

航空宇宙技術研究所報告

TECHNICAL REPORT OF NATIONAL AEROSPACE LABORATORY

TR-1277

並列計算機システム用ジョブスケジューラ

末 松 和 代

1995年9月

航空宇宙技術研究所

NATIONAL AEROSPACE LABORATORY

目 次

1. はじめに	1
2. NWT システムのジョブスケジューラに関する予備的考察	2
2.1 NWT システムの現状	2
2.2 NWT システムのジョブスケジューラに要求される機能	5
2.3 ジョブの状態と遷移	6
3. NWT システムのジョブスケジューリング方式	7
3.1 ジョブスケジューラの機能	7
3.2 ジョブとキューの定義	9
3.3 ジョブ起動条件	9
3.4 ジョブスケジューリング手順	12
4. シミュレーション実験におけるワークロード	12
4.1 ジョブの属性	12
4.2 ジョブ到着時間間隔	14
5. シミュレーション結果の評価・検討	15
5.1 基本実験	15
5.2 負荷実験	19
5.3 スワップアウト効果実験	22
5.4 起動優先ジョブのための PE リザーブ効果実験	23
5.5 プレステージ効果実験	23
5.6 PE 占有時間精度実験	25
5.7 シミュレーションプログラム	26
6. おわりに	26
謝 辞	27
参考文献	27

並列計算機システム用ジョブスケジューラ*

末 松 和 代^{*1}

Investigation on the Job Scheduler for Parallel Computer Systems

Kazuyo SUEMATSU

ABSTRACT

The Numerical Wind Tunnel, NWT, is a CFD-oriented vector parallel computer system with distributed memory. The development of a job scheduler is mandatory in order to operate the system effectively. Proposed herein is a scheduling algorithm which enhances the utilization ratio of processor elements, execution of jobs on the optimum number of processor elements, guarantees the turn-around time of jobs and operation of privileged jobs. The major point of the algorithm is to control job initiation and execution according to waiting time and demand for system resources of jobs. The effectiveness of the proposed scheduling algorithm is verified by various numerical experiments employing a scheduler simulator.

Key Word ; job scheduler, parallel computer system, waiting time, demand for system resources, scheduler simulator

概 要

数値風洞 (NWT システム) は、計算空気力学プログラムの超高速処理を目的とした分散メモリ型の並列計算機システムである。このような並列計算機システムを効率的に運用するためにはジョブスケジューラの開発が必要である。本稿では、①要素計算機 (PE) の高効率利用、②要求台数の PE の割り当て、③適切なターンアラウンドタイムの保証、および④特権ジョブの優先処理を実現することの可能なスケジューリングアルゴリズムについて提案する。本アルゴリズムの特徴は、ジョブの待ち時間およびシステム資源要求量に応じてジョブの起動や実行を制御している点にある。なお、スケジューリングアルゴリズムの有効性については、スケジューリングシミュレータを使用した種々の実験によって実証した。

1. はじめに

1980年代末、航技研では、近年急速に発展している計算流体力学 (以下 CFD と略記する) が今後の我が国の航空宇宙技術の研究開発の基盤技術における中核としての役割を果たすためには、メモリ量が32GB以上、CFDプログラムの実効処理速度が Fujitsu-VP400 の100倍以上の計算機が必須であり、その実現においてはMIMD (複数命令列、複数データ列) 方式のメモリ分散型の並列計算機システムを採用すべきであると指摘されていた^{1),2)}。これに基づき、

上記性能をもつ並列計算機システムの実現を目指して、富士通 (株) との間で共同研究を開始し、ハードウェアの構築、ハードウェア性能を十二分に発揮させるためのオペレーティングシステム機能、プログラム開発環境をより良くするための言語処理ソフトウェア機能等の検討・設計・開発を行った³⁾⁻⁹⁾。さらに、1993年2月には共同研究の成果とも言える数値風洞 (以下 NWT システムと略記する) を導入し、新計算機システムとしての運用を開始した。本稿では、新旧の計算機システムを総称して数値シミュレータシステム (以下 NS システムと略記する) と呼び、区別が必

* 平成7年2月23日受付 (received 23 February 1995)

*1 数理解析部 (Computational Sciences Division)

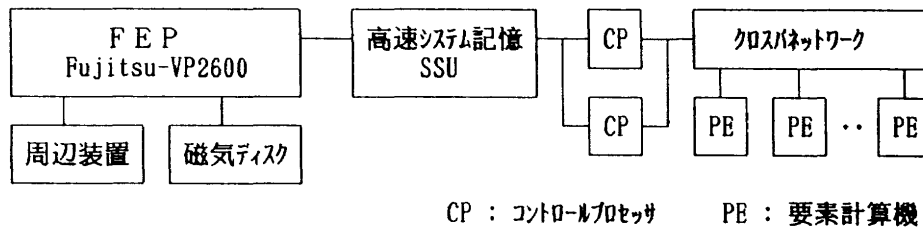


図 1.1 NWT システムのハードウェア構成図

要な場合には、従来の計算機システムを第一期 NS システム、新計算機システムを第二期 NS システムと呼称する。

図 1.1 に NWT システムのハードウェア構成図を示す。NWT システムは、演算処理用の 140 台の要素計算機（以下 PE と略記する）、ジョブおよび入出力制御用の 2 台のコントロールプロセッサ（以下 CP と略記する）、PE-CP 間および PE-PE 間接続用のクロスバネットワーク、ジョブ実行に必要なファイル転送用の高速システム記憶（以下 SSU と略記する）から構成されており、第一期 NS システムで使用していた Fujitsu-VP2600 をフロントエンドプロセッサ（以下 FEP と略記する）として使用している。

NWT システムは、複数ジョブの同時実行が可能であるため、PE の有効利用を図れるスケジューラが必要であったが、並列計算機用に開発されたそのようなスケジューラはなかった。また、並列計算機であるにも関わらず一般の汎用計算機用のスケジューラを使用した場合には PE 稼働率の低下が明白であった（第 2.1 (6) 項に詳述する）ため、新たなアルゴリズムの検討を開始した。その結果、下記特徴を持つジョブスケジューリングアルゴリズムを考案することができた。

- ①任意の PE 台数でのジョブ実行
- ②PE の高効率利用
- ③システム資源要求量に応じた適切なジョブターンアラウンドタイムの保証
- ④緊急ジョブ、特権ジョブの優先処理

なお、本アルゴリズムでは、ジョブ実行時の PE 占有時間を CPU 要求時間から予測する方式を採用しているが、これにより、③項を実現する際のシステム資源利用率の大幅低下防止ばかりでなく、実行中ジョブの終了を待つためにシステム停止を予定時刻に行えないような事態を回避することも可能となった。

これらの機能は、航技研の計算機システムだけでなく、一般の並列計算機システムの運用においても十分有効な機能である。

本稿では、上記スケジューリングアルゴリズムの機能およびスケジューリングシミュレータを使用したアルゴリズムの有効性実証のための実験結果について述べる。なお、NWT システムでは、1994 年 10 月から本アルゴリズムを組

み込み、新方式での運用を開始したところである。

<略号表>

- CFD : 計算流体力学
 CP : 数値風洞のコントロールプロセッサ
 FEP : フロントエンドプロセッサ
 NS : 航技研計算機システム
 (Numerical Simulator system)
 NWT : 数値風洞 (Numerical Wind Tunnel)
 PE : 数値風洞の要素計算機
 SHP : NS システムのベクトルプロセッサ
 SSU : 数値風洞の高速システム記憶

2. NWT システムのジョブスケジューラに関する予備的考察

本章では、NWT システムの現状について簡単に説明を行った後、NWT システムのジョブスケジューラに要求される機能およびジョブの投入から実行終了までの状態遷移について述べる。

2.1 NWT システムの現状

本項では、NWT システムの現状として、スケジューリングに関係のある SSU の設置目的、PE の負荷分散、PE の割り当て方法、ジョブのシステム資源制限項目、大規模ジョブのスワップアウト^{注1)}および現在使用中のジョブスケジューリング方式について説明する。

(1) SSU の設置目的

NWT システムでは、ジョブ実行時の入出力データおよびスワップアウトデータを一時的に保存するために、FEP と CP の間に SSU を設置している（図 1.1 参照）。SSU を有効に使用することができれば、以下に示すように、ジョブ実行時の入出力待ち時間の軽減、スワップアウト、スワップイン^{注2)}処理時間の短縮化が図れる。

最初に入出力待ち時間について説明する。NWT で実行

(注1) 実行中のジョブの情報を継続実行可能な形でディスク等に退避すること。NWT システムの場合の退避先は SSU であり、SSU 容量が不足した場合には FEP 配下の磁気ディスクを使用する。

(注2) スワップアウト中のジョブの情報を PE に復元すること。

する実行形式プログラムおよびデータ等はすべてFEP配下の磁気ディスクにあるため、ジョブ実行時には使用するすべてのデータをFEPからPEへ転送し、ジョブ出力時には出力データをPEからFEPに転送しなければならない。しかし、図2.1に示すように、PE-CP間の転送速度が最大421MB/Sであるのに対し、磁気ディスクは最大4.5MB/S（並列入出力方式により4ファイルに並列に入出力を行った場合でも18MB/S）であり、FEP-PE間のデータ転送では転送速度の最も低い磁気ディスクの速度で抑えられる。しかし、SSUからCPへの転送速度は最大1250MB/Sであるため、ジョブ実行に先立ってあらかじめSSUにデータを転送しておく（この機能をプレステージ機能と呼ぶ。）ことができれば、ジョブ開始時のデータ転送は最大421MB/Sとなり、ジョブの入力待ち時間を大幅に短縮することが可能となる。また、出力処理ではPEからSSUへのデータ転送とSSUからFEPへのデータ転送は独立に処理可能であり、SSUへの書き込みが終了した時点でプログラムは次の処理に移ることが可能となるため、SSU容量が十分にあれば出力も最大421MB/Sとなり、ジョブの出力待ち時間も大幅に軽減できる。

続いて、スワップアウト・スワップイン処理時間である

が、SSUが設置されていない場合にはスワップアウトデータはFEP配下の磁気ディスクに格納することになるため、処理速度は磁気ディスクの転送速度となる。しかし、SSUが使用できる場合にはスワップアウト時はPEの情報をSSUに格納し、スワップイン時はSSUの情報をPEに復元するため、入出力と同様に高速に処理が行える。

しかし、これらの処理でSSU容量が不足した場合にはすべて磁気ディスクとPE間の転送となるため所期の目的が果たせなくなり、入出力処理等に長時間かかることになる。したがって、SSU容量はこれらの目的を達成できる程度に十分にあることが前提となる。また、スワップアウトやプレステージの頻発によりSSUが占有されると同様の結果となるため、NWTシステムでは、スワップアウト量やプレステージ量を制限すべきである。

(2) PEの負荷分散

NWTシステムでは、複数PEを使用するジョブの場合、PEの使用効率を高め、ジョブの処理時間を最短にするためにPEの負荷をできるだけ均等に分散させることがユーザに義務づけられている。ここで、4種類の処理A、B、C、Dから構成されるジョブの負荷分散の例を図2.2に示す。なお、処理Aには8時間、処理B、Cにはそれぞれ2

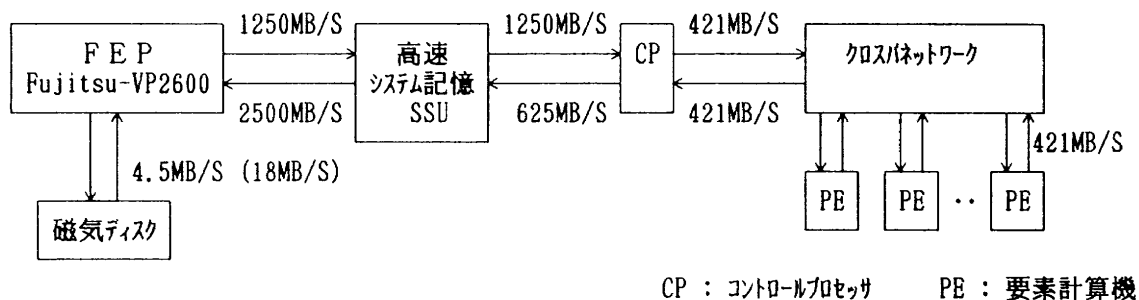


図2.1 データ転送能力（理論性能）

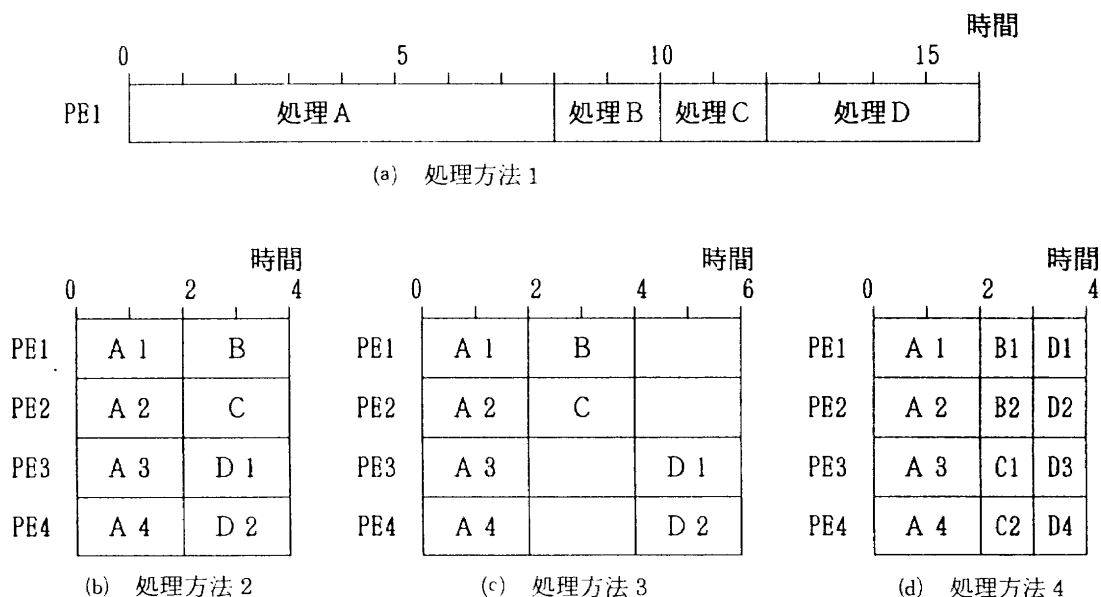


図2.2 並列計算機使用時のプログラムの負荷分散の例

時間、処理 D には 4 時間かかるものとする。

(a)はこのジョブを 1PE で実行した場合を示している。当然のことながら、全処理を行うのに 16 時間かかる。(b)はこのジョブを 4PE で実行した場合を示している。ここで処理 A1, A2, A3, A4 は処理 A を、処理 D1, D2 は処理 D を、それぞれ並列実行可能な処理に均等分割したものである。その結果、ジョブの実行時間は 4 時間となる。ただし、この場合には、処理 B, C, D1, D2 は処理 A の結果を使用して並列に実行可能でなければならない。もし、処理 D が処理 B, C の結果を使用するような場合には、処理 D の開始は処理 B, C の終了を待つため、(c)のように処理しなければならない実行時間は 6 時間となってしまう。この場合は、処理 B は B1, B2 に、処理 C は C1, C2 に、処理 D は D1, D2, D3, D4 に均等分割すれば(d)のような処理となり実行時間は 4 時間となる。一般的に、入出力処理は 1 PE で行うことが多いため、同図の例のように完全な負荷分散は不可能であるが、できるだけ(b)または(d)のように各 PE の処理を分割することが望ましい。このような処理を PE の負荷分散処理と呼ぶ。

(3) PE の割り当て方法

NWT システムでは、ジョブ実行時の PE の割り当て方法は、「マルチジョブ割り当て」と「シングルジョブ割り当て」というふたつの方法が選択可能であった。前者は、PE の有効利用の観点からメモリに余裕のある PE を複数のジョブに割り当てる方法であり、後者は単一ジョブに割り当てる方法である。NWT システムの検討開始に当たり、どちらの方式を採用するかを検討した結果、以下の結論を得た。

- ①「シングルジョブ割り当て」の場合は、「マルチジョブ割り当て」に比べて入出力による PE のアイドルの確率が高いと考えられるが、どちらの場合でも入出力によるアイドルは極めて少ないと考えられる。
- ②「マルチジョブ割り当て」の場合には、ジョブの実行は同一 PE を使用している他のジョブの実行に影響されるため、ジョブ終了の見通しが立ちにくい。
- ③複数 PE を使用するジョブでは、PE の有効利用を図るために PE の負荷分散が重要な要素となる^{(6)~(8)}が、「マルチジョブ割り当て」によりジョブの実行が他のジョブに影響される場合には、PE の負荷分散効果が期待できない。

本システムは CFD プログラムの高速実行を目的として開発するという前提に立っているため、上記 3 つの理由により「シングルジョブ割り当て方式」を採用することに決定した。なお、この決定は、NWT システムの開発方針となった。

(4) ジョブのシステム資源制限項目

一般的に、どの計算機システムでも、システム資源使用量の公平化とプログラム処理異常時の対処のためにシステ

ム資源を制限している。NWT システムの場合もシステム資源の制限は必要であり、投入されるジョブに対して PE 使用台数、CPU 使用時間、経過時間（NWT システム内でジョブが実行を開始してから実行を終了するまでの時間）に関する制限値を設定し、ジョブ実行中に CPU 使用時間や経過時間の制限値を越えた場合にジョブを打ち切る機能がある。ただし、NWT システムで複数 PE を使用する場合には各 PE で異なる処理を行わせることも可能であるため、PE 毎の CPU 使用時間が異なることが予想される。すなわち、CPU 使用時間は、全 PE が使用した総 CPU 使用時間、平均 CPU 使用時間、最大 CPU 使用時間の 3 種類が考えられる。このうち、総 CPU 使用時間、平均 CPU 使用時間は、個々の PE の情報だけでは求めることができないため、制限値を越えた時点でジョブを打ち切るためには常に各 PE の情報を収集するための PE 間通信等が頻発する。したがって、PE 間通信等のオーバーヘッドの増加防止を図るためには総 CPU 使用時間や平均 CPU 使用時間は採用すべきでない。

また、ジョブは終了直前に計算結果を出力することが多いが、この場合には全ての出力データが PE から SSU を経由して FEP に転送された時点でジョブが終了する。すなわち、CPU の処理が全て終了し、出力結果が SSU まで出力された場合でも、FEP への転送中に経過時間切れとなった時にはジョブ処理は中断され出力結果の保証はできない。

以上に述べた理由により、NWT システムの運用ではシステム資源の主たる制限項目として、PE 使用台数、最大 CPU 時間を使用し、プログラム処理異常時の対処を可能にするために大きめに設定した経過時間を補助的に使用するべきであると結論した。したがって、スケジューラにおけるジョブスケジュールの際にもジョブで設定した経過時間は使用できない。

(5) 大規模ジョブのスワップアウト

複数 PE を使用するジョブの場合、処理時間を最短にするために、ユーザは(2)項で述べた PE の負荷分散が義務づけられている。したがって、このように複数 PE を使用するジョブの PE が 1 台でもスワップアウトされるとその分だけ全 PE の処理が遅れることになる。言い換えると、スワップアウトに伴う PE 利用率の低下の割合は規模が大きいジョブをスワップアウトした場合ほど顕著となるため、大規模ジョブのスワップアウトを規制する必要がある。

(6) 現在のスケジューリング方式

現在、並列計算機用に考案されたスケジューリングアルゴリズムの実用例がないため、NWT システムでは、通常の汎用計算機用のスケジューリングアルゴリズム⁽¹⁰⁾を使用している。このアルゴリズムは、CPU 時間、メモリ量等の制限が異なる数個のジョブクラスを定義し、それぞれのジ

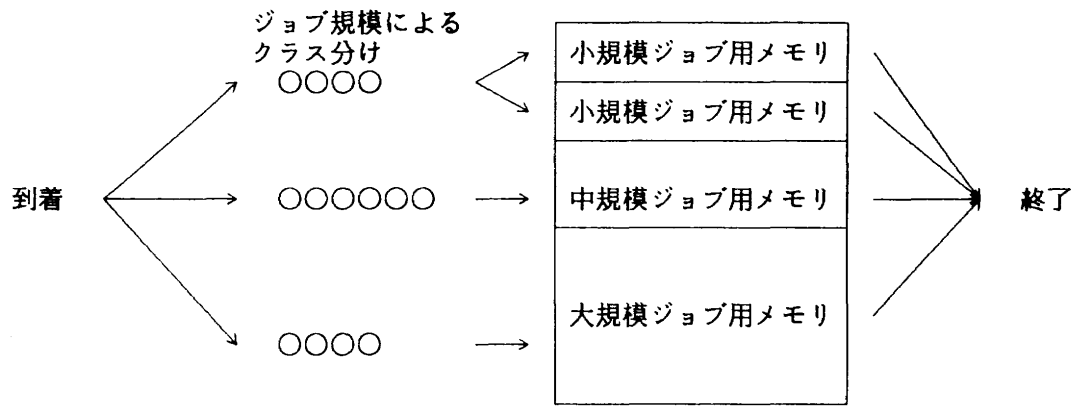


図 2.3 ジョブクラス別多重度制御方式の例

ジョブクラスの多重度（同時に実行可能なジョブ数）の指定に基づいてジョブクラス内ジョブの実行を制御する方式である。本稿では、この方式をジョブクラス別多重度制御方式と呼んでいる。以下に、本方式の特徴と NWT で使用した場合の問題点を説明する。

(a) ジョブクラス別多重度制御方式の特徴

一般的な汎用計算機の運用では、CPU の高効率利用と公平なジョブのターンアラウンドタイムが重要課題である。航技研では、以下に述べるジョブクラス別多重度制御方式を使用することにより、汎用計算機の効率的運用を実現している。

図 2.3 に汎用計算機でのジョブクラス別多重度制御方式例を示す。例ではジョブの規模により、大規模、中規模、小規模という 3 つのジョブクラスを定義し、多重度をそれぞれ 1, 1, 2 としている。この方式は、同一ジョブクラスに投入されたジョブは先着順に処理することを原則としているため、大規模ジョブクラスおよび中規模ジョブクラスのジョブはそれぞれ 1 本、小規模ジョブのみ 2 本まで先着順に処理される。当然のことであるが、各ジョブクラスのジョブを指定多重度で実行するためには、各ジョブクラスのメモリ制限値と多重度の積を加算した値はシステムで使用可能なメモリ量以下でなければならない。この方式の特徴を下記に示す。

- ① 同一ジョブクラス内では先着順に処理されるため投入されているジョブのターンアラウンドタイムの予測が比較的簡単に行える。
- ② 大規模クラスのジョブでも一定時間待てば処理される。
- ③ ジョブクラス制限値に満たないジョブを実行すると未使用のメモリ領域が発生するが、システム内で複数のジョブが実行されている限り CPU の高効率利用が可能である。

(b) NWT での問題点

表 2.1 に、NWT システムでのジョブクラス別多重度制御方式例を示す。ここでは、システムの総 PE 台数が 140

表 2.1 多重度設定例

	クラス別設定値		クラス別使用 PE 台数
	PE 台数	多重度	
クラス 1	1	6	6
クラス 2	2	3	6
クラス 3	4	2	8
クラス 4	10	1	10
クラス 5	30	1	30
クラス 6	80	1	80
合 計	—	14	140

台であるにも関わらず、各ジョブクラスの PE 制限台数と多重度の積を加算した値はシステムで使用可能な PE 台数以下という制限からジョブクラス数は 6 しかない。この場合以下の問題点が発生する。

- ① プログラムにはそれぞれ効率実行するための最適な実行 PE 台数があるにもかかわらず、使用可能な PE 台数が限定される。
- ② プログラムをジョブクラスの制限台数に合わせて実行すると、各 PE のロードバランスの調整を最適にできないために PE の使用効率が低下する。
- ③ 最適 PE 台数でないジョブクラスでプログラムを最適 PE 台数で実行すると、実行中常に空き PE が発生し、PE の使用効率が大幅に低下する。
- ④ 混雑時でもジョブのないジョブクラスがあると未使用の PE が発生する。

したがって、PE の有効利用を図る観点からは、ジョブクラス別多重度制御方式は採用できない。

2.2 NWT システムのジョブスケジューラに要求される機能

ジョブスケジューラに要求される機能は、システム資源の効率的運用機能、適切なジョブターンアラウンドタイム

保証機能、システム運用管理機能に大別することができる。本項では、第一期 NS システムのジョブスケジューラ¹²⁾の特徴をこれらの要求機能毎に示し、それに対応させて NWT システムのジョブスケジューラに必要な機能について述べる。

昭和62年2月に導入された第一期 NS システムはフロントエンドプロセッサに Fujitsu-M780 を、バックエンドプロセッサにベクトル計算機である Fujitsu-VP400 と VP200 (以下これらのベクトル計算機を SHP と略記する) を配置した疎結合の複合計算機システムである。なお、VP200 は平成3年1月にシステム内の位置づけは同じままで VP2600 に代替された。このシステムでは、ベクトル計算機である SHP の有効利用を図るために、実行ステップのみ SHP で処理し、実行ステップ以外の SHP ジョブ、会話型処理等はすべて FEP で処理した。また、ジョブに対するシステム資源の制限は主としてメモリ量と CPU 時間で行った。

(1) システム資源の効率的運用機能

汎用計算機システムである第一期 NS システムでは、システム資源、特に CPU の効率的運用が最重要課題であり、この課題を達成するために実行可能ジョブの確保機能およびジョブクラス別多重度制御機能を使用した。

実行可能ジョブの確保機能は、FEP 内の SHP ジョブを優先的に実行して SHP の実行可能ジョブを一定数確保する機能であり、SHP ジョブのサービスステップ処理待ちによる SHP の空き状態の回避に有効である。ジョブクラス別多重度制御機能の実現方式については第2.1(6)項で述べたが、CPU の高効率利用に有効な機能である。

これに対し、NWT システムでは各 PE の有効利用を図ることが最重要課題である。しかし、汎用計算機用のスケジューリングアルゴリズムであるジョブクラス別多重度制御機能では PE の有効利用が図れないため、システムの有効利用の観点から以下の機能が必要となる。

- ① 使用可能な PE 使用台数を限定しない。
- ② 未使用の PE を極力発生しない割り当て方式を採用する。
- ③ ジョブの流れを適切に制御して PE の利用効率をできるだけ向上させる。

なお、PE の利用効率向上のための機能として、NWT システムでは第2.1(1)項に述べたプレステージ機能が提供されている。PE の有効利用を図るためにはこの機能を積極的に利用する必要がある。

(2) 適切なジョブターンアラウンドタイム保証機能

適切なジョブターンアラウンドタイムを保証することは、研究者が研究計画を遂行するために重要である。ジョブクラス別多重度制御機能は、ジョブクラス多重度分のジョブの実行が保証されるため、ジョブの適切なターンアラウン

ドタイムの保証にも有効な機能である。また、ジョブクラス毎にすべてのジョブに起動^{注3)}優先度、メモリ優先度、実行優先度を設定(運用時間帯ごとに設定可能)し、これらの優先度によりジョブの実行を制御しているため、システム資源が競合した場合にもジョブクラス毎のジョブのターンアラウンドタイムを相対的に保証することが可能である。

NWT システムでも適切なジョブターンアラウンドタイムの保証は重要な課題である。その実現のためには、ジョブの属性やシステム資源要求量に応じたジョブ起動制御機能だけでなく、起動できずに不当に長時間待たされているジョブに対しては実行中ジョブのスワップアウト処理による要求 PE 台数の確保や他ジョブの起動抑止による早期起動を可能とする実行制御機能等が必要となる。

(3) システム運用管理機能

第一期 NS システムの特権ジョブは、システムが原因となる事故ジョブの処理が大半を占めており、発生件数が限られること、個別対応が必要なことからマニュアル操作での高優先度設定機能により対処していた。

これに対し NWT システムでは、HOPE^{注4)}シミュレーションのようなプロジェクトジョブを緊急に処理する必要性が考えられていた。これらのジョブは、必ずしも実行中のジョブを中断させてまで即刻起動する必要があるとは考えられないが、使用 PE 台数が30台程度であり他のシステムで実行することが不可能に近いこと、比較的高頻度で投入されると予想されることから、システム管理者によるマニュアル操作なしで早い時期に起動できるようにする必要があった。そのためには、特権ジョブの概念を取り入れ、特権ジョブを非特権ジョブより優先して起動、実行するためのジョブ起動方式や実行制御方式を組み込むことが必要となる。

2.3 ジョブの状態と遷移

図2.4にジョブ状態遷移図を示す。同図において、スケジューラが意識すべきジョブの状態を①～⑦、起動待ちジョブのデータの状態をⅠ～Ⅱに示す。また、a～eの状態遷移はスケジューラが決定することを示す。本項では、同図を基に、それぞれのジョブ状態と遷移について述べる。

① 起動待ち状態

システムに到着し、起動を待っている状態を起動待ち状態と呼ぶ。また、起動待ち状態のデータは、プレステージが完了しているか否かによりプレステージデータまたはプレステージ待ちデータと呼ぶ。なお、プレステージおよび

(注3) 実行待ちのジョブの実行を決定すること。また、そのために必要となる処理を行うこと。

(注4) H-II orbiting plane の略語。航技研と宇宙開発事業団の共同開発プロジェクト

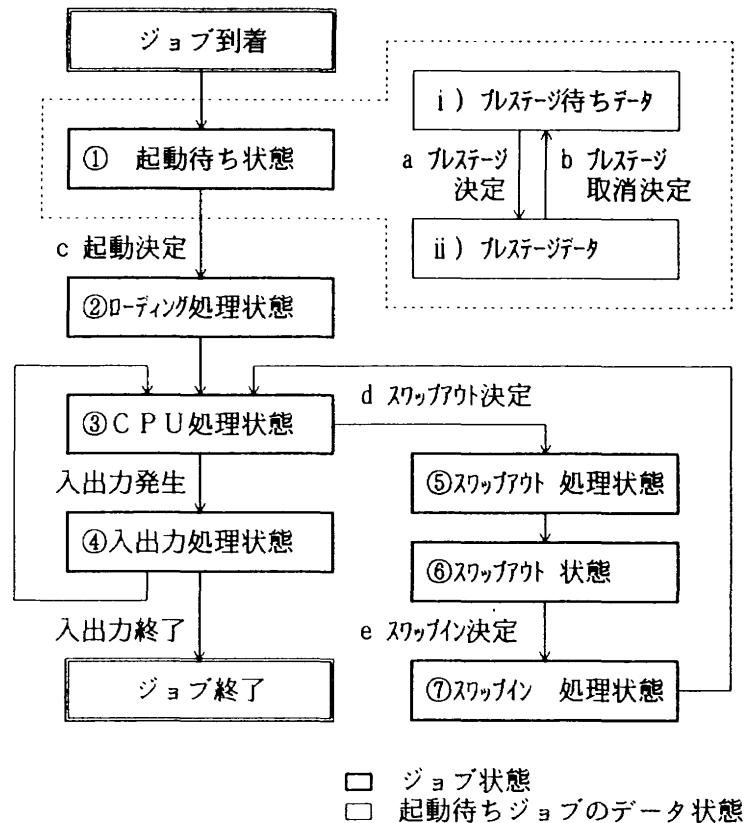


図 2.4 ジョブ状態遷移図

プレステージ取消の決定はスケジューラが行う。

② ローディング処理状態

スケジューラにより起動が決定され、実行開始に必要なデータを PE に転送中の状態をローディング処理状態と呼ぶ。転送データの内、プレステージデータは SSU から PE に、プレステージ待ちデータは FEP から SSU を経由して PE に転送される。

③ CPU 処理状態

ローディング処理状態、入出力処理状態、あるいはスワップイン処理状態でのそれぞれの処理が終了し、CPU 処理中の状態を CPU 処理状態と呼ぶ。なお、CPU 処理によりすべての処理が終わったジョブ、あるいは予定 CPU 時間を使い切ったジョブは終了する。

④ 入出力処理状態

CPU 処理中に入出力要求が発生し、入出力処理中の状態を入出力処理状態と呼ぶ。なお、入出力によりすべての処理が終わったジョブは終了する。

⑤ スワップアウト処理状態

ジョブ実行中にスケジューラによりスワップアウトが決定され、スワップアウト処理中の状態をスワップアウト処理状態と呼ぶ。ジョブが CPU 処理状態の場合にはスワップアウトの決定と同時にスワップアウト処理状態となるが、入力処理状態の場合には入出力処理が終わり CPU 処理状態となった時点でスワップアウト処理が開始される。

⑥ スワップアウト状態

スワップアウト処理が終了し、スワップイン待ちの状態をスワップアウト状態と呼ぶ。

⑦ スワップイン処理状態

スワップアウト状態のジョブのスワップインが決定され、スワップイン処理中の状態をスワップイン処理状態と呼ぶ。この状態のジョブは、スワップイン処理が終了した時点で CPU 処理状態となる。

3. NWT システムのジョブスケジューリング方式

前章で、NWT システムのスケジューラに要求される機能について述べた。本章では、これらの要求事項を満たし、NWT システムを効率的に運用するために考案したスケジューリング方式について述べる。

3.1 ジョブスケジューラの機能

(1) ジョブの分類とジョブのキュー管理機能

システム内のジョブを異なる優先度レベルで処理するために、優先度の低い方から順に α_1 ジョブ、 α_2 ジョブ、 α_3 ジョブに分類する。これらのジョブは、実行待ちジョブ、実行中ジョブごとに図 3.1 に示すジョブ管理キューで管理する。

実行待ちジョブ管理キューおよび実行中ジョブ管理キュー

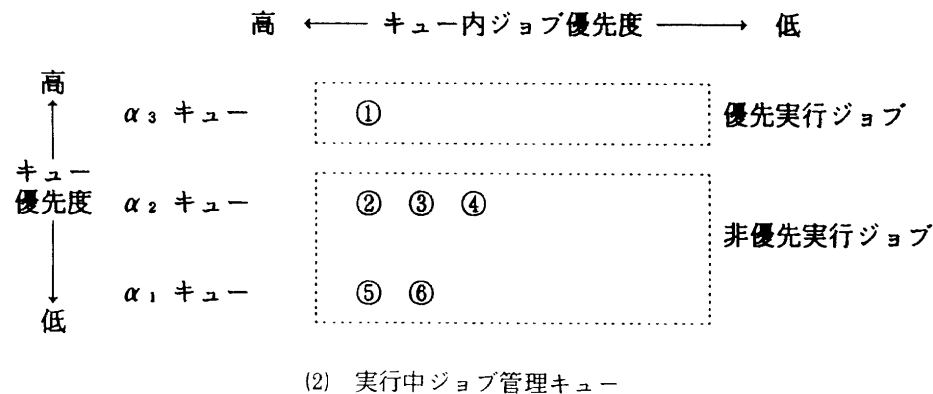
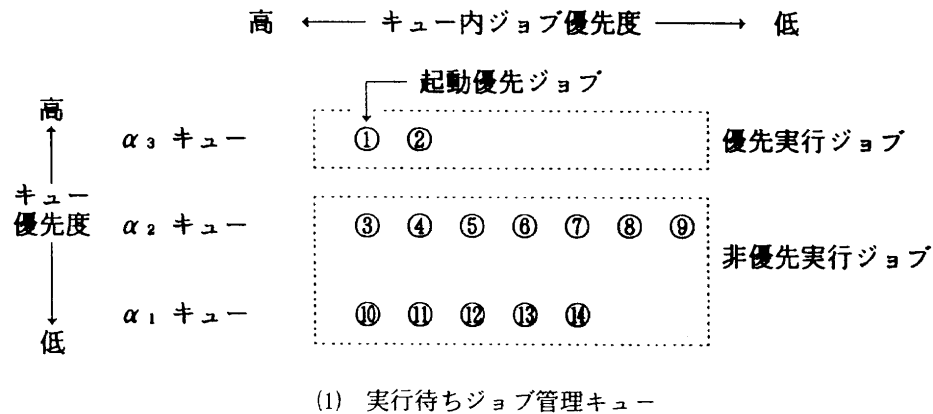


図 3.1 ジョブ管理キュー

ーには α_1 キュー、 α_2 キュー、 α_3 キューと呼ぶ3本のキューを設け、それぞれ α_1 ジョブ、 α_2 ジョブ、 α_3 ジョブを優先度の高い順に管理する。図中の丸付き数字は各キューにつながれているジョブを示す。また、中の数字は実行待ちジョブ管理キューまたは実行中ジョブ管理キュー内の相対的ジョブ優先順位を示す。

なお、待ち時間が長くなったジョブを最優先に処理するためには α_3 ジョブの処理方式を変える必要がある。そこで、 α_3 ジョブを「優先実行ジョブ」、 α_2 ジョブおよび α_1 ジョブを「非優先実行ジョブ」と呼び、ジョブ起動時、PEの割り当て時等において優先的な処理方式を設定する。ただし、 α_3 ジョブの起動時には強硬な手段を用いる必要があるため、すべての α_3 ジョブを対象とすると混乱を生じる恐れがある。したがって、特にその先頭ジョブを「起動優先ジョブ」と呼び、優先的な起動方式を設定する。

(2) ジョブ優先度の動的制御機能

ジョブ属性、PE 要求台数、CPU 要求時間、待ち時間からジョブの所属キューとキュー内優先度を決定するための関数をあらかじめ定義しておく。本稿では、システム資源要求量と待ち時間によって優先度が決定するように第3.2項(1)、(2)に定義する関数を使用しているが、関数を変えることにより、到着順の処理や特定規模のジョブの優先処理等、任意の設定が可能である。

起動ジョブを選択する際には、すべてのジョブがその時点での優先度で処理されるように、所属キューおよび優先度を見直し、必要に応じてジョブを適当なキューへつなぎ直す。

(3) ジョブの起動機能

割り当て可能な空き PE がある場合には、実行待ちジョブ管理キュー内のジョブの起動条件（第3.3項に述べる）を相対的ジョブ優先順位の高いジョブから順に調べ、起動条件を満たしているジョブをすべて起動する。処理中に起動条件を満たさないジョブがあっても、割り当て可能な空き PE がある限り後続のジョブの起動を試みることにより PE の使用効率の向上を図る。

(4) 優先実行ジョブに対する PE の優先割り当て機能

PE の割り当ては、実行中の優先実行ジョブ、実行待ちの優先実行ジョブ、実行中の非優先実行ジョブ、実行待ちの非優先実行ジョブの順に優先する。したがって、スワップイン、スワップアウト、起動のための PE の割り当てはこの優先度にしたがって処理する。しかし、PE 要求台数が極端に多い実行待ちジョブは優先度を高くするだけでは実行不可能な場合も考えられるため、 α_3 キューの先頭ジョブである起動優先ジョブは以下に述べる特別な PE の割り当て方式を採用する。

起動優先ジョブ起動時の要求 PE 台数の割り当ては、空

き PE だけでなく、実行中の非優先実行ジョブのスワップアウトによる空き PE の割り当ても可能とする。それでも起動できない場合には、他のすべてのジョブの起動を抑止した場合のジョブ開始時刻と終了時刻を予測し、その時間帯に確実に実行するために要求台数の PE の割り当てを予約する。起動優先ジョブに対するこの処理を「PE リザーブ」と呼ぶ。

(5) スワップアウト機能

スワップアウト処理では、実行中の非優先実行ジョブの中のジョブ優先度の低いジョブの PE から順に取り上げる。ただし、下記の 3 項目の上限をあらかじめ設定しておく、これらの範囲内でのみスワップアウトを行うものとする。

- ①スワップアウト可能なジョブの規模（要求 PE 台数）
- ②1つのジョブを起動する際のスワップアウト PE 台数
- ③システム全体での総スワップアウト PE 台数

ここで、①は PE の利用率の大幅低下を引き起こす恐れのある大規模ジョブをスワップアウト対象から外すためのパラメータであり、②、③は、スワップアウトの頻発による SSU の専有やシステムオーバヘッドの増大を防止するためのパラメータである。

(6) プレステージ機能

プレステージ処理では SSU の無駄な専有を防止するために、次に起動されるジョブとその起動時間、プレステージに要する時間を予測し、起動直前にプレステージを終了させることが最も望ましい。しかし、ジョブの打ち切りは CPU 時間で行う方針であり、ジョブは終了予定時刻以前に異常終了する場合もあることから、次に起動されるジョブの正確な予測は不可能である。そこで起動待ちジョブの中の優先度の高い順に一定数のジョブをプレステージする方式とし、システム状況の変化によりプレステージすべきジョブが変更になった場合には、それに伴ってプレステージジョブも入れ換えるものとする。

3.2 ジョブとキューの定義

(1) キュー遷移条件

本アルゴリズムでは、待ち時間が一定値に達したジョブの優先処理を行うために、以下の遷移条件を満たしたジョブを α_2 キューおよび α_3 キューへ遷移させる。

α_2 キューへの遷移条件： $WT_j \geq WT_{j_1}$

α_3 キューへの遷移条件： $WT_j \geq WT_{j_2}$

ここで、 WT_j はジョブが投入されてから現在に至るまでの待ち時間であり、 WT_{j_1} 、 WT_{j_2} はジョブのシステム資源要求量や待ち時間から任意に決定することができるが、NWT システムでは以下のように定義する。

a. 一般ジョブの場合

$$\begin{aligned} WT_{j_1} &= [\text{CPU}_j \times \sqrt{\text{PE}_j}] \times \text{FAC}_1 \\ WT_{j_2} &= [\text{CPU}_j \times \sqrt{\text{PE}_j}] \times \text{FAC}_2 \end{aligned}$$

ここで、 PE_j はジョブ j の要求 PE 台数を、 CPU_j は要求 CPU 時間を表す。また、 FAC_1 、 FAC_2 はシステムの混雑度に対処するためのファクタであり、

$$0 < \text{FAC}_1 \leq \text{FAC}_2 < \infty$$

である。NWT システムでは、「大規模 CFD プログラムの高速実行」という設置目的により、PE 台数の多いジョブを優遇している。

本稿では、キュー遷移条件を決定する要素の中の FAC_1 、 FAC_2 をキュー遷移パラメータと呼ぶ。

b. 特権ジョブの場合

$$\begin{aligned} WT_{j_1} &= -\text{WTP} \\ WT_{j_2} &= -\text{WTP} \end{aligned}$$

ここで、 WTP は運用パラメータであり、 $-\text{WTP} < 0$ とする。したがって、先に述べた遷移条件から、投入された特権ジョブは、ただちに α_3 キューにつながる。

(2) ジョブのキュー内優先度

α_1 ジョブ、 α_2 ジョブ、 α_3 ジョブのキュー内優先度 PR_1 、 PR_2 、 PR_3 はそれぞれ以下の関数に従って決定する。なお、関数の値の小さいもの程優先度は高いものとする。

$$\begin{aligned} \text{PR}_1 &= WT_{j_1} - WT_j \\ \text{PR}_2 &= WT_{j_2} - WT_j \\ \text{PR}_3 &= WT_{j_3} - WT_j \end{aligned}$$

ここで、前項の遷移条件から、 α_1 ジョブおよび α_2 ジョブの優先度 PR_1 、 PR_2 は常に正となり、 α_3 ジョブの優先度 PR_3 は常に負となる。また、 $\text{PR}_1 \neq \text{PR}_2$ であることから、 α_2 キューへ遷移したジョブの優先度はすでに存在する α_2 ジョブの優先度より高くなることのあるのに対し、 $\text{PR}_2 = \text{PR}_3$ であるため、 α_3 キューへ遷移したジョブの優先度はすでに存在する α_3 ジョブの優先度より高くなることはない。

また、 α_3 キューにつながれた特権ジョブの優先度 PR_3 は、

$$\text{PR}_3 = -\text{WTP} - \text{WT}_j$$

となることから、 WTP を一般ジョブで予想される最大経過時間より十分大きく設定することにより、一般ジョブより常に優先的に処理することが可能となる。

3.3 ジョブ起動条件

ジョブを起動するためには、要求台数の PE の割り当て

が可能であることが必須であるが、本スケジューラでは効率的、効果的なスケジュールを実現するために、以下の起動条件を使用する。

a. 起動優先ジョブの起動条件

条件 1 : $PE_j \leq free$
 条件 2 : $(PE_j - free) \leq AVSWAP$,
 ただし, $AVSWAP = \max(SWAP1, SWAP2, SWAP3 - SWAP)$

b. 非起動優先ジョブの起動条件

条件 1 : $PE_j \leq free$

ここで使用されている記号の定義は、以下の通りである。

PE_j : 起動対象ジョブの要求 PE 台数

$free$: 起動対象ジョブが使用可能な空き PE 台数

$SWAP$: スワップアウト中の PE 台数

$SWMAX$: スワップアウト可能なジョブの規模 (単位は PE 台数) の上限

$SWAP1$: 現在スワップアウト可能なジョブ (実行中の α_1 ジョブ, α_2 ジョブのうち要求 PE 台数が $SWMAX$ 以下のジョブ) の総 PE 使用台数

$SWAP2$: 一つのジョブを起動する際のスワップアウト PE 台数の上限

$SWAP3$: システム全体での総スワップアウト PE 台数の上限

条件 1 は、空き PE だけでジョブを実行するための条件

である。この場合には、要求台数の PE を割り当ててジョブを実行する。

条件 2 は、空き状態の PE が要求台数に満たないときに、不足分の PE をスワップアウトにより生成するための条件であり、起動優先ジョブの場合のみの条件である。この条件を満足した場合には、スワップアウト可能なジョブの中の実行優先度の低いジョブから順に $(PE_j - free)$ 台数分の PE をスワップアウトし、既存の空き PE とスワップアウトによって生じた空き PE を割り当ててジョブを実行する。

なお、起動条件で使用した「起動対象ジョブが使用可能な空き PE 台数 ($free$)」の算出方法について以下に説明する。

起動優先ジョブの場合、他の実行待ちジョブより優先的に起動する。したがって、空き PE 台数をすべて使用することができるが、使用可能 PE 台数が変化する場合には、その時刻と変化量を考慮しないとシステム停止等の処理が遅延する恐れがある。さらに、非起動優先ジョブの場合には、起動優先ジョブのための PE リザーブ台数をも考慮しなければならない。

ジョブを新たに起動しないと仮定すると、システムで使用可能な空き PE 台数は、①ジョブの終了、②使用可能最大 PE 台数の変更、③PE リザーブの開始、④PE リザーブの終了により変化する。①および②を考慮した場合の例を図 3.2 に示す。図では縦軸に PE 台数を、横軸に時刻をとっている。

説明に先立ち、以下の記号を定義する。

$jmax$: 現在実行中のジョブ数

PE_j : j 番目に終了予定のジョブの使用 PE 台数

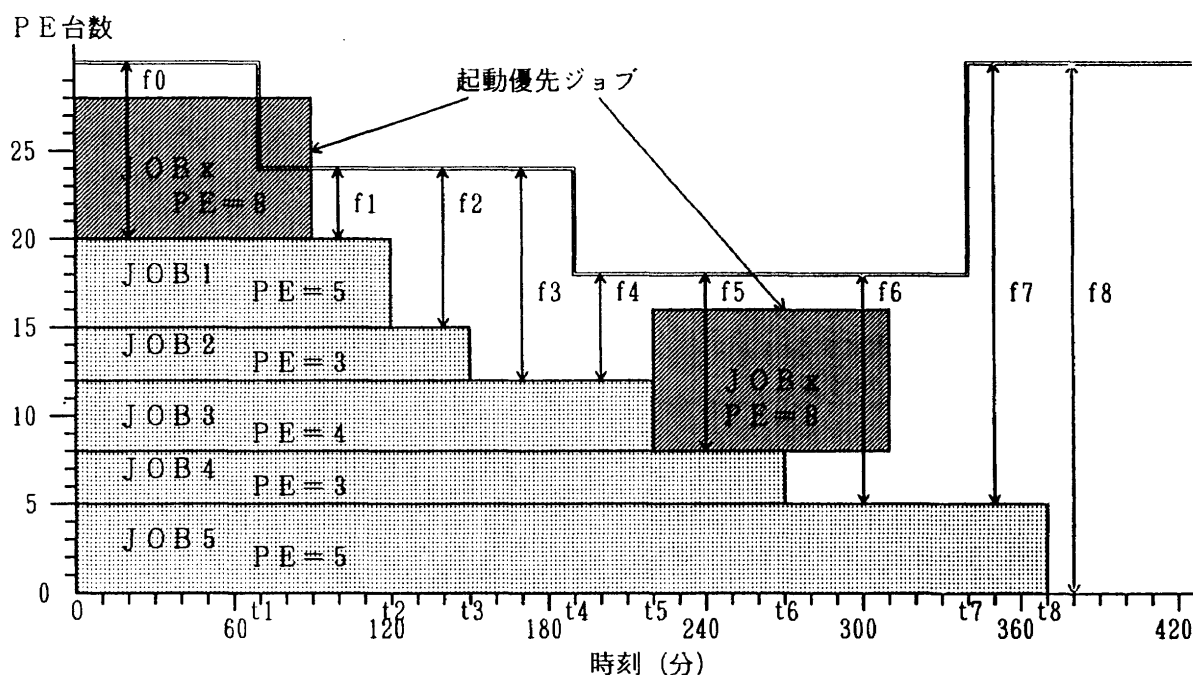


図 3.2 ジョブ起動時の使用可能 PE 台数説明図

(j=1~jmax)

nmax : 空き PE 台数変更事象の数

t0 : 現在の時刻

tn : n 番目に発生する空き PE 台数変更事象の発生時刻 (n=1~nmax)

fn : 時刻 tn における空き PE 台数 (n=0~nmax)

FPEn : 時刻 tn における空き PE 台数の変化量 (n=1~nmax)

PEmax : 時刻 t0 におけるシステムで使用可能な最大 PE 台数

この場合、空き PE 台数の変化事象はジョブの終了および図中の時刻 t1, t4, t7 で示されるような利用可能な PE 台数の変更となる。なお、時刻 tn における空き PE 台数 fn は次式で表される。

$$f_n = \begin{cases} PEmax - \sum_{i=1}^{j_{max}} FPE_n & (n=0 \text{ のとき}) \\ f_{n-1} + \sum_{i=1}^n FPE_n & (n=1 \sim nmax \text{ のとき}) \end{cases}$$

同図では、起動優先ジョブを時刻 t0 から開始した場合と時刻 t5 から開始した場合の例が示されているが、free=「現在の空き PE 台数」と考えて時刻 t1 までに終了しないジョブを t0 に起動した場合には、時刻 t1 に行く予定の使用可能 PE 台数の削減処理が行えないことを示してい

る。すなわち、使用可能最大 PE 台数が変化する場合には、free=「時刻 t0 から時刻 t0+ELAPStj までの間の最小空き PE 台数」としなければならない。ここで、ELAPStj は「起動対象ジョブの PE 占有時間」を表す。なお、この場合の各時刻における空き PE 台数は、表 3.1 に示す表を使用して算出することができる。例では、時刻 t5 に起動可能と判定できるので、t5 を PE リザーブ開始時刻、t5+90 分を PE リザーブ終了時刻として PE をリザーブする。この場合、PE リザーブ開始およびリザーブ終了をシステムの使用可能 PE 台数の変更事象と同様に考えて、表 3.1 を表 3.2 のように変更することにより、PE リザーブ有りの場合にも対応可能となる。

(4) 補足

(1), (2) 項で述べた起動の可否や起動予定時刻の算出には起動対象ジョブの PE 占有時間および実行中ジョブの残り PE 占有時間が必要となる。しかし、NWT システムのジョブ打ち切りのための経過時間は、第 2.1 (4) 項に述べたように、ジョブの出力を保証するために大きめにセットせざるを得ないため、PE 占有時間としては利用できない。仮にユーザの予想する経過時間を知る方法があったとしてもジョブが途中で異常終了する可能性もあり、スケジューラが各ジョブの実 PE 占有時間を正確に知ることは不可能である。したがって、ジョブの PE 占有時間は要求 CPU 時間か

表 3.1 空き PE 台数算出表(1)

n	時刻 (分)	FPEn	fn	備考
0	0	—	10	ジョブ終了
1	70	-6	4	PE 台数削減
2	120	5	9	ジョブ終了
3	150	3	12	ジョブ終了
4	190	-6	6	PE 台数削減
5	220	4	10	ジョブ終了
6	270	3	13	ジョブ終了
7	340	12	25	PE 台数増加
8	370	5	30	ジョブ終了

表 3.2 空き PE 台数算出表(2)

n	時刻 (分)	FPEn	fn	備考
0	0	—	10	ジョブ終了
1	70	-6	4	PE 台数削減
2	120	5	9	ジョブ終了
3	150	3	12	ジョブ終了
4	190	-6	6	PE 台数削減
5	220	4	10	ジョブ終了
6	220	-8	2	PE リザーブ開始
7	270	3	5	ジョブ終了
8	310	8	13	PE リザーブ終了
9	340	12	25	PE 台数増加
10	370	5	30	ジョブ終了

ら予想するものとし、実 PE 占有時間と予想 PE 占有時間の差の程度によりどの程度スケジューリングに影響を及ぼすのかも合わせて検討する。

3.4 ジョブスケジューリング手順

スケジューラの行うべき処理を、キュー登録処理、キュー登録削除処理、キュー遷移処理、起動処理、プレステージ処理に大別し、以下に各々の処理内容を説明する。なお、ジョブ到着時にはキュー登録削除処理を除くすべての処理を、ジョブ終了時およびジョブキャンセル時にはキュー登録処理を除くすべての処理を行う。

(1) ジョブ登録処理

到着したジョブの所属キューおよびキュー内優先度を決定し、決定されたキューに優先度順につなぐ。

(2) ジョブ登録削除処理

ジョブ終了時には実行中ジョブの中から、ジョブキャンセル時には実行待ちジョブまたは実行中ジョブの中から該当するジョブを探し、削除する。

(3) ジョブのキュー遷移処理

- a. α_1 キューの実行待ちジョブを優先度順に調べ、 α_2 キューへの遷移条件を満たしているジョブがあればつなぎ替える。
- b. α_2 キューの実行待ちジョブを優先度順に調べ、 α_3 キューへの遷移条件を満たしているジョブがあればつなぎ替える。
- c. α_1 キューの実行中ジョブを優先度順に調べ、 α_2 キューへの遷移条件を満たしているジョブがあればつなぎ替える。
- d. α_2 キューの実行中ジョブを優先度順に調べ、 α_3 キューへの遷移条件を満たしているジョブがあればつなぎ替える。

(4) 起動処理

本処理は、以下の手順で行う。なお、システムに割り当て可能な空き PE がなくなった時点で本処理を中止する。

- a. スワップアウト中の α_3 ジョブを実行優先順位の高いものから順に調べ、スワップイン可能なものをスワップインする。 α_3 ジョブに対するスワップアウトは行わないが、スワップアウト後に α_3 ジョブとなる場合があるため、この処理が必要となる。
- b. α_3 キューの実行待ちジョブを実行優先順位の高いものから順に調べ、第 3.3 項 a に述べた「起動優先ジョブの起動条件」を満たすジョブが続く限り順次起動する。「起動優先ジョブの起動条件」を満たさないジョブがあった場合には、そのジョブのための PE リザーブ処理を行い、以下の処理へ移る。
- c. α_3 キューの後続のジョブを調べ、第 3.3 項 b に述べた「非起動優先ジョブの起動条件」を満たすジョブが

あれば順次起動する。

- d. スワップアウトされている α_2 ジョブ、 α_1 ジョブがあれば実行優先順位の高いものから順次スワップイン可能なものをスワップインする。
- e. α_2 キュー、 α_1 キューの実行待ちジョブを実行優先順位の高いものから順次調べ、「非起動優先ジョブの起動条件」を満たすジョブがあれば順次起動する。

(5) プレステージ処理

実行待ちジョブの中の優先度の高いジョブから一定数のジョブをプレステージする。そのために、プレステージすべきジョブ名と現在プレステージされているジョブ名を調べて比較し、新たにプレステージすべきジョブ、プレステージを取り消すべきジョブを決定する。決定に従い、プレステージ処理、プレステージの取消し処理を行う。

4. シミュレーション実験におけるワークロード

本章では、第 3 章で述べたスケジューリングアルゴリズムの効果を検証するために行うシミュレーション実験で使用するワークロードについて説明する。なお、ワークロードは、ジョブの属性とジョブの到着率から決定する。

4.1 ジョブの属性

NWT システムは、CFD プログラムの高速実行を目的としている。そのため、シミュレーションでも、航技研の CFD プログラムの実行を前提とする。

以下に、CFD プログラムの特徴とシミュレーションで使用するジョブ属性の決定方法について述べる。

(1) 航技研における CFD プログラムの特徴

第一期 NS システムで処理されていた一般的な CFD プログラムの特徴^{注5)}を以下に列挙する。

- 格子点数 10 万点当たりの計算に 3600 秒程度^{注6)}の CPU 時間を必要とする。
- 格子点当たりのメモリ量は 200 B から 400 B 程度のジョブが多い。
- 定常状態の計算を行うジョブ（定常ジョブと呼ぶ）はジョブ実行開始直後の入力と終了直前の出力があり、計算途中でのデータ入出力はほとんどない。これに対し、非定常の計算を行うジョブ（非定常ジョブと呼ぶ）は一定間隔で出力があり、CPU 使用時間は出力回数に応じて回

(注5) 本スケジューリングアルゴリズムの検討および実験を行った時点では NWT システムはまだ導入されていなかったため、第一期 NS システムでのデータを使用した。

(注6) これは VP-400 を使用した場合の値である。PE の処理能力は VP-400 の処理能力より勝っていると予測されたが、本実験を行った時点では明確な実測値が得られていなかったためこの値をそのまま使用した。

数倍に増える。

- 運用上の統計データによれば、計算結果の出力量はメモリ量の10%程度であり、入力量はメモリ量の18%程度である。入力量は出力量より格子データの分増える。
- 計算プログラムのデバッグ処理は本計算の約1/20程度のCPU時間を要する。
- 実在気体の計算は、一般的な完全気体の計算に比べてメモリ量、CPU時間ともに2倍程度必要である。

また、CFDプログラムをNWTシステムで処理する場合には、

- メモリ量の増加率（メモリオーバーヘッド率と呼ぶ）は1.2から2程度である。

と予想されていた。並列化に伴うメモリ量増加の主な原因は、データ領域の重複である。NWTシステムはメモリ分散型の並列計算機であり、各PE内のメモリの情報はクロスネットワークを介して他のPEで参照することが可能であるが、他のPEの情報を参照するためには双方のPEに領域を重複してとる必要があるため、メモリ量が増加する。

なお、非定常ジョブは定常ジョブを連続して行うこととはほぼ同等であり、入力が無い分だけSSUの負荷が軽くなるため、スケジューリングアルゴリズムを検証するためのシミュレーションでは非定常ジョブを考慮する必要がないと考える。

以上の結果、シミュレーション用ジョブ属性は、メモリ量、格子点数当たりのデータ量、メモリオーバーヘッド率、ジョブファクタ（メモリ量、CPU時間に対するファクタで、実在気体を扱うジョブと完全気体を扱うジョブを発生させるために使用する）、CPUファクタ（CPU時間に対するファクタで、デバッグジョブを発生させるために使用する）の5つの要素を基にして決定することとした。この5つの

要素をジョブ属性の基本要素と呼ぶ。

(2) ジョブ属性の基本要素決定方法

NWTシステムではジョブ属性に未知の要素が多かったため、ジョブ属性の基本要素はそれぞれ独立に決定した。また、属性を決定する際の乱数の発生には、発生データの最小値および最大値が規定でき、パラメータの指定により平均値の異なるデータを発生することができる β 分布¹³⁾を使用した。

実験ではジョブファクタ、CPUファクタは固定し、デバッグジョブの発生率は実験の目的に応じて設定する。また、メモリ量、格子点数当たりのデータ量、メモリオーバーヘッド率は、 β 分布のパラメータをそれぞれ3種類、2種類、2種類用意した。したがって、1つの実験に対して12通りのワークロードを使用したシミュレーションが可能となる。そのため、実験結果の比較においては、メモリ量のタイプを i ($i=1\sim3$)、格子点数当たりのデータ量のタイプを j ($j=1\sim2$)、メモリオーバーヘッドのタイプを k ($k=1\sim2$)とした場合、使用したワークロードを $i/j/k$ のように表す。また、12通りの結果をメモリ量のタイプごとに集計した場合には、 $i/*/*$ ($i=1\sim3$)のように表す。

以下に、基本要素の決定方法を示す。ここで、各項目の後ろの[]内の記号は、各項目のデータの略号として使用する。

① メモリ量 [RMEM] (単位 MB)

ジョブのメモリ量は、80 MB から30000 MB の範囲の β 分布で決定する。最大値の30000 MBはPE 120台を使用するジョブを想定したものである。なお、 β 分布のパラメータはタイプ1 (0.5, 8.0)、タイプ2 (0.5, 5.0)、タイプ3 (0.5, 2.0)の3種類使用し（図4.1(a)参照）、メモリ量の平均値がそれぞれ1840, 2800, 6064の場合の実験を行う。

② 格子点当たりのデータ量 [DGRID] (単位 B)

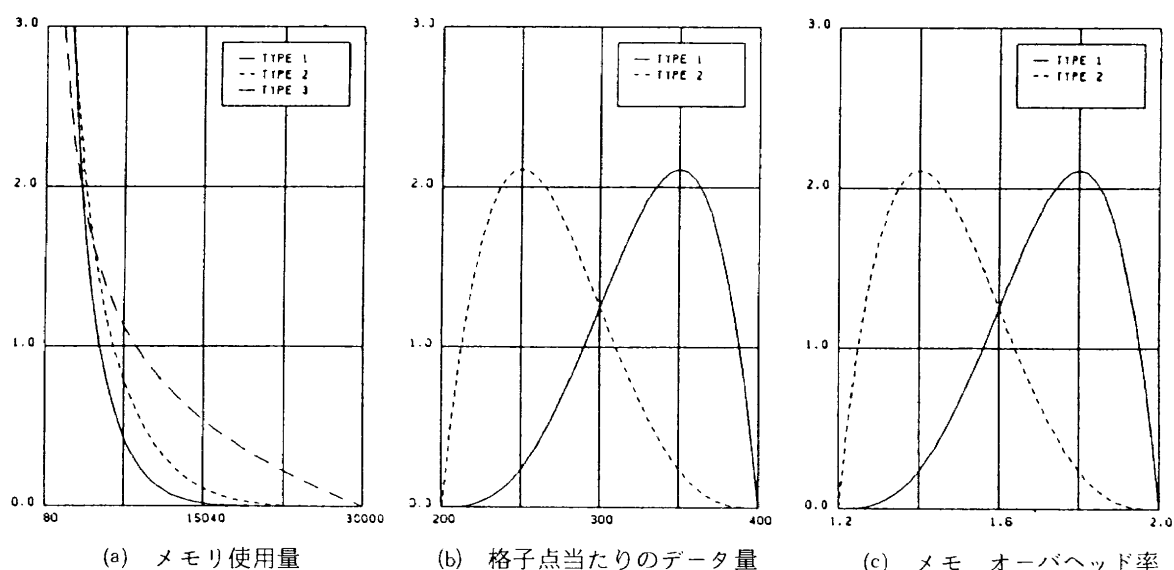


図 4.1 ジョブ属性決定に使用する β 分布

ジョブの格子点当たりのデータ量は、200Bから400Bの範囲の β 分布で決定する。なお、 β 分布のパラメータはタイプ1 (4.0, 2.0), タイプ2 (2.0, 4.0) の2種類使用し(図4.1(b)参照), 格子点当たりのデータ量の平均値がそれぞれ333.3, 266.6の場合の実験を行う。

③ メモリオーバーヘッド率 [OH]

ジョブのメモリオーバーヘッド率は、1.2から2.0の範囲の β 分布で決定する。なお、 β 分布のパラメータはタイプ1 (4.0, 2.0), タイプ2 (2.0, 4.0) の2種類使用し(図4.1(c)参照), メモリオーバーヘッド率の平均値がそれぞれ1.734, 1.466の場合の実験を行う。

④ ジョブファクタ [FACT]

完全気体のジョブと実在気体のジョブを1:1の割合で発生させ、以下の値を設定する。

$$\text{FACT} = \begin{cases} 1.0 & (\text{完全気体の場合}) \\ 2.0 & (\text{実在気体の場合}) \end{cases}$$

⑤ CPU ファクタ [CFACT]

デバッグジョブを設定された割合で発生させ、以下の値を設定する。

$$\text{CFACT} = \begin{cases} 1 & (\text{非デバッグジョブの場合}) \\ 1/20 & (\text{デバッグジョブの場合}) \end{cases}$$

(3) 基本要素以外のジョブ属性の決定方法

以下に、基本要素以外のジョブ属性の決定方法について述べる。

①使用PE台数 [PE]

$$\text{PE} = \text{RMEM}/250.$$

②PE当たりのメモリ量 (単位 MB)

$$\text{PEMEM} = \text{RMEM}/\text{PE}$$

③基本メモリ量 [BMEM] (単位 MB)

$$\text{BMEM} = \text{RMEM}/\text{OH}$$

④計算格子点数 [GRID]

$$\text{GRID} = \text{BMEM} * 1000000 / (\text{DGRID} * \text{FACT})$$

⑤PE当たりの計算格子点数 [PEGRID]

$$\text{PEGRID} = \text{GRID}/\text{PE}$$

⑥予定CPU使用時間 [ACPU] (単位 ms)

$$\text{ACPU} = (\text{PEGRID}/100000.) * 3600. * \text{FACT}$$

⑦CPU使用時間 [CPU] (単位 ms)

アベンドジョブを1割の割合で発生させ、以下の値を設定する。

$$\text{CPUT} = \begin{cases} \text{ACPU} & (\text{非アベンドジョブの場合}) \\ \text{ACPU} * \text{AVRATE} & (\text{アベンドジョブの場合}) \end{cases}$$

ここで、AVRATEは0.0~1.0の一樣乱数により決定する。

⑧入力データ量 [INDATA] (単位 MB)

$$\text{INDATA} = \text{RMEM} * 0.18$$

⑨出力データ量 [OUTDATA] (単位 MB)

$$\text{OUTDATA} = \text{RMEM} * 0.10$$

表4.1は、12種類のワークロードを使用してそれぞれ1000件のジョブを発生し、メモリ量の発生タイプ別にそれぞれの平均値を求めた結果である。

4.2 ジョブ到着時間間隔

システムの負荷は、システムに到着するジョブの平均システム資源要求量とジョブの平均到着時間間隔により決定される。したがって、システムに到着するジョブの属性が決定している場合には、ジョブの到着間隔によりシステムの負荷状態が決まる。

いま、システムの処理能力相当の負荷状態を負荷率1.0, システムの処理能力の x 倍の負荷状態を負荷率 x とする。ジョブの平均CPU使用時間を cpu , 平均使用PE台数を pe , 使用可能PE台数を maxpe としたとき、負荷率 x 場合の到着時間間隔 $\text{ARRIV}(x)$ は以下のようにして求めることができる。

$$\text{ARRIV}(x) = \frac{\text{cpu} \times \text{pe}}{\text{maxpe} \times x}$$

一般的に、負荷状態の異なるシミュレーション実験は負荷率 x を変えることにより行うことができる。しかし、本アルゴリズムの場合には、ジョブ優先度の決定に待ち時間を使用しているため、処理能力以上の負荷を想定した場合には、単に負荷率を上げただけでは除々にジョブが蓄積し優先度設定の意味がなくなってしまうこととなる。実際のシステムの場合、混雑時にはジョブの投入を控えるユーザがいるためシステム内ジョブ数が一定値を越えると到着間隔

表4.1 ワークロード別ジョブ属性

	単位	最小値	最大値	1/*/*	2/*/*	3/*/*
PE台数	台	1	120	7.9	11.8	24.8
主記憶量	MB	80	29760	1846	2821	6050
予定CPU時間	秒	3920	35184	15930	16637	17620
実CPU時間	秒	8	35184	15090	15712	16821
入力データ量	MB	7	3856	212	323	692
出力データ量	MB	4	2142	118	179	384
格子点数	千点	54	90672	3012	4654	10020

は長くなるか到着しないようになるのが普通である。そこで、システム内の全ジョブを処理するのに要する時間の最大値（最大積算時間と呼ぶ）と最小値（最小積算時間と呼ぶ）を定義し、最大積算時間に達した場合にはジョブの到着をストップし、最小積算時間以下になった場合にジョブの到着を開始する方式を採用して、より現実的な状況を作ることにした。ジョブの最大積算時間を DMAX, 最小積算時間を DMIN とした場合、システム内のジョブの最大件数 JMAX, 最小件数 JMIN は以下のようにして求める。

$$JMIN = \frac{DMIN}{ARRIV(1.0)}$$

$$JMAX = \frac{DMAX}{ARRIV(1.0)}$$

なお、シミュレーション実験では、最小積算時間を72時間（3日）、最大積算時間を120時間（5日）としている。

ここで、ディバグジョブ発生率の変化に伴い、発生するジョブ属性と到着率がどのように変化するかを簡単に説明する。第4.1(1)項に述べた理由により、ディバグジョブのCPU時間は本来のジョブのCPU時間の1/20とする。したがって、ディバグジョブの発生率0.0の場合の平均CPU時間を1とすると、発生率0.5のときの平均CPU時間は以下に示すように0.525、発生率0.8のときの平均CPU時間は0.240となる。したがって、ディバグジョブ発生率0.0の場合と同等の負荷状態にするためには、到着時間間隔をそれぞれ0.525倍、0.240倍とする必要がある。

$$1.0 \times 0.5 + 0.05 \times 0.5 = 0.525$$

$$1.0 \times 0.2 + 0.05 \times 0.8 = 0.240$$

シミュレーション実験では、上記のようにディバグジョブ発生率を考慮して到着時間間隔を決定する。

5. シミュレーション結果の評価・検討

本章では、スケジューリングアルゴリズムの有効性を実証するために、以下に示す6種類のシミュレーション実験を行い、実験結果の評価・検討を行う。

(1) 基本実験

本スケジューリングアルゴリズムを使用した場合のシステム資源使用状況およびジョブ処理状況の傾向を調べる。また、ジョブクラス別多重度制御方式の場合と比較してシステム資源使用状況がどの程度改善されるかを予測する。

(2) 負荷実験

負荷率が変化した場合のシステム効率、ジョブ処理状況への影響を調べ、特定の条件のジョブにのみ影響を及ぼすことがないことを確認する。

(3) スワップアウト効果実験

起動優先ジョブを起動させる際のスワップアウトジョブ数に対するシステム効率、ジョブ処理状況への影響を調べる。また、SSUの有効利用を図るため、スワップアウトの

適正条件を予測する。

(4) PE リザーブ効果実験

起動優先ジョブを最優先で起動するためのPEリザーブは、すべてのジョブの起動を抑止した場合と比較してどの程度システム効率に効果があるかを調べる。

(5) プレステージ効果実験

ジョブ開始時の入力待ち時間を削減するための手段であるプレステージ処理のシステム効率、ジョブ処理状況への影響を調べる。また、SSUの有効利用を図るため、プレステージの適正条件を予測する。

(6) PE 占有時間精度実験

ジョブの予想PE占有時間は、要求CPU時間に一定のファクタを掛けて求めるが、ファクタの設定値がシステム効率にどの程度の影響を及ぼすかを調べる。また、ファクタの適正値を予測する。

なお、本稿では、システム結果はデータの最終値で比較するため、システムが定常状態にまで達しており、結果が振動していないことが必要条件となる。このような状態を安定状態と呼び、以下の安定条件で判定する。

本スケジューラは、使用PE台数、CPU要求時間および待ち時間によりジョブの起動からジョブ終了までの流れを制御しているため、システムが安定状態に達した時には経過時間に対する実CPU時間の割合（これをジョブ実行率と定義する）の平均値も安定することが期待できる。そこで、ジョブを使用PE台数が10台以下のジョブ（JOB1と呼ぶ）、11～40台のジョブ（JOB2と呼ぶ）、41台以上のジョブ（JOB3と呼ぶ）という3つのグループに分類し、JOB1とJOB2の平均ジョブ実行率は少数点以下2桁以上の値が連続10回以上変動がなく、発生件数の少ないJOB3の平均ジョブ実行率は少数点以下1桁以上の値が連続10回以上変動がないことをシミュレーションの安定条件とする。なお、実験結果の比較の都合上、シミュレーション時間はすべて50日分とし、安定条件を満たしているか否かは別に確認する。

5.1 基本実験

基本実験では、本スケジューリングアルゴリズムを使用した場合のシステム資源使用状況およびジョブ処理状況の傾向を見るために、通常の状態を想定した負荷率0.8の実験、混雑時を想定した負荷率1.0の実験、システムの極限状態を想定した負荷率2.0の実験を行った。

(1) 実験条件

基本実験の実験条件として以下のようにパラメータを設定する。

- 負荷率：0.8, 1.0および2.0
- ディバグ率：0.0, 0.5および0.8
- α_2 キューへの遷移パラメータ (FAC1)：0.25

- α_3 キューへの遷移パラメータ (FAC2): 3.0
- プレステージジョブ数: 3
- スワップアウトパラメータ: (4, 10)
- FEP 配下の磁気ディスクの入出力速度: 18MB/S
- PE 占有時間ファクタ: 1.0

ここで、上記パラメータについて、簡単に説明する。

- ①FAC1, FAC2 の設定により、 α_2 キュー、 α_3 キューへの遷移条件である WT_{j1} 、 WT_{j2} は、以下ようになる。

$$WT_{j1} = [\text{予定 CPU 時間} \times \sqrt{\text{PE 台数}}] \times 0.25$$

$$WT_{j2} = [\text{予定 CPU 時間} \times \sqrt{\text{PE 台数}}] \times 3.0$$

予定 CPU 時間を 1 時間と仮定すると、 α_2 キューへの遷移時間は PE 台数が 4 台のジョブは 0.5 時間、PE 台数が 36 台のジョブは 1.5 時間、PE 数が 100 台のジョブは 2.5 時間となる。同様に、 α_3 キューへの遷移時間は PE 台数が 4 台のジョブは 6 時間、PE 台数が 36 台のジョブは 18 時間、PE 数が 100 台のジョブは 30 時間となる。

- ②スワップアウトパラメータの (4, 10) は、1 つのジョブを起動する際の最大スワップアウト可能 PE 台数が 4 であり、システムの最大スワップアウト可能 PE 台数が 10 であることを表す。また、プレステージ処理では、起動待ちジョブの中の優先度の最も高いジョブから指定された数のジョブを常時プレステージする。

- ③PE 占有時間ファクタは、スケジューラがジョブの要求 CPU 時間からジョブの PE 占有時間を予測する際に使用する。この場合には、予定 PE 占有時間 = CPU 時間 \times 1.0 となる。

(2) 実験結果の評価・検討

本実験では、ディバグ率が 0.0, 0.5, 0.8 の場合についてそれぞれ 12 のワークロードを使用した実験を行った。ここで、SSU 容量は十分にあることを前提とし、どの程度の SSU 容量が必要であるかを調べるという観点で実験した。

なお、割り当ては 8MB 単位で行った。また、ジョブ起動後のプログラムおよび入力データの転送時間は PE の割り当て後に行われるため、PE 使用時間の中に含まれるべき種類のものである。しかし、これらの値は PE の効率的使用の観点からはできるだけ短いことが望まれること、これらの時間を PE 使用時間に含めた場合には PE が有効に使用されている時間の評価が難しくなることから、ジョブ起動後のプログラムおよび入力データの転送時間は PE 使用時間に含めないことにした。したがって、実システムの平均 PE 使用率に比ベシミュレーションの平均 PE 使用率は実システムの計算方法に比べて低い値を示すことになる。なお、平均 PE 使用台数および平均 SSU 使用量は、平均使用量を 1 時間単位に集計して求めている。

(a) PE 使用状況

図 5.1 は、基本実験結果の平均 PE 使用台数をワークロードのメモリ量発生タイプごとに平均し、棒グラフで示したものである。

この図から、一般的な傾向として、負荷率と平均 PE 使用台数の関係、ジョブ規模と平均 PE 使用台数の関係を挙げるができる。

低負荷を想定した負荷率 0.8 の場合の平均 PE 使用台数は最も高い場合でも 120 台程度であるのに対し、負荷率が 1.0 の場合には 125 台以上、負荷率が 2.0 の場合には大半が 130 台以上となっており、負荷率が高くなるにつれて、平均 PE 使用台数も確実に高くなっている。しかし、例外として、ディバグ率 0.0 の場合には、負荷率 1.0 の時より負荷率 2.0 の時の平均 PE 使用台数が若干減少している。これは、ディバグ率が低いために、他の場合に比べてシステム内のジョブは大規模、長時間ジョブが多くなり、最優先ジョブのための PE リザーブにより起動できないジョブが多かったものと考えられる。

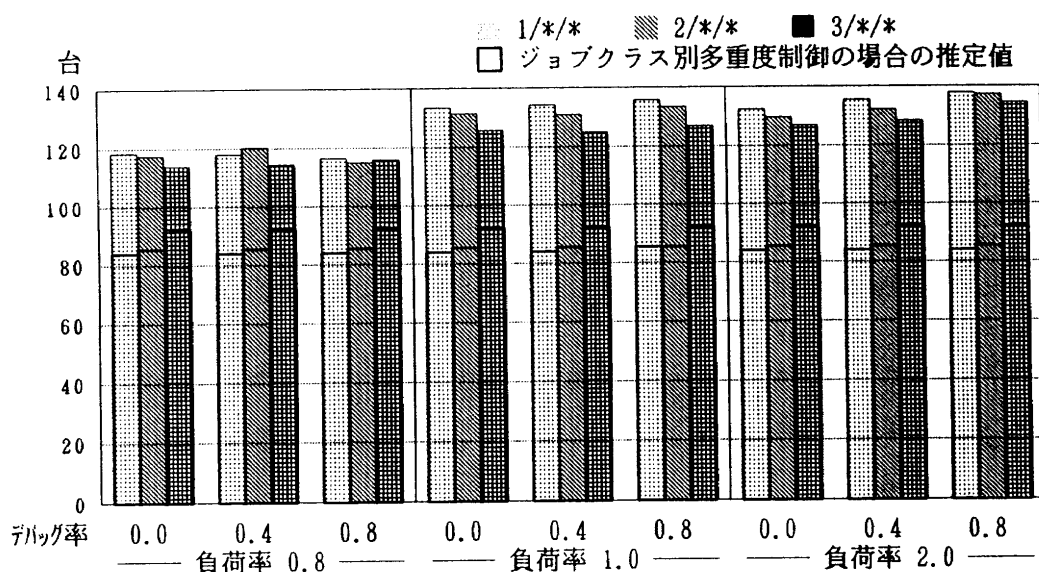


図 5.1 基本実験における平均 PE 使用台数

また、ジョブの規模が大きい場合平均 PE 使用台数が若干の低下が見られるが、これは、実行待ちジョブが多い場合、規模の小さいジョブの方が起動し易いことに起因しており、負荷率 0.8 の場合のように実行待ちとなるジョブが少ない場合にはその影響がほとんど現れていない。

なお、図 5.1 の太枠で囲まれたグラフは、ジョブクラス別多重制御方式でジョブを処理した場合の推定平均 PE 使用台数であり、以下のように算出した。

ジョブクラス多重度制御方式では、多重度分のジョブが実行されていた時には運用 PE 台数がすべて使用されるように設定することにより混雑時には PE の使用率を 100% にすることができるが、ジョブクラスの設定 PE 台数に満たないジョブを実行した時には空き PE が発生する。そのため、推定平均 PE 使用台数 U_{pe} は、運用 PE 台数と PE 有効使用率 E_{rate} （割り当てられた PE が有効に使用されている割合）の積から算出することができる。以下に具体的な算出式を示す。

$$U_{pe} = \text{運用 PE 台数} \times E_{rate}$$

ここで、

$$E_{rate} = \frac{\sum_j (\text{使用 PE 台数} \times \text{経過時間})}{\sum_j (\text{占有 PE 台数} \times \text{経過時間})}$$

である。なお、算出式の中の $\sum_j (...)$ は、すべてのジョブ情報の加算を意味しており、基本実験で得られた課金情報を使用して求めた。

この結果、推定平均 PE 使用台数を求める際の条件を満たしている負荷率 1.0、負荷率 2.0 の場合には、実験結果における平均 PE 使用台数は、ジョブクラス別多重度制御方式の場合の平均 PE 使用台数の推定値より平均 43.90 台分勝っていた。平均 PE 使用台数はシステムのジョブ処理量に直接影響を及ぼすためシステム効率の重要な指標となるが、上記結果を見る限りジョブクラス別多重度制御方式に比べて本アルゴリズムの優位性は明瞭である。

(b) SSU 使用状況

SSU は主にプレステージ、入出力、スワップアウトに使

用されるため、実験結果もこれらの使用量について評価する。

基本実験の平均プレステージ使用量を図 5.2 に示す。プレステージ使用量は、プレステージジョブ数を 3 としていることから負荷率よりもワークロードによる平均 PE 台数の違いの影響が大きく表れている。プレステージ量を負荷率で比較すると、負荷率 1.0 の場合が最も多く、ついで負荷率 2.0、負荷率 0.8 の順になっているが、これはプレステージされたジョブの規模の違いと考えられる。負荷率 0.8 ではジョブの競合が少ないためプレステージされるジョブは大規模のものに限定されないが、負荷率 1.0 では処理されにくい大規模ジョブがプレステージされ、負荷率 2.0 では全体的に処理が滞ったため比較的小規模なジョブもプレステージされたものと考えられる。図 5.3 に平均スワップアウト使用量を示す。負荷が低い場合には α_3 キューまで進むジョブが少ないためスワップアウト使用量は極めて少ないが、負荷率が高くなるに従って単純に増加している訳でもない。これは、ひとつのジョブを起動する際のスワップアウト PE 台数の制限を 4 としたために、比較的小規模なジョブが起動優先ジョブとならない限りスワップアウトの条件が満たされないためと考えられる。このことは、小規模ジョブの発生率が高い $1/*/*$ のスワップアウト使用量が相対的に多いことから裏付けられる。

図 5.4 に平均入出力使用量を示す。入出力使用量はプレステージ使用量に比べて極端に少ない。これは、ジョブの CPU 時間がジョブクラス制限値という形で制限されていないために、CPU 時間に対する相対的な入出力量が減少したためと考えられる。このことは、デバッグ率が高くなるのに伴い入出力量が相対的に増加していることから明らかである。したがって、適切な SSU 容量を見積もる際には、入出力量はプレステージ量やスワップアウト量程厳密に行う必要はないが、実運用において CPU 時間を制限する場合には増加することを考慮すべきである。

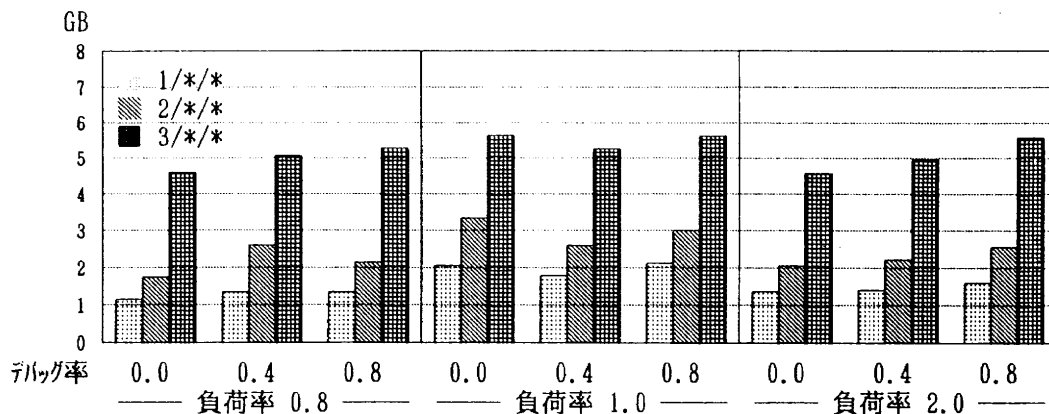


図 5.2 基本実験における平均プレステージ使用量

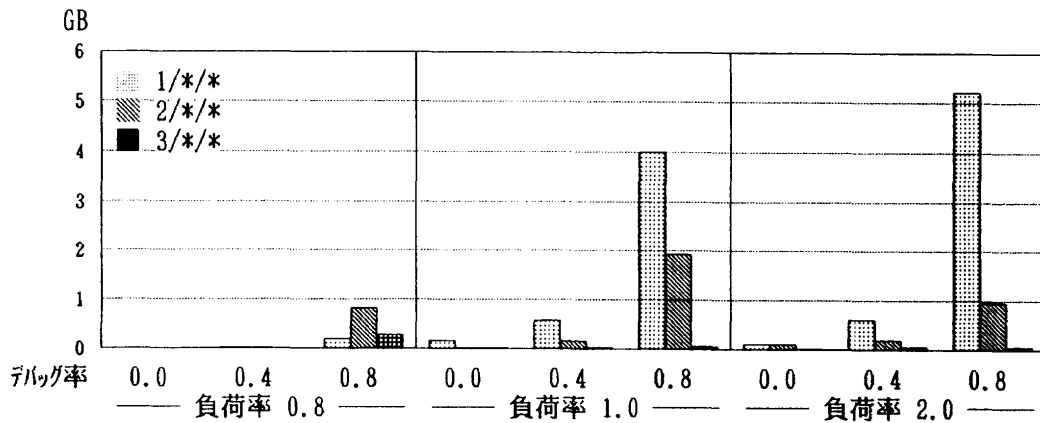


図 5.3 基本実験における平均スワップアウト使用量

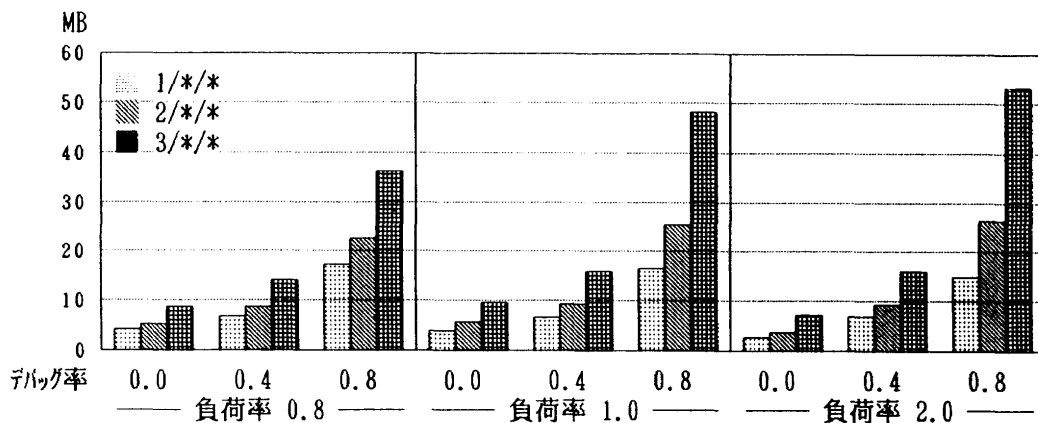


図 5.4 基本実験における平均入出力使用量

(c) ジョブ処理状況

図 5.5 に基本実験結果のジョブ実行率を、図 5.6 に α_3 キューまで遷移したジョブの割合 (以後、 α_3 キューへの遷移率と呼ぶ) を示す。

本章の冒頭で述べたように、ジョブ実行率は経過時間に対する実 CPU の割合であるため、値が 1 に近い待ち時間が少なくジョブは良好に処理されていると考えることができる。負荷率 0.8 の場合のジョブ実行率は、JOB1 は 0.89 以上、JOB2 は 0.62 以上、JOB3 は 0.25 以上を示しており、単純計算すると、PE 要求台数が 41 台以上である JOB3 でさえ CPU 時間の 4 倍程度で処理されていることになる。これに対し、負荷率 1.0 の場合のジョブ実行率は、JOB1 は 0.53 以上、JOB2 は 0.17 以上、JOB3 は 0.05 以上を示しており、JOB1 のジョブは CPU 時間の 2 倍程度、JOB2 のジョブは CPU 時間の 5 倍程度、JOB3 の処理時間は CPU 時間の 20 倍程度かかっていることになる。また、負荷率 2.0 の場合には、JOB1 は 0.14 以上、JOB2 は 0.03 以上、JOB3 は 0.01 以上となり、ジョブの処理状況は全般的に悪化している。

また、 α_3 キューへの遷移率は、負荷の上昇に伴って高くなる傾向があり、ジョブの規模が小さいワークロードの場

合程この傾向が顕著に現れている。この率が極端に高くなった場合には、優先方式の異なる 3 本のキューを用いている意味がなくなる上、起動優先ジョブの処理が増加するためシステムオーバーヘッドも増加する。したがって、 α_3 キューへの遷移率は一定値以下に抑えることが必要となる。

一般的に、ジョブ実行率は負荷が低い場合には高く、負荷が高い場合には低くなるのは当然の結果であるため、これらの数値から単純にスケジューラの良否を判定することはできない。そこで、負荷率の変化に伴いジョブ処理状況がどの程度変化するかを、第 5.2 項の実験により詳細に調べる必要がある。

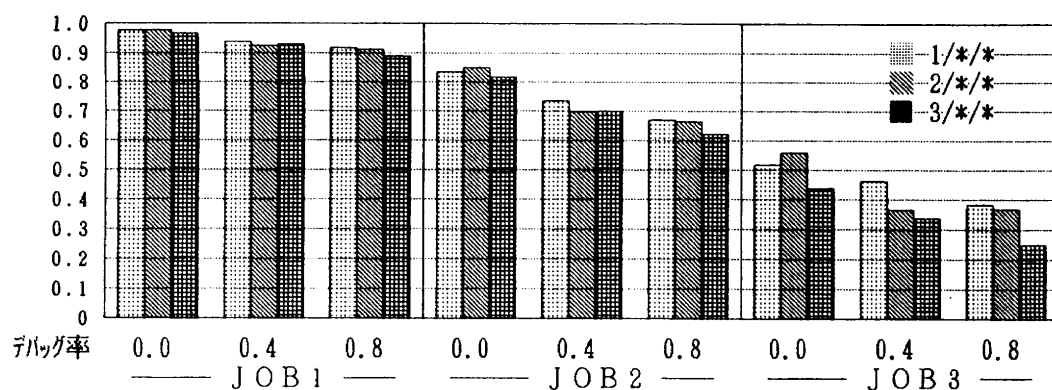
(d) 安定条件

シミュレーションの安定条件に関しては、負荷率 0.8 の場合には 1 日から 41 日程度、負荷率 1.0 の場合には 2 日から 22 日程度、負荷率 2.0 の場合には 1 日から 44 日程度ですべて条件を満たしたことを確認した。

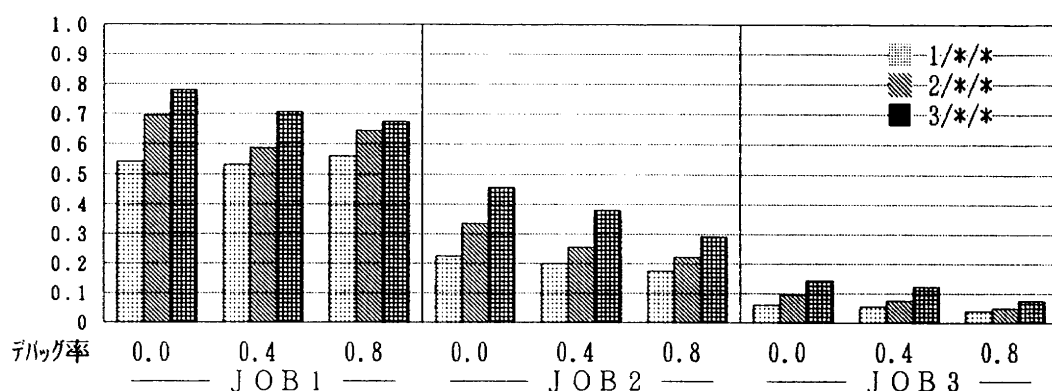
(3) 基本実験から得られた結論

本実験結果から以下の事柄が確認された。

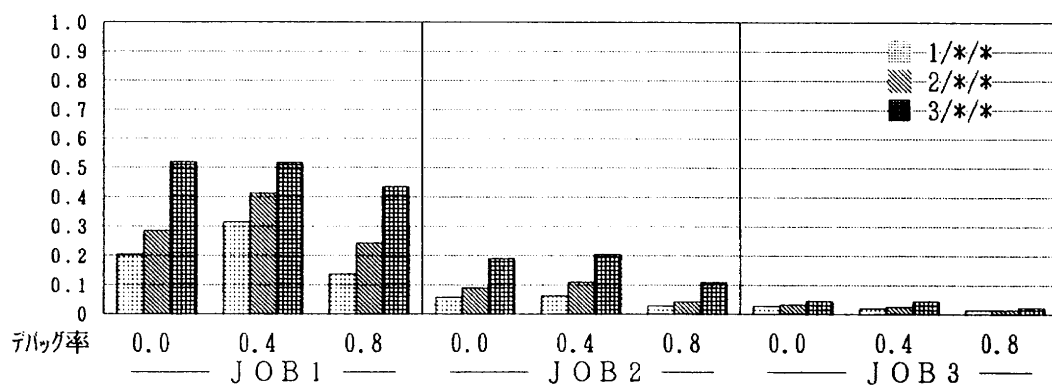
- ① PE の有効利用に関して、本方式はジョブクラス別多重度制御方式に比べて顕著な優位性がある。
- ② 負荷が十分にあれば、ワークロードの規模に関係なく



(1) 負荷率 0.8



(2) 負荷率 1.0



(3) 負荷率 2.0

図 5.5 基本実験におけるジョブ実行率

PE の高利用が可能であると考えられる。

- ③ α_3 キューへの遷移率を一定に抑えるために、混雑時には α_2 キューや α_3 キューへの遷移パラメータを大きめに設定する必要がある。

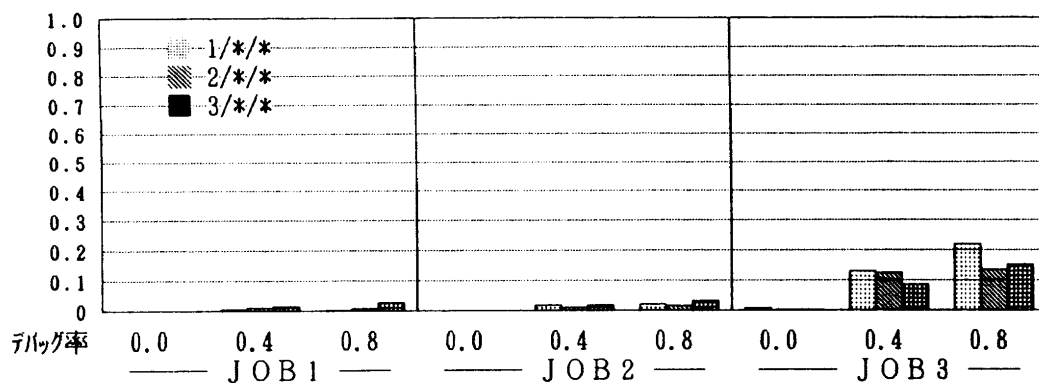
5.2 負荷実験

ほとんど待つこと無しにジョブが処理されていたシステムでも、投入ジョブ数の急激な増加や大規模ジョブ、長時間ジョブの増加に伴う負荷率の上昇により、経過時間が大幅に増加することがある。負荷が増加すればジョブの平均

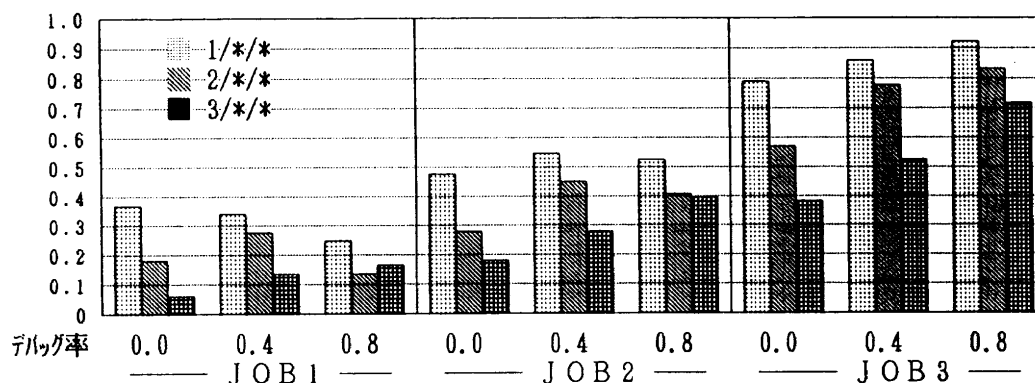
経過時間が長くなるのは当然の事であるが、特定の条件を満たすジョブの経過時間のみが大幅に増加するようなスケジューラは好ましいものではない。このような観点から、考案したスケジューリングアルゴリズムが負荷率の変化にどのように影響されるのか、その影響はワークロードによりどの程度異なるのかを確認する。

(1) 実験条件

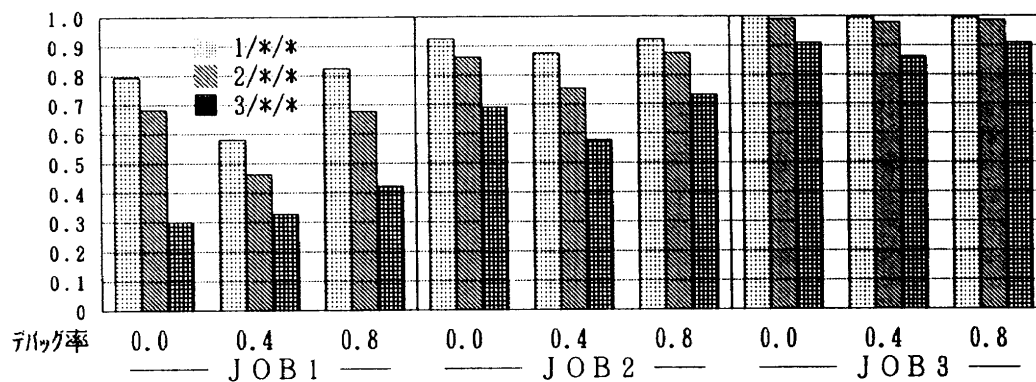
基本実験は負荷率 0.8, 1.0, 2.0 で行った。これに対し本実験では、2 種類のワークロードを使用して 0.8 から 0.2 刻みで 2.0 までの負荷率での実験を行い、その影響を比較



(1) 負荷率 0.8



(2) 負荷率 1.0



(3) 負荷率 2.0

図 5.6 基本実験における α_3 キュー遷移率

する。なお、ワークロードは平均 PE 要求台数および平均 CPU 要求時間が最も小さい 1/1/1 と最も大きい 3/2/2 を使用し、その効果をさらに顕著な形で実行するためにワークロード 1/1/1 はデバッグ率 0.8 で、ワークロード 3/2/2 はデバッグ率 0.0 で行うことにした。したがって、負荷実験の実行条件として以下のようにパラメータを設定した。

- 負荷率 : 0.8, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0
- デバッグ率 : 0.8 (1/1/1 の場合) または 0.0 (3/2/2 の場合)
- α_2 キューへの遷移パラメータ (FAC1) : 0.25

- α_3 キューへの遷移パラメータ (FAC2) : 3.0
- プレステージジョブ数 : 3
- スワップアウトパラメータ : (4, 10)
- FEP 配下の磁気ディスクの入出力速度 : 18 MB/S
- PE 占有時間ファクタ : 1.0

(2) 実験結果の評価・検討

(a) PE 使用状況

図 5.7(1) に負荷実験結果の平均 PE 使用台数を示す。ワークロード 1/1/1 を使用した実験では、負荷率が高くなるに従って僅かずつではあるが平均 PE 使用台数が増加して

いる。これに対し、ワークロード 3/2/2 を使用した実験では、低下する場合もある。これは、起動優先ジョブのための PE リザーブの影響であり、ジョブの規模が比較的小さい場合には問題がないが、規模が大きい場合にはジョブのターンアラウンドタイムの保証と引き換えに PE 稼働率が若干低下することを示している。

(b) ジョブ処理状況

図 5.7(2), (3) にワークロードが 1/1/1, 3/2/2 の場合のジョブ実行率を示す。ワークロードが 1/1/1 の場合も 3/2/2 の場合も、負荷率の上昇に伴い JOB1, JOB2, JOB3 のジョブ実行率は同様に低下していることが判る。

図 5.8(1), (2) にワークロードが 1/1/1, 3/2/2 の場合の α_3 キュー遷移率を示す。負荷の上昇に伴い α_3 キューに達したジョブは予想通り増加しているが、ワークロード 3/2/2 の場合にも JOB1 の値が負荷率 1.4 を境に減少している。これは、ジョブの規模が大きく、負荷が高いため PE 要求台数の大きいジョブが起動されずに大量に残っているため、規模の小さいジョブが起動されやすくなったためと考えられる。また、ワークロードが 1/1/1 の場合にはワークロード 3/2/2 の場合に比べて平均 PE 要求台数、平均 CPU 要求時間が小さいため、発生するジョブ数は相対的に多い。そのため、JOB1 で α_3 キューに達しているジョブの割合は、ワ

ークロードが 1/1/1 の場合には負荷率が 1.4 程度でも 80% 以上であるのに対し、ワークロード 3/2/2 の場合には負荷率 2.0 でも 20% 以下となっている。

(c) 安定条件

シミュレーションの安定条件に関しては、ワークロードが 1/1/1 の場合には 1 日から 6 日程度で、ワークロードが 3/2/2 の場合には 6 日から 47 日程度ですべて条件を満たしたことを確認した。この時間から見てもワークロードが 1/1/1 の場合のジョブ処理状況は安定しているが、3/2/2 の場合には起動優先ジョブのための PE リザーブの影響によってジョブ処理状況が少し不安定になっていることが判る。

(3) 負荷実験から得られた結論

負荷実験結果の検討から以下の事柄が確認された。

- ① ジョブの規模が小さいワークロードの場合には負荷率に応じて PE の利用率も高い値を示している。一方、ジョブの規模が大きいワークロードの場合には、起動優先ジョブのための PE リザーブによりジョブのターンアラウンドタイムと引き換えに平均 PE 使用台数の低下が見られるが、低下の程度は小さく、許容範囲内であると考えられる。
- ② 負荷率の上昇によりジョブ実行率は低下するが、ジョブの規模等に係わらず全体的に低下しており、本アルゴリ

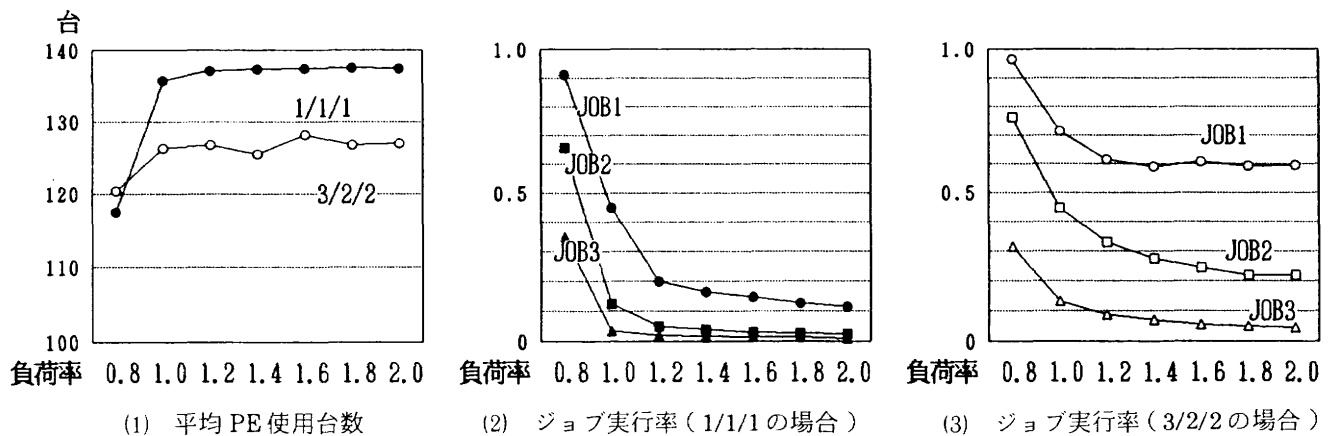


図 5.7 負荷実験結果

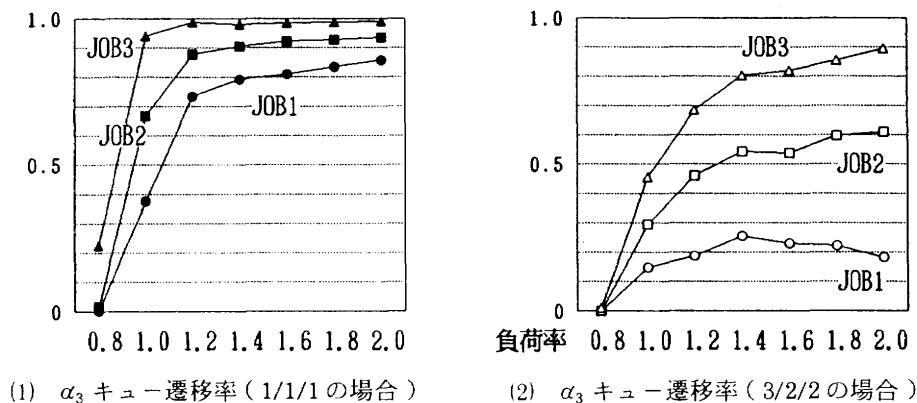


図 5.8 負荷実験結果

ズムは負荷率の変化にも柔軟に対処可能であると判断できる。

- ③ α_3 キューのジョブは優先実行ジョブであるにも係わらず、大半のジョブが優先実行ジョブとなった場合には優先実行の効果が期待できない。そのため、 α_3 キューに達するジョブ数があまり増加しないように、ジョブのキュー遷移パラメータを調節することが望まれる。しかし、負荷率が同じでもジョブ数が多い場合には α_3 キューに達するジョブ数が増加するため、キュー遷移パラメータは、負荷率、ジョブ数等を考慮して決定しなければならない。

5.3 スワップアウト効果実験

本スケジューラでは、適切なジョブターンアラウンドタイムを保証するために、起動優先ジョブの起動時にスワップアウトによる起動を試みている。しかし、スワップアウトによる起動はシステム資源の有効利用の観点からは好ましいものではない。そこで、スワップアウトによりシステム状況が悪化することはないのか、ジョブのターンアラウンドタイムの目的は達成されているのかをスワップアウト効果実験により確認する。

(1) 実験条件

基本実験はスワップアウトパラメータ=(4, 10)で行った。これに対し、本実験ではスワップアウトパラメータ=(0, 0), (4, 10), (8, 20) の場合について行い、その効果を比較する。なお、負荷が極端に高い場合には実行ジョブの大半も α_3 ジョブとなるためにスワップアウト可能なジョブが少なく、効果が期待できない。ディバグジョブが多い場合にもジョブ発生件数が多くなり、ジョブが α_3 ジョブとなる確率が高くなる。そこで、本実験は負荷率1.0、ディバグ率0.0の条件で行うことにした。したがって、スワップアウト効果実験の実行条件として以下のようにパラメータを設定した。なお、ワークロードは全12種類を使用する。

- 負荷率 : 1.0

- ディバグ率 : 0.0
 - α_2 キューへの遷移パラメータ (FAC1) : 0.25
 - α_3 キューへの遷移パラメータ (FAC2) : 3.0
 - プレステージジョブ数 : 3
 - スワップアウトパラメータ : (0, 0), (4, 10) および (8, 20)
 - FEP 配下の磁気ディスクの入出力速度 : 18MB/S
 - PE 占有時間ファクタ : 1.0
- (2) 実験結果の評価・検討
- (a) PE 使用状況

図5.9(1)にスワップアウト効果実験結果の平均 PE 使用台数を示す。平均 PE 使用台数は、最大スワップアウト PE 台数が大きくなるに従い、若干ではあるが向上している。この結果からは、スワップアウトによる顕著な悪影響を見ることはできない。

(b) SSU 使用状況

図5.9(2)に本実験結果の平均スワップアウト使用量を示す。使用量は最大スワップアウト PE 台数が0の場合には当然のことながら0であり、最大値が大きくなるに従って増加している。この結果からも、スワップアウト PE 台数の制限が SSU 使用量を制限するために有効であることが判る。

(c) ジョブ処理状況

図5.10に本実験結果のジョブ実行率を示す。どのジョブの実行率も最大スワップアウト PE 台数が大きくなるに従って僅かではあるが向上している。

(d) 安定条件

シミュレーションの安定条件に関しては、6日から21日程度ですべて条件を満たしたことを確認した。

(3) スワップアウト効果実験から得られた結論

スワップアウト効果実験結果の検討から、以下の事柄が確認された。

- ①最大スワップアウト PE 台数を大きく設定しても、ジョ

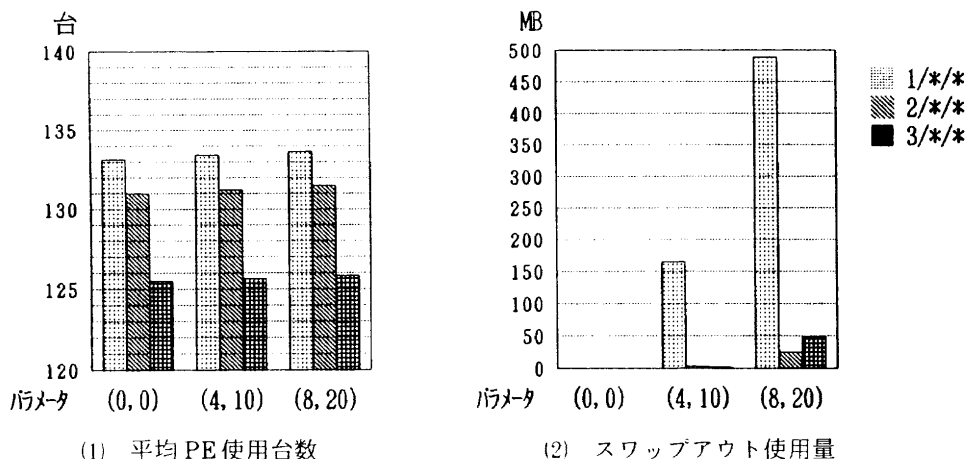


図5.9 スワップアウト効果実験のシステム資源使用量

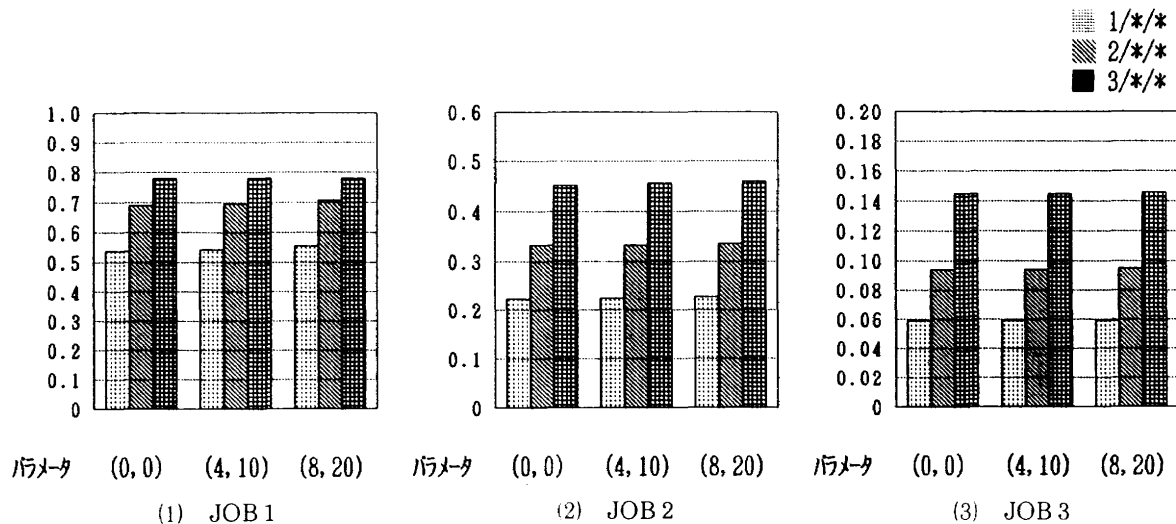


図5.10 スワップアウト効果実験のジョブ実行率

ブ実行率は僅かしか向上しないため、劇的な効果は期待できない。

- ②スワップアウトに一定の制限を設けているため、スワップアウトによる PE の利用率、ジョブ実行率の顕著な低下は見られない。

5.4 起動優先ジョブのための PE リザーブ効果実験

起動優先ジョブが起動できない場合に他のすべてのジョブの起動を抑止する方式と起動優先ジョブの実行に必要な PE のみをリザーブする方式の実験を行い、PE リザーブ方式の優位性を確認する。

(1) 実験条件

PE のリザーブ効果を見るためには、システム状況はある程度混雑している必要があり、ジョブの規模は大きい方が影響が出易いと考えられる。そのため、負荷率は 1.0 および 2.0 とし、デバッグ率は 0.0 として実験する。以下に示すように、PE リザーブ効果実験の実行条件のパラメータを設定する。なお、ワークロードは全 12 種類を使用する。

- 負荷率：1.0 および 2.0
- デバッグ率：0.0
- α_2 キューへの遷移パラメータ (FAC1)：0.25
- α_3 キューへの遷移パラメータ (FAC2)：3.0
- プレステージジョブ数：3
- スワップアウトパラメータ：(4, 10)
- FEP 配下の磁気ディスクの入出力速度：18 MB/S
- PE 占有時間ファクタ：1.0

(2) 実験結果の評価・検討

(a) PE 使用状況

図 5.11(1)に負荷率 1.0 の場合の平均 PE 使用台数を、図 5.12(1)に負荷率 2.0 の場合の平均 PE 使用台数を示す。負荷率 1.0 の場合、ワークロードパラメータが 1/*/* および

2/*/* の場合には PE リザーブ効果は少ないが、発生ジョブの規模が大きい 3/*/* のときには 5 台分以上の効果を示している。負荷率 2.0 の場合には負荷率 1.0 の場合より効果は顕著になっており、特に 3/*/* のときには 11 台分以上の効果を示している。

(b) ジョブ処理状況

図 5.11(2)~(4)に負荷率 1.0 の場合のジョブ実行率を、図 5.12(2)~(4)に負荷率 2.0 の場合のジョブ実行率を示す。

負荷率 1.0 の実験では、ワークロードパラメータが 1/*/*、2/*/* の場合には、全体的なジョブ実行率の効果が僅かながら見られる。また、3/*/* の場合には待たされている小規模ジョブの起動の機会が PE リザーブにより増え、その結果 JOB1 の実行率が大きく向上した反面、JOB3 の実行率が僅かに低下している。

負荷率 2.0 の実験でも、負荷率 1.0 と同様の結果となっているが、ワークロードパラメータが 3/*/* の場合には、JOB3 だけでなく JOB2 も僅かながら実行率が低下している。

(c) 安定条件

シミュレーションの安定条件に関しては、6 日から 42 日程度ですべて条件を満たしたことを確認した。

(3) PE リザーブ効果実験から得られた結論

PE リザーブ効果実験結果の検討から、以下の事柄が確認された。

- ①PE リザーブの平均 PE 使用台数に対する効果は、発生するジョブの規模が大きく、負荷率が高い場合ほど顕著になる。
- ②PE リザーブのジョブ実行率に対する効果は、小規模ジョブに有効である。

5.5 プレステージ効果実験

プレステージは、PE を有効に使用するための手段とし

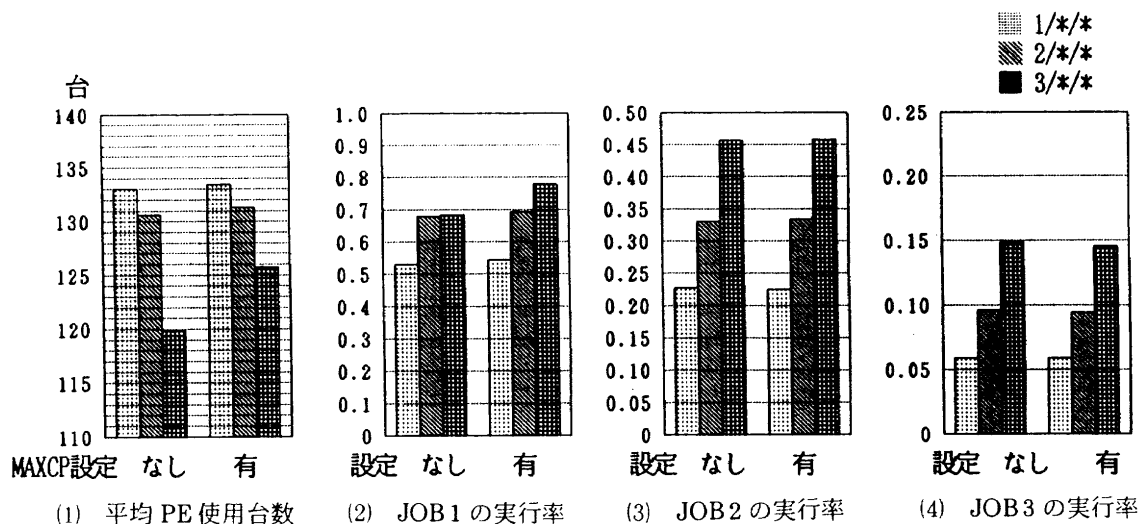


図5.11 PE リザーブ効果実験結果（負荷率1.0の場合）

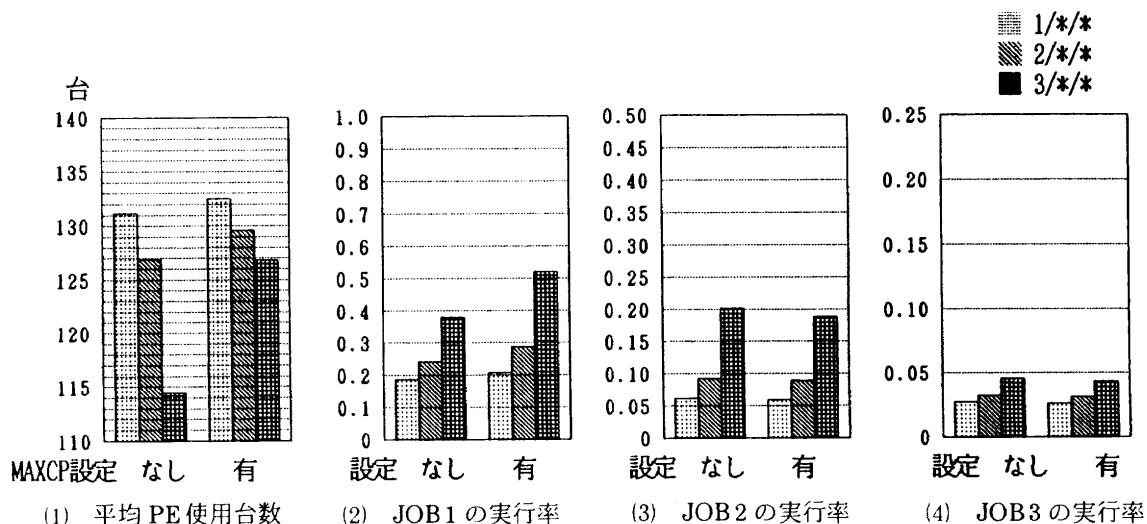


図5.12 PE リザーブ効果実験結果（負荷率2.0の場合）

て期待されている。本項では、最大プレステージジョブ数を変化させた実験を行い、プレステージの効果がどの程度期待できるものかを予測する。

(1) 実験条件

FEP のディスク-SSU 間、SSU-PE 間のファイル転送速度の差が大きい場合、ジョブの入出力量が相対的に高い場合（ディバグジョブが多い場合、ジョブの CPU 時間が制限されている場合などがこれに相当する）にはプレステージの効果が期待できる。ジョブの入出力量が少ない場合には、ファイルの転送速度が多少低くても PE の利用率に悪影響を及ぼすことはないが、ジョブの入出力が多くファイルの転送速度の差が大きい場合には、転送速度の差はジョブが起動してから実際に開始されるまでに要する時間を決定するため、PE の実質的な利用率に直接の影響を及ぼす。そこで、FEP のディスク-SSU 間のファイル転送速度を 3 MB/S と仮定して実験を行うことにした。この

値は並列入出力を行わない場合の 4.5 MB/S から決定した推定実効速度であり、今までの経験に基づき推定した。また、ジョブ入出力を相対的に多い条件で行うためにディバグ率は 0.8 とする。

以上の結果、実験条件として以下のようにパラメータを設定する。なお、本実験はプレステージジョブ数を変化させるためワークロード当たりの実験数が多くなる。そこで、代表的なワークロードとして 1/1/1 の 1 種類を使用する。

- 負荷率：1.0
- ディバグ率：0.8
- α_2 キューへの遷移パラメータ (FAC1)：0.25
- α_3 キューへの遷移パラメータ (FAC2)：3.0
- プレステージジョブ数：0, 1, 2, 3, 4, および 5
- スワップアウトパラメータ：(4, 10)
- FEP 配下の磁気ディスクの入出力速度：3 MB/S
- PE 占有時間ファクタ：1.0

(2) 実験結果の評価・検討

(a) PE 使用状況

図5.13(1)にプレステージ効果実験結果の平均 PE 使用台数を示す。プレステージジョブ数の増加に伴って PE 使用台数は高くなっているが、ジョブ数3の近辺からはほぼ横這い状態となっている。この場合プレステージは PE 4 台分程度の効果となっているが、先に述べたようにジョブの入出力量が相対的に高い場合にはそれに伴う効果が期待できる。

(b) SSU 使用状況

図5.13(2)にプレステージ効果実験結果の平均プレステージ使用量を示す。プレステージ使用量はプレステージジョブ数にはほぼ比例している。

(c) ジョブ処理状況

図5.11(3)にプレステージ効果実験結果のジョブ実行率を示す。プレステージジョブ数が高くなるに従いジョブ実行率も高くなっているが、効果の最も顕著な JOB1 の場合には、プレステージジョブ数3の近辺からはほぼ横這い状態となっており平均 PE 使用台数の結果と一致している。

(d) 安定条件

シミュレーションの安定条件に関しては、2日～4日半程度ですべて条件を満たしたことが確認された。

(3) プレステージ効果実験から得られた結論

プレステージ効果実験結果の検討から、以下の事柄が確認された。

- ①本実験のように入出力量がかなり少ない場合でも、プレステージにより4台分程度の PE 利用率の向上が期待できる。したがって、ジョブの CPU 時間制限を行った場合やジョブの平均 CPU 時間が短い場合には、プレステージの効果は極めて高いと予想される。
- ②SSU の容量に余裕があれば、常時3ジョブ分程度のプレステージを行うことが望ましい。

5.6 PE 占有時間精度実験

本アルゴリズムでは、要求 CPU 時間に一定のファクタを掛けて PE 占有時間を予測するが、ファクタの設定値によりシステム効率がどの程度影響を及ぼされるのか、また適切な設定値はある程度予測できるのかを調べる。

(1) 実験条件

本実験では、ワークロードはジョブ規模およびジョブ時間の平均が最も小さい1/1/1を使用し、ディバグ率も0.8で行うことにした。そして、ファクタ値を0.8～1.4まで0.1刻みに変化させた実験を行うことにした。

以上の結果、実験条件として以下のようにパラメータを設定した。

- 負荷率：2.0
- ディバグ率：0.8
- α キューへの遷移パラメータ (FAC1)：0.25
- α_3 キューへの遷移パラメータ (FAC2)：3.0
- プレステージジョブ数：3
- スワップアウトパラメータ：(4, 10)
- FEP 配下の磁気ディスクの入出力速度：3 MB/S
- PE 占有時間ファクタ：0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4

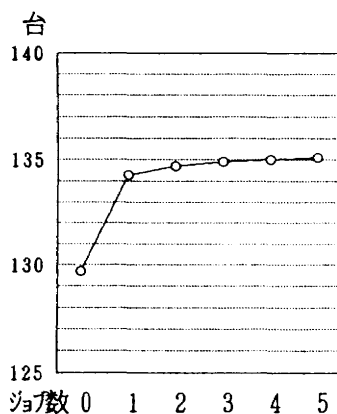
(2) 実験結果の評価・検討

(a) PE 使用状況

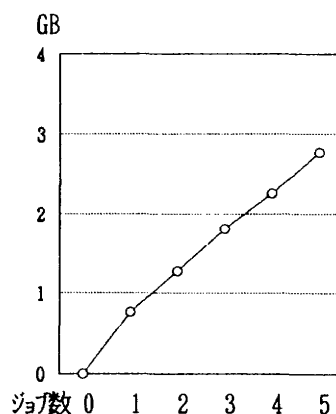
図5.14(1)に PE 占有時間精度実験結果の平均 PE 使用台数を、同図(3)に予想 PE 占有時間内に終了したジョブの割合を示す。同図(3)から、本実験での PE 占有時間ファクタの適正值は1.0から1.1の間にあったと考えられるが、PE 占有時間ファクタが1.0に近づくにつれて PE 使用台数は高くなっており、それ以上ではほぼ横這い状態となっている。

(b) ジョブ処理状況

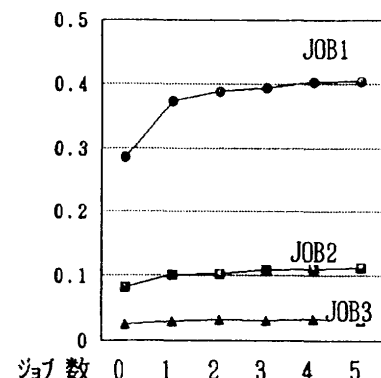
図5.14(2)に PE 占有時間精度実験結果のジョブ実行率を示す。PE 占有時間ファクタの適切な設定は、小規模ジョブの実行状況の改善に有効であることが判る。



(1) 平均 PE 使用台数



(2) 平均プレステージ量



(3) JOB 実行率

図5.13 プレステージ効果実験結果

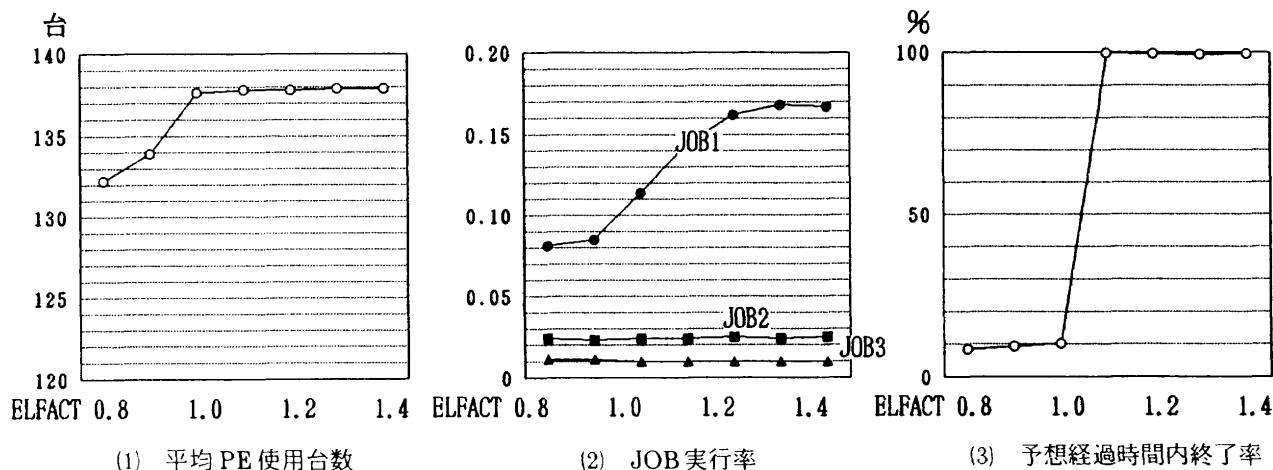


図5.14 経過時間精度実験結果

(c) 安定条件

シミュレーションの安定条件に関しては、1日～1日半程度ですべて条件を満たしたことが確認された。

(3) PE 占有時間から得られた結論

PE 占有時間精度実験結果の検討から、以下の事柄が確認された。

- ① PE 占有時間ファクタを適切に設定することにより PE6 台分程度のシステム効率の向上が期待できる。
- ② PE 占有時間ファクタの設定はあまり厳密に考える必要はなく、適切な設定値が判らない場合には大きめに設定の方が効果的である。

5.7 シミュレーションプログラム

本章で使用した実験結果は、スケジューリングアルゴリズムを評価するために作成したシミュレーションプログラムを使用している。このプログラムは FORTRAN で記述した約2000ステップのプログラムであり、事象を制御するためのメインプログラムと12のサブルーチンから構成されている。代表的なサブルーチンには、ジョブスケジューリングプログラム、ジョブ属性の決定および発生ジョブや実行終了ジョブの統計情報の収集を行うジョブ属性決定プログラム、入出力所要時間の計算および SSU 使用統計情報の収集を行う SSU 管理プログラムがある。

本プログラムを使用して行ったシミュレーション実験の所要 CPU 時間は処理するイベント数に依存するためジョブの処理件数が増えるに従って増加する傾向にある。また、シミュレーションは VP-2600 を使用した TSS 処理で実行したが、基本実験の CPU 使用時間は負荷率 0.8 の場合で 0.459 秒～4.412 秒、負荷率 1.0 の場合で 0.529 秒～6.635 秒、負荷率 2.0 の場合で 0.629 秒～60.139 秒であり、50 日分のシミュレーション実験であるにも係わらず極めて短時間で処理していることから、本シミュレーションプログラムの制御方式がシステム評価手法として極めて有効であったこと

が判る。なお、ジョブ実行率等の計算はシミュレーション結果から出力したジョブ情報（ジョブごとのシステム資源要求量、ジョブ投入時刻、起動時刻、終了時刻等の情報）を使用して別ジョブで処理しており、シミュレーションの所要時間には含まれていない。

6. おわりに

NWT システムは、140 台の PE から構成されている並列計算機システムである。各 PE はピーク性能 1.7 GFLOPS という高性能のベクトル計算機であるため、運用を行う上では PE の有効利用を図ることのできるスケジューラが必須となる。しかし、このような規模の並列計算機システムは世界でも実運用の例がないこと、汎用計算機のスケジューリングアルゴリズムをそのまま適用できないことから、新たなスケジューラを開発する必要があった。

並列計算機システム用スケジューラの開発では、PE の有効利用および適切なジョブターンアラウンドタイムの保証を第一優先で考えた。その結果、無駄な空き PE を発生しない PE 割り当て方式と待ち時間に基づくジョブ実行制御方式を組み込んだスケジューラを考案した。また、NWT システムでの検証実験を最小限度に抑えるために、検証実験に先立って、スケジューリングパラメータの効果、スケジューラの有効性に関する実験をソフトウェアシミュレータを使用して行った。その結果、以下の事柄が確認できた。

- ① 本スケジューラの処理能力はジョブクラス別多重度制御方式の場合と比較して PE 40 台分程度向上する。
- ② 負荷率の変化にも対応可能であり、ジョブの規模等に関係なく公平に処理することができる。
- ③ 負荷が高くなると平均 PE 使用台数も上昇する。まれに起動優先ジョブのために低下する場合もあるが、低下の度合いは極めて少ない。
- ④ スワップアウトはターンアラウンドタイムの適正化に有効に作用する。また、スワップアウトに一定の制限を設

けることにより、システム資源利用率やジョブ実行率の低下は阻止できる。

⑤起動優先ジョブのための PE リザーブは、他のジョブの起動を完全に抑止する場合に比べ最大 PE 10 台分程度の処理能力の向上に匹敵する。

⑥常時 3 ジョブ分程度をプレステージしておくことにより、PE 4 台分程度処理能力が向上する。

上記の事柄を総合した結果、本スケジューラの有効性は極めて顕著であり、NWT 用のスケジューラとして組みむべきであると結論した。

<謝辞>

本スケジューリングアルゴリズムの検討・開発にあたり、財団法人高度情報科学技術研究機構副理事長の三好甫氏から有益なる助言をいただいた。この場を借りて感謝の意を表する。

参 考 文 献

- 1) 三好 甫 ; CFD の推進に必要な計算機性能, NAL SP-13, pp.1~26, 1990年 9 月
- 2) 三好 甫 ; 航技研超高速数値風洞 (UHSNWT) の構想, NAL TR-1108, 1991年 5 月
- 3) 三好 甫 ; 数値風洞 ; 要求要件と概略, NAL SP-16, pp.91~98, 1991年12月
- 4) 三好 甫, 吉岡義朗, 池田正幸, 高村守幸 ; 数値風洞のハードウェア, NAL SP-16, pp.99~106, 1991年12月

- 5) 福田正大, 末松和代, 土屋雅子, 大空 瞭, 工内 隆, 坂本喜則 ; 数値風洞のオペレーティングシステム, NAL SP-16, pp.107~114, 1991年12月
- 6) 福田正大, 中村 孝, 吉田正廣, 岡田 信, 中村修一 ; 数値風洞の言語ソフトウェア, NAL SP-16, pp.115~122, 1991年12月
- 7) 吉田正廣, 中村 孝, 福田正大, 中村修一, 岡田 信 ; NWT 向け並列 Fortran プログラミングについて, NAL SP-19, pp.247~252, 1992年12月
- 8) 中村 孝, 吉田正廣, 福田正大, 中村修一, 村瀬丈夫, 松崎達哉 ; NWT 並列 Fortran に基づく並列評価, NAL SP-19, pp.253~258, 1992年12月
- 9) 土屋雅子, 末松和代, 末松俊二, 畠間晴夫, 森重博司, 山口 靖, 軽部行洋 ; NWT のユーザインタフェース実現方式, NAL SP-19, pp.259~264, 1992年12月
- 10) 情報処理ハンドブック, 情報処理学会, 1989年11月
- 11) 中村 孝, 吉田正廣, 福田正大, 村瀬丈夫, 松崎達哉 ; CFD プログラムによる NWT の性能評価, NAL SP-22, pp.179~184, 1994年 3 月
- 12) 土屋雅子 ; NS システム用ジョブ・ジョブステップ・スケジューラの開発, NAL TR-977, 1988年 6 月
- 13) JUHNSON KGIZ ; Distributions in Statistics ; Continuous Univariate Distribution-2

航空宇宙技術研究所報告1277号

平成7年9月発行

発行所 航空宇宙技術研究所
東京都調布市深大寺東町7丁目44番地1
電話三鷹(0422) 47-5911(大代表) ㊞182
印刷所 株式会社三興印刷
東京都新宿区西早稲田 2-1-18

Printed in Japan