

TR-414

GHC処理系における負荷分散方式の検討

安里 彰、岸本 光弘、

小沢 年弘、細井 聰_(吉士)

July, 1988

©1988, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03) 456 3191-5
Telex ICOT 132964

Institute for New Generation Computer Technology

GHC 处理系における負荷分散方式の検討

A Study of Load Balancing Method on GHC

安里 彰
Akira ASATO

岸本 光弘
Mitsuhiko KISHIMOTO

小沢 幸弘
Toshihiro OZAWA

坂井 雄
Akira SAKAI

林 耕司
Kouji YAMASHI

藤田 彰
Akira FUJI

富士通株式会社

FUJITSU LIMITED

1.はじめに

我々は、第五世代コンピュータプロジェクトの一環として、並列論理マシン PIM の開発を行っている。PIM は 100 万台規模の要素プロセッサを、共有バス／共有メモリによる密結合されたクラスタとクラスタ間ネットワーク結合の二階層で構成されたシステムである。PIM は、並列論理型言語 GHC (1) をベースに設計された KSL を拠言語としている。我々は PIM 上での KSL の実現に先立ち、共有メモリによる密結合並列計算機である Symmetry 上に FGHC の処理系を開発した。

処理系作成の第一の目的は、実際に並列動作する処理系を解析して様々な基礎データを採取し、効率の良い並列処理方式を確立するとともに、PIM 設計へのフィードバックすることである。更なる目的として GHC の言語機能の充実や並列データ構造の確立といった、アプリケーションレベルの研究への展開を目指している。

本稿の前半では我々の作成した処理系について概要を述べ、負荷分散の方法について説明する。

並列にプログラムを実行する場合、各プロセッサ間の負荷分散を上手に行なうことが重要なポイントとなる。負荷分散の手法はいろいろ考えられるが、それらを全て処理系に組み込んで比較を試るのは非常に時間を要し、得策ではない。そこで我々は、GHC の動作に基づいた簡単なモデルを用いて、負荷分散方式の検討を行うことにした。

本稿の後半では、まず上記モデルを提案し、それを用いたシミュレーションの方法を説明する。そして、我々が採用した方式を含む 3 種類の動的負荷分散方式についてのシミュレーション結果を述べる。更に、静的な情報に基づいて負荷分散を行った場合の性能向上について考察する。

2. 処理系の概要

2.1 GHC

GHC (Guarded Horn Clauses) は並列論理型プログラミング言語の一つであり、プログラムは次のようなガード付 Horn 節の集合である。

$H := G_1, \dots, G_m | B_1, \dots, B_n$

ヘッド部 ガード部 ボディ部

FGHC はガード部に組み込み述語のみを許すという制限を GHC に加えたものである。GHC の実行はゴールリダクションによって進行する。

2.2 ソフトウェア構成

本処理系は、FGHC ソースプログラムを中間コード (KSL の抽象命令である KSL1B 命令) に変換するコンパイラ、アセンブラーおよび、中間コードを解釈実行するエミュレータより構成される。エミュレータは UNIX のプロセスで、複数が同時に動作する。

2.3 処理系の仕様

・ゴール

ゴールはゴール・レコード (G.R.) というブロックで表現する。各プロセスはゴール・レコードをキュー管理し、順にリダクション処理を行う。(図 1) ゴール・レコードは専用の領域に割り当てる。

・データ

データはタグ付きで表現される。型は、アトム、リスト等の 6 種類をサポートしている。これらのデータはヒープ領域に割り当てる。(図 1)

・メモリ管理

ヒープ領域、ゴール・レコード領域とともに共有メモリ上にとられるが、両者は別個に管理される。ゴール・レコードはフリーリスト管理を行う。また、これらの領域はプロセス（以後 P Eと書く）毎に割り当て、競合を抑えている。ある P Eのフリーリストが尽きたと、それを他の P Eに通知し、フリーなゴール・レコードの再配分を行う。また、コピー方式の G Cを実現している。ただしこの G Cは並列に実行するものではなく、各 P Eが順に自分のルートからの G Cを行う。

・メタ機能

ゴールの成功、失敗をメタレベルで扱えるように、メタコールをサポートしている。あるメタコールの下のゴールは双方向リンクにより、リング状に結ばれている。

・サスペンション方式

サスペンションは、多重待ちの効率を多少犠牲にしても单一待ちを高速に処理できる方式を工夫した。

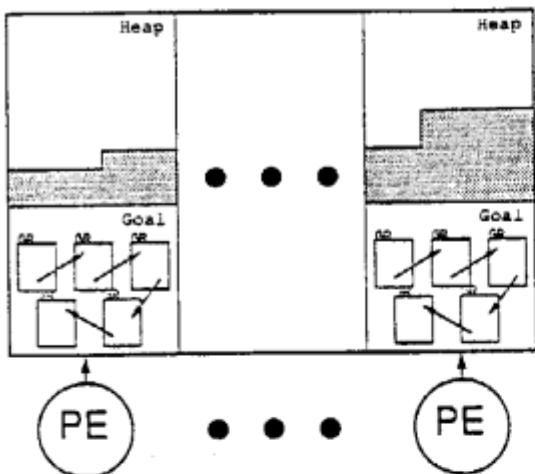


図1 共有メモリ

3. 負荷分散方式

並列システムにおいては負荷分散方式は重要なポイントとなる。負荷分散には動的負荷分散と静的負荷分散がある。前者はプログラム実行中の各プロセッサの負荷に応じて行うもので、後者はプログラム自身に何らかの情報を含ませるものである。我々の処理系では動的負荷分散のみを採用している。以下、我々の方法を説明する。

各 P Eは独自のスケジューリング・キュー（レディ・キュー）を持ち、そこにはゴールがつながっている。

また、それ以外にエキストラ・キューと呼ぶキューを設けている。エキストラ・キューは全ての P Eからアクセス可能である。各 P Eは普段は自分のキューからゴールを得てリダクション処理を行い、生成したゴールは自分のキューにつなぐが、アイドルな P Eができたり、自分のキューが長くなると、生成したゴールをエキストラ・キューにつなぐようになる。また、アイドルな P Eはエキストラ・キューからゴールを得る。この方式を用いれば、アイドルな P Eが生じても、余分なゴールがある限りは速やかにゴールを得られるので、稼働率を高く保つことができる。しかし、アイドルな P Eが頻繁に現れるとエキストラ・キューへのアクセス競合によって、稼働率が落ちる可能性もある。図2に本方式を示す。

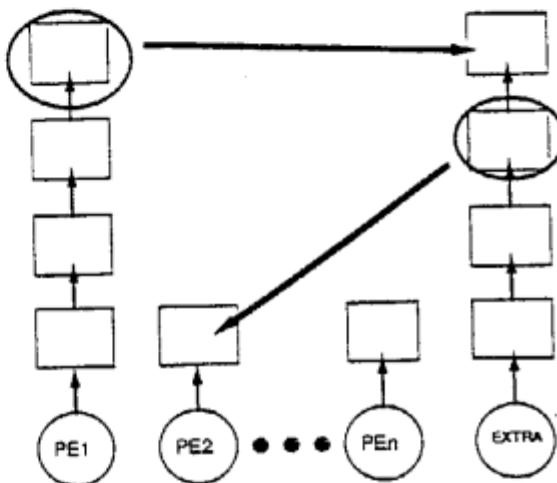


図2 エキストラ・キュー方式

4. モデルによるシミュレーション

4.1 重み分割モデル

負荷分散には多くの方法が考えられる〔2, 3〕が、それらの方法を試験するために処理系自体をインプリメントし直すのは、あまり経済的とは言えない。

そこで我々はGHCの動作に基づく、「重み分割モデル」と呼ぶ簡単なモデルを用いて、負荷分散方式の検討を行うことにした。それは次のようなものである。ゴールに重みという属性を与える。ある重みを持ったゴールはリダクションによって1だけ重みを消費し、いくつかのゴールを生成し、自分は消滅する。新しく生成されたゴール群の重みの合計は、はじめのゴールの重みから1を減じたものになっている。図3に示す。ゴールを何個生成するか、重みをどのように分割するかといった点を工夫することにより、現実のGHCプログラムの挙動を模倣することができる。

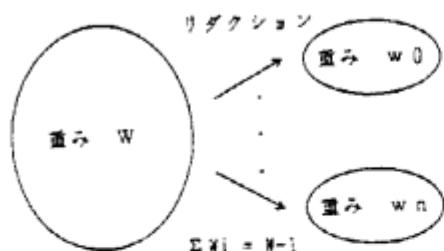


図3 重みの分配

4.2 シミュレーションの方法

本章では上記モデルを用いたシミュレーションについて述べる。複数のPEがあり、各PEは自分のキューリーを持ち、キューリーにはゴールが運ばれる。PEはキューリーの中から1個のゴールを選び、実行（リダクション）する。その結果、新たなゴールが生成され、どこかのキューリーにつなげられる。選ばれたゴールが必ず実行できる保証はなくサスペンドすることもある。また、キューリーが空になってしまったPEは何らかの手段でゴールを得ようとする。キューリーは排他制御の対象であり、同時に1個のキューリーには1個のPEしかアクセスできない。

シミュレーションは、適当な初期状態から各PEがゴールの獲得、実行、生成、分散を、ゴールがなくなまるまで繰り返す様子を1クロックずつトレースすることで行う。PEの数はパラメタとして可変にする。実行、キューリング、サスペンド等に要するクロック数は我々の処理系での値に基づいて決めた定数を用いる。測定項目は稼働率およびサスペンド回数である。

生成したゴールの分散方法、アイドルなPEのゴール獲得方法、実行するゴールの選択方法を規定するのが、負荷分散方式である。その際、各キューリーの長さや、アイドルなPEの数などを基準にするのが、動的負荷分散であり、ゴール自身の持つ情報をもとにするのが、静的負荷分散である。

4.3 プログラムのタイプ

シミュレーションに用いるプログラムのタイプは、実際のGHCプログラムを意識して、典型的な動作をするものとして、TREE型、FILTER型の2タイプを採用した。それぞれの動作を説明する。

① TREE型

あるゴールは重みがなくならない限り2個の子ゴールを生成する。同じ現ゴールから生成された兄弟ゴール間で優先順位をつけることができる。優先度の高いゴールが実行されるまでは、優先度の低いものは実行できず、もしスケジューリングされればサスペンドする。

現の持つ優先度を子が引き継ぐこともできる。その場合は、現より低い優先度のゴール（子から見れば叔父にあたる）は、子ゴールが全て終了するまで、実行が待たされる。図4のような3通りのパターン（A：兄弟間に優先順位がない場合、B：兄弟で優先度が違うが現の優先度は譲り受けない場合、C：現の優先度を譲り受けする場合）が等確率で起こるようにインプリメントした。

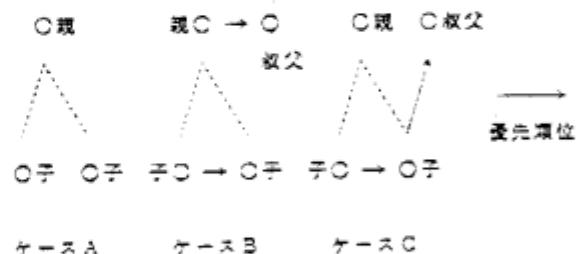


図4 子ゴール生成パターン

② FILTER型

図5のように、優先順位の定まった、一定の幅のゴール群を想定する。図で最も左のゴールは無条件で実行可能であるが、それ以外のゴールは自分の左のゴールが終わるまでは実行できない。各ゴールは1個の子ゴールを生成し、優先順位は子供に引き継がれる。幅は3:2とした。

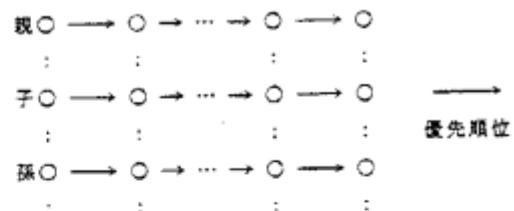


図5 FILTER型

5. シミュレーション(1) — 動的負荷分散

本章では、動的負荷分散方式を変化させてシミュレーションを行った結果を示し、考察を加える。

5.1 各方式の紹介

3通りの方式についてシミュレーションを行った。

① CQ (Common Queue)

キューリーが唯一存在し、全てのゴールはそこにつながる。各PEはゴールを生成すると、一つは自分が統けて実行するために残し、あとは全てキューリングする。

ゴールがある限り、速やかにスケジュールされるので、高い稼働率が期待できるが、キューへのアクセスが集中すると、アクセス競合が生じやすい。

② RQ (Request)

(2)で提案されていた方式。各PEが自分のキューを持つ。アイドルになったPEは、一番長いキューにメッセージを送る。メッセージを受けたPEは送り手のPEのキューに自分のキューからゴールを分け与える。この方式は、自分以外のキューへアクセスする場合でも、メッセージを受けたPEに選定されるので、キュー上のアクセス競合が生じないという長所を持つが、アイドルになったPEがゴールを得るまでにある程度の時間を要するため、稼働率の点では他方式に一歩遅れる。

③ EQ (Extra Queue)

我々が処理系で採用した方式である。第3章で説明したように、各PEのキュー以外にエキストラ・キューを1本設けてスケジューリングを行う。

5.2 シミュレーションの条件

シミュレーションは、PEの台数を4個から4個刻みで24個まで変化させ、稼働率とサスペンド数を測定し

た。リダクションに要したクロック数をT、全体のクロック数をTとすれば、稼働率は次の式で与えられる。

$$\text{稼働率 (\%)} = \frac{T}{t} / T$$

また、PEは常にキューの先頭のゴールを選ぶようにした。その他の条件は以下の通りである。

$$\text{重みの初期値} = 5000 \text{ (TREE)}$$

$$500 * 32 \text{ (FILTER)}$$

$$\text{リダクションに要するクロック数} = 10 + 10n$$

(nは生成するゴールの数)

$$\text{サスペンド処理のクロック数} = 10$$

$$\text{キュー操作に要するクロック数} = 1$$

5.3 シミュレーション結果

測定結果を示す。

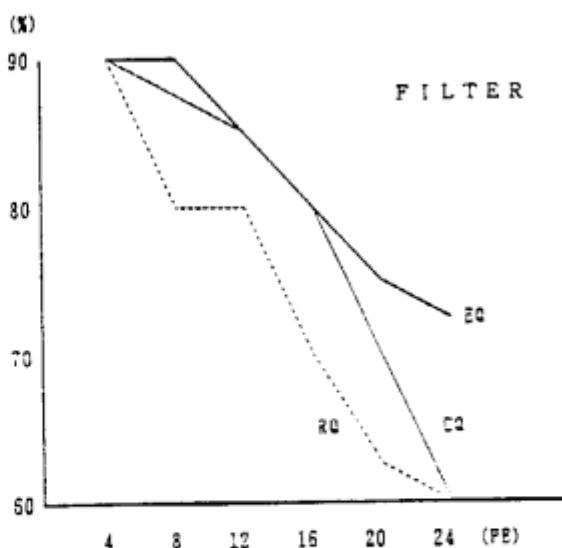
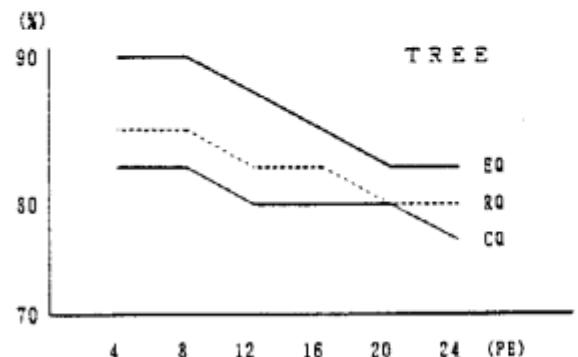


図6 稼働率 (動的負荷分散)

PE \	4	8	12	16	20	24
TREE CQ	1035	1148	1204	1242	1185	1235
EQ	704	692	691	723	785	848
RQ	714	714	692	682	700	724
FILTER CQ	386	1295	1814	2461	3542	4893
EQ	510	703	1690	2471	3644	4657
RQ	237	533	1282	1382	1517	1526

表1 サスペンド数 (動的負荷分散)

どちらのタイプのプログラムにおいても、稼働率はEQが最も優れていることがわかる。また、EQとRQのサスペンド数を比べると、TREEではほぼ等しく、FILTERではRQの方が少ない。にも関わらずEQの稼働率が高いのは、アイドル率の差が効いているからである。CQでは1本のキューへのアクセスが集中するため、本来の稼働率の高さが活かせていないと考えられる。

6. シミュレーション (2) — リザーブ優先方式

6.1 リザーブ優先方式

リザーブされたゴールは絶対実行可能である。

このことから、ゴールを選ぶ基準として、リジュームされたゴールを優先的に選ぶ方式が考えられる。これをRF (Resume First) と名付け、キーの先頭を選ぶ方式DF (Depth First)との比較を試みた。動的負荷分散方式はEQを用い、他の諸条件は前章と同様にした。

6.2 シミュレーション結果

測定結果を示す。

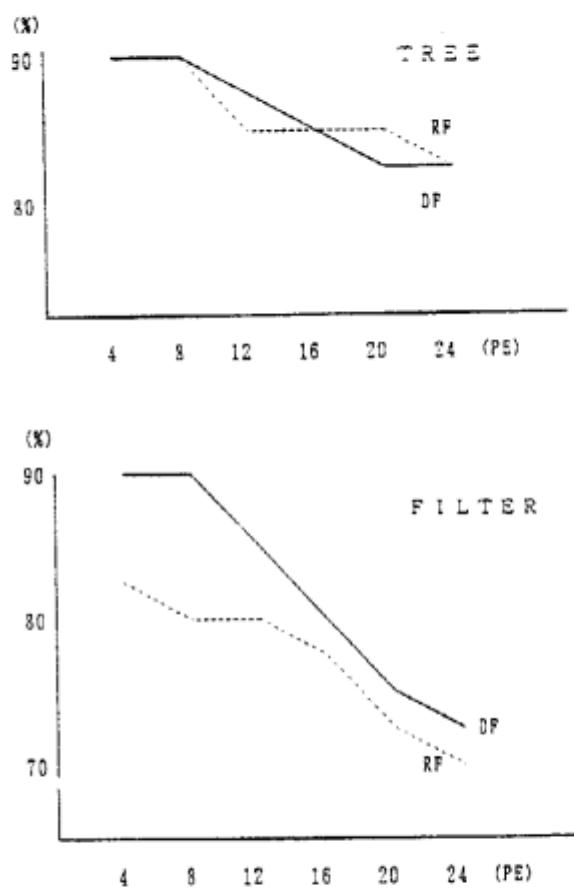


図7 様例実験 (リジューム優先方式)

PE	4	8	12	16	20	24
TREE DF	704	692	691	723	785	848
RF	667	706	734	705	756	824
FILTER DF	510	703	1690	2471	3644	4657
RF	6255	5919	5960	6828	7306	7452

表2 サスペンド数 (リジューム優先方式)

TREEの場合と様例実験、サスペンド数とも両方式で差立った差異はみられない。FILTERではRFの方が、様例実験、サスペンド数の双方で劣っている。結論としては、リジュームしたゴールを優先的にスケジュールしても意味がないと言わざるを得ない。これは、サスペンドしたゴールの子供は親同様にサスペンドすることが多いため、優先的に実行しても仕方ないからと考えられる。特にFILTER型ではその傾向が強い。

7. シミュレーション (3) — 静的負荷分散

我々の処理系では静的な負荷分散は現在行っていないが、処理系の性能を向上させるためには不可欠だと考えている。ここでは、静的負荷分散の目的をサスペンド数の削減に置き、ゴール間の優先度に関する情報を、静的情報と定義する。一般に静的情報を実行以前に得るには、コンパイラの高機能化、ソース上で明示的な記述法の開発などが考えられるが、ここではそういった方法論には触れずに、静的情報がある程度得られた場合の性能向上についてシミュレーションした結果を示す。

7.1 静的情報

① TREE型の場合

静的情報として、ゴールに自然数で表す順位属性を持たせた。兄弟間に優先順位が生じた場合、優先度の高い方から0, 1, ... の順に順位属性を与える。そして、キーの中で順位属性が最小のものをスケジューリングするようとする。また、親の優先度を受け継いだ場合に、順位属性の付け直しを行う方式と行わない方式を区別した。前者をS1、後者をS2と呼ぶ。

② FILTER型の場合

全てのゴールに(i, j)のような2次元座標を振った。ここでiは図7の左から順に、jは上から順に0, 1, ...とした。iが小さいほど優先度は高い。また、子供を生成する度にjは増えていく。座標(i, 0, 0)のゴールは、i = 0なら無条件で動作可能で、それ以外の時は、(i - 1, j)のゴールの実行が終わるまではサスペンドする。ゴールを選ぶ際には、i + jが最小なものを選ぶ方式、iが最小なものを選ぶ方式を採用した。前者をS3、後者をS4と呼ぶ。

7.2 シミュレーション条件

シミュレーションは、前節で述べた方式とDFを用いた。動的負荷分散方式はEQに統一し、その他の測定条件は動的負荷分散の場合と同様にした。

7.3 結果と考察

測定結果を示す。

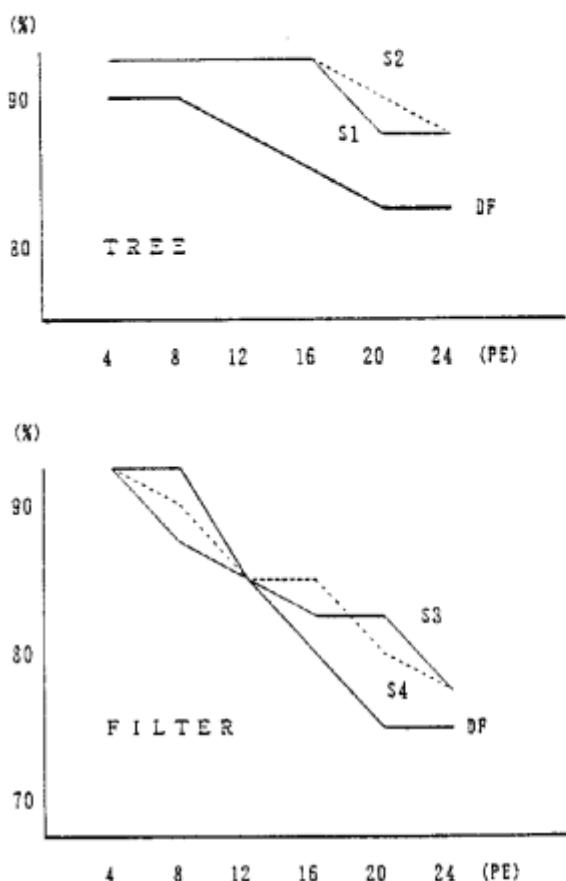


図8 様偏率(静的負荷分散)

PE \	4	8	12	16	20	24
TREE DF	704	692	691	723	785	848
S1	55	62	48	55	94	103
S2	26	29	33	54	62	94
FILTER DF	510	703	1690	2471	3644	4657
S3	671	1916	2498	3038	3381	5014
S4	288	958	1225	1712	2692	3431

表3 ナスペンド数(静的負荷分散)

TREEでは静的負荷分散によって様偏率が高まり、ナスペンドも大幅に抑制されているのがわかる。一方 FILTERではあまり効果が上がっていない。またS1とS2、S3とS4ではS1、S3の方がそれぞれ情報量が多いにも関わらず、むしろ逆効果となっている。TREE型の場合にS2程度の静的情報を取り入れ、FILTER型では単にDepth Firstでスケジューリングするのが良いことが分かった。

8.まとめ

我々の作成したFGHC処理系を紹介した。また、負荷分散方式をシミュレーションするためのモデルを提案し、各種シミュレーションの結果を述べた。シミュレーションの結果、次のような知見を得た。

- 動的負荷分散ではEQが優れている。ただしサスペンド数はRQの方が抑制されている。
- リザーブ優先方式は効果がない。
- 静的負荷分散はTREE型プログラムに有效で、静的情報としてはS2程度で良い。しかし、FILTER型プログラムでは静的負荷分散は効果がなく、単なるDepth First方式で良い。

今後も、今回の結果をふまえて、より良い並列システムの研究を進めていくつもりである。

謝辞 日頃御指導下さる林部長、服部室長および研究室諸兄に感謝致します。

参考文献

- (1) 清水他、"並列論理型言語GHCとその応用" 共立出版
- (2) M.Sato et al, Evaluation of the K1 Parallel System on a Shared Memory Multiprocessor, IFIP WG 10.3 Working Conference on Parallel Processing in PISA ITALY, April 1988
- (3) D.L.Bauer et al, Adaptive Load Sharing Homogeneous Distributed Systems, IEEE Trans. on Software Engineering, MAY 1986