

# 仮想計算機システムの運転方式について

加藤 保夫      川崎 隆二

NTTソフトウェア研究所

仮想計算機システムの運転性を改善する方法として、限定優先方式による運転方式について提案する。

仮想計算機システム上で並走する多数の仮想計算機の中で特定の仮想計算機、例えばオンラインシステムを搭載した仮想計算機にCPU時間を優先して割り付けることにより、他仮想計算機に比べスループットを向上させることができる。既存のFIFO制御を用いて優先制御を行う場合、CPU負荷の大きい仮想計算機が優先指定されると優先度の低い仮想計算機の沈み込みという問題が発生する。これらの問題点を改善するため、限定優先方式ではCPU時間を取り上げる機能と優先度を制限する機能を設けた。

本報告では限定優先制御方式について紹介するとともに、本方式の評価結果を示す。

## An Operating Method for the Virtual Machine System

Yasuo KATO and Ryuji KAWASAKI

NTT Software Laboratories, NTT

1-9-1 Kohnan Minato-ku Tokyo 108 Japan

In this paper a new virtual machine system operating methodology is proposed. The method is named "Limited Priority Scheduling Method (LPSM)".

When various operating systems run parallelly on a virtual machine system, previously it was difficult to get high throughput for individual operating system. Traditional FIFO control method causes "under schedule problem", in which low priority operating system is suppressed to get CPU resource by high priority operating system. The LPSM resolved this problem by controlling each operating system priority.

This paper presents the mechanism of LPSM and its evaluation.

## 1. はじめに

従来仮想計算機システム（VMシステム）は、VMシステムが実現する仮想計算機（VM）の性能的な問題から主に試験システムを搭載していた。しかしVMを高速に動作させるためのハードウェアアシスト機構を備えた計算機の開発により、VMの性能は飛躍的に向上してきている。一方ハードウェアの性能も毎年向上してきており、VMシステムを搭載しても性能的な問題はすくなくなくなっている。このような状況から従来プログラムの試験のみに利用していたVMシステムをバッチサービス用更にはオンラインサービス用のシステムを搭載し利用できるようになった。

VMシステム上に試験システム、バッチシステム、及びオンラインシステムを混在して搭載し運転する場合、搭載システムの特性上の問題の1つとして試験システムのプログラムバグによるループでCPU時間が寡占され、他のVMの沈み込みが発生することが考えられる。一方VMシステムの運転方式として、搭載するシステムのサービス条件から特定VMに優先的にCPU時間を割り当て、当該VMのスループットを向上させる方式が考えられる。これにより例えばオンラインシステムの端末での応答時間を改善することができる。この場合当該VMの負荷が高くなった時他の優先されていないVMが沈み込まないように、各VMにCPU時間を割り付ける方法が必要となる。

本論文では、上記のような各種システムが混在して搭載したVMシステムの運転性向上のためのスケジュール方式として、VMの沈み込み等の問題を解決した限定優先制御方式について一案を紹介し、その評価結果を示す。本限定優先方式を用いることにより特定の仮想計算機のスループットの向上を図ると共に、仮想計算機の沈み込みも防止することができ、仮想計算機システムの運転性向上を図ることができる。

## 2. VMの利用形態モデルについて

### 2.1 搭載システムの特性とサービス条件

検討を進めるにあたりVMシステムに搭載するシステムを試験システム、バッチシステム、オンラインシステムに分類し、各システムの特性とサービス条件をCPU使用率の観点から表1のように整理した。

表1 搭載システムの特性とサービス条件

搭載システム	搭載システムの特性とサービス条件
試験システム	作成したプログラムの試験を行うためのシステムである。実行するジョブとしては、試験環境の作成、試験プログラムの入替え、試験プログラムの走行、試験結果の出力がある。通常状態ではCPU要求率は低いが、試験中であることからトレース等のCPU時間を必要とするデバッグ機能の利用、プログラムバグによるループ状態があり一時的にCPU時間を寡占することがある。CPU時間の割当の観点からは優先してCPU時間を割り当てる必要性はない。また他のシステムがCPU時間を必要としている場合は当該システムのCPU時間を犠牲にして良い。
バッチシステム	バッチ処理サービスを提供するシステムである。長時間のジョブ、CPU使用率の高いジョブを走行させることから、CPU使用率は中程度から高程度となる。CPU時間の割当の観点からは優先してCPU時間を割り当てる必要性は少ない。しかし他の優先度の高いシステムがCPU時間を必要としても、当該システムにはある一定のCPU時間を確保する必要がある。
オンラインシステム	オンラインサービスを提供するシステムである。同時接続の端末台数とそのトラフィックによりCPU使用率は変わるが、CPU使用率は中程度と考えられる。オンラインサービスであることから端末に対する応答時間が短い必要がある。即ち当該システムに対して優先的にCPU時間を割り当てる必要がある。

### 2.2 VMスケジュール方式の考え方

VMスケジュール方式の検討に当たり、前述の搭載システムの特性とサービス条件から以下の設計条件を設定した。

- ① オンラインシステムへの適用が可能なこと。
- ② 各搭載システムのCPU使用率を容易に制御可能なこと。
- ③ 搭載システムの沈み込みが回避可能なこと。

VMシステムでは搭載システム対応にVMを割り当てている。搭載システムの沈み込みとは、VMが要求するCPU時間に対して、VMシステムのスケジュール機能が当該VMに対してCPU時間を割り当てないために搭載システムの実行が非常に遅くなる現象である。例えば特定VMが優先的にスケジュールされる運転方式の場合、当該VMのCPU使用率が非常に高くなると優先度の低いVMの実行が遅くなる。

当該問題はCPU負荷の制御の問題として捉えることができることから、オペレーティングシステム(OS)のジョブスケジュールについても類似の問題がある。OSのジョブは一定の時間で生成/消滅を繰返しており、ジョブ内でのCPU使用率の変動が大きくないことから、OSはジョブの開始を抑制することにより負荷制御することができる。しかしVMシステムでは、以下の理由からVMの起動を押し止動的に負荷制御することは困難である。

- ① オンラインシステム/バッチシステム用のVMではVMシステムの運転開始/終了と同期して起動/停止を行う。
- ② 試験システム用のVMでは利用者が利用するかなり長い時間を単位として起動/停止を行う。
- ③ 搭載システムによるが、CPU使用率の変動が大きい。

### 3. VMスケジュール方式案

VMに対するCPU時間のスケジュール方式として以下の方式案が考えられる。

#### (1) 平等方式 (FIFO方式)

レディ状態の全VMの実行をFIFOにより制御する方式であり、各VMの動作特性が同一の場合各VMのCPU使用率は平等になる。本方式の適用方法としては、全VMに試験システムを搭載している場合のように各VMへ平等にCPU時間を割り当てる場合が考えられる。VMの沈み込みという問題は発生しないが特定VMを優先して制御することはできない。

#### (2) クラス優先方式 (多段ラウンドロビン方式)

OSのジョブ管理等で広く実施されている方法であり、各VMをクラスで分類し、レディ状態のVMのうち優先度の高いクラスのVMを実行し、同一クラス内のVM間はFIFOにより制御する方式である。特定VMを優先して実行する方法として有効であるが、優先度の高いクラスのVMのCPU使用率が高いとVMの沈み込みの問題が発生する。また優先度の高いクラスのVMがウェイト状態になると優先度の低いクラスのVMが走行するので、当優先度の低いクラスのVMがタイムスライスまでCPU時間を使用する場合、優先度の高いクラスのVMはCPU時間を汎山使用することはできない。従って端末への応答時間の条件がきびしいオンラインシステムを搭載するVMシステムには不適当である。

#### (3) 限定優先方式

本論文において提案する方式であり、クラス優先方式に対して応答時間を改善するためCPU時間を取り上げる機能と、VMの沈み込みを回避するため優先度を限定する機能を付加した方式である。優先度の限定はクラス毎設けた優先カウンタにより制御する。優先カウンタの最大値は優先回数(N)と非優先回数(M)を合計した値である。なお優先回数の割合(=  $N / (N + M)$ )を優先率という。スケジュールはクラス優先方式と同様に優先度の高いクラスのVMから検索し、当該クラスの優先カウンタが優先範囲の場合スケジュールする。非優先範囲の場合はスケジュールしない。さらに優先度の低いクラスのVMが走行中に優先度の高いクラスのVMがレディになった時、優先度の高いクラスの優先カウンタが優先範囲であるなら走行中の優先度の低いクラスのVMを停止させ優先度の高いクラスのVMを走行させる。これをCPU時間の取り上げという。優先範囲とは優先カウンタの値が1~Nのときであり、非優先範囲とは  $N + 1 \sim N + M$  のときである。優先カウンタは、①当該クラスのVMがスケジュールされた時、②下位のクラスのVMがスケジュールされた時、更新される。優先度の高いクラスの優先率を1にすることにより当該クラスのVMはレディになった時に即時に走行することになり、応答時間を改善することができる。VMの沈み込みは優先率を0から1の範囲で設定することにより防ぐことができる。各クラスのVMが常にレディ状態の場合の限定優先方式のスケジューリング例を表2に示す。

各スケジュール方式の比較結果を表3に示す。

表2 限定優先方式のスケジューリング例

フェーズ	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
クラス1	1	2	3	1	2	3	1	2	3	1	2	3	1	2	3	1	2	3
クラス2			1			2			3			1			2			3
クラス3									1									2

フェーズ	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36
クラス1	1	2	3	1	2	3	1	2	3	1	2	3	1	2	3	1	2	3
クラス2			1			2			3			1			2			3
クラス3									3									4

(凡例) 数字：優先カウンタ      網掛け有：該クラス走行      網掛け無：該クラス非走行

(条件)

クラス	優先回数	非優先回数
クラス1	2	1
クラス2	2	1
クラス3	3	1

表3 スケジュール方式の比較結果

要求条件	スケジュール方式		
	(1)平等	(2)クラス優先	(3)限定優先
オンラインシステムへの適用	×	△	○
CPU使用率の制御の容易性	×	△	○
VMの沈み込みの発生	○	×	△

#### 4. 評価方法と結果

各方式についてシミュレーションにより評価した。

##### (1) 評価モデル

各VMはスケジューリング時に指数分布に従ってCPU時間を使用し、VMモニタによりタイムスライス切れとなったときには直ちにレディ状態となり、次にスケジューリング時はまた指数分布に従ったCPU時間を使用する。タイムスライス切れとなる前にCPU時間を放棄した場合は、指数分布に従った時間経過後にレディ状態になる。CPU時間を放棄するのは周辺装置との入出力処理のためであるが、ここでは周辺装置との入出力時の競合は無いと仮定した。

##### (2) VMシステムの動作特性の測定法

各VMが動作するときには、VMの走行時間に応じて①VM走行に伴うVMモニタの処理分、②VMモニタのプロセス処理分、③割込処理の非タスク処理分の各オーバーヘッドが発生する。①については当該VMのCPU時間としてカウントした。②については全てのVMより高い優先クラス(クラス0)で動作する独立した処理としてカウントした。③についてはVM走行と同期して余分に時間が経過したとして処理した。

##### (3) シミュレーション条件

VMの優先クラスは、オンライン用、バッチ用及びデバグ用の3クラスとし、表4のシミュレーション条件を仮定した。ここでCPU要求率は当該VMが走行する時に必要とするCPU時間の割合であり、競合がない時はCPU要求率がCPU利用率となる。クラス1はオンライン用であることから優先率を大きくし、クラス2とクラス3はケース1のCPU要求率を考慮して決定した。ケース1は定常時である。ケース2はトラブル等によりデバグ用VMのVM05が高負荷時であり、CPU要求率が100%となり、VMシステムが過負荷となった状態である。

##### (4) シミュレーション結果

以下の場合についてシミュレーションを行い応答時間とCPU取上率を求めた。

- 図1 === ケース1、クラス1のCPU要求率を変化
- 図2 === ケース2、クラス1のCPU要求率を変化
- 図3 === ケース1、クラス1の優先率を変化、クラス1のCPU要求率=20%
- 図4 === ケース1、クラス1の優先率を変化、クラス1のCPU要求率=80%
- 図5 === ケース2、クラス1の優先率を変化、クラス1のCPU要求率=20%
- 図6 === ケース2、クラス1の優先率を変化、クラス1のCPU要求率=80%

応答時間はCPU使用率をもとに算出した。表3のCPU要求率はCPU使用率の期待値であり競合のない状態でのCPU使用率である。VMシステム上で走行した場合、他VMとの競合及びVMモニタとの競合のためCPU使用率はCPU要求率より小さくなる。すなわち競合のない状態での応答時間を1とすれば競合時の応答時間は(CPU要求率/CPU使用率)となる。

CPU取上率とは限定優先方式において発生する状態であり、全スケジューリング回数に対するCPU時間の取り上げによりスケジューリングされた割合である。CPU取上率が高いということは、当該クラスが優先されて実行している指標になる。

表4 シミュレーション条件

優先クラス	VM識別詞	優先制御		CPU要求率(%)		記 事
		優先回数 (N)	非優先回数 (M)	ケース1 (定常時)	ケース2 (高負荷時)	
クラス0 (VMモニタ用)	-	1	1	-	-	
クラス1 (オンライン用)	VM00	9	1	変化	変化	図1、図2
		変化		20	20	図3、図5
		変化		80	80	図4、図6
クラス2 (バッチ用)	VM01	5	2	35	35	
クラス3 (デバグ用)	VM02	10	0	5	100	115
	VM03			5		
	VM04			5		
	VM05			5		

## 5. 考察

### (1) 応答時間の改善効果について

図1、図2の結果から限定優先方式では他の方式に比べクラス1の応答時間が期待値に近いことから、オンラインシステムの応答時間の改善に限定優先方式は有効である。特にクラス3のCPU要求率が高くなった時には有効である。

### (2) CPU使用率の制御について

平等方式と比較した場合、限定優先方式ではクラス1のCPU要求率が増加するに従ってクラス2と3の応答時間が非常に大きくなる。これはクラス1がCPU時間を取り上げたためである。またクラス優先方式と限定優先方式では、クラス2の応答時間が大きく異なる。これは限定優先方式では優先率によりクラス2の走行が制約されているためである。従って各クラスのCPU使用率が優先率により制御できることが分かる。

### (3) VMの沈み込みについて

図1と図2において、限定優先方式のクラス3の応答時間はクラス優先方式より僅かに大きい。すなわち限定優先方式がクラス3の沈み込みが大きいことになる。原因としては、①クラス1の優先率が大きすぎる、②CPU取上機能が影響しクラス2と3のCPU時間がクラス1に取り上げられたためと考えられる。

### (4) CPU取上率について

図1と図2からクラス1のCPU要求率が増加するに従って、CPU取上率は減少している。これはクラス1のCPU要求率の増加に連れてクラス1のCPU使用率が増加すると他のクラスのCPU使用率が減少し、

取り上げ対象の優先クラスが走行しなくなるためである。従ってCPU要求率が100%近くでは急に低下することになる。

図3から図6においては優先率の増加とともにCPU取上率が増加する。CPU取上率の増加は応答時間の減少と対応しており、応答時間の改善にCPU取上機能が有効であることが分かる。優先率が1の場合のCPU取上率は、自分より下位のクラスが走行している確率になる。例えば図3では確率0.55に対し取上率は0.56となっている。また図5では確率0.88(=(35+115)/(20+35+115))に対し取上率は0.87となっている。

#### (5) 優先率の設定について

図3から図6の優先率を変化させた場合の応答時間から、優先率はCPU要求率を考慮して決定すれば良いことが分かる。図3と図4のケース1の定常時は優先率を変えてもクラス1の応答時間の変化は少ない。これはクラス1以外のCPU負荷が少ないためである。しかし図5と図6のケース2ではVM05の影響により優先率に依存して応答時間が大きく変化する。応答時間が1.3となるには、図5のクラス1のCPU要求率が20%の場合は優先率0.4、図6のクラス1のCPU要求率が80%の場合は優先率0.7である。優先率を大きくすればクラス1の応答時間は改善できるが、図4と図6からクラス1のCPU要求率が高いと他のクラスの応答時間がかなり大きくなることより、クラス1の優先率を適当な値にすべきである。

VMシステムに各種のサービスを搭載した場合、オンラインシステムの応答時間を最も速くするには、オンラインシステムの優先率を1にすればよい。しかし他のシステムの応答時間は犠牲になることから、オンラインシステムの応答時間が例えば1.3倍程度になることが許容されるなら、優先率は0.7程度でよい(図6)。これにより他のシステムの応答時間もある程度保証できると考えられる。

また優先すべきシステムがなく全システムを同程度の応答時間にするには、各システムのCPU要求率の割合で優先率を決定する。

## 6. おわりに

今回の評価により限定優先方式についてある程度の解析結果が得られた。即ち優先率を設定することにより、①VMシステムのオンラインシステムへの適応、②CPU使用率の制御、が可能であることが確認できた。しかしVMの沈み込みの制御については十分な評価が得られなかった。この点についてはCPU取上機能を除いた限定優先方式について測定し評価する必要があると考えられる。

またシミュレーション結果として、①変化が異様な部分及び②CPU要求率に対して僅かにCPU使用率が大きいという箇所が発見された。この点について問題の解析が進んでいないが、より精度の高いシミュレーションを行い評価結果を検証する予定である。

さらに優先率/CPU要求率に対するCPU使用率は、VMモニタのタイムスライス値、VM上のプログラムのラン時間/ウェイト時間とその分布特性により変化する。今後はこれらの影響についても含めて評価することとしたい。

### 【参考文献】

- [1]加藤、川崎、「仮想計算機システムにおける優先制御方式について」、  
情処第36回全国大会(1988)、PP353-354
- [2]池ヶ谷、他、「Mシリーズ仮想計算機機構(VM/EX機構)の開発-処理方式」、  
情処第36回全国大会(1988)、PP229-230
- [3]池ヶ谷、他、「Mシリーズ仮想計算機機構(VM/EX機構)の開発-評価」、  
情処第36回全国大会(1988)、PP231-232
- [4]小柳津、他、「DIPS-11/5Eシリーズの実用化」、  
NTT研究実用化報告(1987)、Vol. 36、No. 1、PP49-56
- [5]塩川、他、「DIPS-11/5Eシリーズ論理装置」、  
NTT研究実用化報告(1987)、Vol. 36、No. 1、PP57-65

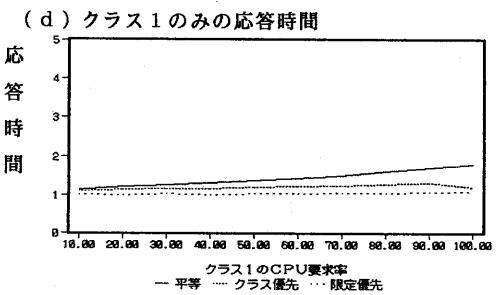
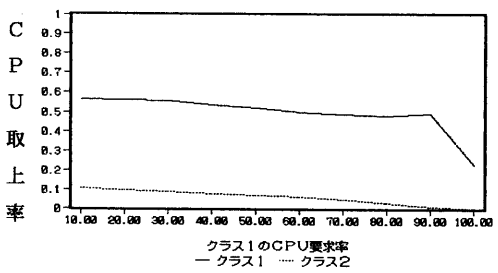
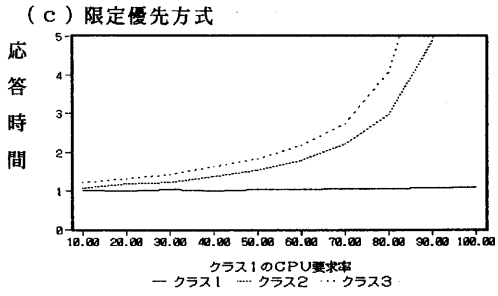
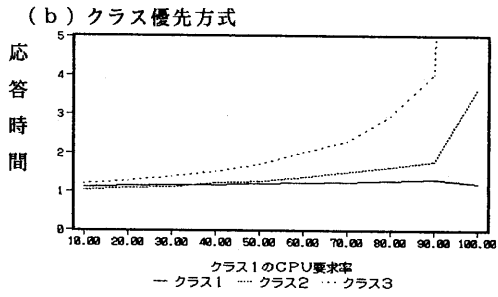
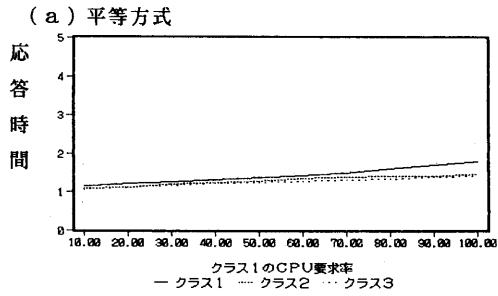


図1 ケース1、クラス1のCPU要求率変化

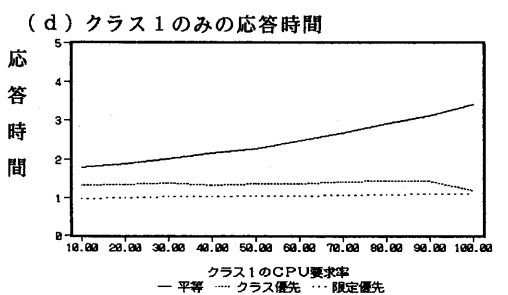
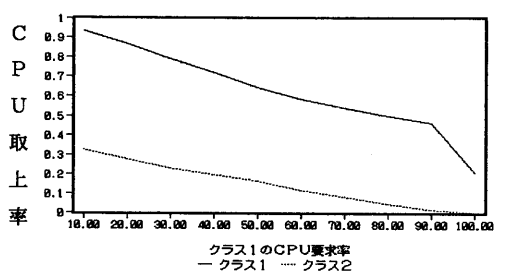
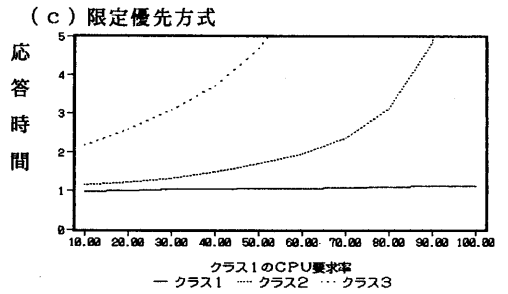
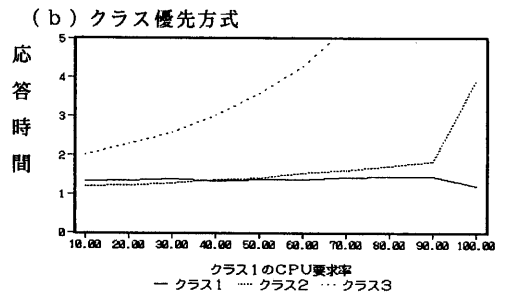
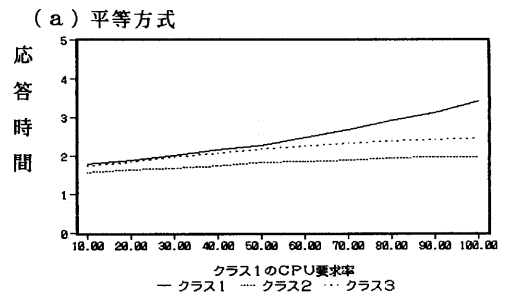


図2 ケース2、クラス1のCPU要求率変化

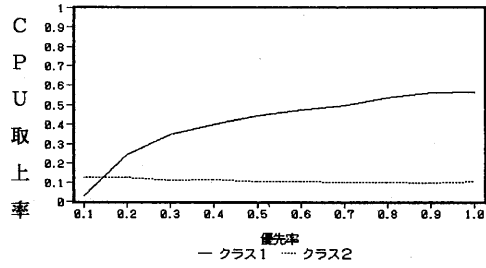
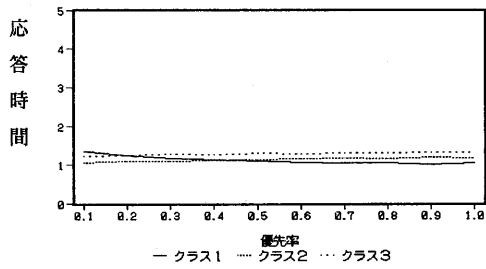


図3 ケース1、クラス1の優先率を変化、  
クラス1のCPU要求率=20%

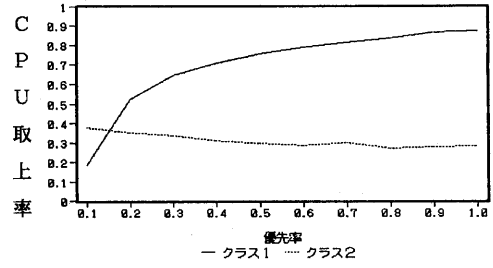
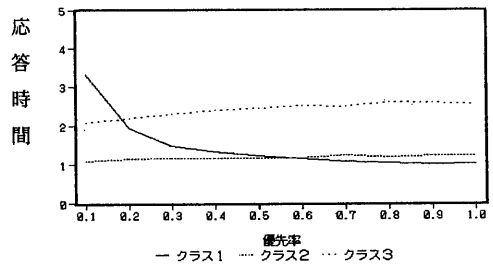


図5 ケース2、クラス1の優先率を変化、  
クラス1のCPU要求率=20%

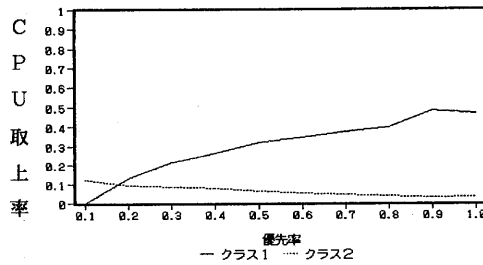
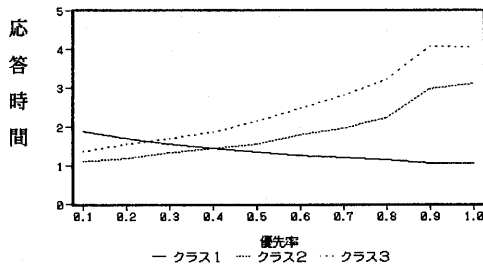


図4 ケース1、クラス1の優先率を変化、  
クラス1のCPU要求率=80%

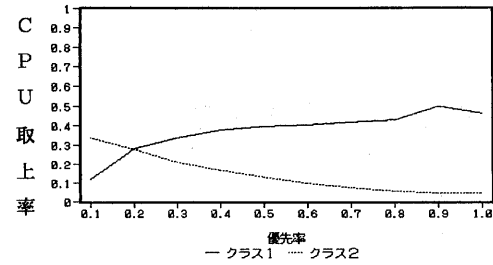
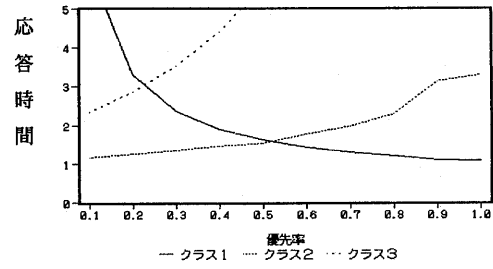


図6 ケース2、クラス1の優先率を変化、  
クラス1のCPU要求率=80%