

高並列計算機CAP-IIの ブロードキャスト・ネットワーク

加藤 定幸 清水 俊幸 堀江 健志 石畑 宏明

(株)富士通研究所

本論文では、高並列計算機CAP-IIのブロードキャスト・ネットワークについて述べる。ブロードキャスト・ネットワークはホスト計算機とセル(プロセッサ)を接続するネットワークで、主にホストとセルのインターフェースに用いられる。ブロードキャスト・ネットワークはホスト-セル間で、2次元配列を高速に交換するスキヤット、ギャザの2つの転送モードを持っている。本論文ではこれらのブロードキャストネットワークの機能と実現方法について述べる。

Broadcast-network of Highly Parallel Processor CAP-II

Sadayuki Kato Toshikyuki Shimizu Takeshi Horie Hiroaki Ishihata

FUJITSU LABORATORIES LTD.

1015, kamikodanaka, Nakahara-ku, Kawasaki 211, Japan

This paper presents a Broadcast-network of highly parallel processor CAP-II. Broadcast-network is a network which connect between host and cells (processor), and it has "scatter" and "gather" transfer mode of Broadcast-network which exchange 2-dimensional array between host and cells rapidly. This paper presents function of Broadcast-network and its hardware design.

1. はじめに

我々は、数値計算と映像生成の高速実行を目的とした、分散メモリ型の高並列計算機CAP-IIを開発している。本論文では、CAP-IIのホスト・セル(プロセッサ)間の通信路であるBroadcast network (Bネット)について述べる。

CAP-IIは64~1024台のセルを持っているので、ホスト計算機はこれらのセルと効率よく通信を行うことが出来なくてはならない。ホスト・セルの通信にはホストからセルへのプログラム、データの転送、セルからホストへの計算結果の転送がある。ホストからの転送のうちプログラムや各セル共通のデータはブロードキャストによって高速に転送可能であるが、セル毎に異なるデータを送る場合あるいは、セルからホストに計算結果を吸い上げる場合は、転送するセルの選択あるいは、転送するデータの選択といったオーバーヘッドがある。しかも、このオーバーヘッドはセルの台数が増加するにつれ大きくなる。

この問題を解決するため、我々はBネットの転送モードとして、従来からあるブロードキャストに加えて、セルに異なるデータを転送するスキヤッタ(分配)、セルのデータをホストに効率よく吸い上げるギャザ(収集)を設けることにした。

本論文では、まず、第2章でBネットの機能、特にスキヤッタ、ギャザについてのべる。第3章ではCAP-IIのコマンドバスのハードウェアについて、スキヤッタ・ギャザの実現法を中心に述べる。第4章では、スキヤッタ、ギャザの両機能について、従来の転送方式との比較を行い、本方式の有効性について検討する。

2. Broadcast network (Bネット)の仕様

2.1 Bネットの機能

CAP-IIのBネットはホスト計算機と全てのセルを接続した共通バスである。BネットはCAP-IIの1000台のセルとホスト間で効率よく通信する事を目的としたネットワークである。並列計算機のセル(プロセッサ)間の通信については広く検討されているが、セルとホスト計算機の通信も以下の理由により高性能化が不可欠である。

CAP-IIにおいては、個々のセルの処理能力の向上および、セル台数の増加によって、システム全体で処理可能なデータの量が増大する。ホスト—セル間通信はこのデータを供給するために高いバンド幅が要求される。また、システムを構成するセルの数が増大するにつれて通信時のセル切替えのオーバーヘッドが増大する。

また、ブロードキャストで全てのセルに同じデータを送るときはホストは1度のデータ転送でデータを送る事ができるが、ホストのデータを分割してセルに送る場合、あるいはセルの計算結果をホスト計算機に集める場合には、受信するセルあるいはデータを送りだすセルの切替えのオーバーヘッドが入る。この切替えの回数はセルの台数に比例して増加し、1回の切替え時間もセルの増加に

よるハードウェア規模の増大にともなって長くなる傾向がある。そのため、長さがnのデータをp台のセルに等分に分割して送るときの転送時間Tは以下ようになる。

$$T = T_c \times p + T_t \times n \quad \text{————— (1)}$$

T_c :セル切替え時間

T_t :1ワードの転送時間

従ってセル台数の増加にともなって、実際のデータの転送時間よりもセルの切替えのオーバーヘッドが大きく影響してくることになる。この問題を解決するためにCAP-IIのBネットには通信切り換え時のオーバーヘッドが少なくなるような転送方式が必要になってくる。

Bネットは論理的にはセルとホストを接続した均一な共通バスに見える。

セルとホストはBネット上では対等であり、セルからセルへのブロードキャストも可能になっている。そこで、ホストとセルをまとめてノードと呼ぶことにする。Bネットには同時に複数のノードよりデータを送りだすことは出来ない。ブロードキャストの権利を獲得したノードをマスタ、その他をスレーブと呼ぶ。また、各ノードにはそれぞれIDがふられている。このIDは各セルで異なっているが複数のノードに共通のIDを与えることも可能である。共通のIDをもったノードを同じグループに属しているノードと呼んでいる。

表1に、Bネット上での転送モードについて説明する。マスタはデータを送りだす際に、1まとまりのデータ(パケットと呼ぶ)のまえにヘッダをつける。ヘッダには転送のモードを示すビットならびにこのデータを受信するグループのIDを示すビットが含まれている。

表1 Bネットの機能

ブロードキャスト	そのパケットを全てのスレーブが受信する。ホストから全てのセルに共通のデータを送ったりプログラムを送るのに使用する
グループキャスト	ヘッダのIDビットと各スレーブのIDを比較して等しいスレーブだけがパケットを受信する
スキヤッタ	2次元配列を分割し、異なるセルに送る転送モード。
ギャザ	セル上の2次元配列をホスト上に集めて、1つの2次元配列を作る転送モード。
その他	①ブートストラップコードをセルに送る。 ②セルにIDを付ける。

2.2 スキャタとギャザ

CAP-Ⅱは主に数値計算を高速実行するために、作られた並列計算機であるが、このような用途では2次元配列のデータを扱うことが多い。また、この配列をセルに送る場合、すべてのセルに同じ配列を送る場合だけでなく、図1のように、配列を分割してセルに割り当てることも多い。

従来のネットワークでこのような事を行う場合、まずデータをホスト計算機上でブロックに切り分けて、それぞれのブロックを送り先のセルを指定して転送していた。したがって、前述のようにデータ転送に要する時間には切りわけの時間と転送先のセルの切替え時間がオーバーヘッドとして加わっていた。しかし、Bネットでは受信側に送られてきたデータが自分が受信するブロックに入っているかどうかを判定するハードウェアをもうけブロックの切りわけを受信側で行っている。この判定の時間がデータの転送速度より早ければ、送信側（マスタ）はブロードキャストと同等の速度でデータを転送することができる。この機能をスキャタとよんでいる。

逆にセルの2次元配列をデータをホストに収集する転送方式としてギャザがある。このモードでは前述のブロードキャストやスキャタとは反対にスレーブからマスタに向かってデータが送られる。各セルは他のセルとホストの通信状況を監視しホスト上の配列のどの部分が転送されているか判定して、自分がデータを送りだすタイミングを決めている。同時に多数のセルがデータを送りだすためにセル間の順序の制御が必要になるが、この制御に要する時間はマスタを決定する調停時間と異なり、データ転送時間と等しいので、ブロードキャストと等しい転送速度でデータを収集することが出来る。

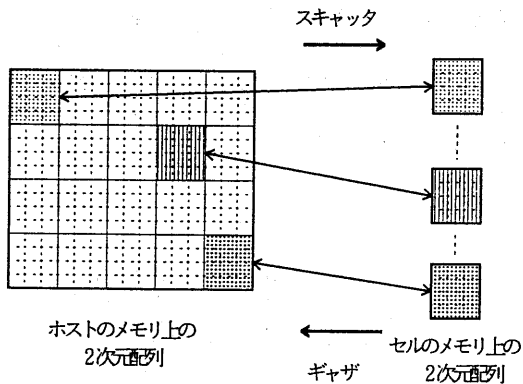


図1. 2次元配列の交換

スキャタ、ギャザ時のブロックのパターンには図2に示すようにX方向Y方向のストライプ状のパターンの組み合わせである。このパターンはそれぞれのスレーブごとに設定し、そのパラメータはX、Y方向それぞれの幅と間隔と原点からのオフセットの6種類である。この

パターンはスレーブごとに独立しているので、1部のスレーブにはブロック状に他のスレーブにはストライプ状に配列を切り分けることも可能である。また、スキャタ時には隣あったセルにオーバーラップしたブロックに切り分けることもできる。

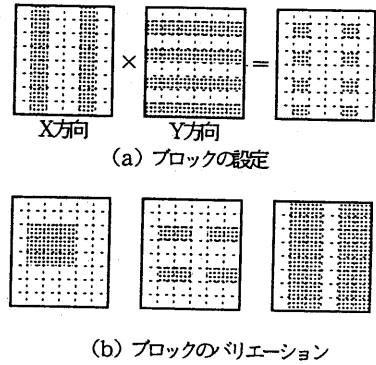


図2. 配列の切りわけ例

3. Bネットのハードウェア

3.1 ネットワークの構造

Bネットは1024台までのセル構成に対応可能でなくてはならない。そこで我々は、Bネットの構成として図3に示すようなリング型のネットワークと階層型のネットワークを組み合わせた構成を取ることにした。バックプレーン間の接続はリング型のネットワークを用い配線を簡略化している。またバックプレーン内ではツリー型のネットワークを用いている。

リング型のバスを用いるとセル台数の増大に比例して伝搬遅延（レイテンシ）が増加する欠点があるが、ツリー状のネットワークと組み合わせて伝搬遅延の増加を抑えている。コマンドバスではブロードキャスト、スキャタ、ギャザなどまとまったデータが流れる事が多く、転送性能はスループットによって決まるので、リング型のバスを用いても性能はそれほど低下しない。

マスタとなるノードを決めるために、ノード間で、コマンドバスの転送権のアービトレーションを行っている。アービトレーションの回路もバスの構造と同様に階層構造となっており、調停に時間を要するが、前にも述べたようにコマンドバスは大きなデータの転送を主に考えているので、アービトレーションの高速化は特に考慮しなかった。

Bネットではホストあるいは任意のセルがマスタになることが出来るので、マスタの位置に応じてネットの構造を変化させる必要がある。具体的にはマスタのあるツリー型ネットへのリングの分岐点の入力側でリングネットを切断する。ツリー型の部分では基本的にはリング

からセルへのブロードキャストになるようにバッファの向きを決めているが、マスタの存在するツリーではマスタからリングとの接続点に向かう経路上のバッファの転送方向を反転してデータをリングへおくりだす。これによってBネットワーク全体はマスタをルートにしたツリー状のネットワークになる。このネットワークはマスタとスレーブの距離が不均一であり、また、MIMD型計算機であるので各スレーブが同期して動作している保証もない。そこで、データの受信が遅れたスレーブにあってもデータの取りこぼしがないようにハンドシェイクを行っている。ただし全スレーブでハンドシェイクを行うと1ワードの転送に時間が掛かるので各バッファ内に2段のFIFOをもうけツリーの各枝毎にハンドシェイクを行っている。したがって、ネットのサイズによらず転送速度は一定である。

ギャザの時はマスタに向かってスレーブからデータが送られるので、ツリー型ネットでは転送方向を全て反転

表2. Bネットワークのハードウェア諸元

トポロジ	リング+ツリー
セル構成	64~1024台
伝播遅延	0.8 μ s(64台) ~ 5.6 μ s(1024台)
転送速度	50MB/秒 (32ビット幅)
専用LSI	CMOS 1.2 μ m ゲートアレイ 33000 ゲート

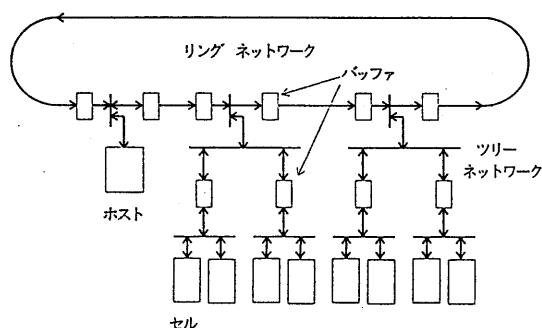


図3 ブロードネットの構造

し、リングの部分ではマスタのあるツリー型ネットへのリングの分岐点の出力側でリングネットを切断する。しかし、複数のバッファが同時にデータを出力することは出来ないため、あとで述べるようにバッファの出力制御を行っている。

また、ネットワーク上のバッファにデータが残っているままで、マスタを切り換えてしまうとデータが正しくスレーブに伝わらない。そこで、データが全てのスレーブに送られた事を検出するために、データ転送終了後、マスタはコマンドバスに転送の終了を示す特殊なデータを送りだす。このデータが全てのセルに行き渡った事をCAP-IIのSネット（同期機構）で検出して、マスタを開放する。ギャザ時はマスタがあらかじめ決められた数のデータを受けたことで終了を判定する。

3.2 専用LSI

BネットワークとセルあるいはホストとのインターフェースにはBIFと呼ばれる専用のLSIをもちいる。BIFの内部構成は図4の様になっている。また、階層バスの各バッファもこのLSIをデータのバッファとして使用している。

- ヘッダデコーダ: ヘッダのデコードを行い、送られてきたデータが自分がこのパケットを受信するのかどうか、また、スキヤッタで受信するのかどうかを判定する。また、グループキャストのためにホストから送られてきたIDを記憶する機能をもつ。

- 送受信バッファ: Bネットワークとのデータの転送を行うバッファ、この部分はBネットワークの各バッファでも使用できるようになっている。

- Bネットワークコントローラ: コマンドバスの獲得、および、スキヤッタ・ギャザ時のデータの転送タイミングの制御を行う。スキヤッタ、ギャザのパラメータはセルのCPUによって設定される。スキヤッタ時は送られてきたデータが自分が受信するデータで無いときは受信バッファをクリアし、受信するときはFIFOに書き込む。ギャザ時はデータを送るときはFIFOのデータを送信バッファにおくり、それ以外の期間はダミーを送る。

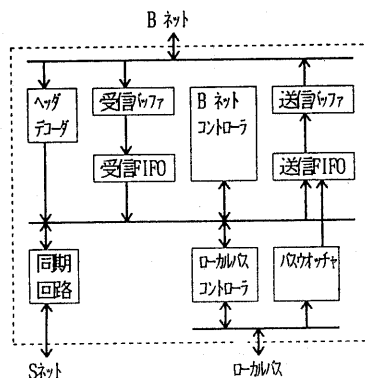


図4 ブロードネット専用LSI

- 送受信FIFO: B ネットとセルのローカルバスの転送速度の差を吸収する8ワードのFIFO
- ローカルバスコントローラ: セルのローカルバスとBIF内部のレジスタあるは、FIFOとのデータ転送の制御をおこなう。
- 同期回路: CAP-IIのSネットの制御を行う。
- バスウォッチャ: セルのローカルバスのパフォーマンス、エラーの監視と記録を行う。

3.3 ギャザの実現

スキケットについては4.1のブロードキャスト機能と4.2のBIF内の選択回路によって実現できる。しかし、ギャザではデータを出力するスレーブを高速で切り換える必要がある、また、MIMD型の計算機では各セルの動作にばらつきがあり、同じタイミングでデータを送りだす保証がないためにあるセルが遅れた時でも、マスタに届くデータの順序が正しい事を保証する必要がある。その為に全スレーブのBIFを同期させる必要があるが1データ毎に全セルで同期を取っていたのではデータ転送速度が同期の速度で制限されてしまう。そこで、各スレーブの同期にはデータの転送プロトコルを同期に利用し、各合流点で独立に同期をとることで、ブロード時の転送速度でギャザを出来るようにした。

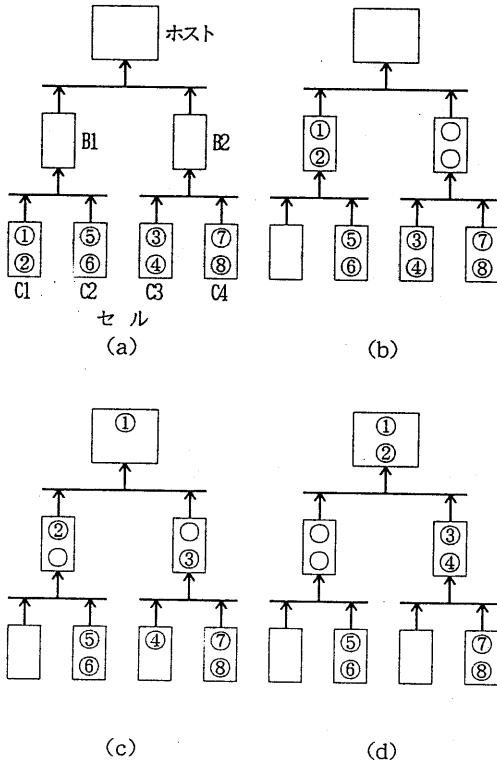


図5 ギャザ時の同期

各セルはパラメータに従ってデータを送りだすが、各スレーブが送りだすデータはホスト上の配列のイメージの中に自分が送りだすデータを埋め込んだデータの列で、自分が送りださない部分のデータはダミーのデータを埋めておく。そのために、データ線にはそのデータがダミーである事を表すフラグをもうける。スレーブから送られてきたデータは各合流点でマージされて送られていく、この時ダミーのデータ同志はダミーのデータを、ダミーのデータでないデータが来たときはそのデータを送るようにする。また、各分岐点で複数のスレーブあるいはバッファの出力を制御は、自分がダミーでないデータを送りだすときに出力をオンにすることによって行う。またフラグのビットは他のデータとは別にオアをとって上位のバッファにわたしている。

以上の動作を図4で説明する。この例では4台のセルがスレーブとなってマスタであるホストにデータをおくる。各セルが送るデータは(a)のセルの箱に入っているデータで、それぞれ番号の順番でホストに転送される。初めにC1以外のセルはダミーのデータを送りだす(b)は、バッファB1, B2に2つのデータが送られた状態である。B2はC3, C4ともにダミーを送りだしているのダミーのデータ(図中の○)が入っているが、B1にはC1のデータが入っている。このあとB1のデータがホストに送られると、B2のなかのダミーも消えるのでその開いた分にC3からのデータが送られる(c)(d)。また、データを送り終わったC1もダミーのデータを送り続ける。

4. 性能予測

Bネットのスキケット、ギャザの機能について、従来方式のブロードキャスト(グループキャスト)で同様の事を行った時との性能比較を行った。

ブロードキャストで同様の処理を行う場合は、ブロックの切りだしや統合をホスト計算機で行わなければならない。しかし、CAP-IIではMSC(DMAC)にメモリ上のとびとびの領域のデータを転送する機能があり、データの送受信にこの機能を利用するとホスト上でのブロックの切りだしを行わなくてもデータを転送できる。ここでは、ブロックの切りだしをMSCで転送時に行った場合(ブロードキャスト①)と、ブロックの切りだしをホスト上でソフトで行った場合(ブロードキャスト②)について検討を行う。

また、転送時の性能にはBネットよりもセルのローカルバスの転送速度(40MB/秒)が支配的であるので、式(1)での転送速度としては、ローカルバスの転送速度をもちいる。また、スキケット、ギャザ時のBIFの設定時間はレジスタへのロード時間のみ、として10メモリサイクル1.6マイクロ秒、MSCの設定時間も同様に1.6マイクロ秒とした、また、ソフトでのブロックの切りだしの時間も1ワード当たり、2メモリサイクル0.32マイクロ秒とした。

また、スキケット、ギャザ時の転送時は初期設定を最初に1度だけやれば良いので、転送時間は、

$$T = T_c + T_t \times n \quad \text{————— (2)}$$

また、ブロードキャスト②ではソフトで切りだしを行う時間が加わるので、

$$T = T_c \times p + (T_t + T_w) \times n \quad \text{————— (3)}$$

T_w:1 ワードあたりのデータ
切りだし時間

100 ワードのデータを転送したときの転送時間について性能予測をした結果を、表 3 にスカッタについて、表 4 にギャザについてしめす。

表 3

セル台数	スカッタ新方式	ブロードキャスト①	ブロードキャスト②
64	1.04	1.13	4.40
1024	1.04	2.66	5.93

表 4

セル台数	ギャザ新方式	ブロードキャスト①	ブロードキャスト②
64	1.04	1.13	4.40
1024	1.04	12.49	15.76

(単位ミリ秒)

本方式では、セル切替えのオーバーヘッドがないので、台数が増えても転送時間は増加しないが、従来方式では台数が増加するにつれて、転送時間が増大しているのがわかる。特に、ギャザを従来のブロードキャストで実現した場合、転送の切替えをホストの指示によっておこなっている。したがって、セルの切替え時間はホストの切替え命令がセルに届く時間と、その命令を受けてからデータがホストへ届く時間の和なので、B ネットの伝播遅延が支配的になっている。そのため、台数が増えたときの性能の低下が著しい。

5. おわりに

高並列計算機CAP-IIのブロードキャストネットワークについて述べた。本ネットワークによって、ホストとセルの間で多量のデータを高速で交換することが可能になる。

現在、我々はCAP-IIのハードウェアの開発をおこなっている。今後は、システムを完成させ実機上での性能評価を行う予定である。

参考文献

- 1) 石畑他, 高並列計算機CAP-IIの構成とメモリシステム, 本研究会投稿予定
- 2) 清水他, 高並列計算機CAP-IIのメッセージコントローラ, 本研究会投稿予定
- 3) 堀江他, 高並列計算機CAP-IIのルーティングコントローラ, 本研究会投稿予定