

ハイパクロスバ・ネットワークにおける並列ソート処理*

1T-5

板倉憲一、朴泰祐、中村宏、中澤喜三郎†

筑波大学 電子・情報工学系‡

{itakura,taisuke,nakamura,nakazawa}@arch.is.tsukuba.ac.jp

1. はじめに

MIMD方式の大規模並列システムにおいてデータベース処理のような非数値処理を行なう場合、プロセッサ(PU)間の通信は不均一、広範囲なものになる場合が多い。

ハイパクロスバ・ネットワーク(HXB)[1]はクロスバスイッチを用いた多段のネットワークである。特に3次元のHXBは比較的小規模なクロスバスイッチを用いて平均距離及び最大距離が短く様々な通信パターンに対応し数千ノード程度の並列処理システムに有効である。

本研究ではHXBの特性を確認するために、非数値処理の例として並列ソート処理の性能評価を行なう。並列ソートのアルゴリズムとしては代表的なbubble、bitonic、radixの3種類を用いる[2][3]。また、他のネットワークアーキテクチャとしてHyper Cube(HC)、Torus Mesh(TM)との比較も行う。

2. ハイパクロスバ・ネットワーク

HXBはHCを拡張したBase m n -cubeをさらに拡張したネットワークである。 n 次元のHCは 2^n 個のPUを結合するが、Base m n -cubeは m^n 個のPUを結合する。各PUは n 本の腕を出し、 m 進数で表記されたPU番号が1桁だけ違う m 台のPUにクロスバスイッチを介してつながる。ハードウェア量はPU数 $P = m^n$ に対して $m \times m$ のクロスバが $m^{n-1}n$ 個必要なのでクロスバのスイッチ数は $mP \log_m P$ 個となる。

HXBは各次元に対して独立な m をとれるようにBase m n -cubeを拡張したものである。3次元のHXBでは、 x, y, z 方向のPU数をそれぞれ X, Y, Z 個とすると、全PU数は $X \cdot Y \cdot Z$ となる。3次元のHXBは3次元のHCのリンク上にPUを全 X (または Y, Z)個配置し、それらの間をクロスバ結合したトポロジになっているので、HCの特性を包含しており同様に3次元TMの特性も包含している。クロスバのスイッチ数は $XYZ(X+Y+Z)$ 個となる。

3. 並列ソートアルゴリズム

ここでは、3つの並列ソートアルゴリズム、bubble、bitonic、radixについてネットワークの転送パターンを中心に述べる。各アルゴリズムの処理時間はstep数、step当たりの通信時間、step当たりのPU処理時間から求まり、bubbleは通信時間が小さく、bitonicはstep数がやや小さい。radixはstep数が極めて小さいが通信時間やPU処理時間が大きいという特性を持つ。PUの台数を P 、総データ数を N 、PU当たりのデータ数を

$n(=N/P)$ とする。処理の前提としては各プロセッサの持っているデータ数 n は等しく、ソート後も全プロセッサは同じデータ数 n を持つものとする。また、bubble及びbitonicのアルゴリズムでは初期状態としてPU内のデータがソートされている必要がある。これはPU毎に逐次クイックソートによって一回だけ行なう。また、各アルゴリズムの処理時間のオーダを表1に示す。

3.1 bubble sort (odd-even sort)

並列化bubble sortは、全プロセッサを1次元方向に並べ、隣同士でデータの交換を行う。ネットワークの転送パターンは1次元の隣接転送のみでありどのようなネットワークでも無衝突転送が可能である。処理step数は P であり、各stepは2フェーズからなる。最初のフェーズで各PUは自分の $n/2$ のデータをマージソートし、次のフェーズでは隣接PU間で $n/2$ のデータを交換し合いマージソートをする。

3.2 bitonic sort

並列化bitonic sortで p 個のPUから成るbitonic列をソートするためには、1回のshuffle転送と1回のbutterfly転送($\log_2 p$ 回のexchange)と $\log_2 p + 1$ 回のマージソートを行なう。処理全体のステップ数は $\log_2 P$ である。butterfly転送のパターンはHC及びHXBでは無衝突転送が可能である。shuffle転送のパターンはHCであってもHXBでもそのまま転送するとネットワーク中の経路で衝突する。しかし、HXBでは最大でも無衝突転送の約2倍の時間で転送は完了する。このことは計算機シミュレーションで確認した。TMではshuffle転送はもちろん、butterfly転送においても経路上で衝突し、回線の仮想化をしていないとデッドロックを起こす。

3.3 radix sort

並列化radix sortはマスタ・スレーブ方式で行なう。マスタは各スレーブから送られてくる部分インデックスを集め、グローバルインデックスを作成する。各スレーブはこれを元に適当な相手とデータ交換を行なう。従って転送パターンは部分インデックスのための P 対1のcollection、グローバルインデックスのための1対 P のbroadcast、各スレーブが持っているデータを担当スレーブに転送する P 対 P のcomplete exchangeがある。HXB上で各PUのデータ転送量を D としたときの転送時間のオーダは以下ようになる。collection転送は全PUが一斉に転送を行ない、マスタは来た順に処理を行なえば $O(P \cdot D)$ の転送時間となる。HXBではbroadcast転送が容易に行なえる[1]なのでその転送時間は $O(D)$ 、complete exchangeは各PUが相対的に位置の等しいPUを相手に送信すると無衝突転送を P 回行なうので $O(P \cdot D)$ の転送時間となる。また、PUの内部

*Parallel Sorting on Hyper-Crossbar Network

†Ken'ich ITAKURA, Taisuke BOKU, Hiroshi NAKAMURA, Kisaburo NAKAZAWA

‡Insitute of Information Sciences and Electronic, University of Tsukuba

処理は基数で分類した数え上げ ($O(n)$) とグローバルインデックス作成 ($O(P \cdot 2^r)$) である。

4. 評価と考察

ここでは HXB の並列計算機を想定する。まず PVM 上で並列プログラムを作成し正当性を確認した後、1つのPUに対応するプログラムからネットワークに関する部分を取り除きワークステーションで実行させてCPU時間を測定する。ネットワークの転送時間についてはシミュレータ [4] を用いて解析を行なう。ネットワークの性能は転送立ち上げ over head を 5μsec、throughput を 200Mbyte/sec とした。データは 1M 個の整数 (4byte) とする。PU 台数を増やした時の PU 処理時間を図 1 に、転送時間を図 2 にそれぞれ示す。

まず、PU 内部処理時間を見ると bubble sort はプロセス台数が増えると n が減るので処理時間は減少するが絶対的な時間は遅い。bitonic sort は全体的に良い性能を示す。これは step が $O((\log_2 P)^2)$ であるためである。radix sort については今回の実装方法では PU 処理時間のうち 2^r に依存する部分が大いので n がある程度大きくないと不利である。この点については改良の余地があると考えられる。 n が P に対して十分大きい時には処理時間は $O(1/r)$ となる。また、 P 少ない時の bitonic sort 処理時間のほとんどは初期のクイックソートの時間である。

データ転送については bitonic sort の毎回のメッセージ長が N-half を越えているので転送時間は回数ではなく転送データ量に依存する。これに対して radix sort は短いメッセージを多く転送するので回数に依存する。PU 数が 500 以上の場合には bitonic sort は転送メッセージ量が減ることで速くなるが、radix sort は転送回数が増えることで遅くなっている。ネットワークの性能によって全体の時間は変わってくるが多くのパラメータで実験したところ PU 数が 500 以上の時には bitonic sort がもっとも良い性能を示すことがわかった。

ネットワーク別に考えると HXB は bitonic sort の shuffle 転送を最も効率良く行なえる。これと同様のことを 1024PU の HC で実現するためには WormHole のようなルーティングが必要で大規模システムではハードウェア的に実現は難しい。TM と比べると bubble sort の隣接転送以外は HXB の方が高い性能を示すので [1][4]、bitonic sort と radix sort は HXB の方が転送時間が短くなる。

5. おわりに

本研究では HXB を用いた並列計算機上での並列ソートの処理性能を評価した。用いたアルゴリズムは bubble sort, bitonic sort, radix sort である。問題規模を一定とし PU 台数が 32 から 1024 の範囲で PU の内部処理時間とネットワークの転送時間に分けて処理時間を求めた。PU 数が数十台の規模の時には radix sort が有効であるが、500 台以上の規模では bitonic sort が有効である。また HXB では bitonic sort で現れる shuffle 転送と butterfly 転送を効率的に処理できるので、500 台以上の PU を用いたマシンで並列ソート処理を行なう場合に極めて有効なネットワークトポロジであることが分かった。

謝辞

本研究に当たり、貴重な御意見をいただいた筑波大学西川博昭助教授に深く感謝します。

参考文献

- [1] 齊藤 哲也 他「超並列計算機のネットワークの実現可能性と性能評価」、情処研報 92-ARC-95-4
- [2] 岡田 英明 他「超並列計算機における各種ソートアルゴリズムの実装と評価」、情処研報 92-ARC-97-12
- [3] Selim G. Akl 著、阿江忠、山下雅史、相原玲二 訳「並列ソーティングアルゴリズム」、啓学出版
- [4] 石橋 規子 他「大規模並列処理ネットワークにおけるランダム転送性能の評価」、情処全大 46 回

表 1:各アルゴリズムの処理時間のオーダー

	step 数	転送時間	PU 処理時間
bubble	P	$O(n)$	$O(n)$
bitonic	$1/2((\log_2 P)^2) + \log_2 P$	$O(n) + O(\log_2 n)$	$O(\log_2 n)$
radix	d/r	$O(P \cdot R) + O(P)$	$O(n) + O(P \cdot R)$

d :key の bit 数、 r :基数の bit 数、 R :基数 (2^r)
bubble, bitonic は 1 回の内部ソート $O(n \log n)$ を行なう。

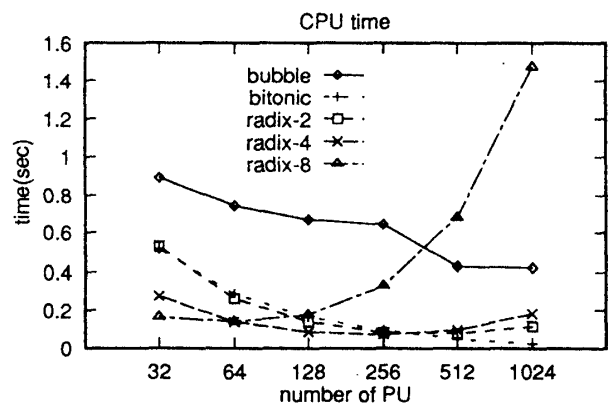


図 1:総データ数固定の PU 処理時間

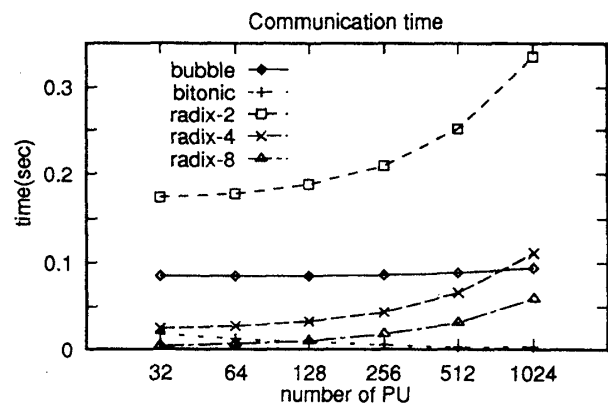


図 2:総データ数固定の通信時間