Development of Parallel Computer "GIGABOX" for Realtime Applications

技術本部	堀		慶	* 3				
	濱	村	秀	彦*2	村	田	英	明*2
エレクトロニクス事業部	溝	П	正	信*1	笠	井	茂	雄*1

プラント制御,監視システムなどのリアムタイムアプリケーションにおいて、コンピュータの高速化は重要な課題であり、 近年、その有効な手段として並列化が注目されている。当社では、並列コンピュータ"ギガボックス"を開発した。ギガボック スは、RISC型プロセッサ"i860 XP"を最大 64 個まで搭載し、9 600 MOPS(=10<sup>6</sup>命令/s)の処理能力を有する。加えて、 100 MByte/sの高速データ転送バス"カスタムバス"により高速入出力が可能である。さらに、ギガボックス上で動作する "負荷分散ソフトウェア"も併せて完成した。ギガボックスは当社製品"X線CT スキャナシステム"で実用化され、処理時間 の大幅短縮に寄与した。

Faster computers are required for such applications as controlling and monitoring plant. To meet this demand, we have developed a parallel computer, the GIGABOX, with up to 64 i860 XP processors, which demonstrates 9 600 MOPS (Million Operations Per Second) processing performance. It also enables fast I/O with a high-speed data transfer bus, the Custom -Bus. In addition, software which can dynamically distribute its load to the parallel processors has been implemented. The GIGABOX, used in an X-ray CT scanner system, has dramatically reduced the processing time.

# 1.まえがき

CPU チップの高速化には目を見張るものがあり、15年前に主 流 であった 8086 プロセッサではわずか1 MIPS (Million Instructions Per Second) ほどであった処理速度が、現在の i860 XP プロセッサでは 150 MOPS [Million Operations Per Second (=MIPS+MFLOPS)<sup>#</sup>] にまで引き上げられた.それにより、 エンジニアリング ワークステーション、パーソナルコンピュー タなどの演算速度は飛躍的に上昇した.しかし、新型のコンピュ ータが発表されても、たちまちアプリケーションはそれを上回る 演算能力を要求するまで成長し、"より早いコンピュータ"を求 める傾向は慢性的に続いている.その結果、コンピュータの並列 化の研究が進み、現在では、実用可能な技術として着実に産業界 に根付きつつある.

(注) Million FLoating-point Operations Per Second 当社においても、プラント制御システムを始めとする各種制御, 監視、シミュレーションシステムにおける処理の高速化、特にリ アルタイム化を実現するため、高速コンピュータのニーズが高ま っている。そこで、これらのアプリケーションに最適な並列コン ピュータ"ギガボックス"を開発した。ギガボックスの特徴は、 i860 XP プロセッサを最大 64 個搭載し、最大 9 600 MOPS の処 理能力を有することと、アドレス変換機構を備え、100 MByte/s の高速データ転送が可能な"カスタムバス"を装備したことであ る。ギガボックスは、本報で紹介する当社製品"X 線 CT スキ ャナシステム"で実用化されたのをはじめ、幾つかのシステムへ の適用が実現しつつある。

本報では、ギガボックスのハードウェア構成、負荷分散ソフト ウェアの概要、及び適用例としての"高速 X 線 CT スキャナシ



ステム"について述べる。

## 2. 並列コンピュータ"ギガボックス"

ギガボックスは、50 MHz で動作するインテル社製の RISC (Reduced Instruction Set Computer)型プロセッサ "i860 XP" を最大 64 個搭載した並列コンピュータシステムである。ギガボ ックスは、図1に示すように、最大 16 枚までの i860 XP 並列演 算ボード (各ボードに 4 プロセッサ搭載) "MHI 860 XP"と、 それらを管理する1枚の CPU ボードが、VME のシャーシに格 納された構成をとる。その処理能力は最大で 9 600 MOPS とな る.

以下に、ギガボックスの核となる演算ボード MHI 860 XP、及 び高速入出力を可能とする高速データ転送バス"カスタムバス" について述べる。

## 2.1 "MHI 860 XP" のハードウェア

i860 XP 並列演算ボード "MHI 860 XP"は、VME 9 U サイ ズのボード上に四つの PE (Processing Element) が搭載されて いる. ボードの写真を図 2 に、ハードウェアブロック図を図 3 に



図2 i860 XP 並列演算ボード "MHI 860 XP" の写真 MHI 860 XP : i860 XP parallel processing board



1860 XP 亚列演算ボード MHI 860 XP のハードワェ 構成を示す. Block diagram of MHI 860 XP

示す. PEは, i860 XP プロセッサ, メモリ, VME バスインタ フェースなどから構成される演算処理の最小ユニットであり, 表1 に示す性能を有する.

ギガボックスでは、アプリケーションの要求する処理能力に合わせたシステム構築が、MHI 860 XP の増設のみで実現され、非常にコストパフォーマンスに優れる.

## 2.2 カスタムバス

並列コンピュータを制御,監視システムなどへ適用する際,入 出力の速度はシステムの性能を大きく左右する.並列コンピュー タのアプリケーションは通常,(1)多量の物理データを処理し, 少量の結果を出力するアプリケーション(X線CTなど),(2) 少量の論理データから,多量の物理データを生成するアプリケー ション(画像合成)のいずれかに分類される.そのどちらの場合

表1 PEの仕様 Specifications of PE

opecifications of TE					
CPU	i 860 XP(50 MHz 動作)				
ローカルメモリ	SRAM 4 MB EPROM 128 KB				
共有メモリ	1 MB				
処理能力	150 MOPS (50 M I P S + 100 MFLOPS)				
VME バスインタフェース	VME Revision C.1 準拠				

# 表2 カスタムバスの仕様

Specifications of Custom-Bus					
データ転送速度	100 MB/s				
データ幅	32 bit				
最大ノード数	80ノード				
バス信号レベル	NTL (NMOS Transistor Logic)				

でも多量のデータ入出力が伴い,その速度によりシステムのスル ープットが決定される.

ギガボックスの入出力用高速データ転送バスとしては、100 MByte/sの転送能力を有する"カスタムバス"を開発し装備し た. このカスタムバスは、VMEのP2バックブレーンを伝送路 とし、並列コンピュータシステムのプロセッサ-入出力間、及び プロセッサ-プロセッサ間のデータ転送をサポートする. その性 能を表2に示す.

カスタムバスには、高速性を更に高めるメカニズムである"ア ドレス変換機構"が装備されている。カスタムバスにおけるデー タ転送では、送信局は送信先のアドレスをデータに付加して送出 し、受信局はアドレス変換機構を用いてそのアドレスを解析し、 データ受信を実行するか否かを決定する。よって、アドレス変換 機構の設定によって、(1)送信局から単一受信局へのデータ転送、 (2)送信局から複数受信局へのデータ転送が自由に選択できる。 その結果、後者では、最大100 MByte/sのN(受信局数)倍の 高速転送能力が得られるのである。

## 3. 負荷分散ソフトウェア

#### 3.1 静的負荷分散及び動的負荷分散

並列コンピュータでは,並列化の対象となるジョブを複数のタ スクに分割し,これらを各プロセッサに割当てて処理を行う。

ここで、並列処理の開始前に各タスクの処理時間が確定してい れば、ジョブ全体の処理時間が短くなるように、あらかじめ各プ ロセッサにタスクを割当てることができる.これを静的負荷分散 と呼ぶ.

しかしながら、並列処理の実行時に各タクスの処理時間が変動 する場合、最短時間になるような負荷分散を、並列処理の開始前 には行うことができない.また静的負荷分散を行うと、実行時に プロセッサのアイドル状態が発生しプロセッサ稼働率が低下する.

このような場合には、実行時にプロセッサ稼働率が向上するように負荷分散を行うことが、高速化に有効である.これを動的負荷分散と呼ぶ.

既存の商用並列コンピュータでは、タスク割当て・実行機能と プロセッサ間通信機能のみが提供されている.したがって、並列 プログラム作成者自身が、プロセッサ稼働率を向上させるための 動的負荷分散を並列プログラムの記述に含める必要があった.こ のことは、記述の複雑化と工数の増加を招いていた.



ギガボックスにおいては,

- ●負荷変動のあるアプリケーションの高速化
- ●並列プログラミングの容易化

の両立を目的として、負荷分散スケジューラ並びに負荷分散ライ ブラリを開発し標準装備した.

3.2 負荷分散スケジューラによる動的負荷分散の実現

ギガボックスに実装した負荷分散スケジューラは、仮想プロセ ッサと実プロセッサのマッピングの概念に基づいて構築しており、 以下のように動的負荷分散を実現した(図4).

並列プログラム作成者は、仮想プロセッサに対して並列処理の 各タスクを割当てるような並列プログラムを記述する.ここで仮 想プロセッサとは、タスクの割当て対象として、負荷分散スケジ ューラが提供するものである.仮想プロセッサは、負荷分散ライ ブラリを通じて無制限に得ることができ、1個の仮想プロセッサ に対して1個のタスクを割当てることができる.したがって並列 処理の開始前に、すべてのタスクを仮想プロセッサへ割当てるこ とができる.この仮想プロセッサに対して、実際に並列計算機で 利用可能なプロセッサを実プロセッサと呼ぶこととする.

並列処理を開始すると、負荷分散スケジューラは仮想プロセッ サと実プロセッサのマッピングを行う.実プロセッサにマッピン グされた仮想プロセッサは、割当てられたタスクを実行する.実 行が終了すると、マッピングが解かれ、実プロセッサはアイドル 状態となる.この実プロセッサは、実行可能であるにもかかわら ずいまだタスクを実行していない仮想プロセッサと直ちにマッピ ングされる.

このように,実プロセッサのアイドル状態にある時間を少なく するよう次々とマッピングすることにより,動的負荷分散を実現 した.

また,並列プログラム作成者にとっては,仮想プロセッサへの 動的負荷分散を行うだけであり,その記述は容易である.

負荷分散スケジューラ及び負荷分散ライブラリは、いずれもC 言語で記述しており、実装に必要なプロセッサ間通信については、 PVM 仕様<sup>(1)</sup>準拠のプロセッサ間通信ライブラリを利用した。 PVM 仕様のプロセッサ間通信ライブラリは、事実上の標準とし て既に多くの並列コンピュータシステムで稼働しているため、本 スケジューラ並びにライブラリを移植することは容易である。

## 3.3 適用例による評価

モンテカルロ法による解析計算の一つに適用し,負荷分散スケ ジューラの評価を行った.

表3 処理時間及びプロセッサ稼働率の比較例 Example of comparison of processing time and operation rate

		静的負荷分散	動的負荷分散
処理時間	(s)	48.4	21.9
平均プロセッサ稼働率	(%)	46	95



(b) スケジューラ動的負荷分散

図 5

プロセッサ状態の履歴を表示し

プロセッサ状態履歴の視覚化 た画面のコピー。 Visualization of processor histories

モンテカルロ法計算では,確率的試行の繰返し計算が処理の大 半を占める.各繰返し処理単位は互いに独立で依存関係が存在し ないため、タスクとして分割し並列処理を行うことができる.こ のとき、各タスクの処理時間は確率変数の値によって異なるため、 静的負荷分散を行った場合には、プロセッサのアイドル状態が発 生しプロセッサ稼働率が低下する.したがって、動的負荷分散を 行いプロセッサ稼働率を向上させることが、全体の処理時間の短 縮に不可欠となる.

1000個の試行計算を対象に、ギガボックスの12個のプロセッ サを用いて、静的負荷分散と動的負荷分散のそれぞれで並列処理 を行った.処理時間並びに平均プロセッサ稼働率の比較を表3に 示す.表3では、動的負荷分散の方が処理時間が短く、平均プロ セッサ稼働率が大きいことが分かる.

また,視覚化ツール HeNCE<sup>(2)</sup>による静的負荷分散と動的負荷 分散のプロセッサ状態履歴を図5に示す.図5では,負荷分散ス ケジューラが,プロセッサのアイドル状態の発生を抑えるように



成を示す. Structure of high-speed X-ray CT scanner

動的負荷分散を行っていることが分かる.

## 4. X線 CT 画像処理の高速化への適用

#### 4.1 高速 X 線 CT スキャナシステム

1971年に頭部診断医療用として、X線を用いた EMI 社のスキャナのプロトタイプが CT スキャナとして登場し、急速な技術進歩と相まって広範に普及した。現在、医療分野ばかりでなく、非接触・非破壊的に内部が観察できる利点から、工業分野での非破壊検査や計測に、X線や y線の CT スキャナが利用されている.

しかしながら、これらの CT スキャナは線源と検出器を機械的 に駆動する方式であるため、一断面画像を得るために必要な撮像 時間が数 s 以上で、流れの時間変動のような高速現象の計測には 適用できなかった。当社の高砂研究所では、気液二相流の時間変 動の動的な計測を主目的に、撮像時間が 0.5 ms の世界最高速の X 線 CT スキャナを開発した. 撮像時間の大幅な短縮によって、 写真撮影でシャッタ時間を短くするのと同様に、ぼけの少ない鮮 明な断面画像を得ることができるようになった. 加えて、ギガボ ックスの利用によって、データの高速処理が可能となった.

この CT スキャナは、高速化の障害となっている機械的駆動部 分をなくし、電気的に X 線の発生部を高速で切換える方式を採 用している.図6に、高速 X 線 CT スキャナの構成を示す.X 線発生管から射出され、測定対象を透過した扇状 X 線、すなわ ち投影データは、向かい側に置かれた検出器によって検知される。 検出器からの計測信号は X 線発生用の制御バルス信号と同期し てギガボックスを内蔵したデータ収録装置に転送されて、コンボ リューションバックプロジェクション法と呼ばれている再構成手 法によって画像再構成が行われる。再構成画像から、流れの瞬時 の密度分布や相分布が得られる.

図7は、空気と水の混合流すなわち気液二相流の計測例で、測 定面での空気と水の瞬時の分布を表す二次元断面画像が、時間的 に連続に得られている。また、これらの二次元断面画像を重ね合 わせて三次元立体像を構成することができ、時間変動する流動状



図7 気液二相流の計測例 気液二相流の二次元断面図(左),及び合成 された三次元立体図(右). Example of gas-liquid two phase flows

況を立体的に把握できる.図7の例では、気液二相流においてス ラグ流と呼ばれている流れの状況が把握されている.

## 4.2 画像再構成の並列処理化

撮像時間が0.5 msの高速X線CTスキャナからは、毎秒2000枚の断面画像を構成できる多量の投影データが得られる. したがって、このシステムでの計測時間を10sとすると、断面 画像20000枚分の再構成処理が必要となる.この画像再構成処 理を、画像サイズ256×256(画素)で、EWS(SUN SPARC Station 5)を用いて実行しようとすると、10s分のデータ処理 に約5hかかることになり、実用上、処理時間が問題となる.

これに対しギガボックスを用いて処理の高速化を実現した. 一 断面画像分の投影データは,相互に依存関係がないため,各プロ セッサに一断面画像ごとの再構成処理を順次割当てる方法で並列 化を行った,その結果,プロセッサを 64 個用いた場合,処理時 間が 7 min 程度まで短縮されることとなった.

## 5. む す び

コンピュータが高速化されれば、制御,監視システムの精度が 向上するのみならず,新たなるアプリケーションの創造も期待で きる.そして,並列コンピュータはその旗手となるべき存在であ る.

まだ実用化の域に入ったばかりの並列コンピュータは、無限の 発展性を秘めている。当社においても、今後、ギガボックスのプ ロセッサ間通信の強化、ソフトウェアの拡充などを図り、より高 速で、信頼性が高く、操作性に優れた並列コンピュータへと改良 するため、鋭意努力する所存である。

#### 参考文献

- Geist, G. A., Beguelin. A., Dongarra, J. J., Manchek, R. and Sunderam, V. S., PVM 3 User's guide and reference manual, Technical Report ORNL/TM-12187, Oak Ridge National Laboratory, September 1994
- (2) Beguelin, A., Dongarra, J. J., Geist, G. A., Manchek, R. and Sunderam, V. S., Graphical development tools for network-based concurrent supercomputing, In Proceedings of Supercomputing 91, Albuquerque, 1991 p.435-444