔽する。しかし、preload 命令と演算命令の間の距離（両命令の間にスケジューリング可能な命令数）はレジスタ数によって制限される。通常の RISC プロセッサのたかだか 32 本の浮動小数点レジスタでは主記憶アクセス・レイテンシを隠蔽するのは難しい。そこで、CP-PACS の PVP-SW では浮動小数点レジスタ数を 128 本に拡張することで長い主記憶アクセス・レイテンシを隠蔽する。また、preload 命令と poststore 命令はすべてのレジスタを指定できるが、命令アーキテクチャレベルの上位互換性を保つために、通常の命令はレジスタ・ウィンドウを介して、拡張したレジスタをアクセスする。

擬似ベクトル処理はループ処理におけるベクトル演算を構成する命令を、スーパスカラ方式によって順次並列発行することによって実現される。ループの 1 イタレーション内では実際の演算と、将来必要とされるデータを先のウィンドウのレジスタへ preload する処理が並行して行われる。さらに、スーパスカラ方式によって preload/poststore 命令と浮動小数点演算命令の 2 命令同時発行が可能であり、ベクトル処理を効率的に行う。CPU は動作周波数が 150 MHz であり、1 命令で浮動小数点の multiply と add を行う複合命令によって理論ピーク性能は 300 MFLOP/s である。

現在の CP-PACS は 2048 台の PU を 1024 台、512 台、256 台×2 の 4 パーティションに分割しており、今回の性能測定には 256 台 (4×8×8) のパーティションを使用した。

### 3. NPB Kernel CG の並列化

本章では NPB version 1 Kernel CG<sup>7)</sup> のアルゴリズムを示し、CP-PACS での並列化について述べる。

Kernel CG は正値対称な大規模疎行列の最小固有値を CG (Conjugate Gradient) 法によって求めるベンチマークである。Class B の問題サイズでは疎行列のサイズは 75000×75000、非零要素率は 0.24% であり 1 行に平均 183 個ある。以下に逐次版 Kernel CG のメインループを示す。

```
do k=1,KERNITR
  call matvec(a, colidx, rowstr, p, q)
  alpha=rho/dotpro(p, q)
  z(1:N)=z(1:N)+alpha*p(1:N)
  rho0=rho
  r(1:N)=r(1:N)-alpha*q(1:N)
  rho=dotpro(r, r)
  beta=rho/rho0
  p(1:N)=r(1:N)+beta*p(1:N)
enddo
```

ここで、 $p$ ,  $q$ ,  $z$ ,  $r$  は  $N$  ( $=75000$ ) 要素のベクトルであり、matvec と dotpro はそれぞれ行列・ベクトル積とベクトルの内積を計算するサブルーチンである。内積計算 (dotpro) といわゆる DAXPY の計算は 3 次元 PU 空間を仮想的に展開した 1 次元 PU 空間に各ベクトルをブロック分割し、データパラレル処理によって行う。これにより、すべての PU がデータ転送を行うことなしに並列処理を行える。なお、内積計算では PU 内の計算が終了した後に Butterfly Summation<sup>11)</sup> を行い、全 PU が同じ計算結果を保持する。

メインループにおいて最も処理時間のかかる部分は行列・ベクトルの積 (matvec) である。matvec の処理は仮想的な 2 次元 PU 平面 ( $P = P_x \times P_y$ ,  $P$  は PU 台数) に行列をブロック・ブロック分割する。これは、単に仮想 1 次元 PU に対して行ブロックマッピング ( $P_x = 1$ ) や列ブロックマッピング ( $P_y = 1$ ) を行った場合も包含しており、この  $P_x$ ,  $P_y$  の組合せによりデータ転送時間の最適化が行えると考えられる<sup>12)</sup>。

なお、各 PU では部分疎行列を値 (a) と桁位置 (colidx) と各行のスタートポイント (rowstr) によって表現する。以下に matvec の並列処理を示す。図 2 は例として 8 PU ( $P_x = 2$ ,  $P_y = 4$ ) による matvec の処理過程におけるデータ転送および演算を示している。図中、斜線部は 5 番の PU が処理するデータを表している。なお、行列  $A$  は簡単のため密行列表現されているが、実際には疎行列である。

(1) Collection: 被乗数ベクトル  $p$  に関し、各 PU

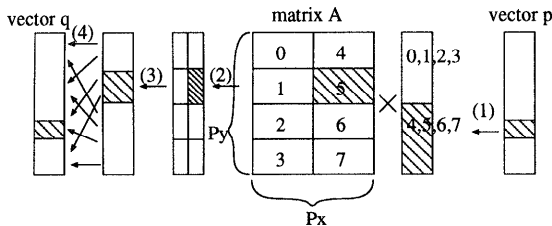


図2 matvec の処理とデータ転送パターン

Fig. 2 Data processing and transfer pattern in "matvec".

はその保持する部分行列の列に対応する  $N/P_x$  要素を必要とする。しかし、初期状態として  $p$  は  $N/P$  に分割されているために  $P_y$  台の PU からなる  $P_x$  個の各グループ内で  $p$  の部分ベクトルを集め、 $N/P_x$  要素のベクトルを作る必要がある。ここでは、この処理をベクトルの Collection と呼ぶ。ベクトルの Collection の方法としては  $P_y$  回のマルチキャスト、 $P_y - 1$  回の「バケツリレー」的な転送などがあるが、広範囲な転送が得意な HXB では各 PU の転送回数が少ない butterfly 転送によって行うのが効率的である。転送パターンは Butterfly Summation と同じであるが、Butterfly Summation がリダクション演算であるために毎回のデータ転送量が一定なのに対して、Butterfly Collection ではデータ転送量が 2 倍ずつ増加する点が異なる。

- (2) **Multiply**: 次に部分行列と部分ベクトルの乗算を行う。この乗算結果はベクトル  $q$  に保持されるが、 $q$  の 1 要素の計算は以下のようなリストベクトル処理になる。

$$q = q + a[i] \times p[\text{colidx}[i]]$$

浮動小数点データ ( $a, p$ ) のロードはレジスタへの preload 命令を使うコードがコンパイラによって生成される。これに対して  $\text{colidx}$  は整数データなので、そのアクセスは通常のロード命令によって行われるが、整数配列の連続領域に対するアクセスなので 1 次キャッシュの高ヒット率が期待できる。

- (3) **Summation**: Multiply で求めた結果は部分和なので、真の乗算結果を得るために  $P_x$  台の PU からなる  $P_y$  個の各グループ内で Vector Summation を行う。この処理は dotpro の Butterfly Summation と同等のリダクションをベクトルに対して行う。
- (4) **Shuffle**: Summation で求めた部分ベクトル  $q$  はこの後で dotpro や DAXPY の計算に用いられるが、各 PU が保持している部分は他のベクト

表1 matvec のデータ転送  
Table 1 Data transfer in "matvec".

処理パターン	転送回数	総転送量
Collection	$\log_2 P_y$	$N/P \times (P_y - 1)$
Summation	$\log_2 P_x$	$N/P_x \times \log_2 P_x$
Shuffle	1	$N/P$

ルのマッピングと異なってしまっている。そこで、各 PU が保持する部分ベクトル  $q$  の位置を合わせるために Shuffle 転送を 1 回行う。

#### 4. 実行時間予測

matvec の処理時間には用いる 3 次元 PU 空間をどのような仮想 2 次元 PU 平面 ( $P_x \times P_y$ ) に展開するかによって決まる。まず、matvec で行われるデータ転送の回数とデータ転送量を表 1 に示す。定性的には  $P_x$  の増加は Summation 時間の増加を招き、 $P_y$  の増加は Collection 時間の増加を招く。簡単なピンポン転送の実験からデータ転送の立ち上げオーバーヘッドは約 5500 clock<sup>\*</sup>、転送スループットは 2 Byte/clock という基本通信性能が分かっているので全転送時間は以下のようなになる。

$$(\log_2 P_x + \log_2 P_y + 1) \times 5500 + N/P \times (P_y + P_y \log_2 P_x)/2$$

次に、CPU の計算時間を求める。計算時間は Multiply の時間と Summation のベクトル和の時間からなる。疎行列の非零要素の分布がほぼ一様なので 2 次元 PU 平面の構成に関係なく Multiply の演算数は一定であると考えられるが、部分行列とベクトルの積は行ごとに行われ、その要素数は  $P_x$  の増加によって減少する。これは PVP-SW のようなベクトル処理機構にとって不利な条件であり、 $P_x$  の増加にともなって演算時間は増加すると考えられる。しかし、ここではこの増加量の見積りが難しいために、十分なベクトル長があるものと判断し、Multiply の演算時間は一定であると予測しておく。

Summation 時のベクトル和の演算量は  $N/P_y$  なので  $P_x$  の増加 (=  $P_y$  の減少) は演算時間の増加を招く。ベクトル和の部分のみを実測した結果、1 要素あたり 7 clock で演算できることが分かったので、Multiply の処理時間を固定とすると Multiply を除く matvec の処理時間は図 3 のように予想され、PU 台数が 256 と 128 のときは  $P_x = 4$ 、64 以下では  $P_x = 2$  のときが最も高速に処理できると考えられる。

\* プログラミングを容易にするために最速の送信命令は使っていない。

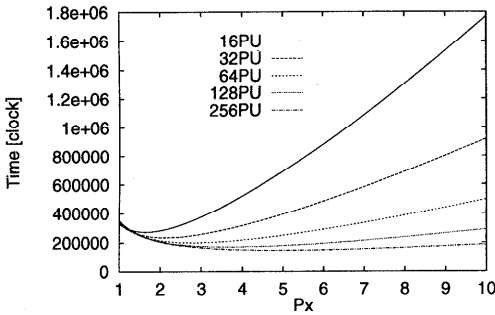


図3 Multiply を除く matvec の予想処理時間

Fig. 3 Estimation time in “matvec” except Multiply.

表2 Kernel CG の処理時間 [sec]  
Table 2 Result of Kernel CG [sec].

#PU	$P_x = 1$	$P_x = 2$	$P_x = 4$	$P_x = 8$
16	106.21	126.41	203.43	252.83
32	55.75	66.18	105.24	129.66
64	30.89	35.07	55.34	69.11
128	19.66	20.54	30.44	39.00
256	13.19	13.00	17.66	19.73

表3 メインループの処理時間の内訳 (256 PU) [clock]

Table 3 Breakdown of elapse time for the main loop (256 PU) [clock].

処理	$P_x = 1$	$P_x = 2$	$P_x = 4$	$P_x = 8$
matvec	804846	797683	1143634	1290212
dotpro	55478	59084	61063	59976
DAXPY	3115	3130	3207	3059

表4 matvec の処理時間の内訳 (256 PU) [clock]

Table 4 Breakdown of elapse time for “matvec” (256 PU) [clock].

処理	$P_x = 1$	$P_x = 2$	$P_x = 4$	$P_x = 8$
Collection	444295	253017	148965	86151
Multiply	351274	490649	873599	1091424
Summation	94	15115	40072	87841
Shuffle	9757	6914	6874	6749

## 5. 性能評価

### 5.1 データ転送時間の最適化

PU 台数が 16 から 256 の場合において、 $P_x$  を 1 から 8 まで変えて Kernel CG の実行時間を測定した結果を表 2 に示す。この結果では、ほぼすべての PU 台数で  $P_x = 1$  のときが最短の処理時間を示し、予想とは異なる結果となった。この原因を究明するために、PU が 256 台の場合について、メインループの処理時間の内訳と matvec の処理時間の内訳を測定した。これらを表 3、表 4 に示す。さらに、データ転送パター

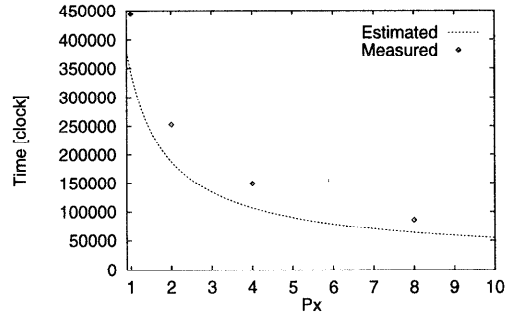


図4 Collection の処理時間 (256 PU)

Fig. 4 Elapse time of “Collection” (256 PU).

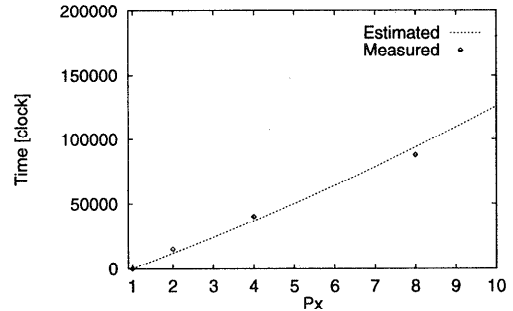


図5 Summation の処理時間 (256 PU)

Fig. 5 Elapse time of “Summation” (256 PU).

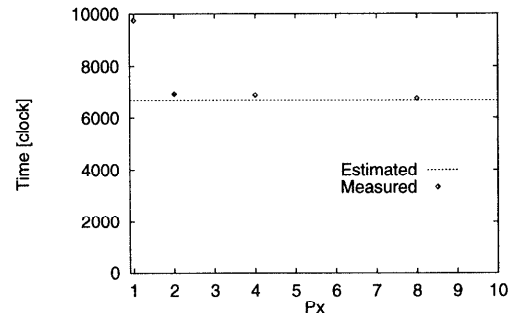


図6 Shuffle の処理時間 (256 PU)

Fig. 6 Elapse time of “Shuffle” (256 PU).

ンごとの予測値と実測値を図 4、図 5、図 6 に示す。

これらの測定には、CPU に内蔵されている、clock に同期してインクリメントされる特別なコントロールレジスタ (すなわちクロックカウンタ) を用いた。これにより、ほぼ clock 周期 (= 6.7 ns) の精度できわめて正確な処理時間の測定が可能である。

表 3 から分かるように、最も処理時間のかかる部分は、予測したとおり matvec である。matvec の処理時間の内訳を表 4 から見ると、Multiply の処理時間は  $P_x$  にともなって増加していることが分かる。これ

は、前章での予測と大きく異なる点であり、表2においてほぼすべてのPU台数の場合に $P_x = 1$ が最適となってしまう原因でもある。一方、matvecにおける各転送パターンごとの予測時間と実測時間を比較すると(図4, 図5, 図6), 予測値と実測値が良い精度で合っていることが分かる。つまり、行列のマッピングによってCollectionとSummationの転送時間がトレードオフしている。ここで、Shuffle転送に関しては多次元のHXBネットワークを用いた場合にネットワーク中で衝突が起こることが分かっているが、予測値は無衝突を仮定しているので若干の誤差が生じている。しかし、Shuffle転送は他の転送に比べて転送回数や転送時間が小さいために、この誤差は全体の予測には影響を及ぼさないと見える。

表4において、Multiplyの処理時間が予想に反して $P_x$ にもなって増加している点に関しては次節で述べる。

## 5.2 CPU 処理時間の最適化

Multiplyの処理時間はPU内での純粋な計算時間のみであり、「演算量が一定なので計算時間が2次元PU構成によらず一定である」と仮定したが、CP-PACSのPUはPVP-SW機構によりベクトルプロセッサと同じような特性を持ち、ループ長が短いループ処理では性能が低下する。コンパイラの生成する擬似ベクトル処理のコードは、一定の長さのベクトル長(CVL: Critical Vector Length)を必要とし、実際のベクトル長がCVL以下であればスカラ処理、CVL以上であれば擬似ベクトル処理を行うようなコードが生成されている。

現在のコンパイラでは、メモリのバンク・コンフリクトなどによる主記憶アクセス・レイテンシの増加に対しても可能な限り対応するために、使用可能な浮動小数点レジスタ(約100本程度)を最大限に利用し、preload命令と実際の演算命令の間に他のデータのpreloadを挿入することによって、主記憶アクセス・レイテンシを極力隠蔽するような擬似ベクトルコードを生成する。このために、一般にCVLは大きくなりMultiplyのコードではCVL=50であった。これに対して、今回の測定に用いたプログラムでは部分行列の1行ごとに積を計算するので、ループ長は1行の非零要素数になる。Kernel CG(Class B)では1行の平均非零要素数は183個であり、平均ベクトル長は $183/P_x$ となる。つまり、コンパイラの生成したCVL=50のコードでは、 $P_x$ の増加にもなってベクトル長が短くなり、通常のload命令によるスカラコードでMultiplyの処理が行われる場合が増加する

ためにMultiplyの処理時間が増加してしまう。

CVL=50では、実際の演算が始まるまでに50要素分のpreloadを必要とし、この時間はベクトル処理でのスタートアップ時間に相当する。しかし、PVP-SWはスーパスカラ処理のプロセッサにデータ供給能力を強化して擬似的にベクトル処理を行うので、このスタートアップ時間を短くしたコードを作成することも可能であり、本質的には短ベクトル処理に対応できる。そこで、我々はMultiplyの部分においてのみアセンブラソースレベルのチューニングを行い、CVL=20のプログラムを作成して評価を行った。このようなアセンブラレベルのチューニングはNPBのルールでは禁止されているが、コンパイラへのディレクティブなどを用いてCVLをFORTRANソースレベルでチューニングすることは容易であり、CP-PACSの本質的な性能評価を行う目的で有効であると考えられる。

表5は作成したCVL=20のアセンブラコードによるMultiplyの時間を示している。 $P_x = 1$ のときはほとんどのベクトル長が50を超えているためにMultiplyの処理時間はCVL=50のコード(表4)とほぼ同じであるが、 $P_x = 2$ では35%、 $P_x = 4$ では39%短縮される。これは、 $P_x$ が大きくなるのにもなって短ベクトルの割合が増えるためである。 $P_x = 8$ では29%の短縮にとどまっているが、これはCVL=20のコードでもスカラ処理になってしまう20以下のベクトル長の割合が増えるためである。このMultiplyを用いた全体の処理時間を表6, 図7に示す。この結果は予想転送時間(図3)に単調増加するMultiplyの処理時間を足した振舞いをし、PU台数が大きいときには $P_x$ にもなうMultiplyの時間増加分が通信時間の変化とほぼ同じオーダになるので、 $P_x$ によって処理時間がトレードオフする。

以上のように、通信時間の最適化とPU内部処理時

表5 matvecの処理時間の内訳(256PU, CVL=20) [clock]  
Table 5 Breakdown of elapse time for "matvec" (256PU, CVL=20) [clock].

処理	$P_x = 1$	$P_x = 2$	$P_x = 4$	$P_x = 8$
Multiply	357818	419084	531332	769355

表6 Kernel CGの処理時間(256PU, CVL=20) [sec]  
Table 6 Result of Kernel CG (256PU, CVL=20) [sec].

#PU	$P_x = 1$	$P_x = 2$	$P_x = 4$	$P_x = 8$
16	101.25	110.22	138.96	198.46
32	54.51	58.98	71.04	100.95
64	30.30	31.09	36.93	52.54
128	18.38	17.59	19.92	27.90
256	12.94	11.38	11.86	15.41

表 7 他の並列計算機の Kernel CG の処理時間 [sec]  
Table 7 Results of Kernel CG on other machines [sec].

#PU	16	32	64	128	256
CP-PACS	101.25	54.51	30.30	17.59	11.38
CRAY T3E (Nov 96)	107.1	59.3	33.8	22.9	22.2
IBM RS/6000 SP Wide-node 1 (67 MHz) (Mar 94)	88.4	52.53	33.79	25.44	-
IBM RS/6000 SP P2SC node (120 MHz) (Nov 96)	62.21	36.22	22.71	-	-

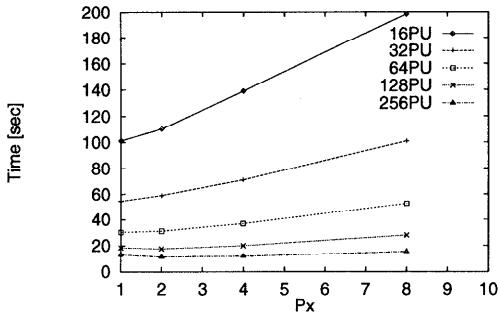


図 7 Kernel CG の処理時間 (CVL=20)  
Fig. 7 Result of Kernel CG (CVL=20).

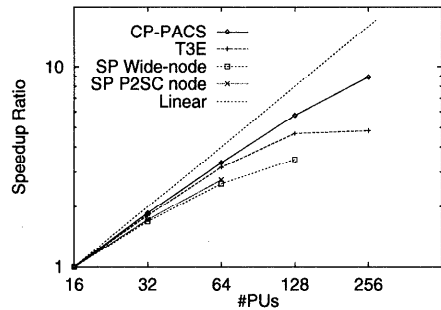


図 8 速度向上率  
Fig. 8 Speed up ratio.

間の最適化がトレードオフする場合、ここに示したような注意深いチューニングが必要であることが分かった。この種の問題は、ハイパフォーマンス・コンピューティングにおいて大きな位置を占める、ベクトル処理機能を持つ PU を用いた並列計算機において、今後重要になっていくものと思われる。

## 6. 他の並列計算機との比較

最後に NPB version 1 の結果<sup>13)</sup>を基に、Cray T3E, IBM RS/6000 SP system と CP-PACS を比較する。各並列計算機の結果を表 7 に、各々の 16 PU の処理時間を 1 とした速度向上率を図 8 に示す。

Cray T3E と比較するとすべての PU 台数において処理時間が短く、特に PU 台数が増えたときに CP-PACS は台数効果を維持していることが分かる。また、RS/6000 SP P2SC に対しては、ノードプロセッサの peak 性能差 (CP-PACS: 300 MFLOP/s, SP P2SC: 480 MFLOP/s) から絶対時間は CP-PACS が劣る。しかし、SP P2SC では 64 PU の段階で性能が頭打ちになり始めているのに対して、CP-PACS は台数効果を維持していることが分かる。T3E や SP P2SC での測定に用いられた並列アルゴリズムは公開されていないが、PU 台数の増加によってデータ転送の粒度が細かくなることが予想される。このような細粒度のデータ転送に対して、CP-PACS のリモート DMA による

低オーバーヘッドの通信が有効に働き、CP-PACS では台数効果が維持できているといえる。

## 7. おわりに

本論文では、NPB kernel CG (Class B) によって超並列計算機 CP-PACS の性能評価およびその解析を行った。CPU に組み込まれているいわゆるクロックカウンタを用いることにより、様々な処理時間を細かく測定・解析することができた。解析の結果、データ転送に関しては、基本的な諸元 (転送立ち上げオーバーヘッドとスループット) から予測される時間と実測がほぼ一致することを確認した。このことから、多くの転送パターンを無衝突で扱える HXB ではデータ転送時間の正確な見積りが可能であることが分かった。

Kernel CG は様々なデータ転送パターンを必要とし、その多くが細かいデータ転送となるために並列化効率が悪い。この問題を解決するために行列を 2 次元 PU 平面にブロックブロック分割しデータ転送を最適化する手法を用いた。しかし、このマッピングが PVP-SW による擬似ベクトル処理のベクトル長を短くしてしまい、内部処理時間の増大を招く結果となった。現在のコンパイラはバンクコンフリクトなどによる長い主記憶アクセス・レイテンシにできるだけ対処するようなコードを作成するが、逆にそのようなコードが擬似ベクトル処理をする最短ベクトル長を大きく

してしまう。今回の評価対象である Kernel CG のような短ベクトル処理を多く行う問題では擬似ベクトル処理機構を有効に活用できないことが分かった。これに対して、擬似ベクトル処理の最短ベクトル長 (CVL) を小さくすると、短ベクトル処理においても擬似ベクトル処理が効果的であることが分かった。

一般に超並列計算機では通信時間を減らすためにアルゴリズムの選択や細かいチューニングが行われるが、これに対して HXB は広い選択範囲を与える。さらにハイパフォーマンス・コンピューティングの分野ではベクトル処理は必須であり、この性能は処理のベクトル長に依存する。これに対しては短ベクトルに対しても PVP-SW が小さいオーバーヘッドで処理可能であることを示した。この種のチューニングの問題は、今後の並列ベクトル機において重要になっていくと考えられる。

結果として擬似ベクトル処理を最適化することで、他の並列計算機と比較して CP-PACS は Kernel CG (Class B) 問題に対して高い絶対処理性能と台数効果を持つことを示した。今後は、他の NPB Kernel やアプリケーションに対する評価を行い、CP-PACS の有効性を示したい。

謝辞 超並列計算機 CP-PACS を利用する機会を与えていただいた筑波大学計算物理学研究センター関係各位に感謝します。なお、本研究の一部は創成的基礎研究費 (08NP0401) の補助によるものである。

### 参 考 文 献

- 1) 中澤喜三郎ほか：CP-PACS のアーキテクチャの概要，情報処理学会研究報告，ARC-108-9 (1994)。
- 2) Iwasaki, Y.: Status of the CP-PACS Project, *Proc. Lattice '96, Nucl. Phys. B (Proc. Suppl.)* (1996)。
- 3) Nakamura, H., et al.: A Scalar Architecture for Pseudo Vector Processing based on Slide-Windowed Registers, *Proc. ACM International Conference on Supercomputing '93*, pp.298-307 (July 1993)。
- 4) 位守弘充ほか：スライドウィンドウ方式による擬似ベクトルプロセッサ，情報処理学会論文誌，Vol.34, No.12, pp.2612-2613 (1993)。
- 5) 田中輝雄ほか：識別子を用いたデータ転送方式を基本とする MIMD 型並列計算機アーキテクチャ，並列処理シンポジウム JSPP '89 (1989)。
- 6) 朴 泰祐ほか：ハイパクロスバ・ネットワークの性能評価，信学技報，CPSY-93-40 (1993)。
- 7) Bailey, D., et al.: The NAS Parallel Benchmarks, Technical Report, RNR-94-007, NAS (1994)。

- 8) Bailey, D., et al.: The NAS Parallel Benchmarks 2.0, Technical Report, NAS-95-020, NAS (1995)。
- 9) 松原正純ほか：超並列計算機 CP-PACS のネットワーク性能評価，情報処理学会研究報告，HPC-67-10, pp.55-60 (1997)。
- 10) Lam, M., et al.: The cache performance and optimization of blocked algorithms, *Proc. ASPLOS-IV* (Apr. 1991)。
- 11) 板倉憲一ほか：並列計算機 CP-PACS の CG ベンチマークによる性能評価，情報処理学会研究報告，HPC-63-6, pp.31-36 (1996)。
- 12) Agarwal, R.C., et al.: High-performance parallel implementations of the NAS kernel benchmarks on the IBM SP2, *IBM Systems Journal*, Vol.34, No.2, pp.263-272 (1995)。
- 13) Saim, S., et al.: Parallel Benchmark (Version 1.0) Results 11-96, Technical Report, NAS-96-018, NAS (1996)。

(平成 9 年 11 月 7 日受付)

(平成 10 年 4 月 3 日採録)



板倉 憲一 (学生会員)

昭和 44 年生。平成 5 年筑波大学第三学群情報学類卒業。平成 7 年同大学院工学研究科前期博士課程修了。現在同研究科後期博士課程に在籍中。並列計算機アーキテクチャの研究に従事し、性能評価方法等に興味を持つ。



松原 正純 (学生会員)

昭和 48 年生。平成 8 年筑波大学第三学群情報学類卒業。平成 10 年同大学院工学研究科前期博士課程修了。現在同研究科後期博士課程に在籍中。並列計算機の研究に従事。



朴 泰祐 (正会員)

昭和 59 年慶應義塾大学工学部電気工学科卒業。平成 2 年同大学院理工学研究科電気工学専攻後期博士課程修了。工学博士。昭和 63 年同大学院工学部物理学科助手。平成 4 年筑波大学電子・情報工学系講師，平成 7 年同助教，現在に至る。超並列計算機のアーキテクチャおよび性能評価，ハイパフォーマンス・コンピューティングへの応用の研究に従事。電子情報通信学会会員。





中村 宏 (正会員)

昭和 60 年東京大学工学部電子工学科卒業。平成 2 年同大学院工学系研究科電気工学専攻博士課程修了。工学博士。同年筑波大学電子・情報工学系助手。同講師、同助教授を経て、平成 8 年より東京大学先端科学技術研究センター助教授。計算機アーキテクチャ、ハイパフォーマンスコンピューティング、計算機の上位レベル設計支援の研究に従事。本会平成 5 年度論文賞、平成 6 年度山下記念研究賞各受賞。電子情報通信学会、IEEE、ACM 各会員。



中澤喜三郎 (正会員)

昭和 30 年東京大学工学部応用物理卒業。昭和 35 年同大学院数物系博士課程応用物理修了。同年日立製作所入社。TAC, HITAC 5020, E/F, 8800/8700, M-200 H/280 H, 680 H, S-810 等、超大型コンピュータ・スーパーコンピュータの開発に従事。平成元年より筑波大学電子・情報工学系教授、計算物理学センター向きの超並列処理システム CP-PACS の研究に従事。平成 8 年より電気通信大学情報工学科教授、平成 10 年より明星大学情報学部教授。工学博士。電子情報処理学会、IEEE、ACM 各会員。平成 5 年度本学会論文賞、平成 8 年度本学会功績賞各受賞。