

直接网络中的自适应路由算法分析*

刘燕 杨晓东 王志英

(国防科技大学计算机系 长沙 410073)

摘要 互连网络是大规模并行计算机的重要组成部分,路由算法是其中决定网络性能的重要因素,本文在直接网络结构基础上对路由算法进行讨论,给出了一种分类方法,并着重对采用虫孔路由由开关技术的自适应路由算法进行分析,为进一步的评价和设计新的算法提供了参考。

关键词 大规模并行计算机,直接网络,自适应路由算法,前进型路由算法,后退型路由算法

分类号 TP393, TP301

Adaptive Routing Algorithms in Direct Networks

Liu Yan Yang Xiaodong Wang Zhiying

(Department of Computer Science, NUDT, changsha, 410073)

Abstract Interconnect network is an important part of massively parallel processors (MPP), and routing algorithm constitutes the primary factor influencing on the performance of it. In this paper, we discuss the routing algorithms for direct networks, and study the wormhole-routed adaptive algorithms in detail. Finally, we give some available points to design and evaluate new algorithms.

Key words massively parallel processors, direct networks, adaptive routing algorithms, progressive routing algorithms, backtracking routing algorithms

大规模并行计算机(MPP)被认为是最有潜力达到万亿次计算能力的技术。MPP是把大量的由微处理器、存储器和必要的控制逻辑等组成的处理结点,经一定拓扑结构的高带宽低延迟的互连网络连在一起实现高度并行操作以获得高速度的。其中,互连网络与系统的性能密切相关,且网络占整个机器系统的成本和功耗的很大比例,是MPP中的关键部件。

一个互连网络通常由拓扑结构、开关技术、流量控制和路由算法四方面来表征,其中路由算法描述了信息在网络中如何选取路径,其效率对网络性能起着很关键的作用,因此对其研究极为重要。路由算法可分为确定性和自适应路由算法两种,确定性算法实现简单,已在很多商用MPP中得以实现;而自适应路由算法正成为新一代MPP系统的采用对象,如Cray T-3E系统采用了完全自适应路由算法。本文在直接网络的基础上,对自适应路由算法进行一个系统的分类,然后针对每类算法进行分析,并着重对采取虫孔路由由开关技术的自适应路由算法进行研究和分析。

1 自适应路由算法及分类

自适应路由算法与预先唯一确定路径、不受网络状态影响的确定性路由算法(如E-cube)不同,它对于一对源和目的结点,视网络的工作状态,可有多条选取的路径,因而有避免死锁、提高网络的带宽利用率和网络容错能力的好处;另外,自Dally提出虚通道概念之后,采用虚通道的自适应路由算法实现起来更为经济和灵活,从而使自适应路由算法得到了极大的发展。一个好的路由算法应有三个特

* 国家 863 计划和九五国防预研基金资助项目
1997 年 2 月 27 日收稿
第一作者:刘燕,女,1971 年生,博士生

实现起来更为经济和灵活,从而使自适应路由算法得到了极大的发展。一个好的路由算法应有三个特点:低通讯延迟、高网络吞吐率和VLSI工艺上的易实现性,近年来很多算法又将容错作为一项重要指标,旨在寻求自适应性、性能和容错之间的最佳平衡。自适应路由算法的一种分类如图1所示。在分类第一层,算法分为虫孔路由型和非虫孔路由型,其中前者是指在虫孔路由开关技术下采用的算法,而后者是在其它路由开关技术下(如存储转发、线路交换、虚跨步、流水线路交换PCS等)下采用的路由算法。在分类第二层,算法分为前进型和后退型,前者指“向前”传送信息、无后退回溯能力,后者指算法系统搜索网络、需要时可后退回溯。虫孔路由型由于受其流水特性限制无法回溯,全部为前进型。在分类的第三层,算法分为最小路由算法(路由总是选取从源到目的结点间最短的路径)和非最小路由算法(路由针对网络当前状态,允许信息沿一长路径进行传送)。在分类的最底层,算法分为完全自适应的和部分自适应的路由算法。前者可选取所有路径;而后者由于在路由上加以限制,仅有部分路径可选。

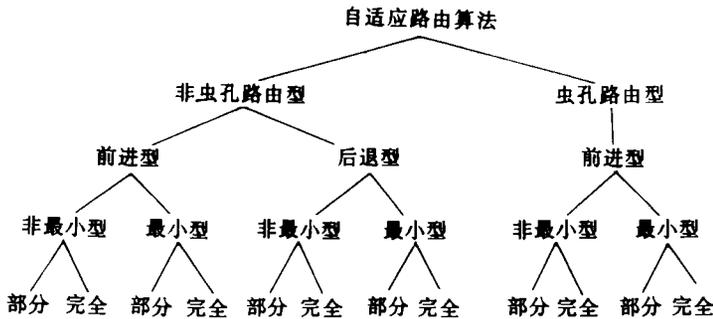


图1 自适应路由算法的一个分类

2 非虫孔路由型算法

在非虫孔路由算法中,前进型与后退型路由算法的区别在于当中间结点没有可选的通道空闲时采取的不同措施。前进型会等待、放弃或绕道到其它路径上去,而后退型本着“寻找其它路径比等待一条路径变为可用要好得多”的原则,会释放一些通道并后退回溯到先前的某结点处再在网络的其它部分寻找路径。在前进型的算法中,较为典型的有采用线路交换的基于已走步数的算法、A1算法、Idle算法等,该类算法类似于后面所述虫孔路由型算法,故在此不作讨论。在后退型算法中,有代表性的有采用最小路径进行深度优先搜索路由的EPB算法、采用最小和非最小路径进行深度优先搜索路由的EMB算法、将二者结合起来的TPB-u算法和采用PCS开关技术且路径中仅允许有小于等于 m 个非最小通道的MB-m算法等。后退算法容错性能较好,但寻径头、延迟及附加成本都相对较大、且需解决活锁问题,复杂性较高,因此需在容错、性能和成本之间进行充分考虑后决定是否采用这种算法。这类算法一般适用于对容错要求较高的系统中。表1给出了对 k 元 n 方体网络,几种后退型路由算法相对于E-cube算法的寻径头大小比较。

3 虫孔路由型算法

由于虫孔路由开关技术的广泛采用,虫孔路由型算法是目前自适应路由算法中研究和采用得最多的一类,在直接网络中尤其如此,最有代表性的为Linder和Harden提出的 $2P_n$ 算法^[1],Glass和Li提出的基于拐弯模型(Turn model)的算法^[2],Chien和Kim提出的平面自适应路由算法^[5],Berman和Gravano提出的 $*$ -channel算法^[6]等,它们分别代表了以下三类很有研究和实现意义的算法。

表 1 几种算法的寻径头比较

算法	类别	寻径头大小
E-cube	确定性算法	$n \log_2 k$
EPB	完全自适应的最小后退算法	$n \log_2 k + 1$
EMB	完全自适应的非最小后退算法	$n \log_2 k + 1$
TPB-u	部分自适应的非最小后退算法	$n \log_2 k + 1$
MB-m	部分自适应的非最小后退算法	$n \log_2 k + n \log_2 m + 1$

3.1 虚通道数大的完全自适应路由算法

这类算法通过对每个物理通道采用大量的虚通道而达到无死锁和完全自适应的双重目的。以 n 维 mesh 结构上的 $2Pn$ 算法为例。算法中为每个物理通道设 2^n 条虚通道, 每条物理通道的虚通道数可用一个 n 位二进制数来表示。当从源 $s = s_{n-1} s_{n-2} \cdots s_1 s_0$ 向目的结点 $d = d_{n-1} d_{n-2} \cdots d_1 d_0$ 传送信息时, 算法为信息设置一个如下的标志 $t = t_{n-1} t_{n-2} \cdots t_1 t_0$:

$$t_i = \begin{cases} 1 & \text{if } s_i < d_i \\ 0 & \text{if } s_i > d_i \\ 0 \text{ or } 1 & \text{if } s_i = d_i \end{cases}$$

当消息从源结点前进了 i 步到达中间结点, 则其占用该中间结点的序号为 t 的虚通道, 将信息传向目的结点。该算法是完全自适应的最小路由算法。

$2Pn$ 类算法的突出优点是自适应度高, 但所需的大量虚通道使硬件增加, 尤其当网络维数较大时成本很高; 其次, 由于该算法每步都基于当前结点可用的局部信息来路由信息, 故所选路径并不一定全局最佳, 如在均匀和热点流量模式下, 其性能甚至还不如 E-cube 算法。该类算法仅适用于低维小规模网络上某些流量模式非均匀的应用。有算法对其进行改进, 如每步基于某类优先信息(如已走步数)来路由信息, 使性能有所提高^[7]。但巨大的虚通道需求量促使更多的研究者转入开发成本量较低的自适应路由算法。

3.2 虚通道数少(或无)的部分自适应路由算法

该类算法是实际实现中采用较多的自适应路由算法, 它在硬件成本与自适应性间进行折衷, 将自适应性限制在某部分, 从而减少了避免死锁所需增加的硬件。

首先以静态限制路由自适应性的基于拐角模型的算法为例。拐角模型不是基于增加物理或虚通道, 而是基于分析网络中信息可转弯的方向和转弯所形成的环路, 通过强制算法在路由消息过程中不出现某些转弯而阻止网络环路的出现, 从而达到避免死锁的目的。以其中自适应性最好的负优先算法为例, 在 2 维 mesh 结构中, 设定西和南方负向、东和北为正向, 该算法禁止从正向到负向的拐弯(即禁止北到西、东到南两个转弯)。当从源结点向目的结点传送信息时, 算法要求首先在西和南方向上自适应地路由(如有必要), 再在东和北两个方向上自适应地路由, 直至传向目的结点。该类算法是部分自适应的非最小路由算法。

基于 Turn model 的算法的突出优点是对硬件逻辑要求较简单、成本较低, 无需增加虚通道即可达到无死锁和部分自适应性。负优先算法对于某些非均匀流量模式性能很好, 如对矩阵转置其性能优于确定性算法 2 倍。但该类算法由于偏重于某些通道, 打破了流量的均衡性, 易使网络过早进入饱和状态, 如对于均匀流量模式其性能甚至还不如确定性路由算法。因此该类算法仅对某些非均匀流量模式的应用性能较好。为使流量均衡, 很多研究者又在对该类算法进行改进, 提高其在其它流量模式下的性能^[4]。

平面自适应路由算法代表了动态限制路由自适应性的部分自适应路由算法。以 n 维 mesh 结构上的算法为例, 算法中为每个物理通道设 3 条虚通道, 将整个网络分成 $A_0 A_1 \cdots A_{n-1}$ 共 n 个自适应平面, 其中每个平面 A_i 仅包含 i 和 $i+1$ 两维。当从源向目的结点路由信息时, 算法按从 A_0 到 A_{n-1} 的顺序进行路由, 在每个平面 A_i 中, 选择第 i 和 $i+1$ 维中任意趋近目标的通道进行自适应路由、直至某中间结点第

i 维上的坐标等于目的结点在该维的坐标。在路由经过所有自适应平面后,信息到达其目的结点。该算法中任何时刻自适应性均限制在两维,从而达到了无死锁和降低网络成本的目的。这类算法仅需固定数目的虚通道,所需硬件成本和复杂性较低,且在流量分配上较 Turn model 均匀,对多数流量模式实际性能均较高,所以适用性较强,但其仍存在自适应性受限的缺点,采用时尚需针对具体应用进行相应的改进。

3.3 虚通道数较少的完全自适应路由算法

为以较少虚通道数达到完全自适应性,很多研究者对死锁的充分必要条件进行研究,并在此基础上提出新的算法,Berman 和 Gravano 等提出的 * -channel 算法是其中最为典型的一个,该算法是基于“通道依赖图是动态非循环的,则可避免死锁”提出来的。以 n 维 mesh 结构上算法为例,算法中为每个物理通道设置 4 条虚通道,并将虚通道分为带 * 号和不带 * 的。并约定:在带 * 号的通道上执行维序的确定性路由,在不带 * 号的通道上执行维序路由所不允许的路由。当从源向目的结点路由信息时,在中间结点处可选择任一空闲的不带 * 号通道,而仅可选择符合维序关系的空闲的带 * 号的通道,由此将信息传向目的结点。不带 * 号的通道提供了自适应性,带 * 号的通道避免了死锁,该算法为完全自适应的最小路由算法。

* -channel 在均匀和某些非均匀的(如位反)流量模式下性能都较好,尤其在均匀模式下其性能远优于确定性算法,这在自适应路由算法中是不多见的。但其开关成本和复杂性相对较高。

表 2 给出了 n 维 mesh 结构上几个典型自适应路由算法的一个比较,其中 h_- 、 h_+ 和 h 定义为,当从源 $(s_{n-1}s_{n-2}\cdots s_1s_0)$ 向目的结点 $(d_{n-1}d_{n-2}\cdots d_1d_0)$ 路由信息时,在某中间结点 $(p_{n-1}p_{n-2}\cdots p_1p_0)$ 处,设满足 $d_i < p_i$ ($0 \leq i \leq n-1$) 的数目为 n_1 ($0 \leq n_1 \leq n$), 则:

$$h_- = n_1$$

$$h_+ = n - n_1$$

设满足 $d_i \neq p_i$ ($0 \leq i \leq n-1$) 的数目为 n_2 ($0 \leq n_2 \leq n$), 则:

$$h = n_2$$

表 2 mesh 结构上几个自适应路由算法比较

算法	类别	共享一物理通道的虚通道数	某时刻路由信息可选的维数
负优先算法	部分自适应非最小路由算法	1	若 $h_- > 0$, 则为 h_- ; 否则为 h_+
$2pn$	完全自适应最小路由算法	2^n	h
平面自适应算法	部分自适应最小路由算法	3	2
* -channel 算法	完全自适应最小路由算法	4	h

通过对大量已有算法和资料进行分析,我们认为:1) 虚通道和限制部分路由技术是避免死锁的有效方法;2) 互连网络中最昂贵的资源为物理通道所需线路,其次为缓冲、开关和其它控制电路。虚通道虽不似物理通道那么昂贵,但其所必须的缓冲队列及相关的控制逻辑仍不可避免地带来额外成本,因而其数目选择上应受成本和复杂性限制,不宜过多;3) 对某些流量模式的应用,自适应性较高的算法性能并不一定优于自适应性较低的算法性能,有些自适应算法的性能甚至低于确定性路由算法,算法的性能与具体的消息流量模式有很大关系,自适应算法设计不当可能会导致硬件成本较高、但性能较差。因此在为一个系统设计路由算法时,应充分考虑硬件成本和实际应用,选择设计性能价格比高、自适应性满足需要、实现尽量简单的无死锁路由算法,在成本、性能和实现复杂性上进行合理的折衷。

4 小结

针对直接网络中的自适应路由算法进行讨论,特别对虫孔路由型的自适应路由算法进行了分析。结果表明:1) 算法的选用和设计在具体应用密切相关(如对容错要求较高的应用,选择后退型算法好;而对延迟要求较高的应用,选择虫孔路由型算法较好);2) 自适应性并非越高越好,这与具体应用的流量模式和算法的流量分配均衡性有关,因此要根据具体应用特点选取实现后性能价格比高的算法。

进一步的工作将通过性能分析和模拟实验对算法进行定性和定量的评价，并据此提出实际性能更高的算法。

参考文献

- 1 Linder D H Harden J C. An adaptive and fault-tolerant wormhole routing strategy for k-ary n-cubes. *IEEE Trans on Computers*, 1991, 402~412
- 2 Glass C J, Ni L. M. The turn model for adaptive routing. In Proc. 19th Int. Symp. Comput. Arch, 1992
- 3 Glass C J, Ni L. M. Fault-tolerant wormhole routing in meshes without virtual channels. *IEEE Trans. on parallel and Distributed systems*, 1996, 7 (6): 620~636
- 4 Upadhyay J H *et al.* Efficient and balanced adaptive routing in two-dimensional meshes. Prasant @ iastate. edu.
- 5 Chien A A, Kim J H Planar-adaptive routing: low cost adaptive networks for multiprocessors. In Proc. 19th Int. Symp. Comput. Arch, 1992
- 6 Berman P E *et al.* Adaptive deadlock-and livelock-free routing with all minimal paths in torus networks. In Proc. 4th. Symp. on Parallel Algorithms and Architeches, 1992
- 7 Boppana R V *et al.* A framework for designing deadlock-free wormhole routing algorithms. *IEEE Trans. parallel and Distributed systems*, 1996, 7 (2): 169~183
- 8 Gaughan P T and Yalamanchili S A family of fault-tolerant routing protocols for direct multiprocessor networks. *IEEE Trans on. T. Gaparallel and Distributed systems*, 1995, 6 (5): 482~497
- 9 Gaughan P T *et al.* Adaptive routing protocols for hypercube interconnection networks, *IEEE Comput. mag*, 1993, 26: 12~23
- 10 Jose Duato *et al.* Highly adaptive wormhole routing algorithms for n-dimensional torus. *DIMACS serials in Discrete Mathematics and Theoretical computer Science*. 1995, 21: 87~103