

細粒度並列処理向け相互結合網 HTN の適応ルーティング

†的山 和也, †三浦 康之

湘南工科大学 工学部情報工学科

● 1. はじめに

近年 LSI の大集積化、三次元 LSI の登場により、LSI 上に多数の PE を配置し、PE 間をネットワーク結合する技術が注目され、さまざまな三次元 LSI 向けの相互結合網が提案されている[3]。これらのルーティングアルゴリズムとして、一般的には固定ルーティングが用いられているが固定ルーティングは適応ルーティングと比べて、チャンネルの利用率が劣り、大きなデータ通信の際に遅延を起ししやすい。我々は、低遅延な適応ルーティングの実装法を研究中である。更なる適応ルーティングアルゴリズムの充実により、HTN の大幅な性能向上が見込める。

本研究では、細粒度並列処理向けの相互結合網 HTN の適応ルーティングを高性能なシミュレーションプログラムを用いて、性能評価、検証を行う。これまでに HTN の固定ルーティングの実装と評価が完了したので、本稿ではそれらの結果を示すと同時に適応ルーティングアルゴリズムの提案を行う。

● 2. HTN(Hierarchical Torus Network)

2.1 構造

HTN[1]は基本モジュール(Basic Module: BM)が $(m \times m \times m)$ の三次元トーラス網の構造をしており、BM を複数個用いて $(n \times n)$ の二次元トーラスと相互結合させた構造をしている。本研究では $m=n=4$ としている。HTN を作成するにあたっては、上位レベルのネットワークとして BM の一番上のレベルに位置する PE のみを使用する。実際の HTN の構造を図 1 に示す。

2.2 HTN の固定ルーティング

HTN の固定ルーティングアルゴリズムを図 2、図 3、図 4 に示す。

固定ルーティングは基本的に最上位レベルから最下位レベルまで、レベル順に転送が行われる。

HTN の固定ルーティングでは、始めに上

位レベルの二次元トーラスのルーティングを行って目的の PE のある BM まで転送を行い、その後 BM 内の転送が行う。

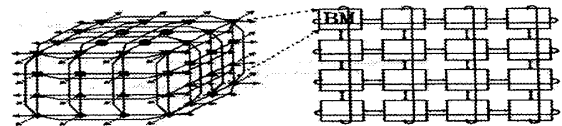


図 1 HTN 構造

```

BM_ROUTING(src, dest){
  if(src-z != dest-z){
    tag = (dest-z - src-z);
    if((tag < -2)||((tag > 0)&&(tag <= 2)))
      return z+;
    else
      return z-;
  }
  else if(src-x != dest-x){
    tag = (dest-x - src-x);
    if((tag < -2)||((tag > 0)&&(tag <= 2)))
      return x+;
    else
      return x-;
  }
  else if(src-y != dest-y){
    tag = (dest-y - src-y);
    if((tag < -2)||((tag > 0)&&(tag <= 2)))
      return y+;
    else
      return y-;
  }
}

```

図 2 BM ルーティング

```

nextroute(src, dest){
  if(src == dest)
    return 0;
  if(src-z != dest-z){
    tag = (dest-z - src-z);
    if((tag < -2)||((tag > 0)&&(tag < 2)))
      if((tag == 2)&&(src-x >= 2))
        if((tag == -2)&&(src-x >= 2))
          if((src-x >= 12)&&(src-x <= 15))
            return z+;
          else
            BM_MESHROUTING(src, (src-x-4) + 12);
      else if(src-z <= 0)
        return z-;
    else
      BM_MESHROUTING(src, src-x-4);
  }
  else if(src-x != dest-x){
    tag = (dest-x - src-x);
    if((tag < -2)||((tag > 0)&&(tag < 2)))
      if((tag == 2)&&(src-x >= 2))
        if((tag == -2)&&(src-x >= 2))
          if((src-x >= 8)||((src-x == 7)
            ||(src-x == 11)||((src-x == 15))){
            return x-;
          }
      else
        BM_MESHROUTING(src, x-1);
    else
      if((src-x >= 0)||((src-x == 4)
        ||(src-x == 8)||((src-x == 12))){
        return x+;
      }
      else
        BM_MESHROUTING(src, x);
    else
      BM_ROUTING(src, dest);
  }
}

```

図 3 上位レベルでのルーティング

```

BM_MESHROUTING(src, dest){
  if((src-z != dest-z){
    tag = (dest-z - src-z);
    if(tag > 0)
      return z+;
    else
      return z-;
  }
  else if(src-x != dest-x){
    tag = (dest-x - src-x);
    if(tag > 0)
      return x+;
    else
      return x-;
  }
  else if(src-y != dest-y){
    tag = (dest-y - src-y);
    if(tag < 0)
      return y-;
    else
      return y+;
  }
}

```

図 4 BM MESH ルーティング

図 2 は BM 内での転送を行うルーティングである。src-z、src-x、src-y はそれぞれ現在地の PE 番号を示す。dest-z、dest-x、dest-y は行き先の PE 番号を示す。'z+'と書かれたものは z 軸の上方向に進むと解釈する。'z-'と書かれたものは z 軸の下方向に進むと解釈する。他も同様である。

図3では上位レベルのネットワークである二次元トーラスでのルーティングである。src-z1、src-x1はそれぞれ上位レベルの二次元トーラスの現在のPE番号を示す。dest-z1、dest-x1は行き先のPE番号を示す。図4は上位レベルでの移動を終えて、目的BMにたどりついた後に行うルーティングである。

このようにBM内だけでの通信と上位レベルでの移動を行う時の2つに分けて行っている。

● 3. 実験環境

3.1. シミュレーションプログラムの構成

本研究で用いたPEの構成はノードプロセッサ、送信用ネットワークインターフェイス、受信用ネットワークインターフェイス、そしてルータとなっている。ルータの中身はFIFO、制御回路、デマルチプレクサ、マルチプレクサ、クロスバスイッチとなっている。これら全部を1つのPEとして扱う。図5はPEの構成図となっている。

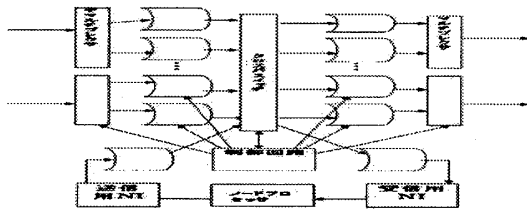


図5 PEの構成

3.2. シミュレーションプログラム環境設定

本プログラムでのPE数は1024となっている。HTNではBM(4×4×4)が16個並んでいる。MESHではBM(32×32)となっている。

● 4. 実行結果

本実験ではパケットの送信はランダムで行っている。HTNの固定ルーティングでシミュレーションした結果を図6に示す。

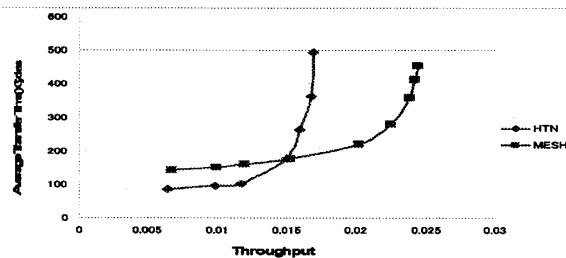


図6 実行結果グラフ

図6の結果よりMESHの固定ルーティングとHTNの固定ルーティングの最大スループットがMESHの方が約1.4倍の結果となっている。HTNのような階層型相互結合網では、上位レベルリンクが混み合うことから、上位レベルリンクに対して適応ルーティングを用いることにより、通信性能の大幅な向上が期待できる。

● 5. HTNの適応ルーティング

5.1 Turnモデル

デッドロックは結合網内のバッファが論理的に循環構造を作るため起こる。

Turnモデル[2]ではパケットがルーティング中に方向を変えるパターンに制限を加え、循環させないようにするものである。このモデルは論理的な循環構造に着目し、結合網に依存しないのが特徴である。これにより、故障地点や混雑地点を迂回する適応ルーティングが可能となる。

5.1.1 NF(North Fast)法

本研究ではHTNの適応ルーティングとしてNF法のHTNへの実装を検討中である。HTNでもっとも混雑が予想されるのが上位レベルの転送部分である。そこで、上位レベル転送についてNF法の実装を考える。図7にNF法のTurnモデルを示す。NF法では図中の上に行く動作に制限を加えることでデッドロックを回避する手法である。NF法の動作は目的PEが、自分よりも下か同じ高さにいる時経路を自由に選べ、目的PEが自分よりも上にいる時は固定ルーティングを行う。現在NF法はBM内リンクのうち上位レベル転送に使う部分についてのみ実装を始めている。そのため、実際にNF法を適用しているのは上位レベルでのMESH網だけである。今後上位レベル部分の実装を終えた後は下位レベルに対しての実装も考えている。



図7 TurnモデルNF法

● 6. 今後の予定

今回の実験により、HTNの固定ルーティングの実装と評価は完了した。今後、NF法用いたルーティングを作成し、動作確認後、性能評価および検証を行う。また、他の適応ルーティングアルゴリズムの実装を考えており、三次元LSI向けの階層型相互結合網DDR(Dynamic Dimension Reversal)法[3]の性能評価、検証を行う予定である。

● 参考文献

[1] M.M. Hafizur Rahman, Hierarchical Interconnection Networks for Massively Parallel Computers
 [2] G. J. Glass and L. M. Ni, "Maximally Fully Adaptive Routing in 2D Meshes," ISCA92, pp. 278-287.
 [3] 三浦康之, 堀口進, 福士将, 細粒度並列処理向け相互結合網 TESH における適応型ルーティングアルゴリズム