

NCモデルに基づくMIMD型並列計算機のための 高性能プロセッサ

A high performance processor for MIMD machines based on NC model

朴 泰佑、野村 茂弘、天野 英晴

Taisuke BOKU, Shigehiro NOMURA, Hideharu AMANO

慶應義塾大学 理工学部 電気工学科

Department of Electrical Engineering, Faculty of Science and Technology, Keio University

1. はじめに

近年の科学技術計算はその応用範囲の広がりや問題領域の拡大によりますます大量の計算を必要としている。このようなニーズに対応するためにスーパーコンピュータは必要不可欠となっている。しかし、既存のスーパーコンピュータのほとんどはパイプライン方式のベクトル計算機であり、問題の系が不規則であったり、計算中に多くの条件分岐が存在するような場合、問題が持つ本来の並列性を十分に生かして効率の高い並列処理を行うことが極めて難しい。

そこで、このような問題に対処するために最近ではMIMD方式の科学技術計算用並列計算機の開発が活発化している。筆者らの研究室で開発中の一般疎行列専用並列計算機(SM)²⁻¹¹もその一つである。(SM)²⁻¹¹ではユーザはNCモデルと呼ばれる並行プロセス・モデルに基づいてプログラミングを行い、計算は任意数の協調プロセスによって実行される。現在のプロトタイプではシステム内のプロセスの実行管理はDIPROSと呼ばれる分散処理系によって行われている。

DIPROSはNCモデルに基づき、プロセス間の通信や切り替えを行うが、現状ではこれら全てをソフトウェアによって行っているため、システムの効率が低い。特に、プロセスの粒度の小さい問題を解く際に効率の低下は著しくなる。

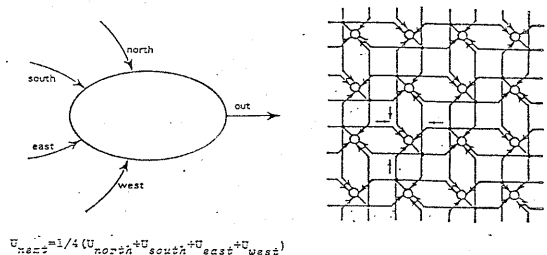
我々が現在開発中の高機能プロセッシング・ユニットは、このようなプロセス間通信(IPC: InterProcess Communication)管理をハードウェアにより行い、また高性能な浮動小数点演算装置(FPP: Floating Point Processor)を効率的に利用するコントローラを備えている。本報告ではこのプロセッシング・ユニットのアーキテクチャを提案する。

2. NCモデル

科学技術計算の分野において、問題の対象となる系が存在し、その形状が終始変化しないような問題は数多い。例えば、有限要素法、偏微分方程式によるポテンシャルの計算、電子回路の過渡解析等がこれにあたる。このような問題では処理を行う複数の節点と、各節点の状態を他の節点に伝える通信網によって問題空間を表現することができる。

我々はこのような系で表わされる問題を効率的に並列処理するための問題記述モデルとしてNCモデルを提案した¹¹。NCモデルでは能動的な処理を行うノード(節点)と、ノード間のデータを転送するネットワーク(通信網)を使って問題を表現する。例えば、2次元のポアソン方程式を差分法によって計算する場合、空間分割された各領域をノードとみなし、隣接するノード間を結ぶ正方形状の結合がネットワークの形状となる。各ノードにおける処理とネットワークを<図1>に示す。

NCモデルにおけるノードとネットワークには次のような特徴がある。



<図1> NCモデル

(1) ノード及びネットワークは処理の実行開始から終了まで存在し、その機能及び形状は変化しない(静的である)。

(2) 各ノードには複数の入力または出力ポートがあり、ノードはこれらのポートを介してネットワークに結合される。

(3) ネットワークは各ポート間を結ぶ結合線により構成される。結合線は1つの出力ポートと1つ以上の入力ポートを結合したものである。

(4) ノード間のデータ転送はポートを介して行われる。データを出力するノードは出力ポートからデータを送る。そのデータは結合線を通して、接続されている全ての入力ポートに送られ、相手側のノードは対応する入力ポートからデータを受け取る。

(5) ノード間の同期は入力ポート側にある単一のバッファによってとられる。データの送信は、結合されている全ての入力ポートのバッファが空いていなければ行うことができ、そうでなければ待たされる。

以上の特徴のうち、(4)のデータ通信方法をマルチキャストと呼ぶ。我々が対象としている問題においては一般的に一つの節点の状態は他の複数の節点に伝えられる必要がある場合が多い。従って、このような問題において、マルチキャストは効率の良い通信方法である。

NCモデルにおけるノードは協調型計算システムにおけるプロセスとみなすことができ、ネットワークはプロセス間の通信経路である。つまり、NCモデルは静的な系を対象とする並行プロセス・モデルである。

3. 並列計算機(SM)²⁻¹¹

NCモデルに基づいて記述されたプログラムを効率的に実行するためには、柔軟な通信形態を持つMIMD型並列計算機が好ましい。本研究室ではNCモデルに基づくアーキテクチャ及びプロセス間通信機能を持つ一般疎行列専用並列計算機(SM)²⁻¹¹⁽²⁾を提案した。

(SM)²⁻¹¹はRSM(Receiver Selectable Multicast)と呼ばれるPU(Processing Unit)間通信機構を持ち、1PUから複数のPUへのマルチキャストを効率よく行う。RSMは一種の仮想共有メモリ・システムであり、全てのPUは1つの仮想アドレス空間を共有するが、全てのアドレスに対して実際のデータ・メモリが存在するわけではない。データを送信するPUはRSM上の1つのアドレスにデータを書き込む。各PU上にはRSMマネージャと呼ばれる選択機構が備わっており、書き込まれたアドレスによってそのデータがそのPUにとって必要であるかどうかを判断し、もし必要であればRSMマネージャ内の受信データ・メ

モリに書き込み、必要でなければ通信に関係しない。これによって、データの必要性が各PU毎に独立に判断され、結果としてマルチキャストが実現される。このように、NCモデルにおけるマルチキャストはRSMを利用して効率的に行われる。

(SM)²⁻¹¹におけるプログラミング言語として、ユーザは並行プロセス記述言語、NCCを利用できる。NCCはNCモデルに基づく言語であり、静的なプロセス及びプロセス間を結合する結合線を通じたマルチキャストをIPCとして提供する。ユーザはNCモデルにおける基本概念をそのまま生かしてプログラミングを行うことが出来る。NCCに関しては次節で詳述する。

NCモデルにおけるノード(プロセス)及び結合線のネットワークが静的であるため、計算実行中のプロセスの生成・消滅や、プロセス間結合線のトポロジーの変更は行われない。

4. NCC(Node oriented Concurrent C)⁽³⁾

NCCはNCモデルに基づく並行プロセス記述言語である。NCCにおいては、NCモデルにおけるノードをプロセスと呼ぶ。プロセスはプロセス・タイプと呼ばれるテンプレートから作り出すことができ、そのプロセスはそのプロセス・タイプの定義部に記述された処理を行う。また、一つのプロセス・タイプから複数のプロセスを生成することができる。

1つのNCCプログラムはプロセス宣言部、プロセス記述部、ネットワーク記述部の3つの部分から成る。プロセス宣言部では必要な全てのプロセスと、そのプロセスのプロセス・タイプを宣言する。プロセス記述部では各プロセス・タイプの処理を記述する。ネットワーク記述部では各プロセスが持つポート(出力ポート、入力ポート)の結合を手続き的に記述し、システム全体のネットワークを構成する。

NCCには一般的なC言語の機能に、IPCを実行する関数と、マルチ受信エントリをサポートするエントリ構文が追加されている。IPC関数及びエントリ構文は以下の通りである。

・IPC(プロセス間通信)関数

```
send(port,data)
```

出力ポートからデータを送信する。

```
receive(port,var)
```

入力ポートからデータを受信し、変数に格納する。

```
probe(port)
```

入力ポートにデータが到着しているかどうかを調べる。

・マルチ受信エントリ構文

```
entry{  
  (port1,var1):  
    statement1;  
    continue または break;  
  (port2,var2):  
    .  
    .  
    .  
}
```

各入力ポートにデータが到着したら対応する変数に格納し、対応する文を実行する。各文の最後がcontinueの場合はさらに別のエントリの受信を待ち、breakの場合はその文を実行した後、まだ未受信のエントリがあっても構文を抜ける。

このように、NCCにおけるIPCはポートを介したデータの送受信によって行われ、各ポートはネットワーク記述部において結合線によって結合される。この時、1つの出力ポートは複数の入力ポートと結合することができ、この結果、データのマルチキャストが実現される。

NCCはNCモデルの基本概念をすべて取り入れた言語で、分布系のシミュレーションや行列計算などの科学技術計算を並列アルゴリズムを使って容易に記述することができる。

5. DIPROS

NCCにおいてユーザは任意個のプロセスを作成することができる。実際の計算はこれらのプロセスを並列実行させることによって行われる。この時、プロセスの数が(SM)²-IIシステム上のPU台数を上回った場合、1つのPU内で複数のプロセスを並行実行させる必要が生じる。また、このような場合、IPCは、PU内ですむ場合もあれば、PU間交信となる場合もある。

そこで、このような実行時におけるプロセスの管理やIPCのサービスを行うために、システム内にオペレーティング・システムが必要となる。(SM)²-IIは、NCCの処理系及び(SM)²-II本体のオペレーティング・システムとしてDIPROS¹⁴⁾と呼ばれるシステム・ソフトウェアを持つ。

各PU上にはLocal Processing System (LPS) が常駐し、それらがRSMを介して交信し合うことによりDIPROSを構成している。

LPSの機能を以下に挙げる。

(1) IPCのサポート

プロセスが出すIPC要求を受け、交信の方法(PU内交信かPU間交信か)を判別してデータの受信バッファへのマルチキャスト(送信の場合)あるいはバッファからのデータの読出し(受信の場合)を行う。

(2) プロセスのスイッチング及びスケジューリング

実行中のプロセスがIPC待ちになった時、プロセスをサスペンドし、別の実行可能プロセスを走らせる。また、サスペンドされているプロセスのIPCが終了後、そのプロセスが再起動されるようにスケジューリングする。

(3) マルチ受信エントリのサポート

マルチ受信エントリ構文を実行するプロセスのエントリを管理し、データの到着順に対応する文が実行されるようにスケジューリングする。

(4) ホスト計算機へのファイル・アクセス及び入出力のサービス

(SM)²-IIはスタンド・アロン型の計算機ではない。そのため、ファイル・アクセスやユーザとの交信は結合されているホスト・マシンを通じて行う。プロセスから出された入出力要求はLPSを介してホスト・マシンに伝えられる。

ホスト計算機上でコンパイルされたNCCプログラムはDIPROSによって(SM)²-IIの各PUにダウン・ロードされ、DIPROSの管理の下で実行される。計算中に発生するIPC要求やファイル・アクセス要求などはすべてLPSによって処理され、他のLPSとの交信またはホスト計算機との交信を行うことによって実行される。

DIPROSは現在、(SM)²-IIのプロトタイプ2号機上に実装され、稼働中である。

6. 高機能プロセッシング・ユニット

NCCによって記述できる問題領域は広く、これまでに種々のアプリケーションが作成されている。我々は様々なアプリケーション・プログラムを実行し、実際にプロセスを実行している時間、即ち有効稼働時間が全計算時間中に占める割合を実測によって求めた。この割合を有効稼働率と呼ぶ。

その結果、有効稼働率は45%から25%と、かなり低い値であることがわかった。DIPROSはプロセス及びIPCに関する全ての処理をソフトウェアによって行っているが、この結果からそのオーバーヘッド

が無視できない大きさであることがわかる。

頻繁にIPCを行うようなプロセスはIPC処理時間だけでなく、プロセスのスイッチング回数及びスケジューリングに要する時間を増大させ、結果としてシステムのオーバヘッドを増加させる。従って、計算時間に対してIPCの頻度が高いプロセス、即ち粒度の小さいプロセスが多いほど、有効稼働率は低くなる。

粒度の小さいプロセスから成る問題を効率的に解くためには、現在のソフトウェアによるプロセス及びIPCの管理ではパフォーマンス不足である。このためには、DIPROSの全部あるいは一部の機能をハードウェア化し、計算(プロセス)の実行とDIPROSの処理を並列に行う必要がある。即ち、1PU内にプロセスを実行するプロセッサと、プロセス・スケジューリングやデータのマルチキャスト、他のPUとの通信などを行うプロセッサの2つを設け、両者を並列実行させる。

メッセージ・バッシングをベースとする計算機ではこのようなオーバヘッドは一般的にかなり大きい。そこで、このようなメッセージ処理を専用のプロセッサを用いて負荷分散する研究が行われている^{[5][6]}。また、外部から到着したメッセージにより自動的にプロセッサ内の処理が行われるようなアーキテクチャも提案されている^[7]。しかし、我々のシステムではシステム全体のメッセージ転送だけではなくPU内におけるメッセージ転送とそれに付随するプロセス・スイッチングも高速に実行する必要がある。このようなハードウェアを持ったプロセッサとしてはTransputer^[8]が存在する。TransputerのIPCはCSP^[9]に基づくランデブーである。しかし、我々が対象としているような問題においては、ランデブーによる一対一通信よりも一枚バッファによるマルチキャストの方がより効率的な処理が行えると考えられる(ランデブーに基づく通信ではIPCのプロック率が高くなってしまう)。従って、NCモデルに基づくハードウェアIPCサーバが必要である。

一方、科学技術計算では大量の浮動小数点演算を行うため、高性能なFPPが利用される。そこで、このようなFPPを有効利用するために、専用のFPPコントローラを設ける。これにより、FPPはある程度独立に演算を行うことができるようになる(この必要性は第8節で述べる)。従って、1つのPUはこれら3つのモジュールから構成される。これらのモジュールをタスク・エンジン、IPCエンジン、FPPエンジンと呼ぶことにする。各モジュールの機能は次の通りである。

タスク・エンジン：NCCプログラムのプロセスを

逐次実行する。プロセスがIPCを行う際にコンテキストをIPCエンジンに渡し、次に実行可能となった別のコンテキストをIPCエンジンから受け取ってそのプロセスを続行する。

IPCエンジン：タスク・エンジンからのIPC要求を受け、PU内あるいはPU間のマルチキャストを実行する。同時にIPCの完了により実行可能となったプロセスをタスク・エンジンが実行できるようにスケジューリングする。

FPPエンジン：タスク・エンジンがプロセスを実行中に発生する浮動小数点演算要求を受け、一連の計算を行う。また、浮動小数点演算が他のプロセスからのデータ待ちである場合、データが到着し次第、タスク・エンジンを介さずIPCエンジンによって直接起動されることもある。

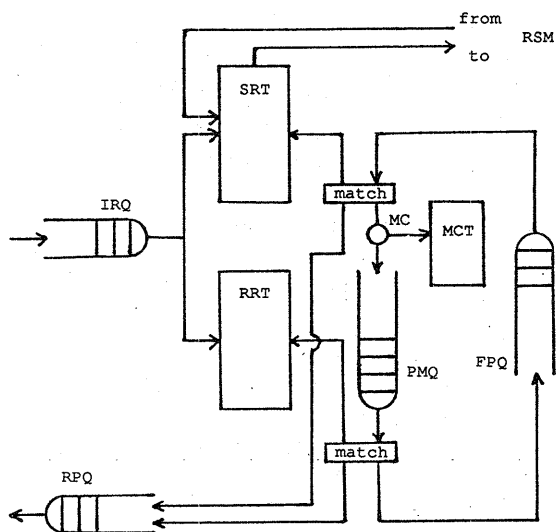
これら3つのモジュールが並列に動作できるようになれば、システムのパフォーマンスは飛躍的に増大すると考えられる。しかし、これらの内、IPCエンジンはその機能の複雑さからハードウェア化するのは非常に難しい。特に、対象が一般的な並行記述言語であればIPCを完全にハードウェア化するのは極めて困難である。しかし、(SM)²-IIで使われているNCCはモデルの単純さと、静的であるという点をうまく生かせば比較的簡単な構造でハードウェア処理できると考えられる。次節ではIPCエンジンの機能と、それらをどのような方法でハードウェア化するかについて述べる。

7. IPCエンジンのアーキテクチャ

IPCエンジンの主な機能は(マルチ受信エントリを含む)IPC要求の処理、プロセスのスケジューリングであるが、これらは密接に関係し合っており、切り離すことはできない。本節ではこれらの機能を総合的に処理するアルゴリズムとその実現に必要なアーキテクチャについて述べる。

IPCエンジン中には2つの流れ、即ちデータの流れとコンテキスト(プロセス)の流れがある。データの流れをバッファの状態遷移でとらえると、バッファは空の状態とデータが入った状態を往復する。また、コンテキストの流れは、送受信を待っている状態と送受信が完了して再実行可能になった状態からなる。これらの状態を管理する機構として次の5つの要素を考える。

- S R T (Send Request Table)
送信要求を出したコンテキストを登録するテーブル
- R R T (Receive Request Table)
受信要求を出したコンテキストを登録するテーブル
- F P Q (Free Port Queue)
対応するバッファが空である出力ポートを保持するキュー
- P M Q (Private Message Queue)
マルチキャストされたデータを受信プロセス別に保持するキュー
- R P Q (Ready Process Queue)
送受信が完了して再実行可能になったプロセスを保持するキュー



<図2> IPCエンジンのブロック・ダイアグラム

<図2>にIPCエンジンのブロック・ダイアグラムを示す。これらのデータ構造を用いて、IPCの実行とプロセス・スケジューリングを行う方法は次の通りである。

まず、初期状態として全ての出力ポートをFPQに入れる(全てのバッファは空である)。到着したIPC要求は送受信の別によりSRTもしくはRRTに登録される。システム内では絶えずスケジューラが動作しており、実行可能なIPCを見つけ出す。スケジューラはFPQを順にスキャンし、対応する送信待ちのプロセスをSRTより見つけ、見つかった場合はデータを全受信プロセス用にコピーしてPMQに入れると同時にそのプロセスをSRTから削除してRPQに入れる。また、スケジューラはPMQを順にスキャンし、対応する受信待ちのプロセスをRRTより見つけ、見つかった場合はそのプロセスをRRTから削除してRPQに入れる。この時、もしそのプロセスがそのマルチキャストにおける最後の受信者であった場合は、そのポートをFPQに入れる。

以上のようにしてデータとプロセスはこれらのデータ構造の中を巡回し、IPC及びスケジューリングが実現される。ここで、送信と受信の処理を分けてるのは、両者の処理に並列性があり同時実行が可能であるためである。従って、IPCエンジンの中はリクエスト受付部、送信側スケジューラ、受信側スケジューラの3つの部分から成り、並列実行される。

以上はPU内のプロセスの間だけで通信が行われる場合である。通信相手がPU外にいる場合はRSMを介した通信が必要となる。この場合、送信に関しては

直接RSMに出力される。出力後のデータのマルチキャストはRSMにより自動的に実行される。また、RSMから受信したデータはPU内でさらにマルチキャストする必要があるため、この場合はデータがあたかもPU内のプロセスから送信されたようにSRTに登録する。その後の処理は送信側スケジューラが行ってくれる。このように、IPCエンジンではデータのマルチキャストを、相手プロセスがPU内にある場合も一貫して処理することができる。

マルチ受信エントリ構文のサポートは、基本的に受信側スケジューラによって行われる。この場合は、まずプロセスがエントリ構文に入る際、全受信エントリに対応する受信要求をRRTに登録する(この時、再実行開始番地も一緒に登録しておく)。後は受信側スケジューラに任せればよいが、問題なのはそのエントリがbreakにより終了している場合である。この解決法として、RRTにはそのエントリがブレイクであるかどうかを一緒に登録しておき、スケジューラが受信可能と判断した際に、そのプロセスに関する他の受信要求を全てRRTから削除する。こうすればその後の受信は成立せず、ブレイクが実現される。

以上のように、IPCエンジンでは3つの独立した機構によりNCCの処理系として必要な全機能を高速に実行することができる。

8. FPPエンジン

科学技術計算用の計算機に使われるFPP (Floating Point Processor) は一般的に数MFLOPSという高い

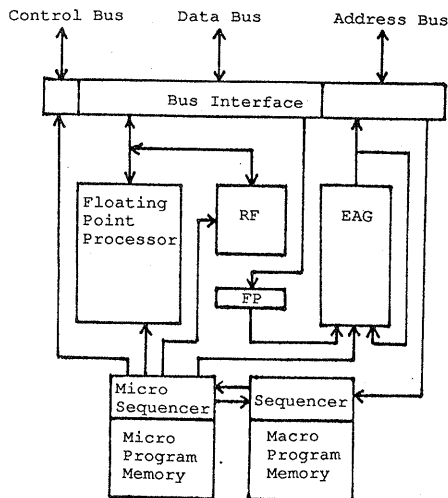
演算能力を持っている。しかし、この値はFPP単体の性能であって、実際の使用状況（例えばメイン・プロセッサによって起動される、データをメモリから読み取る等）においてはその能力を有効に発揮できるとは限らない。たとえ倍精度乗算が200ナノ秒でできたとしても、そのデータを用意し、CPUから起動されて結果を返すまで数マイクロ秒かかっているとは高性能のFPPを使う意味がない。

いかにしてFPPを有効利用するかは、FPPアクセスのためのハードウェア・アーキテクチャとコンパイラのインテリジェンスにかかっている。浮動小数点演算が必要な度に、CPUが命令を解釈して各命令毎にFPPを起動する方式ではFPPの能力を生かしきれない。そこで、まず、ある程度のマクロ命令をプログラムとして実行できるコントローラと、ユーザ・プログラム中の浮動小数点演算部分のみをコンパイルしたマクロ命令シーケンスを用意する。そして、CPU（タスク・エンジン）はどのマクロ命令シーケンスを実行すべきかというトリガだけをコントローラに与えるようにする。FPPエンジンはこのようなFPPとそのコントローラから成る。FPPエンジンのブロック・ダイアグラムを<図3>に示す。

FPPエンジンのコントローラは水平型マイクロ・プログラムに従って一連の命令（マクロ命令）セットを解釈・実行する。また、単純なアドレッシング・モードを用意し、タスク・エンジンに頼らずにFPPエンジンだけでメモリをアクセスできるようにすることで、CPU-FPP間のやりとりを最小限に抑える。コンパイラはプログラム中に浮動小数点演算があると、できるだけ長いマクロ命令シーケンスを生成し（これはFPPエンジンを、できるだけ長期間タスク・エンジンの制御を受けずに動作させ続けるためである）、そのシーケンスの実行をトリガする命令だけをタスク・エンジンのプログラムとする。従って、コンパイラは高機能でなければならず、ある程度のインテリジェンスが要求される。

また、複数のプロセスを並行実行する場合、各プロセスのデータは局所変数に格納されている場合が多い。そこで、現在タスク・エンジンが実行中のプロセス（コンテキスト）のフレーム・ポインタの値をマクロ命令トリガの際に渡し、FPPエンジンではこのポインタを使ってオフセット計算を行うようにする。これにより、FPPエンジンはコンテキストを識別することなしに計算の実行が可能となる。

この方式ではタスク・エンジンとFPPエンジンが並列に処理を行うことができる。但し、コンパイラはタスク・エンジンがFPPエンジンの計算終了を待たなければならない点を検出して、その点で2つのエンジンの同期をとるようなコードを生成しなければなら



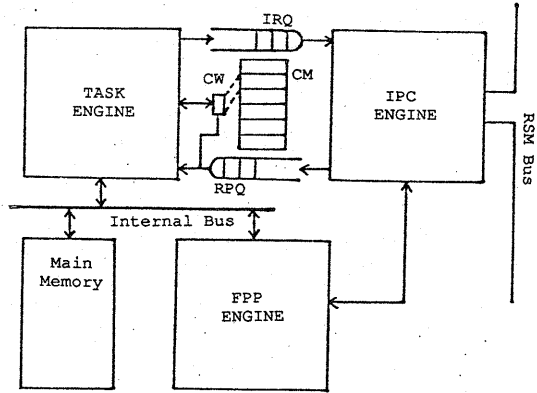
<図3> FPPエンジンのブロック・ダイアグラム

ない。また、2つのエンジンは各々独立にメモリをアクセスするため、FPPエンジンがメモリをアクセスする場合は、タスク・エンジンにバス使用要求を出してバスを譲ってもらわなければならない。この調停作業はある程度の時間を浪費するため、両者のアクセスが頻繁に衝突するようなプログラムの場合、無理に並列実行をせずにタスク・エンジンを止めた方が良い場合も考えられる。この判断はやはりインテリジェンスを持ったコンパイラが行わなければならない。

9. タスク・エンジンとIPCエンジンの結合

以上3つのエンジン、タスク・エンジン、IPCエンジン、FPPエンジンについて述べてきたが、高機能プロセッサにとって重要なのはこれらのエンジンをどのように組合せ、交信させるかである。例えば前節の最後に触れたように、タスク・エンジンとFPPエンジンがメモリを共有すると両者をうまく並列実行させることが難しくなる。この2つのエンジンの関係はどちらもユーザ・プロセスの実行の一部を行う点で共通なので、どちらかをメモリから切り離すわけにはいかない。しかし、例えばタスク・エンジンとIPCエンジンのようにその機能が完全に分散されている場合はその関係に適した結合方式を考えなければならない。

高機能プロセッシング・ユニットの全体構成を<図4>に示す。タスク・エンジンとIPCエンジンの関係は、一種のクライアント-サーバの関係である。即ち、タスク・エンジンがIPC要求を出し、IPCエンジンがその要求を受けると同時に既に要求が満たされたコンテキストをタスク・エンジンに返すことによ



<図4>プロセッシング・ユニットの全体構成

って両者の実行が継続される。この時、両者の動作が非同期に並列実行されるようにするためには、両者を結合するリクエスト及びプライのキューが必要である。これらのキューをIPCリクエスト・キュー及びレディ・プロセス・キューと呼ぶ。この2つのキューによって両者は互いに独立に動作できるようになる。

タスク・エンジン内のプロセスがIPC要求を出した場合、タスク・エンジンはIPCリクエスト・キューに必要な情報を入れる。IPCエンジンは要求をキューから取り出し、必要な処理を行った後、再実行可能なプロセスをスケジューリングしてレディ・プロセス・キューに入れる。

以上の点で特に問題となるのは次の3点である。

(1) サスペンドされているプロセスをどのように管理するか

コンテキストは基本的にタスク・エンジンが管理する。このために、メイン・メモリとは別にコンテキスト・メモリを用意し、プロセスがサスペンドされている間、全ての環境をこのコンテキスト・メモリに保存する。そして、IPCエンジンはそのプロセスのプロセスidによりプロセスを管理し、レディ・プロセス・キューにプロセスidを入れることによりタスク・エンジンにどのコンテキストを再起動すべきかを教える。

また、コンテキスト・メモリへのアクセスを簡単化するためにレディ・プロセス・キューと連動して動作するコンテキスト・ウィンドウを用意する。コンテキスト・ウィンドウはコンテキスト・メモリ中の現在のコンテキストだけをタスク・エンジンに見せ、タスク・エンジンはこのウィンドウを通して環境の保存・復帰を行う。これにより、タスク・エンジンは実質的にどのコンテキストを実行中かをほとんど意識せずに処理を進めることができる。

(2) 送受信データの受渡しをどのような形で行うか
IPCエンジンはメモリに結合されていないため、プロセスの変数に直接アクセスすることができない。そこで、データを送信する場合はタスク・エンジンがデータそのものをIPCリクエスト・キューに入れてIPCエンジンに渡す。また、受信されたデータはIPCエンジンがレディ・プロセス・キューに直接入れ、タスク・エンジンがそのデータを対応する変数領域に書き込むことによって受信が成立する。

この方法は比較的大きなデータをやりとりする場合、効率落ちる可能性があるが、現行のNCCではIPCデータ・タイプとして倍精度浮動小数のみをサポートしているためこの方法でも問題は生じない。

(3) マルチ受信エントリの管理をどのように行うか
マルチ受信エントリの場合、複数のコンテキストが同一プロセスで生じる点が問題である。NCCの文法ではマルチ受信エントリの全てのコンテキストの実行開始環境はエントリ構文に入る直前の状態であると定められている。そこで、タスク・エンジンが各コンテキストを起動する際はコンテキスト・メモリ中の環境をそのまま使い、そのコンテキストが終了した段階で環境をコンテキスト・メモリに戻さなければよい。但し、CPUのレジスタ中でプログラム・カウンタだけは各コンテキストにより異なるため、IPCリクエストを出す際にそのコンテキストの開始番地と一緒にIPCエンジンに渡すようにする。

以上の様にしてタスク・エンジンとIPCエンジンは2つのキューを介して結合するだけで全ての処理を行うことができる。

10. FPPエンジンとIPCエンジンの結合

前節ではタスク・エンジンとIPCエンジンの結合について論じた。FPPエンジンはタスク・エンジンによって起動されるため、IPCエンジンとは直接結合される必要はない。しかし、NCCで記述された各種アプリケーションを調べると、マルチ受信エントリにおいて、受信後に実行される文は多くの場合簡単な(短いという意味ではなく制御のほとんどない)浮動小数点演算のみであった。このような場合、タスク・エンジンが実行する命令はFPPエンジンの対応するマクロ命令シーケンスのトリガという極めて短い処理だけである。しかし、前節の様な構成だと、タスク・エンジンはこれだけのために環境の復帰を行わなければならない、非常に大きなオーバーヘッドとなる。そのコンテキストの環境の中で、実際にFPPエンジンに使用されるのはフレーム・ポインタの値だけであり、その他のレジスタ等は復帰する必要がない。

そこで、このようなFPPエンジンだけを使用するコンテキストに対してはIPCエンジンから直接FPPエンジンを起動することを考える。このような場合、IPCエンジンはマルチ受信エンタリ用のIPCリクエストの際に、FPPのマクロ命令シーケンスの実行開始番地とフレーム・ポインタを受け取っておき、それらを使って直接(タスク・エンジンを介さず)FPPエンジンを起動する。この際、タスク・エンジンとIPCエンジン間でFPPエンジンへのアクセスの衝突が生じるが、これは簡単な調停装置により制御できる。

もし、コンテキストが浮動小数点演算と通常の処理の両方を含む場合はいずれにせよタスク・エンジンを使用しなければならない。この場合は通常の方法で、タスク・エンジンによりコンテキストを起動し、そこからFPPエンジンを起動させる。このどちらの方法を採用するかはコンパイラにより決定される。

以上のように3つのエンジンはそれぞれ適した方法で結合され、各々が並列に実行されることにより極めてパフォーマンスの高い高機能プロセッサが実現できる。

11. まとめ

本報告では、科学技術計算用MIMD型並列計算機のための高機能プロセッシング・ユニット(PU)のアーキテクチャを提案した。このPUは基本的に3つの処理モジュール、タスク・エンジン、IPCエンジン、FPPエンジンから成る。そして、NCモデルに基づくマルチプロセッシングをハードウェア的に処理することができる。

これら3つのモジュール間はキューあるいはバスによって結合され、並列に動作する。特にタスク・エンジンとIPCエンジンは定常的に並列動作することが可能であると考えられる。

当初の我々の目的は現在我々が開発中の一般疎行列専用並列計算機(SM)²⁻¹¹の性能をより高めるためのIPC処理機構のハードウェア化であった。しかし、ここで提案したアーキテクチャは(SM)²⁻¹¹に限らず、NCモデルに基づくMIMD型計算機であれば利用可能である。特にIPCエンジンは、タスク・エンジンのCPUを選ばず、FPPへのダイレクト・アクセス機能を使用しなければ特殊なFPPエンジンも必要としないため、CPUやFPPのバージョン・アップに容易に対応できるような構成になっている。

現在我々は、アーキテクチャの完全な設計を行い、各種パラメータ(内部テーブルの大きさやキューの最大長など)を決定するためにエミュレータを作成中である。このエミュレータは提案したアーキテクチャと

ほぼ同じ構成をしているが、IPCエンジンの部分が汎用のCPUとメモリから構成されている。これは、IPCエンジンの完全なハードウェア化が難しいため、構造やアルゴリズムの詳細を決定し、実際の計算を行ってタスク・エンジンとの並列実行度を調べるためのものである。

エミュレータのIPCエンジンはソフトウェアによって実行されるが、タスク・エンジンと切り離され並列に動作するため、この点だけでも従来の(SM)²⁻¹¹より高速になる。しかし、この形ではプロセスの粒度が小さい場合にタスク・エンジンがプロセス待ちのアイドル状態となり、処理効率の低下は避けられない。やはり、IPC処理をハードウェア化することによって広範囲な粒度に対応できる高機能PUが実現できる。

現在、IPCエンジンの許容プロセス数として100プロセス程度を見込んでいる。RSMは一般的な共有メモリ・システムに比べバス・コンフリクトによる性能低下が少ないため、クラスタ構成をとることによる数千台のPUからなるシステムを構築することができる。このような大規模システムでは数万点以上の節点の計算を行うことも可能である。

エミュレータによる検証が終了しIPCエンジンを含む高機能PUの機能が実証された後は、実際にIPCエンジンのハードウェア化を行う。また、将来的にはIPCエンジンをVLSI化することにより、3チップ構成の高機能PUによる並列計算機を作成することを目標としている。

[参考文献]

- [1] T.Kudoh, H.Amano, T.Boku, "NDL: A language for solving scientific problems on MIMD machines," Proc. of 1st Super Computing Symp., Dec.1985
- [2] H.Amano, T.Boku, T.Kudoh, "(SM)²⁻¹¹: The new version of Sparse Matrix Solving Machine," Proc. of the 12th Int. Symp. on Comp. Arch., Jun.1985
- [3] 高橋範朗, 他「大規模MIMD型並列計算機のための並行プロセス記述言語NCC」第32回情報処理学会全国大会、1986年3月
- [4] T.Boku, H.Amano, T.Kudoh, "DIPROS: A distributed processing system for NDL on (SM)²⁻¹¹," Proc. of the 20th Hawaii Int. Conf. of System Sciences
- [5] U.Ramachandran, M.Solomon, M.Vernon, "Hardware Support for Interprocess Communication," Proc. of the 14th Int. Symp. on Comp. Arch., Jun.1987
- [6] J.W.Wendorf, H.Tokuda, "An Interprocess Communication Processor: Exploiting OS/Application Concurrency," CMU Internal Report, Mar.1987
- [7] W.J.Dally, et.al., "Architecture of a Message-Driven Processor," Proc. of the 14th Int. Symp on Comp. Arch., Jun.1987
- [8] D.May, R.Shepherd, "The Transputer Implementation of OCCAM," Proc. of Int. Conf. of FGCS,
- [9] C.A.R.Hoare, "Communicating Sequential Processes," Comm. of ACM, Aug.1978