

# 1H-4 高並列計算機AP1000のB-Net アーキテクチャと評価

加藤 定幸 石畑 宏明 清水 俊幸 堀江 健志

(株)富士通研究所

## 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AP1000を開発している。

AP1000のように最大1024台のセル(プロセッサ)から構成される並列計算機ではホスト計算機とセル間の通信が問題になる。ホスト計算機が個々のセルと通信するには、通信するセルの切替えが必要になる。AP1000のホスト-セル間ネットワークであるBネットはこの問題を解決するためにスキヤタ、ギャザという転送モードを持っている。

## 2. Bネットの機能

セルから見たAP1000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。同時に複数のセルあるいはホストがBネットにデータを送りだすことはできない。

Bネットでのデータ転送はパケットと呼ばれる単位で転送制御をおこなう。パケットは大別して次の4種類の転送モードをもっている。

- ① ブロードキャスト
- ② グループキャスト
- ③ スキヤタ
- ④ ギャザ

ブロードキャストではネットワークに接続されたすべてのセルがデータを受信し、グループキャストでは1つ以上のセルがデータを受信する。ブロードキャストとグループキャストでは、Bネットに送りだされたデータは対象となるセルに全てのデータが取り込まれる。

スキヤタ・ギャザは1つのパケット内のデータをセルごとに分割して転送したり、それぞれのセルが送りだしたデータをあたかも1つのパケットとしてホストとして取り込んだりする転送モードである。

例えば、図1の様に2次元配列を分割してセルに割り当てて通信を行う場合、従来の方式ではセルごとにグループキャストを繰り返す必要があるが、スキヤタとギャザを使うと1回の転送でデータを送ることができる。

スキヤタ・ギャザ時のデータ分割は図2の示すようなバリエーションがあり、アプリケーションの性質に応じて適当な物を選んでしようとする。図中で同じ網が掛かっている部分ごとにセルに分割して転送する。

また、Bネット上ではホストもセルも同等に扱われるので、例えばセルから他のセルへスキヤタをおこなったりギャザを行うことも可能である。

## 3. Bネットの構造

Bネットは1024台までの拡張性と高速転送を同時に満たすために図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て実現している。

Bネットではホストあるいは任意のセルがデータを送りだすので、送信元の位置に応じてネットの構造を変化させる必要がある。具体的には送信元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、送信元の存在するツリーでは送信元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。

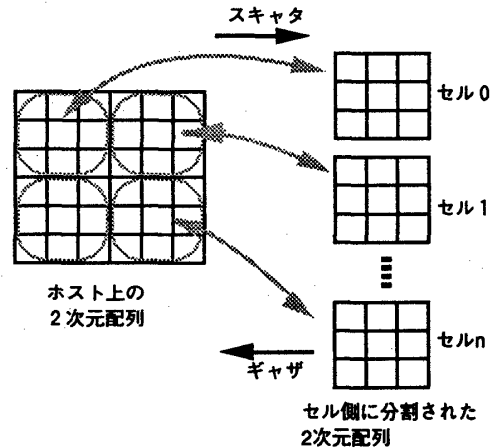


図1 スキヤタとギャザ

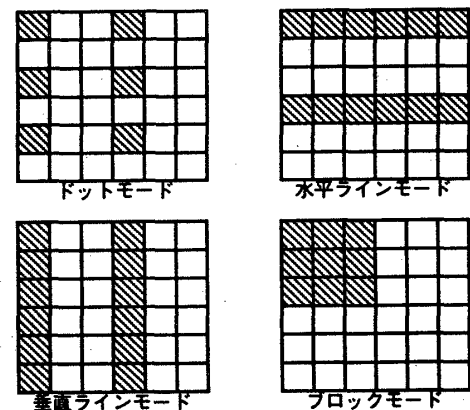


図2 スキヤタ・ギャザのモード

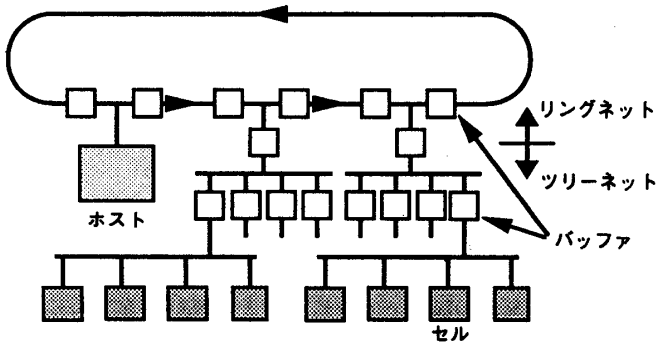


図3 ブロードキャストネットワークの構造

スキヤタ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と等しいので、送信側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。

逆にギャザでは、前述のブロードキャストやスキヤタとは反対に複数の送信元からホスト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出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出すためにセル間の順序の制御が必要になるが、この制御はデータ転送時間と等しい時間でおこなっているので、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。

4. 性能評価

16台のシステムで転送性能の評価をおこなった結果を以下に示す。転送速度は全転送データ数を初期設定の開始から最後のデータの受信までの時間で割ったものである。

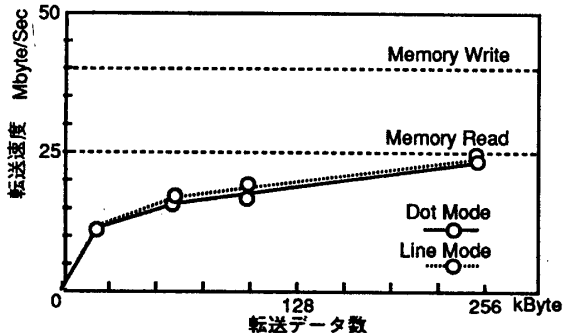


図4 スキヤタの転送速度

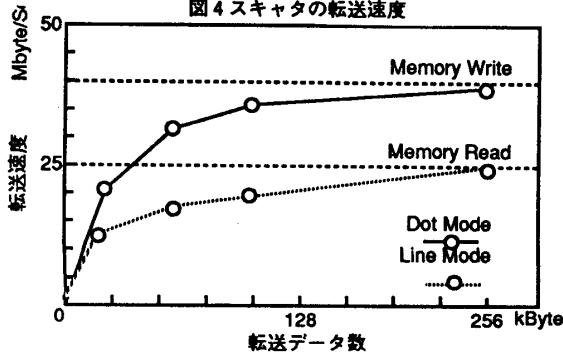


図5 ギャザの転送速度

まず、図4、図5にラインモード、ドットモードでの転送速度を示す。転送データ数が少ないほど、初期設定のオーバーヘッドの影響が大きいと鳴っている。図中の水平な線はセルのローカルバスのメモリリードとメモリアイト時の理論性能をしめしている。通常Bネットの転送レートはメモリリードサイクルによって制限されているが、ドットモードではセルでのメモリアクセスが間欠的になるので転送性能はホスト側のライトサイクルまで転送性能が向上する。

図7、8にラインモードで転送したときに、スキヤタ、ギャザと1つ1つのセルを切り換えて通信する従来方式との比較をおこなった例についてしめす。従来方式ではセルが切り替わるたびに、初期設定をしているので、このオーバーヘッドの為に見かけ上の転送速度が低くなってしまふ。特にホストにデータを集める時には、セルが変わる度にBネットの転送方向を切り換えているために見かけ上の転送速度が大きく低下していることがわかる。

5. まとめ

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。

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. 情処研報 90-ARC-83-37
- 2) 加藤他, 高並列計算機CAP-IIのブロードキャストネットワーク. 情処研報 90-ARC-83-39

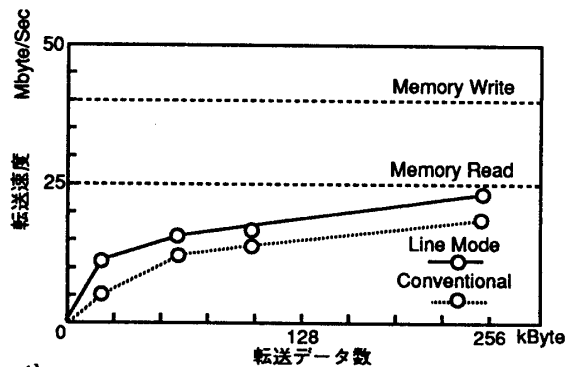


図6 従来方式との比較(スキヤタ)

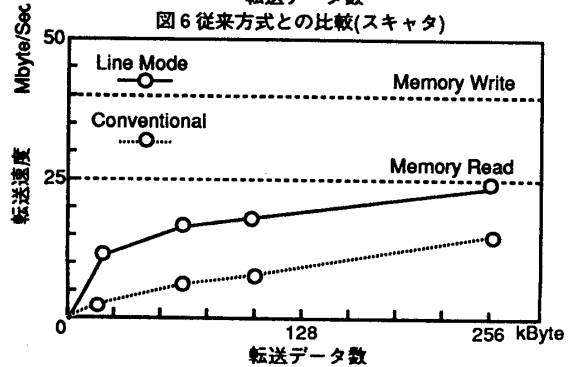


図7 従来方式との比較(ギャザ)